



School of Physics,
Engineering and
Computer Science

MSc Data Science Project

7PAM2002

Department of Physics, Astronomy and Mathematics

Data Science FINAL PROJECT REPORT

Project Title:

Stock Market Analysis

Student Name and SRN:

Ajay Santhosh K V

23024049

Supervisor: William Cooper

Date Submitted: June 23rd, 2025

GitHub: <https://github.com/ajaysanthoshk/Final-Project.git>

Word Count: 5020 (excluding abstract, references and appendices)

DECLARATION STATEMENT

This report is submitted in partial fulfilment of the requirement for the degree of Master of Science **in Data Science** at the University of Hertfordshire.

I have read the detailed guidance to students on academic integrity, misconduct and plagiarism information at [Assessment Offences and Academic Misconduct](#) and understand the University process of dealing with suspected cases of academic misconduct and the possible penalties, which could include failing the project or course.

I certify that the work submitted is my own and that any material derived or quoted from published or unpublished work of other persons has been duly acknowledged. (Ref. UPR AS/C/6.1, section 7 and UPR AS/C/5, section 3.6)

I did not use human participants in my MSc Project.

I hereby give permission for the report to be made available on module websites provided the source is acknowledged.

Student Name printed:Ajay Santhosh K V

Student Name signature



Student SRN number:23024049

UNIVERSITY OF HERTFORDSHIRE

SCHOOL OF PHYSICS, ENGINEERING AND COMPUTER SCIENCE

ABSTRACT

This research paper investigates the performance of three deep learning models—Long Short-Term Memory (LSTM), Gated Recurrent Unit (GRU), and Convolutional Neural Network (CNN)—for stock price prediction using a large-scale dataset of S&P 500 constituents. The dataset comprises over **619,000 daily stock price observations** (including open, high, low, close, and volume) from **2013 to 2018**, enriched with technical indicators such as **moving averages**, **volatility**, and the **Relative Strength Index (RSI)**. The data was pre-processed into time-series sequences and fed into each model for forecasting. Results indicate that GRU model produced the highest accuracy with RMSE = 1.48, MAE = 1.18, and R² = 0.94, outperforming LSTM and CNN. LSTM model came in the second place, yet CNN model performed badly, ostensibly due to CNN only being capable of identifying relatively short-term dependencies. These results confirm that GRU can be utilized for practical and effective financial time-series forecasting.

TABLE OF CONTENTS

Abstract	3
List of figures.....	3
Chapter 1: Introduction	7
1.1 Background of the Study	7
1.2 Problem Statement	7
1.3 Research Objectives	7
1.4 Research questions	8
1.5 Importance of the Research.....	8
1.6 Structure of the Dissertation.....	8
Chapter 2: Literature Review.....	10
2.1 Introduction.....	10
2.2 Overview of Machine Learning in Financial Forecasting	10
2.3 Comparative Understanding of LSTM, GRU, and CNN in Sequential Modelling	11
2.4 Applications of LSTM, GRU, and CNN in Stock Market Forecasting	12
2.5 Challenges and Future Directions in LSTM-Based Stock Market Forecasting	12
2.6 Conclusion	14
Chapter 3: Methodology.....	15
3.1 Introduction.....	15
3.2 Tools and Environment	15
3.3 Research Design	15
3.4 Data Collection and Preprocessing	16
3.5 Ethical Considerations	16
Chapter 4: Implementation and Results.....	19

4.1 Data preparation	19
4.2 Data Exploration and Market Behaviour.....	19
4.3 Model Implementation and Training	31
4.4 Evaluation and Results Interpretation	33
4.5 Conclusion	36
Chapter 5: Discussion.....	37
5.1 Summary of EDA findings and Model Outcomes	37
5.2 Results Comparison with Literature Review	38
5.3 Identifying the Best Model Based on Project Objectives.....	38
5.4 Limitations and Future Work	38
5.5 Conclusion	39
References	40

LIST OF FIGURES

Figure 1: Stock Closing Prices Over Time	19
Figure 2: Distribution of Daily Returns – AAPL	20
Figure 3:30-Day Rolling Volatility.....	21
Figure 4:AAPL - Moving Averages	22
Figure 5: Correlation matrix of top 20 stocks	23
Figure 6: Stock Price Distributions	24
Figure 7: Scatter Plot Matrix	25
Figure 8: Candlestick chart of AAPL.....	26
Figure 9: Bollinger Bands – AAPL	26
Figure 10: Stock Closing Price Correlation	27
Figure 11: High & Low prices Over Time.....	28

Figure 12: Stock's prices from 2017 to 2018	28
Figure 13: Stock Prices in January 2018.....	29
Figure 14: Trading Volume by company wise	30
Figure 15: Predicted vs Actual with Confidence Intervals	33
Figure 16: SHAP Value Importance Across Time Steps.....	35
Figure 17: Feature Importance Value	36

CHAPTER 1: INTRODUCTION

1.1 Background of the Study

Stock market and other financial markets exhibit intricate complexity and incessant alteration as well as unpredictability. Various forecasting methods allow investors as well as market analysts to predict market fluctuations so they can achieve greater returns as well as reduce risks. Nonlinear nature of time series data is a serious complication to conventional forecasting models of ARIMA and GARCH models as stated by (Huimin, et al., 2024). Deep learning and artificial intelligence possess the potential to revolutionize finance in recent times. Recurrent Neural Networks with Long Short-Term Memory (LSTM) segments possess superior ability in anticipating time series data as well as stock market values. Long Short-Term Memory models outshine traditional models such as ARIMA in predicting stock prices in volatile market environments as stated by (Botunac, et al., 2024). The financial sector has embraced GRUs alongside CNNs as newer technology solutions. (Rahman, et al., 2019) documented that GRU models bring computation benefits which enable them to successfully predict market price fluctuations. Time series forecasting benefits from implementations of CNNs as described by (Wibawa, et al., 2022). We employed LSTM and GRU and CNN in this research to determine their individual forecasting abilities for S&P 500 constituents alongside an examination of their potential deployment in trading systems.

1.2 Problem Statement

Research has extensively analysed stock market predictions with machine learning methods but there is a lack of knowledge about deep learning architecture performance comparison on real trading data and existing studies lack connections between predictions and investments, causing practitioners in finance to lose interest.

1.3 Research Objectives

- To analyse the current body of literature on stock market prediction methodologies, highlighting the departure from traditional statistical approaches to advanced deep learning approaches.

- To perform exploratory data analysis (EDA) on selected stocks from the S&P 500 index to further identify features in trends, volatility, volumes, and technical factors.
- To create and implement LSTM, GRU, and CNN models for the prediction of stock prices based on historical time-series data.
- To investigate the interpretability of predictions generated by each model in the context of making financial decisions.
- To provide recommendations by assessing the EDA results and the model's performance to build better trading strategies.

1.4 Research questions

Which deep learning model (LSTM, GRU, CNN) is more capable of predicting a short-term price change in stocks belonging to the S&P 500, and how usable are these predictions as part of a trading strategy?

1.5 Importance of the Research

This research adds to academic research and industry practice in a few ways. Academically, this research provides a thorough geopolitical analysis of three popular deep learning models in financial time series forecasting. Practically, incorporating trading simulations helps to connect model forecast to real-world investment decisions. The results will likely benefit data scientists, financial analysts, fintech startups and investment firms before they create an AI trading tool.

1.6 Structure of the Dissertation

This dissertation is split up as follows:

Chapter 1 introduces the background, problem, research aims and objectives and importance of research.

Chapter 2 literature review brings forward the relevant out of the financial forecasting models pertinent to the dissertation I.e. LSTM and other deep learning models.

Chapter 3 outlines the research methodology (data collection & pre-processing, the design of the model and evaluation metrics).

Chapter 4 discusses the implementation and experimental results.

Chapter 5 discusses the interpretation of model results and comparison to peer-reviewed literature related to model suitability.

Chapter 6 concludes by summarizing the key findings, modelling contributions, and directions for future research

CHAPTER 2: LITERATURE REVIEW

2.1 Introduction

This literature review examines contemporary developments in stock market analysis, concentrating on the effectiveness of deep learning models, especially Long Short-Term Memory (LSTM) networks. The review discusses how these models improve prediction accuracy in an environment with volatility and market complexity; limitations are also noted and challenges in applying these models to financial forecasting are discussed.

2.2 Overview of Machine Learning in Financial Forecasting

Financial forecasting has changed radically through the incorporation of machine learning (ML) that offers the ability to analyse large datasets and recognize patterns that may not be observable with conventional/previous models (Raza, 2023). Stock price volatility together with time series patterns used to be determined by past forecasting systems through ARIMA and GARCH quantitative methods. Financial market non-linearity on a large scale remains unreachable for models which need linearity assumptions as their foundation. Prediction models progressed by incorporating decision trees and support vector machines (SVMs) and random forests from ML models while data extraction occurred according to (Kumar, et al., 2023). Strong neural networks enabled these ML models to extract flexible complex information from complex dataset through deep learning models. Deep neural networks (DNNs) and recurrent neural networks (RNNs) took models to another development stage because they enabled the computational interpretation of temporal relations within specific datasets according to (Goodfellow, et al., 2016). The association of LSTM networks proves to be particularly useful for financial time series forecasting since it maintains extended sequence memory alongside suitable architectural properties (Moghar & Hamiche, 2020). The industry has adopted LSTM networks instead of statistical methods for financial forecasting due to their adaptive properties in data detection and decreased need for human oversight of detailed relationships.

2.3 Comparative Understanding of LSTM, GRU, and CNN in Sequential Modelling

LSTM is a kind of Recurrent neural network proposes to overcome the gradient disappearance problem. i.e., Researches developed the more advanced LSTM networks because of the gradient vanishing issue that initially affected RNN. The gating mechanism of LSTMs allows attainment of better performance for stock price sequence at each and every processing step. The architectural design includes three important gates which function to regulate the data flow across the entire network systems between input and output as well as data storage and release (Houdt, et al., 2020). LSTM has enhanced ability to handle multi-dimensional data effectively through a technique of attention-based LSTM and hierarchical LSTM. Moreover, it can manage both short-term dependencies and long-term dependencies (Lindemann, et al., 2021). LSTM exhibit superior performance in terms of higher accuracy and reduced error rate than GRU while assessing EEG data for predicting future cognitive states (Rivas, et al., 2025). Furthermore, LSTM excels the base RNN model in forecasting price tendencies of financial time series data as per (Al-Selwi, et al., 2024)

In the same vein, Gated Recurrent Units (GRUs) provide a more streamlined and computationally efficient approach than LSTMs. GRUs consolidate the input and forget gate, resulting in a single update gate for the model, which decreases both model complexity and maintains long-term dependencies (Shen, et al., 2018). While GRUs appear to be faster and require less training resources, there are studies that suggest LSTMs perform better on more involved tasks like predicting financial time series (Rivas, et al., 2025)

At this moment. Convolutional Neural Networks (CNNs) are another area of research with regard to time-series data analysis. CNNs have strengths in recognizing local dependencies and are most often used for image processing. Some research supports CNNs being useful in improving local pattern detection in hybrid model performance with either LSTM or GRU in financial applications (Singh, et al., 2023)

2.4 Applications of LSTM, GRU, and CNN in Stock Market Forecasting

LSTM networks have wider application areas in the numerous fields of finance, healthcare, construction and science. For instance, a study of (Malashin, et al., 2024) demonstrates that LSTM efficiency in forecasting time-series trends as well as model building of sequential data for the purpose of comprehension of intricate molecular structure and transient behaviours of polymeric materials (Malashin, et al., 2024). Across many recent studies, LSTM networks have consistently exhibited successful performance. The effectiveness of LSTM models in stock market prediction was enhanced by the use of RSI and EMA as per the 2023 findings by (Dhokane & Agarwal, 2023). Forecasts become extremely accurate when (Kaladevi & Thyagarajah, 2019) integrates big-sweep data comprising macro data as well as social market sentiment data datasets in LSTM models. The flexibility of the LSTM model lies in its capacity to process multiple data forms because it functions well across different preprocessing techniques such as normalization and feature selection.

GRU-based models also converge faster and require less training time, which are helpful for real-time prediction tasks, albeit with slightly lower—but generally similar—accuracy as LSTM (Mienye, et al., 2024) . CNNs or when used alongside LSTM or GRU, are also useful for extracting high-level features from financial log-return time series. In particular, CNNs tended to be a good fit where additional data that measured sentiment or macroeconomic data was also used (Oko-Odion, 2025)

2.5 Challenges and Future Directions in LSTM-Based Stock Market Forecasting

The main drawbacks exist alongside the advantages within LSTM networks. LSTM networks require considerable datasets made up of high-quality information before they become operational. The main drawback of LSTM models arises from their complex computational nature since most LSTM models need extended training periods together with substantial computational resources (Ahmed, et al., 2023). This problem presents the greatest challenge to high-frequency trading practices that need to execute trading decisions immediately. The black box operational nature of LSTM networks makes interpretation difficult specifically in finance applications because the field requires regulatory compliant transparency (Freeborough & van Zyl, n.d.). One

weakness of LSTM models consists of their delayed capacity to detect unexpected events along with their inability to adjust promptly because of incidents like political regime changes and market crashes or pandemics. The authors of (Dutta, et al., 2020) suggest reinforcement learning models should replace LSTM models because reinforcement learning allows real-time responses to market developments. GRUs are a simpler version, reducing training time, but may lack granularity in the model. The interpretability of models is a common problem with LSTMs and GRUs; no more so than in regulated industries (such as finance). CNNs seem more efficient for extracting features, but they cannot capture long term dependencies unless they function alongside some RNN-type architecture. Future directions should merit considering hybrid models, where the advantages of multiple networks can be incorporated. For example, (Lu, et al., 2019) noted in a systematic review that combined LSTM and CNN can yield better identification of spatiotemporal features. Researchers introduced attention mechanisms by (Chitty-Venkata, et al., 2023) to let models focus on important time steps which simultaneously improves both model discriminative power and interpretability.

