# BCSE497J - Project-I

# STOCK MARKET PREDICTION USING LONG SHORT TERM MEMORY  MODEL

| 21BCE2431 | BHAVANA SINGH |
|-----------|---------------|

Under the supervision of

**Dr. Saurabh Agrawal**

Assistant Professor Sr. Grade 1

School of Computer Science and Engineering (SCOPE)

**B.Tech.**

*in*

**Computer Science and Engineering**

**School of Computer Science and Engineering**

**VIT**®
**Vellore Institute of Technology**
(Deemed to be University under section 3 of UGC Act, 1956)

November 2024

# DECLARATION

I hereby declare that the project entitled **Stock Market Prediction Using LSTM Model** submitted by me, for the award of the degree of *Bachelor of Technology in Computer Science and Engineering* to VIT is a record of bonafide work carried out by me under the supervision of **Dr. Saurabh Agrawal.**

I further declare that the work reported in this project has not been submitted and will not be submitted, either in part or in full, for the award of any other degree or diploma in this institute or any other institute or university.

Place : Vellore

Date :

**Signature of the Candidate**

# CERTIFICATE

This is to certify that the project entitled **Stock Market Prediction Using LSTM Model** submitted by **Bhavana Singh ( 21BCE2431 )** , **School of Computer Science and Engineering**, VIT, for the award of the degree of *Bachelor of Technology in Computer Science and Engineering*, is a record of bonafide work carried out by him / her under my supervision during Fall Semester 2024-2025, as per the VIT code of academic and research ethics.

The contents of this report have not been submitted and will not be submitted either in part or in full, for the award of any other degree or diploma in this institute or any other institute or university. The project fulfills the requirements and regulations of the University and in my opinion meets the necessary standards for submission.

Place : Vellore

Date :

**Signature of the Guide**

**Examiner(s)**

**Head of the Department**

**Programme**

**Dr. UMADEVI K S**

**SCOPE**

# ACKNOWLEDGEMENTS

# ABSTRACT

Predicting and analyzing the Indian stock market is a complex task due to the inherent volatility, non-linearity, and dynamic nature of the financial market. Accuracy in forecasting is crucial for investors, as it can minimize risks while maximizing profits. The Long Short-Term Memory (LSTM) model, a subset of Recurrent Neural Networks (RNNs), has emerged as a powerful tool for time series prediction.

In this paper, we explore the application of LSTM for forecasting the Indian stock market. The model's performance is highly dependent on the careful selection of hyperparameters, such as optimizer, batch size, activation function, and the number of hidden layers. We compare stateless and stateful LSTM models and test different hyperparameter configurations to optimize the model for better accuracy. A dataset from the National Stock Exchange (NSE) was used, and the Root Mean Square Error (RMSE) metric was employed to evaluate predictive performance.

Our findings demonstrate that proper tuning of hyperparameters and model configurations significantly enhances prediction accuracy, outperforming traditional methods such as Support Vector Machines (SVM), Naive Bayes (NB), and K-Nearest Neighbors (KNN).

*Keywords:*

 *Long Short-Term Memory (LSTM), Recurrent Neural Network (RNN), Stock Price Prediction, Hyperparameter Tuning, Time Series Forecasting, RMSE, Indian Stock Market.*

# TABLE OF CONTENTS

# LIST OF FIGURES

# 1. <u>INTRODUCTION</u>

## 1.1 Background

Stock price prediction has always been one of the most interesting and challenging topics for researchers, investors, and traders. Accurate forecasts of stock prices may result in gaining an advantage in the capital markets for traders by allowing them to optimize their investment strategy. Conventionally speaking, there have been two main approaches for predicting future price movements: technical analysis and fundamental analysis. While technical analysis, therefore, focuses on the research of patterns within historical data in a bid to identify repeatable trends, fundamental analysis looks at general economic indicators, the health of a company's finances, and market conditions. The forecasting of the stock price, however, has become even more sophisticated with the big data development and, correspondingly, the development of artificial intelligence.

Within the last couple of years, rapid growth in machine learning techniques has considerably enhanced the predictive capability of models of stock prices. Among the most in use, deep learning approaches, such as RNNs and their variant, the Long Short-Term Memory (LSTM) networks, have gained quite considerable popularity. LSTM networks are designed to process data as a sequence; that is why they are particularly effective in time-series forecasting—for example, in stock prices. More crucially, compared to the prior RNNs, which are subject to the problem of vanishing gradients, LSTM is able to capture long-range dependencies; this makes it absolutely unique for processing nonlinear and highly dynamic stock data.

Studies have continuously shown the strength of LSTM over its classical time-series counterparts, such as Auto-Regressive Integrated Moving Average and Multi-Layer Perceptron, in terms of accuracy. Such models have been applied to various markets, including the NYSE, Indian National Stock Exchange, and Turkish BIST100. Furthermore, an important step in tuning performance to show even better performance with an LSTM model is to optimize hyperparameters: variable tuning related to the selection of activation functions, number of batches, and the number of hidden layers within the network structure. Further tuning of these parameters using techniques like metaheuristic optimization algorithms, such as Search Economics and Genetic Algorithms, is an instance of making the predictions more reliable and full of foresight.

Besides historic price data, modern approaches to stock price prediction increasingly make use of alternative data inputs derived from sources such as social media sentiment and financial news. Indeed, there exists research that has established a close relationship between public sentiment, particularly on Twitter, and stock market trends. Due to this added layer of information, the machine learning model will be able to pick out more important insights about market behavior.

## 1.2 Motivation

The unpredictable and volatile nature of stock markets has long posed a formidable challenge for investors, traders, and financial analysts. The task of forecasting stock prices is complex, given that market prices can be influenced by a multitude of factors, including macroeconomic trends, political events, and unexpected global crises. Although there is an abundance of historical data available, it is often difficult to leverage this data effectively to predict stock movements due to the complexities and nuances of market dynamics. Traditional methods such as technical analysis, which focuses on historical price patterns, and fundamental analysis, which evaluates a company's financial health, have been widely used but often fall short when it comes to achieving high accuracy in predicting volatile, time-sensitive stock movements.

Technical analysis, for instance, may fail to capture the full range of factors influencing the stock market, particularly in times of economic turmoil or unusual price swings. Fundamental analysis, on the other hand, typically examines company-specific factors but may overlook the broader macroeconomic and geopolitical influences that can drive market trends. Traditional statistical models like ARIMA (Auto-Regressive Integrated Moving Average) have been used for time-series forecasting; however, they assume linearity and stationarity, which are often unrealistic assumptions in the complex world of financial markets. Similarly, other machine learning models, such as Multilayer Perceptrons (MLPs), are not designed to capture sequential dependencies in data and may struggle with stock price data, which is inherently temporal and highly non-linear.

In recent years, the rapid advancement of artificial intelligence and machine learning has opened up new possibilities in financial forecasting. Techniques from deep learning, particularly Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) networks, have shown promise in handling sequential and temporal data, making them well-suited for time-series analysis. LSTM, in particular, overcomes the limitations of traditional RNNs by effectively addressing the vanishing gradient problem, allowing it to capture long-term dependencies in the data. This capability makes LSTM models ideal for predicting stock prices, which are influenced not only by historical trends but also by real-time events, such as news reports and social sentiment. Unlike traditional models, LSTMs can accommodate both short-term fluctuations and long-term patterns in stock price data.

This research is motivated by the transformative potential of LSTM networks in stock price forecasting. Deep learning models, particularly LSTM, provide a significant advantage in handling complex, high-dimensional data. LSTM's ability to learn patterns in sequential data enables it to outperform traditional models in capturing the nuances of stock market behavior. By integrating additional variables, such as social sentiment and economic indicators, LSTM models can provide a more comprehensive view of the factors influencing stock prices. Furthermore, LSTM networks can be optimized through hyperparameter tuning, which allows for greater model customization and adaptability to specific datasets. Fine-tuning

parameters, such as learning rate, batch size, and the number of LSTM units, can significantly enhance the model's predictive accuracy and robustness.

This project aims to leverage the strengths of LSTM networks in the context of stock price forecasting, with a particular focus on hyperparameter tuning to optimize model performance. The objective is to create a predictive model that not only performs well in terms of accuracy but also provides investors and analysts with reliable insights into market trends. By addressing the limitations of traditional methods and implementing an LSTM-based approach, this project contributes to the growing body of research on AI applications in financial markets. In doing so, it seeks to bridge the gap between theoretical advancements in machine learning and practical applications in stock market analysis.

## 1.3 Scope of the Project

The project focuses on creating an advanced stock price prediction model using LSTM This project focuses on developing a robust stock price prediction model using LSTM networks, with a particular emphasis on the Indian stock market. While global stock markets exhibit some common characteristics, each market has unique attributes and challenges. The Indian stock market, characterized by high volatility, sector-specific risks, and regulatory dynamics, presents an ideal environment to test the capabilities of advanced machine learning models. By focusing on the Indian market, this project aims to provide valuable insights for local investors and analysts, while also developing a framework that could be adapted for other markets.

The model will incorporate a broad spectrum of features beyond historical prices. These features include fundamental data, such as company earnings and balance sheet metrics; macroeconomic indicators, such as inflation rates and interest rates; and technical indicators derived from price and volume data. By integrating these diverse features, the model aims to capture a holistic view of the factors influencing stock price movements. This approach is particularly important in the Indian market, where price movements can be influenced by domestic economic policies, regulatory changes, and global economic trends.

A critical aspect of the project involves systematic data collection and preprocessing. The quality and relevance of input data are essential for building an effective predictive model. In this project, data preprocessing steps will include handling missing values, scaling features to ensure numerical stability, and transforming variables to reduce multicollinearity. For example, certain financial metrics may be highly correlated, leading to issues in model interpretability and overfitting. By addressing these issues through feature selection and transformation techniques, the model can be made more resilient and interpretable.

Hyperparameter tuning will play a crucial role in optimizing the LSTM model architecture. LSTM models, while powerful, require careful configuration to perform optimally. Key hyperparameters, such as the number of LSTM units, the dropout rate, and the learning rate,

will be fine-tuned using techniques like grid search or Bayesian optimization. Tuning these parameters is essential for overcoming challenges like the vanishing gradient problem, which can hinder model performance in long time-series datasets. Additionally, techniques such as dropout regularization will be used to prevent overfitting, ensuring that the model generalizes well to unseen data.

The model's performance will be evaluated using standard metrics such as Root Mean Squared Error (RMSE) and Mean Absolute Percentage Error (MAPE). These metrics provide a quantitative measure of prediction accuracy, allowing for a clear comparison with traditional methods like ARIMA. In addition to these standard metrics, the project will also examine directional accuracy, which assesses the model's ability to predict the correct direction of price movement. This is particularly valuable for investors, as predicting the direction of price changes can be more informative than predicting the exact price.