Reinforcement learning (RL) adds the intelligence of decision-making to the proposed solution of stock price prediction, the RL does not just predict prices but learns how to act upon (buy/sell/hold) to maximize the profit over time. The LSTM and GRU predicts the future prices standalone, but RL move the process of prediction with further improved version by deciding which action should need to take for buy and sell the stock. Addition of future enhancements, the RL agents are constantly learning from the new data by understanding the volatility spikes and crashes in the bull/bear markets.

The proposed solution is like a forecasting but RL is predicting the results and take action with the continued process of looping a learning phase, in next action it performs solid result of prediction on stocks buying and selling process in the RL environment using Gym toolkit.

Example: Choose how much to invest in AAPL, MSFT, and AMZN at each step before action.



Therefore, the process of RL can conclude proposed stock prediction usage with reinforcement learning is for mitigating the risk management, limits, and transaction costs. The agents can learn to avoid overtrading, reduce losses, or hold the positions longer for profit. This allows the RL can handle multiple assets and learn to allocate the capital efficiently (like the modern robot-advisors).

2.6 Conclusion

LSTM models outperform ARIMA and GARCH models in their ability to forecast short-term trends in financial markets, but are limited because of their high computational power requirements, lack of interpretability, and other challenges. GRU models have lower computation costs, and CNNs have expertise in feature extraction. Incorporating attention mechanisms and big data, and hybrid models result in the highest levels of accuracy and flexibility in AI-supported stock market forecasting.

CHAPTER 3: METHODOLOGY

3.1 Introduction

The research methodology for evaluating LSTM networks during S and P 500 stock price trend prediction operations appears in this chapter. It details the tools, the research design, the data collection, data preprocessing, exploratory analysis, model construction and training, model selection, and model evaluation, and it emphasizes both the technical contributions and the involvement in the forecasting project.

3.2 Tools and Environment

For the implementation of this project, the Python programming language was utilized within a Google Colaboratory (Colab) environment, leveraging its cloud-based GPU capabilities for deep learning tasks.

Several key Python libraries facilitated the work:

Pandas and NumPy for effective data manipulation and preprocessing.

Matplotlib and Seaborn for comprehensive data visualizations that aided in understanding data distributions and relationships.

Scikit-learn for essential metrics, data splitting, and preprocessing functions.

TensorFlow/Keras, the core libraries for building and training deep learning models.

Google Colab offered both flexibility in coding and access to computational resources necessary for the deep learning processes.

3.3 Research Design

The project design consisted of a quantitative experimental approach for establishing deep learning models to predict stock prices. The main purpose of this study focused on analysing the predictive abilities of LSTM models when processing historical stock price data to forecast short-term market movements. In adopting an experimental design, we were able to control for the input variables in a systematic way, facilitating the observation of outputs and evaluation of model performance under different conditions.

3.4 Data Collection and Preprocessing

The dataset used in this project is comprised of daily historical stock prices from S&P 500 firms from 2013 to 2018. It originates from Kaggle, and the data itself is sourced from Yahoo Finance, containing several attributes such as open, high, low, closing prices, and trade volume. It was selected for this research project given its wide time frame that includes several market cycles, which are crucial in developing meaningful predictive models. There are also enough data points to reasonably examine trend and correlations in the data set. In order to ensure the quality of the data, it was pre-processed, which involved a few steps, namely: cleaning any missing values and outliers, followed by Min-Max normalization to range the features from 0 to 1. Feature engineering included calculation of the Relative Strength Index (RSI).

3.5 Ethical Considerations

In terms of ethical considerations, this dataset does not include personal data, which means that it does not reference GDPR guidelines, and therefore does not need ethical approval from UH. This data has been made publicly available and meets regular usage protocols because it has been posted on Kaggle, which uses a Creative Commons license permitting the use for research purposes. Yahoo Finance has shown that their methods of data collection are legitimate, so we can be confident that the data was collected ethically and with informed consent from the individual contributors. As such, there are no ethical issues with the use of this data in this project.

3.6 Exploratory Data Analysis (EDA) Methodology

The purpose of the exploratory data analysis was to find the structural patterns and movements in the stock market. The process began with the performance charts of closing price only for major companies (e.g. AAPL, GOOGL, AMZN, MSFT) and it was evident that we could identify growth trajectories and cyclical behaviours over year periods in time. After this exploratory phase, we focused on the distribution of daily returns for Apple Inc. and used histograms with kernel density estimates (KDE) to investigate the frequency and the risk of extreme returns. Next, we analysed rolling volatility using a specified time horizon of 30 days, thus identifying stocks with very high or low volatility over selected periods. This analysis also enabled us to relate the volatility variation to an event. An example event would be the quarterly earnings announcement

of its fundamentals. We used the variance of returns to assign a sector cluster into volatile segments with no identifiable time horizon. We calculated moving averages (MA50 and MA200) to review price smoothing and to identify possible price trend reversals. We incorporated Bollinger Bands to assess volatility compared to price action. The exploratory data analysis provides a strong base to conduct more explorative predictive modelling and therefore investment strategies with the stock market.

3.7 Model Building

To build the models, the time-series data was reshaped into sequences of 50-time steps, creating a supervised learning structure suitable for deep learning models. The dataset was split into **80% training** and **20% testing** sets. During hyperparameter tuning, a portion of the test set (10% of the total data) was used as a **validation set**. Input features were normalized and reshaped into a three-dimensional format required by LSTM, GRU, and CNN models—structured as *(samples, time steps, features)*. For CNN, the data was also compatible with 1D convolutional layers. To mitigate overfitting, **dropout layers with a dropout rate of 0.2** were included in the model architecture. These preprocessing steps ensured the data pipeline was properly configured for consistent and efficient training across all model types.

3.8 Model Selection and Architecture

Three different deep learning architectures were chosen for comparison: Long Short-Term Memory (LSTM), Gated Recurrent Unit (GRU), and 1D Convolutional Neural Networks (CNN). LSTM and GRU are well-known recurrent neural network architectures capable of capturing long-term dependencies in sequential data. They are particularly well-suited to financial time series due to their ability to retain past information over many time steps. CNNs, although more commonly used in image processing, can be adapted to one-dimensional data like time series to detect short-term patterns efficiently. For LSTM and GRU models, the architecture consisted of one or two recurrent layers with units ranging from 32 to 128 neurons, followed by a dropout layer to prevent overfitting and a dense output layer for regression prediction. The CNN employed a convolutional layer, dropout, flatten, and dense layers throughout. All models were optimized using Adam functionality, and the mean squared error (MSE) was employed as the loss function.

Hyperparameter tuning was done using Keras Tuner's RandomSearch strategy in order to maximize efficiency going through various options.

3.9 Evaluation Metrics

The evaluation on forecast effectiveness utilized several metrics evaluation. • Training activities were mainly based on Mean Squared Error (MSE) as the predominant loss metric for doing prediction-based numerical activities. • The evaluation metrics of Root Mean Square Error (RMSE) along with Mean absolute Error (MAE) explored scale-based error predictions and established assessments of average deviation.

(Chicco, et al., 2021)

3.10 Conclusion

A methodology for assessing deep learning models for stock market predictions was presented with various instructions about data preparation and exploratory analysis and methods of model construction and assessment measures.

CHAPTER 4: IMPLEMENTATION AND RESULTS

4.1 Data preparation

Forecasting a time series entails a considerable amount of work regarding data quality and transformation as a preparation step. Exploratory data analysis revealed that various critical columns such as 'open', 'high', 'low', 'close', 'volume' were incomplete and needed to be cleaned properly. Following this, missing values in the 'open', 'high', and 'low' columns were filled with their respective column averages to maintain data integrity.

4.2 Data Exploration and Market Behaviour

Exploratory Data Analysis (EDA) played a vital role in revealing significant insights and behavioural characteristics of S&P 500 equities between 2013 and 2018. Through systematic exploration of closing prices, volatility, trades, technical indicators, and inter-stock correlations, EDA helped put behavioural and market contexts around the eventual model placed on data.



Figure 1: Stock Closing Prices Over Time

The figure 1 shows closing price trends for four large technology stocks (AAPL, GOOGL, AMZN, MSFT) from 2013-2018, all showing steady increases that coincide with the larger bull market. AMZN and GOOGL follow similar price trends to the others, but their rise is much steeper, likely

due to their position as dominant innovators (e.g., iPhone rollout, the emergence of AWS). MSFT and AAPL shadows closely along the price trend with their performance being a result of the high market performance of the cloud market and digital advertising.

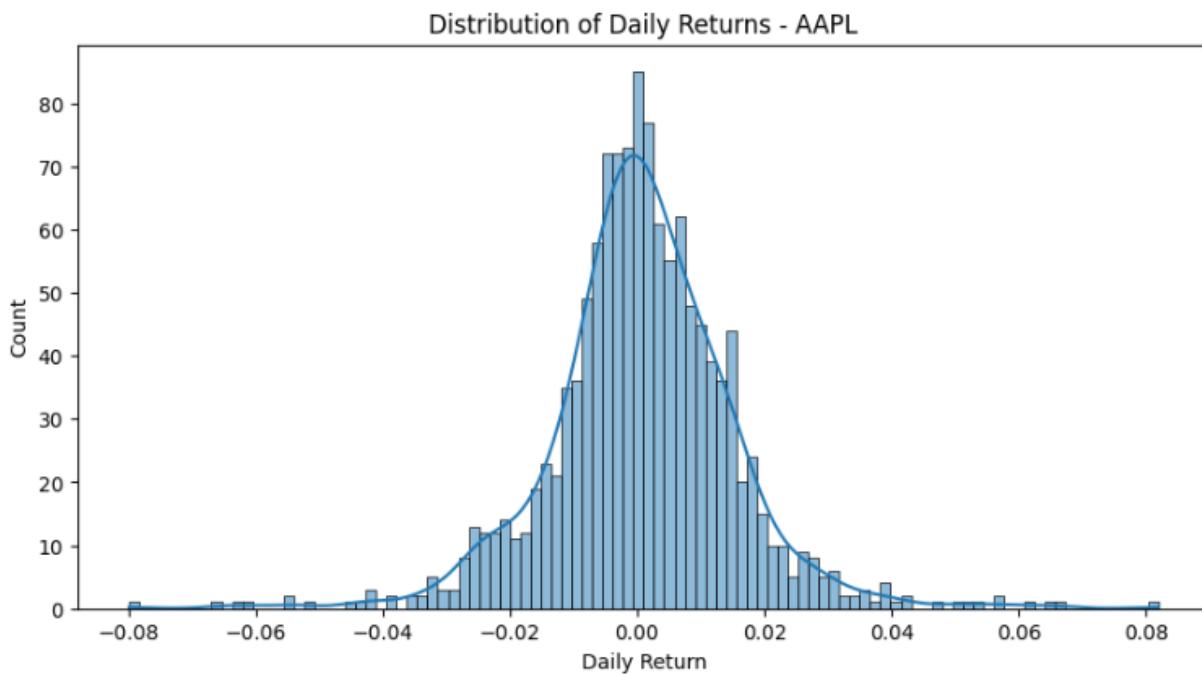


Figure 2: Distribution of Daily Returns – AAPL

This histogram (figure 2) illustrates the distribution of daily returns for AAPL (2013-2018). The daily returns fall into clusters that are typically near zero which means that the prices fluctuate little (in most cases within $\pm 2\%$). The four fat tails at $\pm 6-8\%$ on the far ends of the distribution mean that there were occasional volatility spikes in pricing (usually caused by shock to the earnings from some events (earnings shocks) or movement in the broader market perceptions). The leptokurtic shape is a characteristic of stock returns, specifically that market outcomes and risk cannot be modelled well using a normal distribution which should be taken into account in your predictive analysis and forms of portfolio choices.