Ultimately, the goal of this project is to develop a reliable forecasting tool that can support investors and stakeholders in making informed decisions within a volatile market environment. By focusing on LSTM and hyperparameter tuning, the project aims to create a model that is not only accurate but also interpretable and adaptable to changing market conditions. Additionally, the project seeks to contribute to the field of financial predictive analytics by demonstrating the potential of advanced machine learning models in stock price forecasting. By doing so, it hopes to pave the way for future research and advancements in machine learning applications within the finance sector.

In conclusion, this project will address the complexities of stock price prediction by leveraging the strengths of LSTM networks and incorporating a wide range of features. The systematic approach to data collection, preprocessing, and model optimization will ensure that the model is robust, accurate, and interpretable. This project will not only provide valuable insights for the Indian stock market but also serve as a blueprint for developing advanced predictive models for other financial markets.

# 2. PROJECT DESCRIPTION AND GOALS

## 2.1 Literature Review

The progression of stock price prediction methods reflects a shift from early theoretical approaches to advanced machine learning techniques, each iteration seeking to address the complexities unique to financial markets. Foundational theories in this field, such as the Random Walk Theory and the Efficient Market Hypothesis (EMH), introduced in the 20th century, laid the groundwork for much of modern financial research and have greatly influenced the understanding of stock market behavior.

The **Random Walk Theory** suggests that stock prices move unpredictably, making it impossible to accurately forecast future values based solely on past data. According to this theory, price changes are random and independent, which implies that no consistent patterns exist that could be leveraged for profit. This concept was revolutionary and contentious, challenging traditional assumptions about market predictability. However, critics argued that the Random Walk Theory disregards the impact of external factors—like economic policies or investor behavior—that can lead to identifiable trends. Additionally, the idea of complete randomness fails to account for certain behavioral or psychological patterns that can affect market prices, such as herd behavior or emotional biases, suggesting that prices may not be entirely unpredictable after all.

The **Efficient Market Hypothesis (EMH)**, developed by economist Eugene Fama, provides an alternate view, suggesting that all available information is immediately absorbed into asset prices. EMH asserts three forms of market efficiency: weak, semi-strong, and strong. Each level represents a different degree of information integration. In the weak form, asset prices reflect all past trading data. In the semi-strong form, prices incorporate all publicly available information, and in the strong form, prices reflect both public and insider information. EMH contends that it is impossible to achieve consistently above-average returns in an efficient market since all information is already factored into prices. However, critics of EMH, particularly behavioral economists, argue that inefficiencies exist due to psychological factors—such as overconfidence and panic—that skilled investors can potentially exploit. Despite its limitations, EMH remains a key theory in finance, helping to shape the study of market efficiency and investor behavior.

**Machine learning and statistical models** have since evolved to address the limitations of these early theories. One of the first machine learning approaches applied to stock prediction was the **AutoRegressive Integrated Moving Average (ARIMA)** model. ARIMA operates under the assumption of linear relationships in time-series data, making it suitable for stationary series but limited in handling the nonlinear behavior often seen in financial data. Financial time series, especially during periods of high volatility or economic disruptions, frequently exhibit non-stationary patterns that ARIMA struggles to capture. The linear nature of ARIMA restricts its usefulness in stock prediction, particularly in high-frequency trading and volatile market conditions where complex interactions between market variables can occur.

**Support Vector Machines (SVMs)** have also been explored for stock forecasting, primarily for classification and regression tasks. SVMs show good performance in short-term predictions but often face challenges when applied to large, interdependent, and nonlinear datasets, typical of financial markets. While SVMs utilize kernel functions to map data into higher-dimensional spaces, helping detect more complex relationships, they are inherently limited in capturing sequential dependencies. Time-series data in stock markets depend heavily on historical trends, a dependency SVMs struggle to model accurately over extended periods. Moreover, SVMs lack a mechanism for processing sequential data natively, limiting their effectiveness for long-term forecasting, where past patterns and relationships heavily influence future outcomes.

The advent of **deep learning models**, particularly **Long Short-Term Memory (LSTM)** networks, marks a promising shift in stock prediction capabilities. LSTM networks, a specialized form of Recurrent Neural Networks (RNNs), are designed to handle sequential data and retain information across extended sequences, which sets them apart from traditional models. LSTMs incorporate a gating mechanism with input, forget, and output gates that regulate information flow. This design allows LSTM models to selectively remember relevant information and forget irrelevant data, making them well-suited for time-series applications such as stock price forecasting. Unlike ARIMA, which assumes linearity, LSTM networks excel in capturing both linear and nonlinear relationships within stock data, making them robust in highly volatile and complex market conditions.

LSTM models have been successfully applied not only to individual stock price predictions but also to broader indices like the **S&P 500 and Nifty 50**, where they provide significant improvements in accuracy over conventional methods. By learning and identifying patterns in historical data, LSTMs can detect subtle signals that are easily missed by traditional models. However, despite their success, LSTM models come with challenges, particularly around **hyperparameter tuning**. Model performance in LSTMs is sensitive to hyperparameters such as the number of hidden layers, learning rate, and batch size. Fine-tuning these parameters is essential to avoid issues like overfitting (where the model learns too much from training data and fails to generalize to new data) and underfitting (where the model does not capture the underlying patterns in the data). Finding the right balance requires extensive experimentation, computational resources, and technical expertise.

To optimize LSTM models, researchers have explored **hyperparameter optimization techniques** such as **Search Economics**. This systematic approach to adjusting hyperparameters aims to balance predictive accuracy and computational efficiency. By refining key parameters, Search Economics enhances the accuracy of LSTM models without excessive computational costs, making them feasible for real-time stock forecasting. The success of LSTM networks in stock prediction suggests they hold significant potential to transform financial forecasting by offering a model that captures both past trends and future price movements.

In conclusion, LSTM models represent a substantial advancement in stock price prediction by addressing the limitations of earlier methods like ARIMA and SVM. With improved

hyperparameter tuning and optimization techniques, LSTMs are well-positioned to serve as effective tools for real-time market forecasting, providing traders and investors with valuable insights in a data-driven financial environment.

## 2.2 Research Gap

The implementation of Long Short-Term Memory (LSTM) models in stock market forecasting has shown considerable promise, especially with its capacity to handle sequential and time-dependent data. However, certain research gaps remain, particularly in relation to the Indian stock market and the unique volatility it experienced post-COVID-19. The pandemic brought about unprecedented challenges, disrupting global economies and inducing substantial fluctuations in stock markets. In the Indian context, these fluctuations were heightened by regional factors such as varying lockdown measures, changes in consumer behavior, and economic uncertainty. Traditional forecasting models that relied heavily on historical data struggled to adapt to this new, unpredictable environment. LSTM models, which were initially developed to capture temporal patterns in relatively stable datasets, also encountered challenges in addressing the extreme shifts and volatility introduced by the pandemic.

One of the foremost gaps in research is the **need to validate LSTM models against the volatile market conditions that emerged post-COVID-19**. The pandemic-driven disruptions created nonlinear trends and drastic short-term fluctuations that were difficult for existing models to predict accurately. While LSTM networks have the advantage of retaining long-term dependencies, the unique post-pandemic behavior of markets has not been extensively studied within this framework. This period has highlighted the need for more adaptive models that can respond to sudden shifts in market behavior, making it necessary to explore LSTM's capacity for handling these atypical patterns effectively. Further research is required to determine if and how LSTM models can be modified or combined with other techniques to better capture these nonlinear market dynamics and, importantly, to validate their performance in high-volatility environments.

Another critical area of research involves **comparing LSTM with other hybrid models** such as **Gated Recurrent Unit (GRU)** and **CNN-LSTM hybrids**. GRU, like LSTM, is a type of recurrent neural network but is designed with a simpler architecture that retains many of LSTM's advantages while requiring fewer computational resources. This makes GRU potentially more efficient, especially in high-frequency trading scenarios where real-time predictions are critical. GRU's simpler structure may also contribute to faster convergence, which could be beneficial in handling rapid market changes. However, there is limited research on how GRU performs in comparison to LSTM specifically in post-pandemic environments, where rapid adjustments in model performance are necessary.

The **CNN-LSTM hybrid model** is another innovative approach that has gained attention. By combining convolutional neural networks (CNN) with LSTM, this model leverages CNN's

ability to automatically detect features within the data before passing them to the LSTM for sequential processing. This feature extraction step could be advantageous in stock market forecasting, as CNNs may help the model detect local patterns or "micro-trends" in the stock data that might otherwise be overlooked. For example, CNN layers could identify minor yet consistent patterns in stock price changes that LSTM layers can then use to improve their sequence-based predictions. Although promising, studies comparing the performance of CNN-LSTM hybrids with traditional LSTMs are limited, especially in the context of the Indian stock market, which is known for its diverse economic and regulatory influences.

A significant research gap also exists in **developing multi-target LSTM models capable of forecasting multiple stock price metrics simultaneously**. The majority of existing studies focus on single-target predictions, often forecasting closing prices without considering other valuable metrics such as opening, high, and low prices. In practice, these additional price points can provide a more comprehensive view of a stock's daily behavior, and their prediction could offer investors a fuller picture of market trends. A multi-target LSTM model would be capable of forecasting all these price points in a single framework, potentially increasing the model's usefulness and precision in trading strategies. By incorporating multiple targets, the model could capture the interdependencies between different price metrics—such as the relationship between opening and closing prices—and thereby enhance its predictive power.

Multi-target forecasting not only provides a more holistic understanding of market behavior but also addresses a practical need for more comprehensive predictive insights. For instance, traders often rely on high and low prices to set stop-loss limits and make decisions on entry and exit points, while opening prices can signal the day's initial momentum. Despite these practical applications, there is limited research on multi-target LSTM models, and even fewer studies have evaluated their effectiveness specifically in the context of emerging markets like India. Developing such models could pave the way for more accurate, versatile, and actionable insights for traders and analysts.

To further expand on the LSTM's capabilities, research on **hyperparameter tuning and model optimization** tailored to post-pandemic market conditions could be highly valuable. LSTM models are sensitive to parameters such as the number of layers, batch size, learning rate, and sequence length. Fine-tuning these parameters is essential to achieve optimal performance, but current literature lacks extensive research on how these parameters can be adjusted for high-volatility conditions, such as those experienced during the pandemic. Optimization techniques, such as grid search, random search, or more advanced approaches like Bayesian optimization, could be explored to determine the best configuration for LSTM models in the Indian market context. Additionally, regularization techniques, like dropout and early stopping, could be applied to prevent overfitting, particularly in volatile markets where historical data may not be entirely representative of future trends.