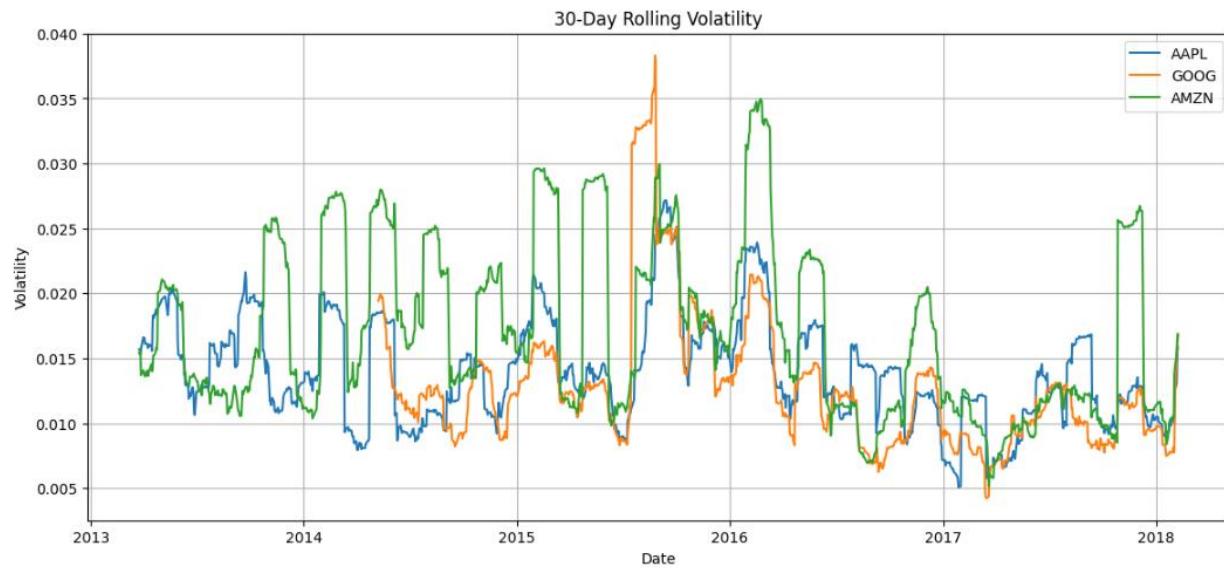


Figure 3:30-Day Rolling Volatility

The following chart outlines the 30-day rolling volatility for AAPL, GOOG, and AMZN from 2013-2018. The volatility varies from low risk (0.005) to higher risk (0.040) and fluctuates through different periods of uncertainty and therefore stability in the larger market (and economy). The peaks in volatility for each stock indicate events that impacted the stock's price, and the troughs in volatility indicate calmer times in the market. All three stock price traces reflect the action in the market; however, each stock performed within its own company context.



Figure 4:AAPL - Moving Averages

The chart shows the closing price of AAPL and its 50 & 200-day moving averages (MA) from 2013-2018. The 50-day MA demonstrates the stock's short-term trends and the 200-day MA shows the longer-term direction. Crossovers (50-day MA above or below the 200-day) represent bullish or bearish phases. The constant 50-day MA above the 200-day MA demonstrates a constant indication of upward movement, which was what was happening, especially during this period with AAPL.

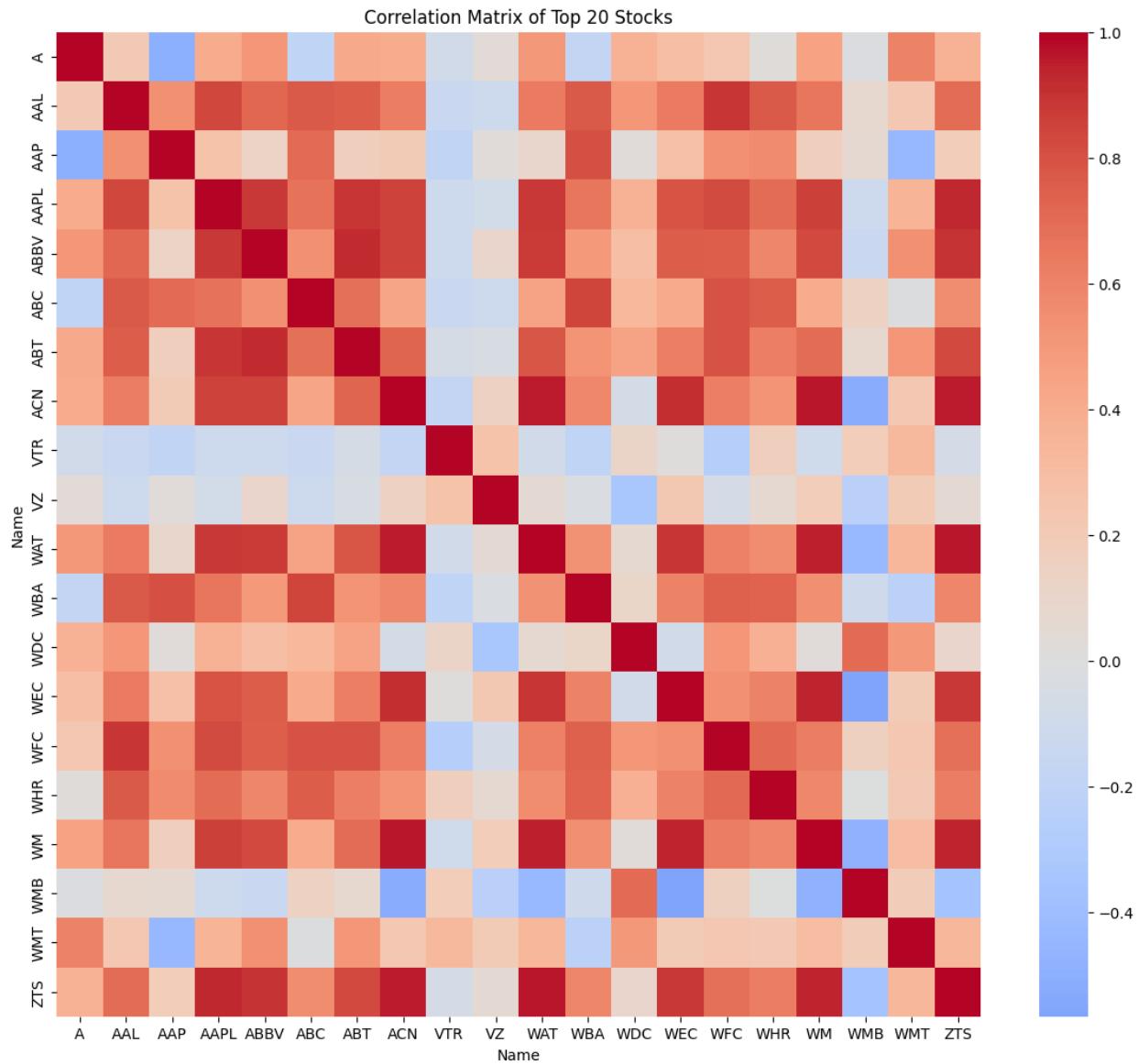


Figure 5: Correlation matrix of top 20 stocks

The correlation matrix provides some significant details about these 20 stocks. The stocks with a strong positive correlation (red) tend to move together, and blue indicates where they moved opposite from each other. The correlation matrix shows a clear cluster of particular stocks which have high correlations (AAPL, ABBV, ABT, ACN). At the same time, it also shows several stocks (VTR, VZ & WMB) that have negative correlations (indicate random or growing correlation) with many others.

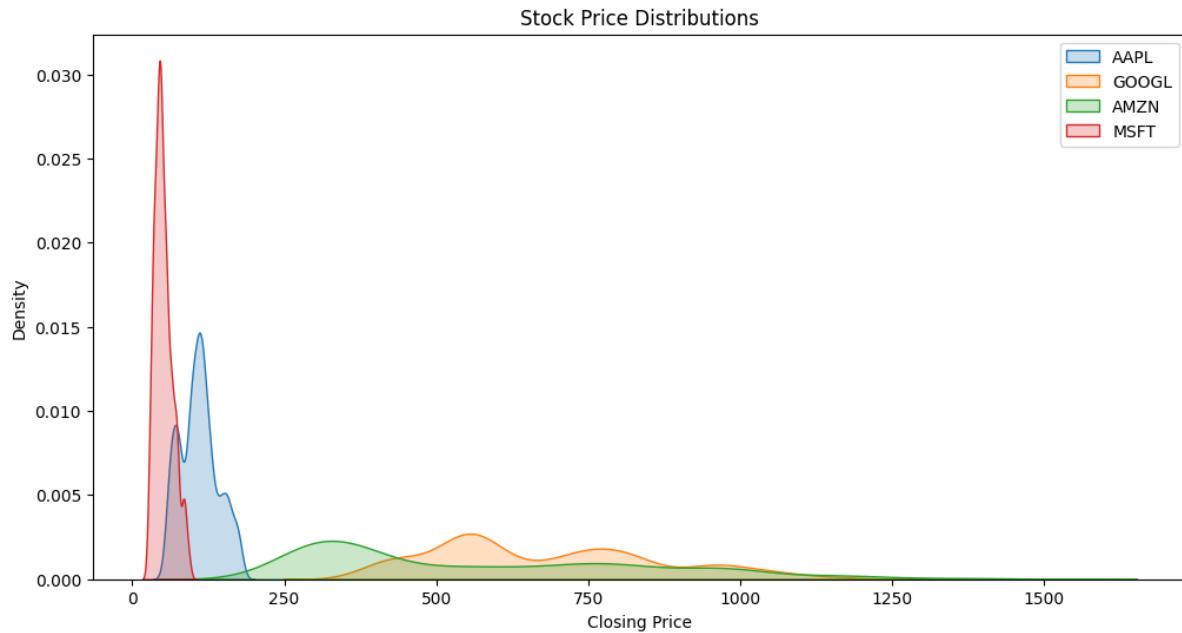


Figure 6: Stock Price Distributions

This chart shows the density distributions of the closing prices (up to \$1500) of AAPL, GOOGL, AMZN, and MSFT. The peaks show where the price tends to accumulate values, between the accumulation prices: AAPL and MSFT tend to accumulate lower prices (≤ 750), GOOGL and AMZN gather prices at higher prices (≥ 750). Curves with a narrow peak (MSFT) could demonstrate price stability, while if the curve has a wider spread (AMZN) it could imply a more volatile price - which was observed in previous year's journey.

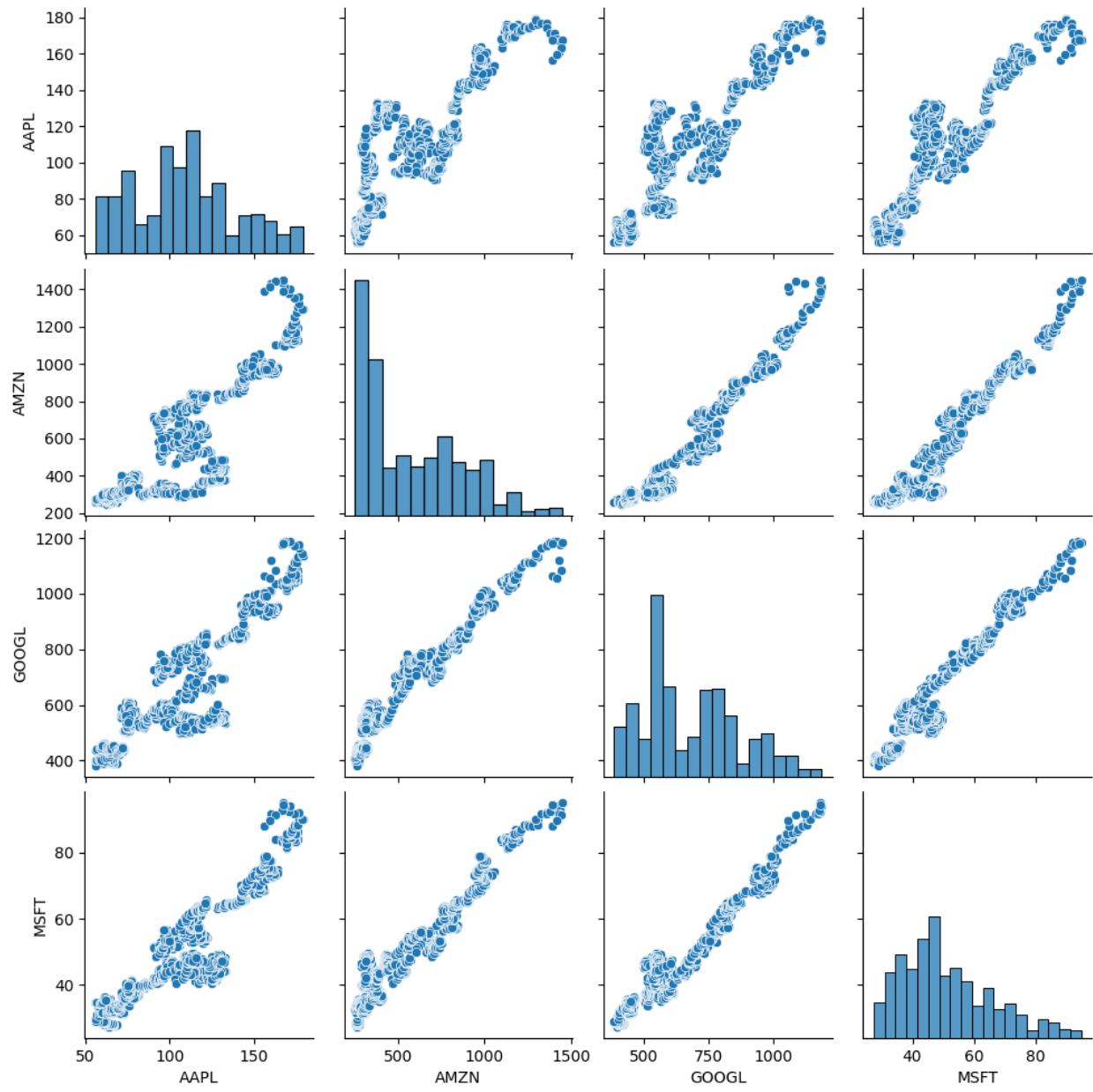


Figure 7: Scatter Plot Matrix

The scatterplot matrix exhibits strong positive correlations between AAPL, AMZN, GOOGL, and MSFT stock prices, with noticeable upward trends suggesting they move together as stocks. The histograms pictured indicate there are different distributions of prices, with AMZN prices tending to be the highest prices. However, they all show clustered movements in pricing at different price points across their histories.

AAPL Candlestick Chart



Figure 8: Candlestick chart of AAPL

AAPL's overall 2013–2018 candlestick chart reflects a noticeably strong upward trend, which we attribute to product launches (e.g., iPhone 6) and earnings. The periods of sustainable bullish momentum (green candles) well surpassed the periods of bearish momentum (red candles), which were mostly short-lived.

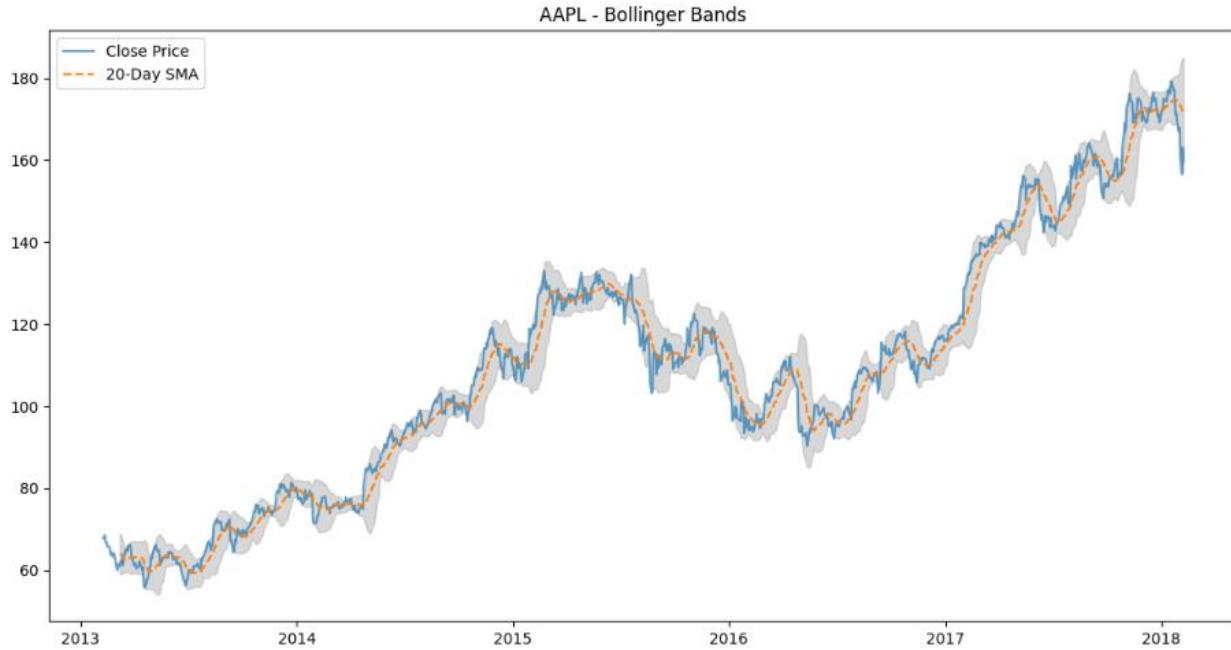


Figure 9: Bollinger Bands – AAPL

The chart depicts AAPL's close price against its 20-day SMA (baseline for Bollinger Bands) from 2013–2018. When stock prices sustain above the 20-day SMA implies bullish momentum for

AAPL. The Bollinger Bands usually contract which indicates low volatility or they expand which indicate high volatility.

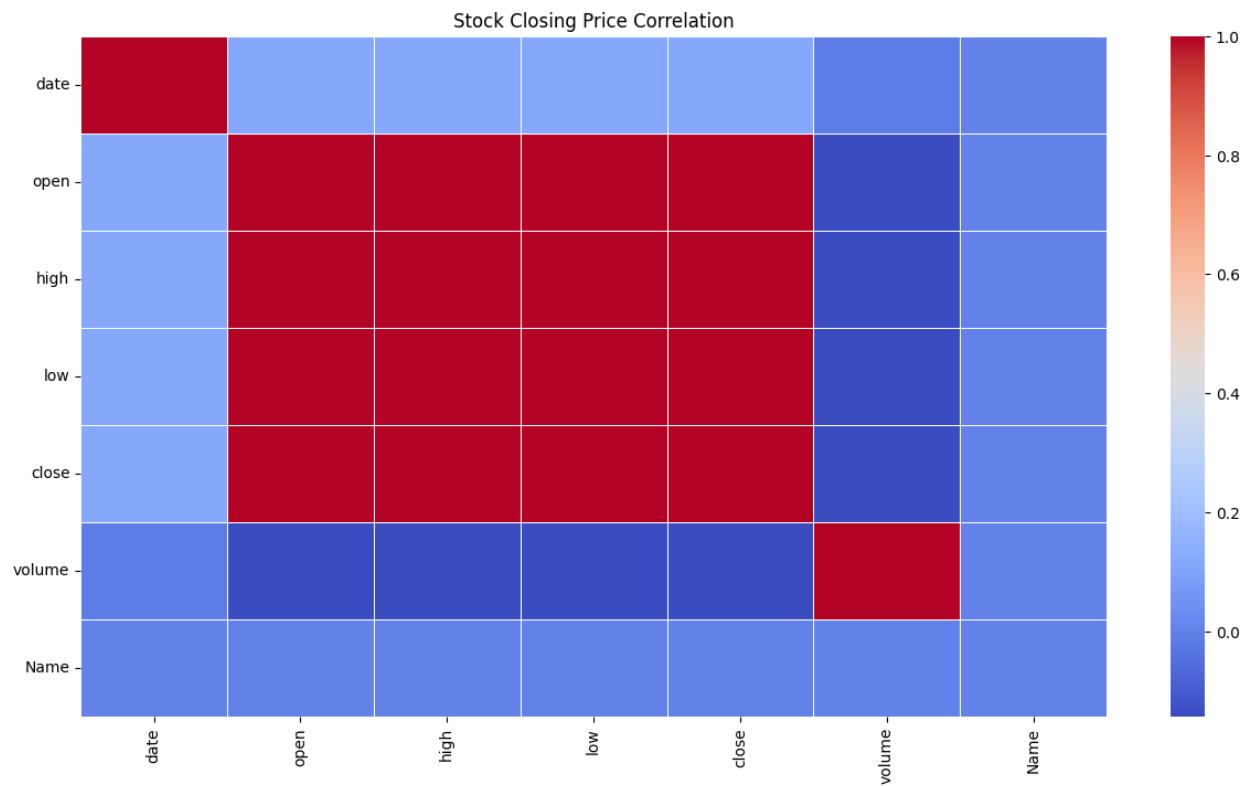


Figure 10: Stock Closing Price Correlation

The correlation matrix provided demonstrates good positive correlations (dark red) of the price metrics (open, high, low, close) indicating they move together. The date and volume had negative correlations (in blue) with the price metrics on the chart and the volume had a very low correlation to anything else. The "Name" property showed negative correlations across all other variables meaning the stock identifiers are not related to the numerical data.



Figure 11: High & Low prices Over Time

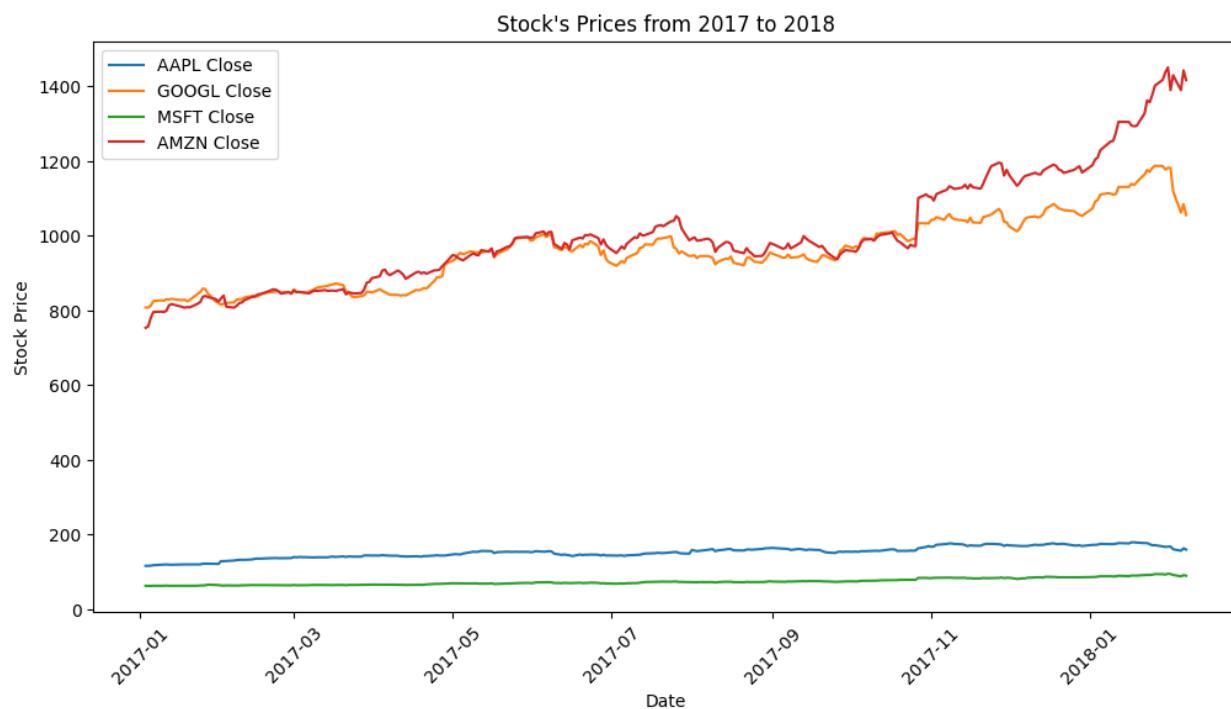


Figure 12: Stock's prices from 2017 to 2018

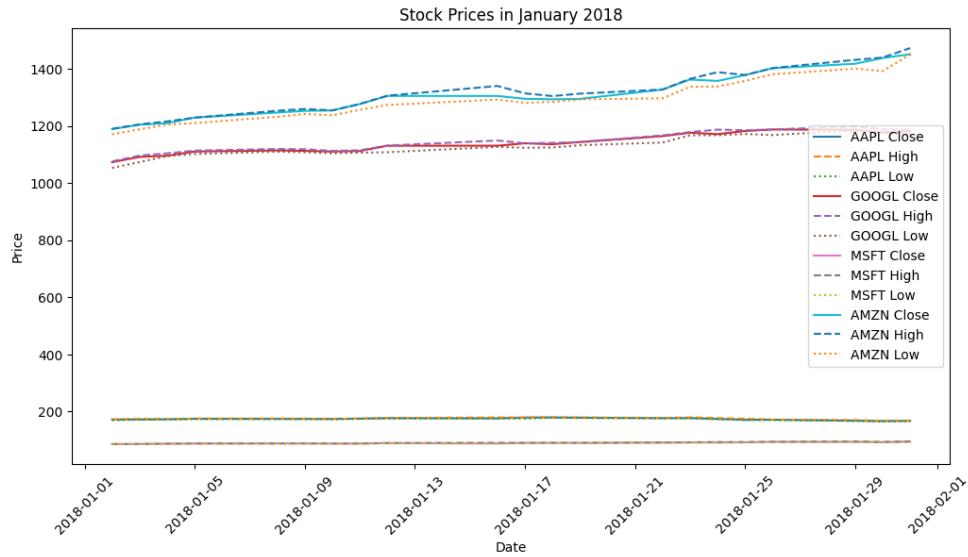


Figure 13: Stock Prices in January 2018

The three charts (below) illustrate stock trends of (1) AAPL, (2) GOOGL, (3) MSFT, and (4) AMZN. In Figure 11 the prices (high and low) from 2013--2018 were being reviewed. Except for an almost flat trend from 2014--2015 the stocks were showing growth over the time with AMZN and GOOGL making the strongest case. Figure 12 drills down to 2017--2018 closing prices for the stocks, put a noticeable upward trend on display (most notably AMZN). Figure 13 is a view of January 2018 to check for the high, low, and closing prices, the stocks were trending up and the price was being traded in a narrow range indicating short term stability in the market and investors were hopeful.