In summary, several key research gaps exist in the implementation of LSTM models for stock market forecasting in the post-pandemic context, particularly within the Indian market. These include the need for validation of LSTM performance in highly volatile environments,

comparative studies involving GRU and CNN-LSTM hybrid models, and the development of multi-target LSTM models that can forecast multiple price metrics. Addressing these gaps could lead to significant advancements in the predictive accuracy and versatility of LSTM-based models in financial forecasting. By tailoring LSTM models to capture post-pandemic market behavior and leveraging alternative architectures, future research can better equip financial analysts and investors with tools to navigate the complexities of modern markets.

## 2.3 Objectives

The primary objectives of the project are as follows:

1. **Predict Future Stock Prices:**
   - Develop an LSTM (Long Short-Term Memory) model to forecast future stock values for Indian companies using historical data.
   - Concentrate on successfully anticipating price values based on key indications, such as stock prices for select companies such as HDB, TATASTEEL, RELIANCE, and others.
   - Data such as historical closing prices, trade volumes, and other technical characteristics are useful for understanding market patterns.

2. **Time-Series Analysis Using LSTM:**
   - Model sequential dependencies in stock data by utilizing LSTM and concentrating on past prices and trading volumes.
   - To obtain the information required for a thorough grasp of market patterns, use yfinance.
   - Manage time-series data well, since accurate future prediction depends on temporal correlations.

3. **Handle Long-Term Dependencies in Stock Data:**
   - Make better predictions of future trends by utilizing the key information that may be retained from prior data by leveraging the unique memory capabilities of LSTM.
   - Better long-term forecast accuracy for stock prices can be achieved by overcoming the vanishing gradients problem that frequently affects regular RNNs.

4. **Improve Decision-Making for Investors:**
   - Identify trends and forecast stock price movements for high-performing firms to give traders and investors important information.
   - Utilizing cutting-edge deep learning algorithms, produce projections that assist with data-driven, strategic investment decisions.

5. **Evaluate Model Performance and Improve It:**

- Metrics like Root Mean Squared Error (RMSE) and Mean Absolute Percentage Error (MAPE) can be used to assess how well the LSTM model is performing.
- In financial applications, where even little variations can have a big financial impact, it is imperative to minimize these inaccuracies in order to achieve the maximum possible accuracy.

6. **Address Research Gaps in Indian Stock Market Prediction:**
   - Fill in the Research Gaps in the Indian Stock Market Forecast by using LSTM models, close the research gap in predicting the Indian stock market after COVID-19.
   - Examine the effectiveness of LSTM as a contemporary AI strategy in comparison to conventional statistical models, offering an understanding of how effectively it adjusts to the unstable market conditions.

## 2.4 Problem Statement

The project aims to address the challenges involved in predicting stock prices within the volatile and complex Indian stock market using modern AI techniques like LSTM. The key challenges include:

1. **Impact of Post-Pandemic Volatility:**
   - Following the COVID-19 pandemic, there was a great deal of volatility in the Indian stock market due to a combination of global disruptions and economic uncertainties.
   - Unreliable predictions have resulted from traditional prediction models' inability to effectively represent these turbulent dynamics.
2. **Limitations of Traditional Methods**:
   - The sequential and non-linear nature of stock data makes it difficult for conventional approaches like ARIMA and SVM to accurately represent it.
   - These techniques frequently fall short of capturing long-term dependencies, which are essential for accurate stock price forecasting.
3. **Research Gaps in Predictive Model Comparisons and Multi-Value Efficiency:**
   - Research on evaluating the real-time performance of various predictive models in the Indian market is scarce, especially when it comes to the models' ability to predict multiple values concurrently (e.g., opening, high, low, and closing values).
   - In order to fill in these research gaps, this project uses LSTM to assess the efficacy and efficiency of the model using data from the Indian stock market.
4. **Solution Approach:**
   - By identifying intricate sequential connections in historical data, LSTM models can be used to forecast changes in stock prices.