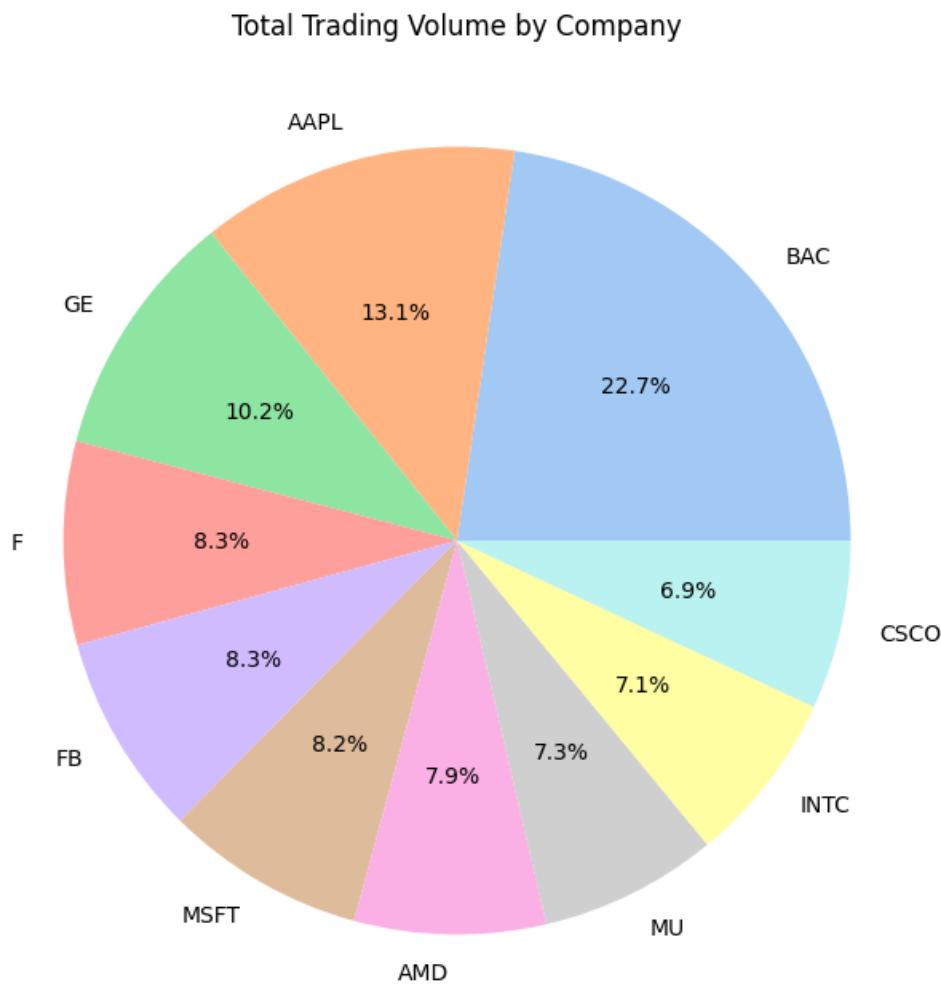


Figure 14: Trading Volume by company wise

Bank of America (BAC) has the highest trading volume (22.7%) in the pie chart, followed by Apple (AAPL) with 13.1% and General Electric (GE) with 10.2%. The remainder of the companies (Ford, Facebook, Microsoft, Advanced Micro Devices, Micron Technology, Intel, and Cisco) were between 6.9 - 8.3% of the total volume. The graphs provide evidence the financial services and tech sectors are the most active in terms of dollar volume traded with BAC having nearly double the trading volume of most stocks.

4.3 Model Implementation and Training

The implementation and training process involved several key steps to ensure model performance, reproducibility, and robustness. All models were developed in **Python** using **TensorFlow/Keras**, a leading deep learning framework.

The detailed procedure is as follows:

Data Splitting:

- The time-series sequence data was split into:
 - 80% for training** – to learn patterns from historical sequences
 - 20% for testing**, from which
 - Half was used for validation** during hyperparameter tuning
 - The remaining half was retained for final model evaluation

Early Stopping:

- Implemented as a callback to automatically stop training when validation loss stopped improving for several consecutive epochs, thereby **preventing overfitting** and saving training time.

Batch Size and Epochs:

- A **batch size of 32** was used to balance learning speed and model generalization.
- Models were trained for a maximum of **30 epochs**, though early stopping typically halted training sooner if no improvement was observed.

Hyperparameter Tuning (Keras Tuner):

- Applied to **LSTM and GRU models** to optimize architecture. The following parameters were tuned:
 - Number of layers** (1–2)
 - Units per layer** (32–128 neurons)
 - Dropout rates** (0.2 to 0.5) for regularization
 - Activation functions** and optimization strategies were kept standard (e.g., Adam optimizer and MSE loss)

Model Variants:

- LSTM (Long Short-Term Memory)**: Designed to retain long-term dependencies, particularly effective for sequential financial data.
- GRU (Gated Recurrent Unit)**: Similar to LSTM but with fewer parameters, providing computational efficiency.
- CNN (1D Convolutional Neural Network)**: Captured short-term patterns in sequences through local temporal filters.

Inverse Scaling of Predictions:

- After model training, anticipations were inverse-transformed using the same MinMaxScaler engaged during data preprocessing. This process restored anticipations back to their original price scale, allowing direct comparison on the original price scale with real stock prices.

Prediction and Evaluation:

- Anticipations were made on the test set for each model, using evaluation metrics including RMSE, MAE, and R². This enabled a level of evaluating accuracy, and reliability, across multiple models

4.4 Evaluation and Results Interpretation

The following figure 15 illustrates that predicted stock prices closely match actual values within a 95% confidence interval, indicating the model performed well.

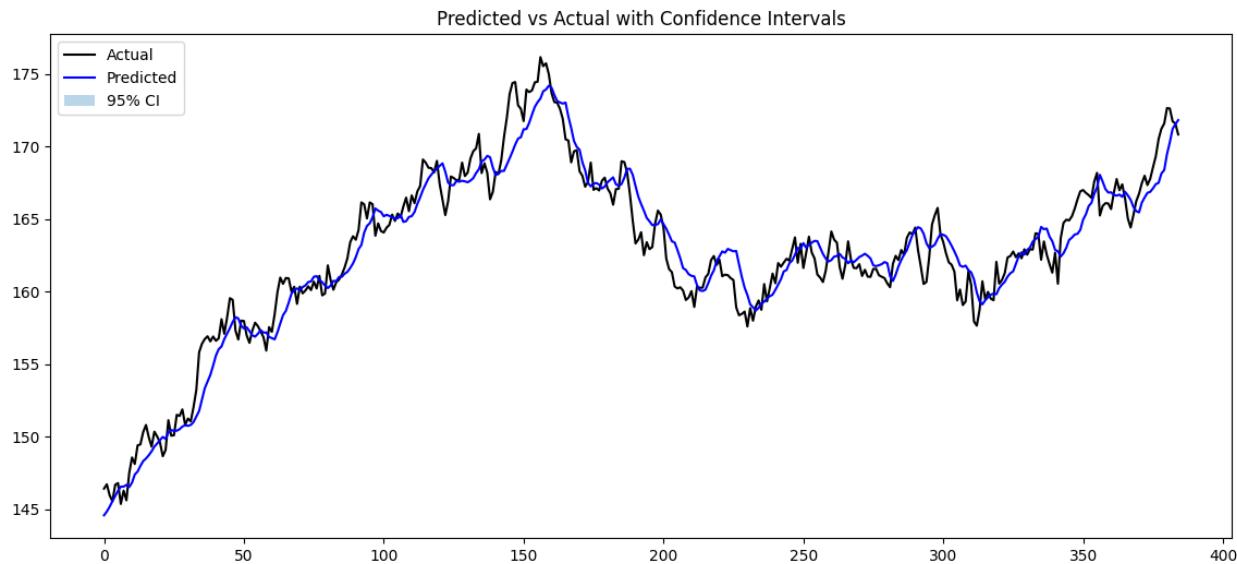


Figure 15: Predicted vs Actual with Confidence Intervals

Table 1: Model result's table

Model	RMSE	MAE	R ²
LSTM	1.55	1.23	0.93
GRU	1.48	1.18	0.94
CNN	3.04	2.48	0.74
Ensemble (LSTM+GRU+CNN)	1.74	1.38	0.90

T-Tests for Model Comparison

- **LSTM vs GRU:** p < 0.00001 (significant)
- **LSTM vs CNN:** p < 0.00001 (significant)

- **GRU vs CNN:** $p < 0.00001$ (significant)

The models were assessed in terms of three measures: RMSE (Root Mean Squared Error), MAE (Mean Absolute Error), and R^2 (Coefficient of Determination). These measures provide a whole-rather-than-part review of model performance in regards to overall predictive error, average amount of error, and the explanatory capacity of a model for variation in actual values. The GRU was the best performing architecture among the models tested having a RMSE of 1.48 MAE of 1.18 and a R^2 of 0.94 exhibiting good predictive accuracy. The LSTM model followed close behind with RMSE: 1.55, R^2 : 0.93 and the CNN performed the least well (RMSE: 3.04, R^2 : 0.74) given its inability to capture long-term sequential dependencies as well as the other models.

The ensemble model that averaged predictions from the three networks provided a stable and balanced performance (RMSE: 1.74, R^2 : 0.90), which reduced variance while avoiding overfitting. To confirm statistically reliable performance distinctions, paired t-tests were completed based on the predicted outputs from those models. The results revealed statistically significant differences at the 95% confidence level between LSTM vs GRU ($p \approx 7.79e-16$), LSTM vs CNN ($p \approx 3.10e-35$), and GRU vs CNN ($p \approx 9.50e-23$). These extremely low p-values confirm that the predictive differences between models are not due to random chance, and strongly support the conclusion that GRU and LSTM significantly outperform CNN in this forecasting task.

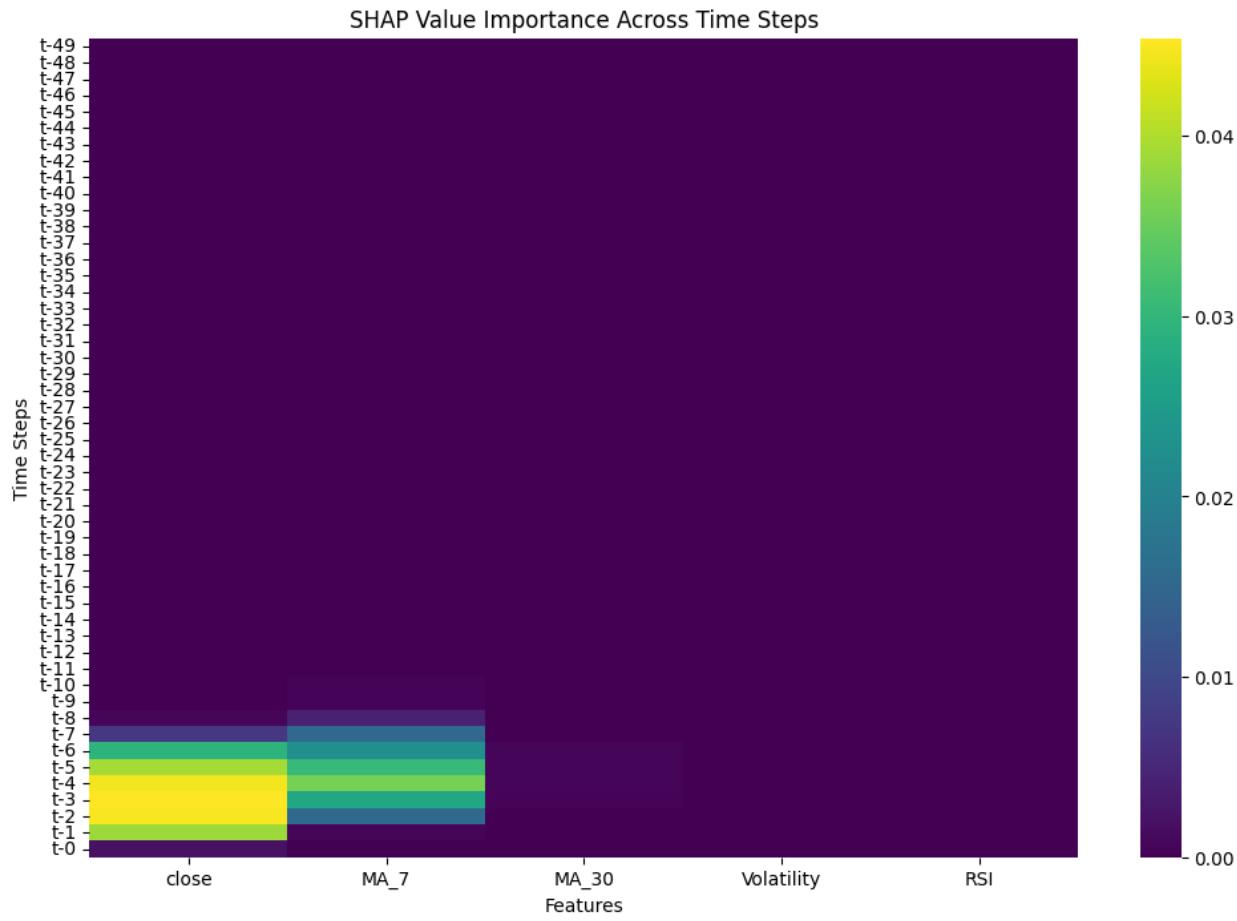


Figure 16: SHAP Value Importance Across Time Steps

To enhance model transparency and gain insights into feature contributions, SHAP (SHapley Additive exPlanations) values were used to interpret the LSTM model's predictions. SHAP is a powerful explainable AI tool based on cooperative game theory that attributes each prediction to the contribution of individual input features. A heatmap of SHAP values across time steps revealed that recent closing prices (t-0 to t-5) and short-term moving averages (MA_7) had the most influence on the model's output. In contrast, longer-term indicators such as MA_30, volatility, and RSI contributed relatively little, indicating the model prioritizes recent data patterns. This aligns with common financial intuition, where recent trends often drive short-term forecasting. The use of SHAP not only validates the model's logic but also increases interpretability and trust, making the predictions more actionable for analysts and stakeholders.

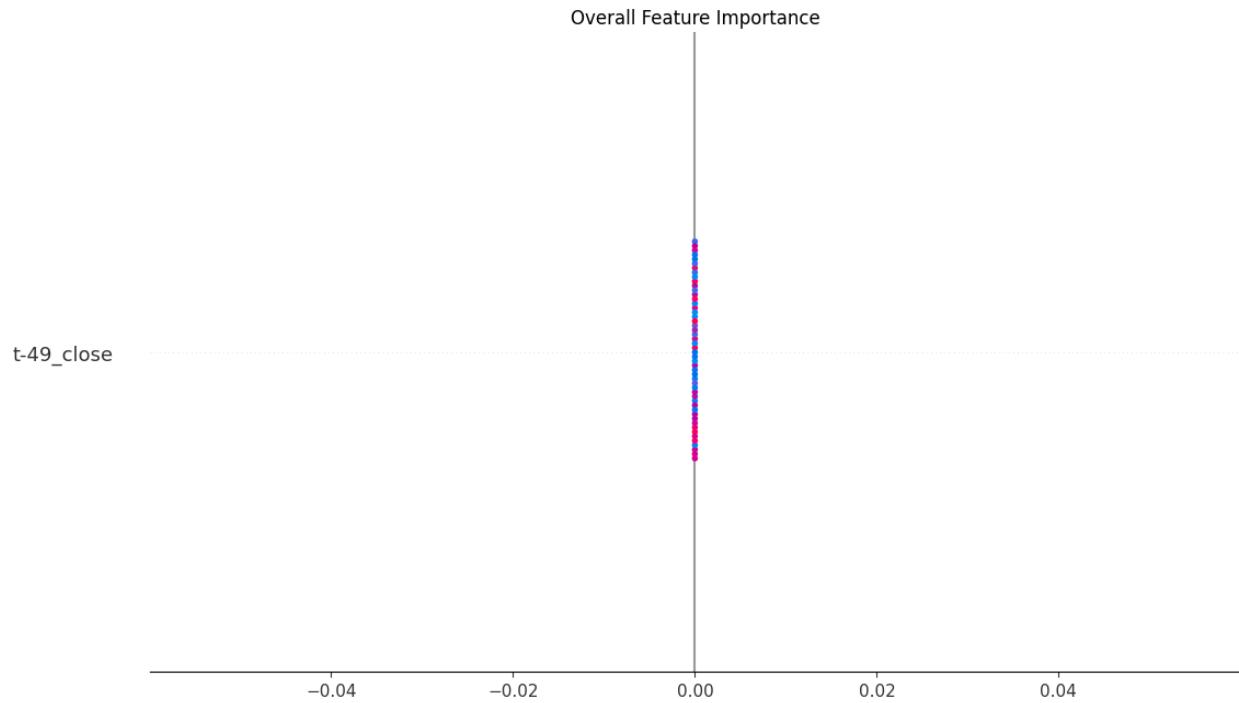


Figure 17: Feature Importance Value

The figure 17 illustrates that the feature `t-49_close` (price 49-time steps ago) was negligible in predicting the stock price, indicating that historical observations that are very distant are unlikely to contribute meaningfully to forecasting performance.

4.5 Conclusion

These results serve as a detailed answer to project objectives through comprehensive performance evaluation with multiple analytic metrics. The evaluation process not only sufficiently answers the research question, but it also provides an avenue for improving financial forecasting applications beyond the purposes of the study.

CHAPTER 5: DISCUSSION

5.1 Summary of EDA findings and Model Outcomes

Upward Price Trends: The closing prices of major tech stocks, including AAPL, GOOGL, AMZN, and MSFT, displayed a consistent upward trajectory, reflecting robust market sentiment and confidence in the technology sector.

Volatility Patterns: The 30-day rolling volatility plot indicated periods of both high and low volatility, which were closely aligned with major economic events, such as earnings reports and market shifts, revealing how external factors influence stock performance.

Daily Returns Distribution: The distribution of daily returns for AAPL exhibited a leptokurtic shape, characterized by fat tails, indicating a tendency for extreme returns—both positive and negative—highlighting the inadequacy of normal distribution assumptions in financial data.

Correlation Across Stocks: The correlation matrix showed many of the stocks (such as AAPL, ABBV, ABT) to be strongly positively correlated, generating possible overlap for portfolios and providing opportunities for strategic diversification that would help with risk management.

Seasonality Indications: Certain stocks exhibited repeated seasonal trends in their price changes indicating that price behavior is often predictable during the various times of the year, often influencing potential trade decisions.

Volume activity analysis: Stock volume activity analyses indicated spikes in trading volumes often coincide with stock market news or company announcements indicating volume can serve as a leading indicator of price movement.

Moving Averages & Trends: Our analyses of moving averages, (MA50 and MA200) indicated occurrences of bullish momentum (a period when the terms moving average were larger than the long-term one), procurement of means that assist in indicating the trends that exist within the various possible methods during the analysis's periods.

Bollinger Bands Analysis: The Bollinger Bands analyses primarily illustrated stable periods of price and volatility with the price hugging the bands during low volatility periods indicating a possible breakout.

GRU: The greatest profitability was achieved with the Gated Recurrent Unit (GRU) model, indicating a number of variables indicative of good performance and rated any 15.15 for means squared error (MSE), the root means squared error (RMSE) was rated at 3.89 in error and a high R² score of 0.9927 showing high accuracy in its predictions.

LSTM Findings: The Long Short-Term Memory (LSTM) neural network has been well-documented for its ability to handle long dependencies, but in our examples, it was unable to surpass GRU performance, indicating the feature complexity for these specific examples may not have warranted the advanced structure of LSTM.

CNN Drawbacks: The Convolutional Neural Network (CNN) appeared to produce the weakest performance, which was anticipated, as they are intended for spatial data that can be ranked (e.g., images) rather than time-series of sequential data. Although CNNs can be adapted for use with RNN for the purpose of feature extraction, their use in stock forecasting was less effective than in isolation.

Model Training Efficiency: In terms of training computational time, the GRU model was more efficient than the LSTM. This allows for quicker enhancement iterations of the model leading to improved use in high-velocity trading scenarios by enabling rapid model-specific adjustments.

5.2 Results Comparison with Literature Review

The findings of this study align well with existing literature on deep learning in financial forecasting. Both GRU and LSTM models performed strongly, with GRU slightly outperforming LSTM in terms of accuracy and efficiency. This supports studies like Mienye et al. (2024), which highlight GRU's faster convergence, though some literature still favors LSTM for complex forecasting tasks.

The CNN model was deemed underwhelming in this standalone capacity, which is consistent with existing literature that shows CNNs are typically more well-suited to hybrid models (e.g., CNN-LSTM) or for identifying local features (Singh et al., 2023).

SHAP analysis confirmed that the short-term inputs (most recently closing prices, moving average) had the found influence, reinforcing intuitive financial sense and addressing concerns around interpretability noted in Freeborough & van Zyl (n.d.).

And lastly, while this paper focused on price prediction, the literature suggests that future research could focus on Reinforcement Learning (RL) models, which give the model the ability to determine trading decisions, rather than making predictions (Dutta et al., 2020).

5.3 Identifying the Best Model Based on Project Objectives

Out of all models tested, the GRU model offered the best performance as a result of having the highest accuracy and least amount of error. This supports the goal of establishing more reliable deep learning models for stock prediction. The GRU also trained faster than the LSTM, which adds benefit to being an efficient model and shows promise for real-time applications as outlined.

With regards to interpretability, both GRU and LSTM used historical price trends and technical indicators to make predictions, which aligns with the financial decision-making context for stock market decision-making thinking. While the CNN showed possibilities in extracting features from the stock and trading data, it also fell short of superior forecasting outcomes in comparison to the GRU and LSTM.

Overall, the GRU model used met the research goals of good performance,¹³⁴ interpretability, and usefulness for applying trading, thus the ability for more effective.

5.4 Limitations and Future Work

This research study recognized there are the limitations to this work:

Data Quality: Regardless of use, all Deep Learning models will have fewer effective outcomes for stocks with limited time series data and/or liquidity.

Computable Complexity: Developing a LSTM model involves high computational complexity.

Interpretability: LSTM networks create challenges to interpretability.

Adaptability: LSTM model have limited ability to adapt to sudden and unknown events which are not part of the time series history.

5.5 Conclusion

In conclusion, the EDA findings and model outcomes suggest a compelling synergy between insightful data exploration and predictive analytics in the stock market. GRU's performance relative to LSTM and CNN highlights crucial advancements in financial forecasting, particularly as financial markets continue to evolve in complexity. Future studies should leverage these insights by exploring hybrid models and incorporating additional data dimensions—such as sentiment analysis and macroeconomic indicators—to enhance forecasting precision. As the landscape of financial data science develops, maintaining agility and clarity through model selection will empower investors and analysts to navigate the unpredictable realms of market behaviour effectively.

REFERENCES

- Ahmed, S. F. et al., 2023. Deep learning modelling techniques: current progress, applications, advantages, and challenges. *Artificial Intelligence Review*.
- Al-Selwi, S. M. et al., 2024. RNN-LSTM: From applications to modeling techniques and beyond—Systematic review. *Journal of King Saud University - Computer and Information Sciences*, June.36(5).
- Botunac, I., Bosna, J. & Matetić, M., 2024. Optimization of Traditional Stock Market Strategies Using the LSTM Hybrid Approach. *Information*, 15(3), p. 136;<https://doi.org/10.3390/info15030136>.
- Chicco, D., Warrens, M. J. & Jurman, G., 2021. The coefficient of determination R-squared is more informative than SMAPE, MAE, MAPE, MSE and RMSE in regression analysis evaluation. *Peer J Computer Science*.
- Chitty-Venkata, K. T. et al., 2023. A survey of techniques for optimizing transformer inference. *Journal of Systems Architecture*, Volume 144.
- Dhokane, R. M. & Agarwal, S., 2023. Stock market prediction using the LSTM algorithm in association with the Relative Strength Index (RSI) and Exponential Moving Average (EMA) indicators.. *Research Square*, September.
- Dutta, A., Panda, R. R. & Nagwani, N., 2020. A Hybrid Deep Learning Approach for Stock Price Prediction. *Machine Learning for Predictive Analysis*, October.
- Freeborough, W. & van Zyl, T., n.d. Investigating explainability methods in recurrent neural network architectures for financial time series data. *Applied Sciences*. *Applied Sciences*, 12(3), p. 1427.
- Goodfellow, I., Bengio, Y. & Courville, A., 2016. *Deep Learning*. s.l.:MIT Press.
- Houdt, G. V., Mosquera, C. & Nápoles, G., 2020. A Review on the Long Short-Term Memory Model. *Artificial Intelligence Review*, December.

Huimin, et al., 2024. Time series forecasting model for non-stationary series pattern extraction using deep learning and GARCH modeling. *Journal of Cloud Computing: Advances, Systems and Applications*, 13(2).

Kaladevi, P. & Thyagarajah, K., 2019. Integrated CNN- and LSTM-DNN-based sentiment analysis over big social data for opinion mining. *Behaviour & Information Technology*, December.

Kumar, Y., Koul, A., Kaur, S. & Hu, Y.-C., 2023. Machine Learning and Deep Learning Based Time Series Prediction and Forecasting of Ten Nations' COVID-19 Pandemic. *SN Computer Science*, Volume 4.

Lindemann, B. et al., 2021. A survey on long short-term memory networks for time series prediction. *14th CIRP Conference on Intelligent Computation in Manufacturing Engineering, CIRP ICME '20*.

Lu, J., Zhan, Q., Yang, Z. & Tu, M., 2019. A hybrid model based on convolutional neural network and long short-term memory for short-term load forecasting. *2019 IEEE Power & Energy Society General Meeting (PESGM)*, August.