- Examine the efficacy and precision of the LSTM model in comparison to conventional prediction models, taking into account the post-COVID-19 dynamics of the Indian market.
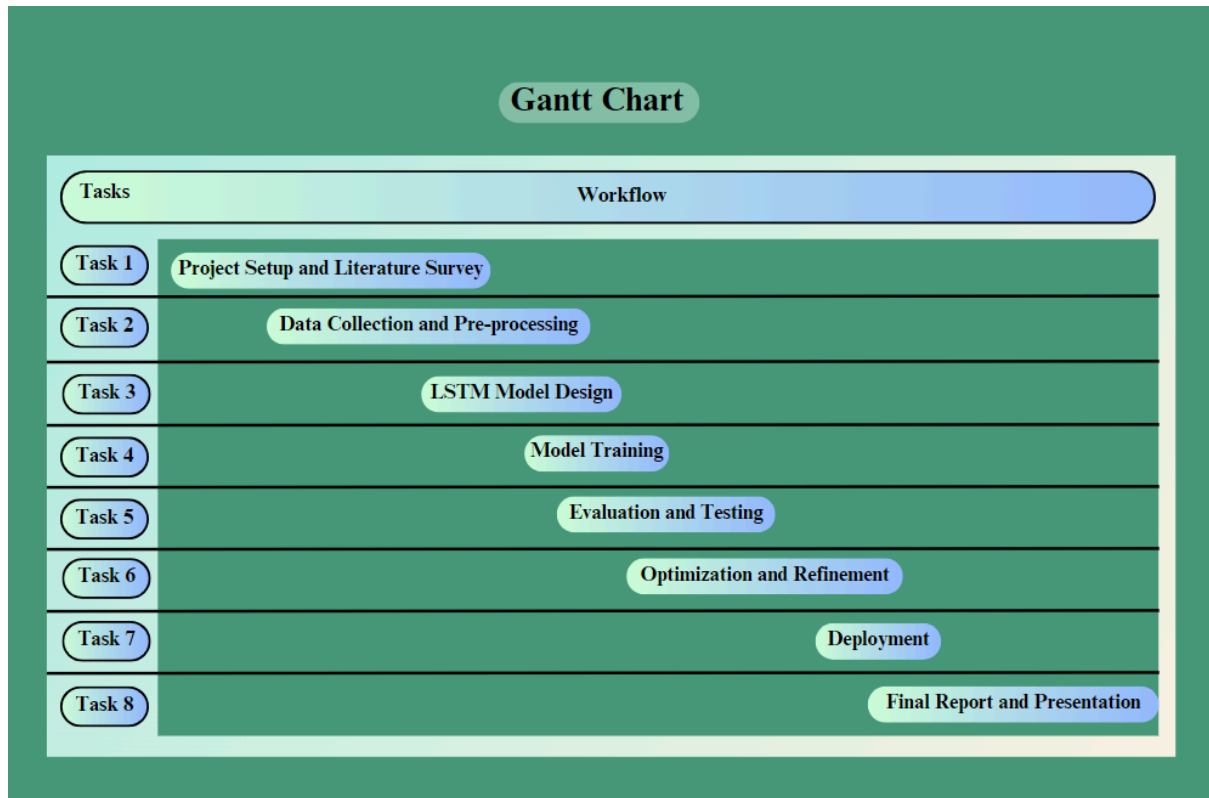
## 2.5 Project Plan



**Fig. 1. Gantt chart**

# 3. <u>TECHNICAL SPECIFICATION</u>

## 3.1 Requirements

### 3.1.1 *Functional*

- **Data Collection:** The system must collect data from various sources, including historical stock prices, macroeconomic indicators, and technical metrics.
- **Data Preprocessing:** The system needs to preprocess the data to ensure clean datasets free from noise and irrelevant information.
- **Feature Selection:** The system should implement effective feature selection to identify the most influential factors affecting stock price behavior.
- **Model Training:** The LSTM model should be trained using historical data to detect significant patterns and relationships in stock price movements.
- **Hyperparameter Tuning:** The system should allow hyperparameter tuning to optimize performance and improve prediction accuracy.
- **Real-time Prediction:** It should provide real-time stock price predictions, offering timely insights for investors and stakeholders.
- **User Interface:** The system must offer a user-friendly interface that allows users to visualize data, review predictions, and adjust model parameters easily.
- **Reporting:** It should generate comprehensive reports summarizing model performance, prediction accuracy, and key insights from the analysis.

### 3.1.2 *Non-Functional*

- **Performance:** The system should ensure minimal latency in processing and analyzing data, delivering timely predictions and insights.
- **Scalability:** It must be scalable to handle increasing amounts of historical data and user demands without compromising performance.
- **Reliability:** The system should maintain high reliability, ensuring minimal downtime and incorporating robust error-handling mechanisms.
- **Security:** It should guarantee data integrity and confidentiality by enforcing strict access controls and encryption techniques.
- **Usability:** The system should be intuitive, easy to navigate, and thoroughly documented for end-users.
- **Maintainability:** It should be designed for easy maintenance and updates, featuring modular components and well-documented code.
- **Compliance:** The system must adhere to relevant financial regulations, ensuring data protection and the ethical use of financial information.

## 3.2 Feasibility Study

### 3.2.1 Technical Feasibility

- **Readily Available Technology**: The project will apply to machine learning and AI technologies, especially the LSTM networks, which are very common in time series prediction. They are perfectly documented and supported by extensive libraries like TensorFlow, Keras, and Scikit-learn. They ensure a firm basis for the predictive model. In addition, the data pre-processing tools and their databases can be used to maintain the stock market data, such as Pandas and MySQL, making the respective technical framework solid.
- **Technical Expertise**: The project requires a team that has extensive expertise in machine learning, AI, and financial data analysis. Deep learning proficiency, especially regarding the training of LSTM models, will become crucial, as predictive algorithms will be developed. Knowledge of the fundamentals of the stock market, macroeconomic indicators, and technical research is also necessary for effective feature selection. Success in either training or hiring such expertise will come to define the project's sustainability and efficiency. These fields have skilled professionals available, which means the project's technical requirements are attainable.
- **Infrastructure**: This should include enough computational resources, such as high-performance servers or cloud computing platforms-including AWS or Google Cloud-to handle large datasets, train complex LSTM models, and perform extensive tuning of hyperparameters. These resources allow for timely and efficient processing and analyses of stock market data. A cloud infrastructure provides scalability and flexibility; therefore, computational demands can be managed with the growth of a dataset.
- **System Integrations**:The system will be integrated well with the existing Data Collection and Financial platforms**.** This system will tap into APIs for stock prices and economic data from third-party providers. Utilizing libraries supported by virtually every IT infrastructure ensures seamless compatibility. Efficient storage and retrieval of data in relational databases and interfacing with cloud services make this prediction model seamless across various financial and trading environments.