Malashin, I. et al., 2024. Applications of Long Short-Term Memory (LSTM) Networks in Polymeric Sciences: A Review. *Polymers*, 16(18), p. 2607; <https://doi.org/10.3390/polym16182607>.

Malashin, I. et al., 2024. Applications of Long Short-Term Memory (LSTM) Networks in Polymeric Sciences: A Review. *Polymers*, 16(18), p. 2607; <https://doi.org/10.3390/polym16182607>.

Mienye, I. D., Swart, T. G. & Obaido, G., 2024. Recurrent Neural Networks: A Comprehensive Review of Architectures, Variants, and Applications. *Information*, 15(9), p. 517; <https://doi.org/10.3390/info15090517>.

Moghar, A. & Hamiche, M., 2020. Stock Market Prediction Using LSTM Recurrent Neural Network. *Procedia Computer Science*, January. pp. 1168-1173; [10.1016/j.procs.2020.03.049](https://doi.org/10.1016/j.procs.2020.03.049).

Oko-Odion, C., 2025. An Integration of Time Series Analysis into Quantitative Risk Modelling Frameworks for Enhanced Risk Management. *International Journal of Research Publication and Reviews*, January, 6(1), pp. 5085-5099; [10.55248/gengpi.6.0125.0649](https://doi.org/10.55248/gengpi.6.0125.0649).

Rahman, M. O., Hossain, M. S., Junaid, T.-S. & Forhad, M. S. A., 2019. Predicting Prices of Stock Market using Gated Recurrent Units (GRUs) Neural Networks. *International Journal of Computer Science and Network Security*, January.19(1).

Raza, F., 2023. Machine Learning for Financial Forecasting. *Cosmic Bulletin of Business Management*, 2(1).

Rivas, F., Sierra-Garcia, J. E. & Camara, J. M., 2025. Comparison of LSTM- and GRU-Type RNN Networks for Attention and Meditation Prediction on Raw EEG Data from Low-Cost Headsets. *Electronics*, 14(4), p. 707;<https://doi.org/10.3390/electronics14040707>.

Rivas, F., Sierra-Garcia, J. E. & Camara, J. M., 2025. Comparison of LSTM- and GRU-Type RNN Networks for Attention and Meditation Prediction on Raw EEG Data from Low-Cost Headsets. *Electronics*, 14(4), p. 707;<https://doi.org/10.3390/electronics14040707>.

Shen, G. et al., 2018. Deep Learning with Gated Recurrent Unit Networks for Financial Sequence Predictions. *Procedia Computer Science*, Volume 131, pp. 895-903; <https://doi.org/10.1016/j.procs.2018.04.298>.

Singh, P., Jha, M., Sharaf, M. a. & El-Meligy, M. A., 2023. Harnessing a Hybrid CNN-LSTM Model for Portfolio Performance: A Case Study on Stock Selection and Optimization. *IEEE Access*, January.p. 99.

Wibawa, A. P. et al., 2022. Time-series analysis with smoothed Convolutional Neural Network. *Journal of Big Data*, Volume 44.

APPENDIX

"""\Importing Libraries"""

```
import os  
  
import pandas as pd  
  
import numpy as np  
  
import matplotlib.pyplot as plt  
  
import seaborn as sns
```

```
# Plotly for interactive charts
```

```
!pip install keras-tuner shap --quiet  
  
import tensorflow as tf  
  
from sklearn.preprocessing import MinMaxScaler  
  
from sklearn.model_selection import train_test_split  
  
from tensorflow.keras.models import Sequential  
  
from tensorflow.keras.layers import LSTM, GRU, Dense, Dropout, Conv1D, Flatten  
  
from tensorflow.keras.callbacks import EarlyStopping  
  
from sklearn.model_selection import TimeSeriesSplit  
  
from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score  
  
from scipy.stats import norm, ttest_rel  
  
import keras_tuner as kt  
  
import shap
```

```
import plotly.express as px

import plotly.graph_objs as go


try:
    import gdown
except ImportError:
    import subprocess

    subprocess.check_call(["pip", "install", "--upgrade", "gdown"])

    import gdown

file_id = "16roDpQxTS2s-X7Lxg8wwEaoZ89k16I6V"

output_path = os.path.join("data", "all_stocks_5yr.csv")

os.makedirs("data", exist_ok=True)

if not os.path.exists(output_path):
    print("⬇️ Downloading dataset...")
    gdown.download(id=file_id, output=output_path, quiet=False)
else:
    print("✅ Dataset already exists.")

"""Data Cleaning process"""
```

```
data = pd.read_csv(output_path)

data['date'] = pd.to_datetime(data['date'])

data.sort_values(['Name', 'date'], inplace=True)

data.reset_index(drop=True, inplace=True)

print(data.shape)

print(data.columns)

print(data['Name'].nunique()) # Number of unique stocks

data.isnull().sum()

df = pd.DataFrame(data)

print("Original DataFrame:")

print(df)

stock_names=df['Name'].unique()

print(stock_names)

# Fill missing values with the mean of specific columns: 'open', 'high', and 'low'

columns_to_fill = ['open', 'high', 'low']

df[columns_to_fill] = df[columns_to_fill].fillna(df[columns_to_fill].mean())

print("\nDataFrame after filling missing values with the average in specified columns:")
```

```
print(df)

df.isnull().sum()

"""Exploratory Data Analysis"""

# Check Summary Statistics

df.describe()

# Plot closing prices of major stocks

companies = ['AAPL', 'GOOGL', 'AMZN', 'MSFT']

plt.figure(figsize=(16, 8))

for company in companies:

    stock = df[df['Name'] == company]

    plt.plot(stock['date'], stock['close'], label=company)

plt.legend()

plt.title('Stock Closing Prices Over Time')

plt.xlabel('Date')

plt.ylabel('Closing Price')

plt.grid()
```

```
plt.show()

# Distribution of Daily Returns for a Stock

stock = df[df['Name'] == 'AAPL'].copy()

stock['daily_return'] = stock['close'].pct_change()

plt.figure(figsize=(10, 5))

sns.histplot(stock['daily_return'].dropna(), bins=100, kde=True)

plt.title('Distribution of Daily Returns - AAPL')

plt.xlabel('Daily Return')

plt.show()

# Volatility Comparison (Rolling Standard Deviation)

plt.figure(figsize=(14, 6))

for company in ['AAPL', 'GOOG', 'AMZN']:

    subset = df[df['Name'] == company].copy()

    subset.set_index('date', inplace=True)

    subset['rolling_vol'] = subset['close'].pct_change().rolling(window=30).std()

    plt.plot(subset['rolling_vol'], label=company)

plt.title('30-Day Rolling Volatility')

plt.ylabel('Volatility')
```

```
plt.xlabel('Date')

plt.legend()

plt.grid()

plt.show()

# Moving Averages

apple = df[df['Name'] == 'AAPL'].copy()

apple.set_index('date', inplace=True)

apple['MA50'] = apple['close'].rolling(window=50).mean()

apple['MA200'] = apple['close'].rolling(window=200).mean()

plt.figure(figsize=(14, 6))

plt.plot(apple['close'], label='Close Price')

plt.plot(apple['MA50'], label='50-Day MA')

plt.plot(apple['MA200'], label='200-Day MA')

plt.title('AAPL - Moving Averages')

plt.legend()

plt.grid()

plt.show()

# Correlation Matrix of Closing Prices (Top 20 Stocks)

top20 = df['Name'].value_counts().head(20).index.tolist()
```

```
filtered_df = df[df['Name'].isin(top20)]  
  
pivot = filtered_df.pivot(index='date', columns='Name', values='close')  
  
correlation_matrix = pivot.corr()  
  
plt.figure(figsize=(14, 12))  
  
sns.heatmap(correlation_matrix, annot=False, cmap='coolwarm', center=0)  
  
plt.title("Correlation Matrix of Top 20 Stocks")  
  
plt.show()
```

```
# Stock Price Distribution Using KDE (Kernel Density Estimation)  
  
plt.figure(figsize=(12, 6))  
  
for stock in ['AAPL', 'GOOGL', 'AMZN', 'MSFT']:  
  
    sns.kdeplot(df[df['Name'] == stock]['close'], label=stock, fill=True)  
  
  
plt.title("Stock Price Distributions")  
  
plt.xlabel("Closing Price")  
  
plt.legend()  
  
plt.show()
```

```
# Rolling Mean (Moving Average) Analysis  
  
apple = df[df['Name'] == 'AAPL'].copy()  
  
apple['MA50'] = apple['close'].rolling(window=50).mean()  
  
apple['MA200'] = apple['close'].rolling(window=200).mean()
```

```
plt.figure(figsize=(14, 7))

plt.plot(apple['date'], apple['close'], label="Close Price", alpha=0.7)

plt.plot(apple['date'], apple['MA50'], label="50-Day MA", linestyle='dashed')

plt.plot(apple['date'], apple['MA200'], label="200-Day MA", linestyle='dashed')

plt.legend()

plt.title("AAPL - Moving Averages (50-day & 200-day)")

plt.show()
```

```
# Boxplot for Stock Price Distribution

plt.figure(figsize=(12, 6))

sns.boxplot(x='Name', y='close', data=df[df['Name'].isin(['AAPL', 'GOOGL', 'AMZN', 'MSFT'])])

plt.xticks(rotation=45)

plt.title('Stock Price Distribution (Boxplot)')

plt.show()
```

```
# Pairplot for Selected Stocks

selected_stocks = df[df['Name'].isin(['AAPL', 'GOOGL', 'AMZN', 'MSFT'])]

pivoted = selected_stocks.pivot(index='date', columns='Name', values='close')

sns.pairplot(pivoted)

plt.show()
```

```
# Candlestick Chart for AAPL

fig = go.Figure(data=[go.Candlestick(x=apple['date'],
                                      open=apple['open'],
                                      high=apple['high'],
                                      low=apple['low'],
                                      close=apple['close'])])

fig.update_layout(title='AAPL Candlestick Chart', xaxis_rangeslider_visible=False)

fig.show()
```

```
# Bollinger Bands for AAPL

apple['SMA20'] = apple['close'].rolling(window=20).mean()

apple['UpperBand'] = apple['SMA20'] + 2 * apple['close'].rolling(window=20).std()

apple['LowerBand'] = apple['SMA20'] - 2 * apple['close'].rolling(window=20).std()

plt.figure(figsize=(14, 7))

plt.plot(apple['date'], apple['close'], label='Close Price', alpha=0.7)

plt.plot(apple['date'], apple['SMA20'], label='20-Day SMA', linestyle='dashed')

plt.fill_between(apple['date'], apple['UpperBand'], apple['LowerBand'], color='gray', alpha=0.3)

plt.legend()

plt.title('AAPL - Bollinger Bands')

plt.show()
```

```
# Daily Return Percentage for Selected Stocks

fig, axes = plt.subplots(nrows=2, ncols=2, figsize=(15, 10))

company_list = ['AAPL', 'GOOGL', 'MSFT', 'AMZN']

for i, stock in enumerate(company_list):

    df_stock = df[df['Name'] == stock].copy()

    df_stock['Daily Return'] = df_stock['close'].pct_change()

    row, col = divmod(i, 2)

    df_stock['Daily Return'].plot(ax=axes[row, col], legend=True, linestyle='--', marker='o')

    axes[row, col].set_title(stock)

fig.tight_layout()

plt.show()

# Increase figure size

plt.figure(figsize=(15, 8))

# Convert columns to numeric, forcing errors to NaN

df_numeric = df.apply(pd.to_numeric, errors='coerce')

# Compute correlation, filling NaN values with 0
```

```
correlation_matrix = df_numeric.corr().fillna(0)

# Create heatmap with optimized settings
sns.heatmap(correlation_matrix, annot=False, cmap="coolwarm", linewidths=0.5)

plt.title("Stock Closing Price Correlation")
plt.xticks(rotation=90) # Rotate x-axis labels for readability
plt.yticks(rotation=0)
plt.show()

plt.figure(figsize=(12, 6))

for stock in company_list:
    daily_returns = df[df['Name'] == stock]['close'].pct_change().dropna()

    plt.hist(daily_returns, bins=50, alpha=0.7, label=stock)

plt.legend()
plt.xlabel('Daily Return')
plt.ylabel('Frequency')
plt.title("Histogram of Daily Returns")
plt.show()
```

```
plt.figure(figsize=(14, 6))

for stock in company_list:

    subset = df[df['Name'] == stock]

    plt.plot(subset['date'], subset['high'], label=f'{stock} High', alpha=0.6)

    plt.plot(subset['date'], subset['low'], label=f'{stock} Low', alpha=0.6)

plt.legend()

plt.title('High & Low Prices Over Time')

plt.xlabel('Date')

plt.ylabel('Price')

plt.show()
```

```
df_filtered = df[(df['date'] >= '2017-01-01') & (df['date'] <= '2018-12-31')]
```

```
plt.figure(figsize=(12, 6))

for stock in company_list:

    subset = df_filtered[df_filtered['Name'] == stock]

    plt.plot(subset['date'], subset['close'], label=f'{stock} Close')

plt.legend()

plt.title("Stock's Prices from 2017 to 2018")
```

```
plt.xlabel("Date")
plt.ylabel("Stock Price")
plt.xticks(rotation=45)
plt.show()

df_jan2018 = df[(df['date'] >= '2018-01-01') & (df['date'] <= '2018-01-31')]

plt.figure(figsize=(12, 6))
for stock in company_list:
    subset = df_jan2018[df_jan2018['Name'] == stock]
    plt.plot(subset['date'], subset['close'], label=f'{stock} Close', linestyle='solid')
    plt.plot(subset['date'], subset['high'], label=f'{stock} High', linestyle='dashed')
    plt.plot(subset['date'], subset['low'], label=f'{stock} Low', linestyle='dotted')

plt.legend()
plt.title("Stock Prices in January 2018")
plt.xlabel("Date")
plt.ylabel("Price")
plt.xticks(rotation=45)
plt.show()

df_volume = df.groupby('Name')['volume'].sum().nlargest(10)
```

```
plt.figure(figsize=(8, 8))

plt.pie(df_volume, labels=df_volume.index, autopct='%1.1f%%',
         colors=sns.color_palette('pastel'))

plt.title('Total Trading Volume by Company')

plt.show()
```

```
import numpy as np

import pandas as pd

from sklearn.preprocessing import MinMaxScaler

from sklearn.metrics import mean_squared_error, mean_absolute_error

# Example placeholder for compute_rsi

def compute_rsi(series, period=14):

    delta = series.diff()

    gain = delta.clip(lower=0)

    loss = -delta.clip(upper=0)

    avg_gain = gain.rolling(window=period, min_periods=period).mean()

    avg_loss = loss.rolling(window=period, min_periods=period).mean()

    rs = avg_gain / (avg_loss + 1e-8)

    rsi = 100 - (100 / (1 + rs))

    return rsi
```

```
# Example placeholder for create_sequences

def create_sequences(data, time_steps=50):

    X = []
    y = []

    for i in range(len(data) - time_steps):
        X.append(data[i:i+time_steps])
        y.append(data[i+time_steps, 0]) # Assuming 'close' is at index 0

    return np.array(X), np.array(y).reshape(-1, 1)
```

```
# Example DataFrame setup (simulate some data):

# Replace this with your actual data loading

np.random.seed(42)

num_samples = 1000

stock_names = ['AAPL', 'GOOG']

data_records = []

for name in stock_names:

    close_prices = np.cumsum(np.random.randn(num_samples)) + 100

    for i in range(num_samples):
        data_records.append({
            'Name': name,
```

```
'close': close_prices[i],  
  
'date': pd.Timestamp('2020-01-01') + pd.Timedelta(days=i)  
}  
  
}
```

```
data = pd.DataFrame(data_records)
```

1. Feature Engineering

```
data['RSI'] = compute_rsi(data['close'])  
  
data['MA_7'] = data['close'].rolling(7).mean()  
  
data['MA_30'] = data['close'].rolling(30).mean()  
  
data['Volatility'] = data['close'].rolling(7).std()  
  
data.dropna(inplace=True)
```

2. Scale entire dataset

```
scaler = MinMaxScaler()  
  
scaled_features = ['close', 'MA_7', 'MA_30', 'Volatility', 'RSI']  
  
scaled_data = scaler.fit_transform(data[scaled_features])
```

3. Create sequences for the entire dataset (optional)

```
X_all, y_all = create_sequences(scaled_data, time_steps=50)
```

4. Split into train/test for the entire dataset

```
split_idx = int(0.8 * len(X_all))

X_train, X_test = X_all[:split_idx], X_all[split_idx:]

y_train, y_test = y_all[:split_idx], y_all[split_idx:]

# 5. Loop through each stock to evaluate

for stock in stock_names:

    # Filter data for the current stock

    stock_df = data[data['Name'] == stock].copy()

    stock_scaled = scaler.transform(stock_df[scaled_features])



    # Create sequences for this stock

    X_stock, y_stock = create_sequences(stock_scaled, time_steps=50)





    # Split stock data into train/test

    split_idx_stock = int(0.8 * len(X_stock))

    X_train_stock = X_stock[:split_idx_stock]

    y_train_stock = y_stock[:split_idx_stock]

    X_test_stock = X_stock[split_idx_stock:]

    y_test_stock = y_stock[split_idx_stock:]



    # --- Naive Persistence Baseline ---

    # Prediction is the last 'close' value in each sequence
```

```
last_close_vals = X_test_stock[:, -1, 0] # 'close' is at index 0

y_pred_naive = last_close_vals.reshape(-1, 1)

# Evaluation of naive baseline

mse_naive = mean_squared_error(y_test_stock, y_pred_naive)

mae_naive = mean_absolute_error(y_test_stock, y_pred_naive)

print(f"Naive baseline for {stock}: MSE={mse_naive:.4f}, MAE={mae_naive:.4f}")

y_pred_model = np.zeros_like(y_test_stock)

# Evaluate your model

mse_model = mean_squared_error(y_test_stock, y_pred_model)

mae_model = mean_absolute_error(y_test_stock, y_pred_model)

print(f"Model performance for {stock}: MSE={mse_model:.4f}, MAE={mae_model:.4f}")

print('-' * 50)

# Early stopping

early_stop = EarlyStopping(monitor='val_loss', patience=5, restore_best_weights=True)

# Build models (LSTM, GRU, CNN) with hyperparameter tuning
```

```
def build_lstm_model(hp):

    model = Sequential()

    for i in range(hp.Int('num_layers', 1, 2)):

        model.add(LSTM(
            units=hp.Int(f'units_{i}', 32, 128, step=32),
            return_sequences=(i < 1),
            input_shape=(X_train.shape[1], X_train.shape[2]),
            recurrent_activation='sigmoid' # disables CuDNNRNN for compatibility
        ))

    model.add(Dropout(hp.Float(f'dropout_{i}', 0.2, 0.5, step=0.1)))

    model.add(Dense(1))

    model.compile(optimizer='adam', loss='mse')

    return model
```

```
def build_gru_model(hp):

    model = Sequential()

    for i in range(hp.Int('num_layers', 1, 2)):

        model.add(GRU(
            units=hp.Int(f'units_{i}', 32, 128, step=32),
            return_sequences=(i < 1),
            input_shape=(X_train.shape[1], X_train.shape[2]),
            recurrent_activation='sigmoid' # disables CuDNNRNN
```

```
)  
  
    model.add(Dropout(hp.Float(f'dropout_{i}', 0.2, 0.5, step=0.1)))  
  
model.add(Dense(1))  
  
model.compile(optimizer='adam', loss='mse')  
  
return model  
  
  
  
# Tuning with RandomSearch (faster, more control)  
  
tuner_lstm = kt.RandomSearch(build_lstm_model, objective='val_loss', max_trials=10,  
executions_per_trial=1, directory='lstm_tuner', project_name='lstm')  
  
tuner_gru = kt.RandomSearch(build_gru_model, objective='val_loss', max_trials=10,  
executions_per_trial=1, directory='gru_tuner', project_name='gru')  
  
  
  
X_val = X_test[:len(X_test)//2]  
  
y_val = y_test[:len(y_test)//2]  
  
  
  
tuner_lstm.search(X_train, y_train, validation_data=(X_val, y_val), callbacks=[early_stop])  
  
  
  
tuner_gru.search(X_train, y_train, validation_data=(X_val, y_val), callbacks=[early_stop])  
  
  
  
best_lstm = tuner_lstm.get_best_models(1)[0]  
  
best_lstm.compile(optimizer='adam', loss='mse')
```

```
best_gru = tuner_gru.get_best_models(1)[0]

best_gru.compile(optimizer='adam', loss='mse')

# CNN

model_cnn = Sequential([
    Conv1D(64, 3, activation='relu', input_shape=(X_train.shape[1], X_train.shape[2])),
    Dropout(0.3),
    Flatten(),
    Dense(64, activation='relu'),
    Dropout(0.3),
    Dense(1)
])

model_cnn.compile(optimizer='adam', loss='mse')

# Train models (use high epochs for better training)

EPOCHS = 30

history_lstm = best_lstm.fit(X_train, y_train, epochs=EPOCHS, batch_size=32,
    validation_data=(X_test, y_test), callbacks=[early_stop], verbose=1)

history_gru = best_gru.fit(X_train, y_train, epochs=EPOCHS, batch_size=32,
    validation_data=(X_test, y_test), callbacks=[early_stop], verbose=1)

history_cnn = model_cnn.fit(X_train, y_train, epochs=EPOCHS, batch_size=32,
    validation_data=(X_test, y_test), callbacks=[early_stop], verbose=1)
```

```
# Predictions

pred_lstm = best_lstm.predict(X_test)

pred_gru = best_gru.predict(X_test)

pred_cnn = model_cnn.predict(X_test)

ensemble_pred = (pred_lstm + pred_gru + pred_cnn) / 3


# Inverse scaling predictions

def inverse_preds(preds):

    return scaler.inverse_transform(np.hstack((preds, np.zeros((len(preds), 4)))))[:, 0]


# Actual prices

y_test_actual = inverse_preds(y_test.reshape(-1, 1))

# Predictions

pred_lstm_actual = inverse_preds(pred_lstm)

pred_gru_actual = inverse_preds(pred_gru)

pred_cnn_actual = inverse_preds(pred_cnn)

ensemble_actual = inverse_preds(ensemble_pred)


# Evaluation

def print_metrics(true, pred, name):
```

```
print(f"{name} -- RMSE: {np.sqrt(mean_squared_error(true, pred)):.2f}, MAE:  
{mean_absolute_error(true, pred):.2f}, R2: {r2_score(true, pred):.2f}")
```

```
print_metrics(y_test_actual, pred_lstm_actual, "LSTM")
```

```
print_metrics(y_test_actual, pred_gru_actual, "GRU")
```

```
print_metrics(y_test_actual, pred_cnn_actual, "CNN")
```

```
print_metrics(y_test_actual, ensemble_actual, "Ensemble")
```

```
# Confidence intervals (uncertainty)
```

```
def get_ci(preds):
```

```
    preds_array = np.array([preds for _ in range(100)])
```

```
    mean = preds_array.mean(axis=0)
```

```
    std = preds_array.std(axis=0)
```

```
    z = norm.ppf(0.975)
```

```
    lower = mean - z * std
```

```
    upper = mean + z * std
```

```
    return inverse_preds(mean), inverse_preds(lower), inverse_preds(upper)
```

```
ensemble_mean, ensemble_low, ensemble_high = get_ci(ensemble_pred)
```

```
# Plot predictions with confidence intervals
```

```
plt.figure(figsize=(14,6))
```

```

plt.plot(y_test_actual, label='Actual', color='black')

plt.plot(ensemble_mean, label='Predicted', color='blue')

plt.fill_between(range(len(ensemble_mean)), ensemble_low, ensemble_high, alpha=0.3,
label='95% CI')

plt.legend()

plt.title("Predicted vs Actual with Confidence Intervals")

plt.show()

```

```

# T-tests between models

print("T-test LSTM vs GRU:", ttest_rel(pred_lstm_actual, pred_gru_actual).pvalue)

print("T-test LSTM vs CNN:", ttest_rel(pred_lstm_actual, pred_cnn_actual).pvalue)

print("T-test GRU vs CNN:", ttest_rel(pred_gru_actual, pred_cnn_actual).pvalue)

```

```
# 1. Select a representative sample from training data
```

```
background = X_train[np.random.choice(X_train.shape[0], 100, replace=False)]
```

```
# 2. Create explainer - Use KernelExplainer as fallback
```

```

def model_predict(X):

    """Reshape data for model prediction"""

    if len(X.shape) == 2:

        X = X.reshape(-1, X_train.shape[1], X_train.shape[2])

    return best_lstm.predict(X, verbose=0)

```

```
explainer = shap.KernelExplainer(model_predict, background.mean(axis=0).reshape(1, -1))

# 3. Calculate SHAP values for test samples

X_shap = X_test[:50].reshape(50, -1) # Flatten time steps and features

shap_values = explainer.shap_values(X_shap)

# 4. Reshape and plot

feature_names = ['close','MA_7','MA_30','Volatility','RSI']

time_step_names = [f't-{i}' for i in range(X_train.shape[1]-1, -1, -1)]

# Aggregate absolute SHAP values across time steps

shap_abs = np.abs(shap_values).mean(axis=0)

shap_abs = shap_abs.reshape(X_train.shape[1], len(feature_names))

# Create heatmap of feature importance over time

plt.figure(figsize=(12, 8))

sns.heatmap(shap_abs,
            xticklabels=feature_names,
            yticklabels=time_step_names,
            cmap='viridis')

plt.title("SHAP Value Importance Across Time Steps")
```

```

plt.xlabel("Features")

plt.ylabel("Time Steps")

plt.show()

# Alternative: Bar plot of overall feature importance

plt.figure(figsize=(10, 5))

shap.summary_plot(shap_values, X_shap,
                  feature_names=[f"t{t}_{f}" for t in time_step_names for f in feature_names],
                  plot_type='bar',
                  show=False)

plt.title("Overall Feature Importance")

plt.tight_layout()

plt.show()

```

Additional EDA plots

AAPL Candlestick Chart