### 3.2.2 Economic Feasibility

- **Cost-Benefit Analysis:** Open-source tools like Python, TensorFlow, Keras, and yfinance offer a cost-effective LSTM-based stock prediction model on a cloud platform, reducing dependency on proprietary tools and enabling accurate forecasts for improved investment choices.

- **Budget:** A detailed budget plan should outline costs for cloud infrastructure, GPU services, model training personnel, maintenance, and real-time prediction enhancements, despite the open-source nature of the software tools.

- **Return on Investment (ROI):** The initiative provides data-driven insights to financial institutions and investors, reducing risks, increasing profitability, and resulting in significant ROI, improved financial performance, and long-term cost savings.

- **Funding:** The project may require additional funding from investors or stakeholders to scale, ensure reliable real-time deployment, and enable feature growth, UI enhancements, and model performance for wider market application.

### 3.2.3 Social Feasibility

- **User Acceptance:** Investors should be able to easily access the system and obtain actionable and unambiguous insights. It seeks broad adoption among individual traders and investors who frequently lack access to specialized tools by offering stock market analysis during periods of extreme volatility.

- **Training and Support:** To guarantee that users are able to comprehend and use the stock market predictions in an efficient manner, sufficient training courses and user support materials should be made available. By doing this, users will be able to optimize the system's worth and make wiser financial choices.

- **Ethical Considerations:** Data privacy should be protected, and model forecasts should be transparent and devoid of bias that could mislead investors. The application of AI and machine learning models, such as LSTM, in financial analysis should abide by ethical standards.

- **Impact on Investor Behavior:** Through lowering the level of uncertainty around stock projections, the project may have an impact on investors' decision-making processes and foster confidence in their financial planning. In order to effectively handle these developments, it is crucial to inform investors about the appropriate application of predictive insights and to promote deliberate decision-making.

## 3.3 System Specification

### 3.3.1 Hardware Specification

To complete the LSTM model for stock market prediction, the following requirements are recommended:

- *Processor :* Intel i7 or AMD Ryzen 7 processor
- *Graphics Processing Unit (GPU) :* A powerful processor is required to meet the needs of deep learning models such as LSTM. training. Since the LSTM model includes matrix processing and takes advantage of GPU parallelism, NVIDIA RTX 3060 or above is recommended. especially when dealing with various stock indices and historical data.

- *Memory (RAM) :* 16 GB(32 GB recommended)
- *Network :* Stable connection.A reliable network connection is required to collect and preprocess the history of online products such as yfinance.
- *Storage :* 500GB SSD

### 3.3.2 Software Specification

- *Operating System*:
  - ➢ Compatible with Windows, macOS, and Linux to support various development and deployment environments.
- *Programming Languages:*
  - ➢ Python: Primary language for predictive modeling and data preprocessing due to its extensive ML libraries.
  - ➢ SQL: Used for efficient database queries and management of historical stock price data.
- *Development Environment:*
  - ➢ Jupyter Notebooks and PyCharm for testing, debugging, and model visualization.
  - ➢ Git and GitHub for version control and collaboration.
- *Libraries & Frameworks:*
  - ➢ TensorFlow and Keras for LSTM model development.
  - ➢ Pandas and NumPy for data manipulation.
  - ➢ Matplotlib and Seaborn for data visualization.
  - ➢ Scikit-learn for feature selection, scaling, and data splitting.
- *Database:*
  - ➢ MySQL or SQLite for storing historical stock data, macroeconomic indicators, and technical metrics.
- *Security:*
  - ➢ SSL encryption for secure data transmission.
  - ➢ OAuth for user authentication and authorization.

# 4. DESIGN APPROACH AND DETAILS

## 4.1 System Architecture



**Fig. 2. System Architecture**

- **Data Source:**
  - ➢ Historical stock data is fetched from external APIs like Yahoo Finance or Alpha Vantage.
  - ➢ The data includes stock prices, volumes, and other key indicators, usually in CSV or JSON formats.
- **Data Preprocessing:**
  - ➢ Raw data is cleaned by handling missing values and removing any noise or anomalies.
  - ➢ Data is normalized to ensure it is scaled correctly, making it easier for the model to process.
  - ➢ The data is split into training and testing sets, ensuring that the model is evaluated on unseen data.
  - ➢ Feature engineering is applied, adding important stock market indicators like moving averages to improve prediction accuracy.
- **LSTM Model:**
  - ➢ The LSTM (Long Short-Term Memory) network is used because of its ability to capture both short-term and long-term patterns in time-series data.

- ➢ The input layer receives sequential stock data (e.g., the last 60 days of stock prices).
  - ➢ The LSTM layers analyze and process time-series data to detect dependencies and trends.
  - ➢ The output layer predicts the next time step's stock price based on historical data.
- **Training Process:**
  - ➢ The model is trained by fitting it to the training data over multiple epochs.
  - ➢ Backpropagation and gradient descent are used to adjust the model weights, minimizing prediction errors.
  - ➢ Hyperparameters like learning rate, batch size, and the number of LSTM layers are tuned for optimal performance.
- **Prediction and Evaluation:**
  - ➢ After training, the model is used to predict future stock prices.
  - ➢ The predictions are evaluated using metrics such as Root Mean Squared Error (RMSE) and Mean Absolute Percentage Error (MAPE), measuring the accuracy of the model.
- **Storage:**
  - ➢ Predicted stock prices are stored in a database along with performance metrics for further analysis and retrieval.
- **Dashboard & Visualization:**
  - ➢ The final output is visualized on a dashboard where users can view actual vs. predicted prices and stock trends.
  - ➢ Interactive charts and graphs allow users to explore the model's predictions in an intuitive and user-friendly interface.

## 4.2 Design

### *4.2.1 Data Flow Diagram*



This flowchart shows the stages involved in developing and deploying a stock price prediction model based on a Long Short-Term Memory (LSTM) network.

**Here is a breakdown of the steps-:**

- **Collect Stock Data (yfinance):** Begin by gathering historical stock data from a reputable source, such as the yfinance library.
- **Preprocess data using pandas and MinMaxScaler:** Clean and prepare the data to use with the LSTM model. This usually includes data cleaning that entails removing missing values, addressing outliers, and sometimes transforming data types.
- Scaling is the process of normalizing data using techniques such as MinMaxScaler to guarantee that all features fall within a comparable range.
- Split the data into training and testing sets to train the LSTM model and then evaluate its performance.
- **Train the LSTM Model (Keras):** Using the Keras library, create and train the LSTM model using the training data. You will define the model's architecture (layers, neurons, etc.) and initiate the training process.
- **Evaluate the model (RMSE, visualization):** Evaluate the model's performance using relevant measures such as Root Mean Squared Error (RMSE). You may also see how the predictions compare to the real data for further study.

- **Deploy Model (web application):** Once you're satisfied with the model's performance, make it available for use. This might include incorporating it into a web application where users can interact with the model and receive stock forecasts.

*4.2.2 Use Case Diagram*



The use case diagram depicts the interactions between a user and a stock prediction system. Here's a breakdown of the main components and their relationships:

**Actors:**

- **User:** The individual who interacts with the system is referred to as the user.
- **Server:** The infrastructure that runs the system and handles user queries.

**Use cases:**

- **Collect and Preprocess Stock Data:** The system collects stock data from various sources, such as yfinance, and prepares it for analysis.
- **View Stock Data and Trends:** The user can visually analyze historical stock data to detect patterns or trends.
- **Scale and Normalize Data:** Data is scaled and normalized to guarantee consistency and avoid issues with varying magnitude values.
- **Split Data into Training and Testing Sets:** Divide the data into training and testing sets so that the prediction model may be trained and evaluated.
- **Train the LSTM Model:** To learn stock price trends, the training data is fed into an LSTM (Long Short-Term Memory) model.

- **Model Evaluation:** To examine the trained model's performance, measures such as RMSE are used and shown.
- **Plot Stock Prices and Moving Averages:** The system allows you to visually represent stock prices and moving averages to provide additional context.
- **Calculate Percentage Change in Stock Prices:** The system estimates percentage changes in stock prices to assist users in understanding price fluctuations.
- **Deploy Model as a Web Application:** The trained model is deployed as a web application, which the user can access via a web interface.
- **Use Predictions to Make Informed Investment Decisions:** Users can utilize the model's predictions to make more educated investing decisions.

**Relationships:**

- **User → Use Cases:** Users initiate most use cases.
- **Use Cases → Server:** The server performs the necessary operations for each use case.
- **Use Cases → User (for visualization):** User interacts with the system to view visualizations and results.

Overall, this graphic shows how a user may use a Stock Prediction System to collect data, train a model, assess its performance, and make investment decisions based on its forecasts.

## 4.2.3 Class Diagram



The figure depicts the classes and relationships in a stock prediction system.

- **StockPredictor** is the key class that controls the entire process. It prepares the stock data with a DataPreprocessor, trains an LSTMModel on the preprocessed data, and then assesses the model's performance using an Evaluator.
- **DataPreprocessor** loads, cleans, normalizes, and generates datasets from raw stock data.
- The **LSTMModel** is a neural network model that learns to predict future stock prices using historical data.
- **Evaluators** assess the model's performance by comparing its forecasts to real stock prices. It then depicts the results.

The arrows between the classes show the relationships:

- StockPredictor uses DataPreprocessor for data preprocessing.
- StockPredictor uses LSTMModel for model training and prediction.
- StockPredictor uses Evaluator for evaluation and visualization.
- DataPreprocessor offers training data to the LSTMModel.
- The LSTMModel offers predictions to the Evaluator.

Overall, this diagram depicts a procedure for creating and deploying a stock prediction system based on deep learning.

## 4.2.4 Sequence Diagram



The sequence diagram illustrates how several components of an LSTM model-based stock price prediction system interact with one another. The steps are broken down as follows:

- **User generates a StockData instance:** The user generates a StockData object by interacting with the system. The retrieval and processing of stock data will be handled by this object.
- **StockData preprocesses and downloads stock data:** The StockData object gets stock price data from a source and formats and cleans the data, among other preparatory operations.
- **Data is divided into training and testing sets by StockData:** The preprocessed data is divided into two sets: a testing set for LSTM model performance evaluation and a training set for LSTM model training.
- **An LSTMModel instance is created by the user:** An LSTMModel object is created by the user. The Long Short-Term Memory (LSTM) model that will be used to forecast stock prices is represented by this item.
- **The model is compiled via LSTMModel:** The particular LSTM architecture selected determines the parameters and structure of the LSTMModel.
- **The model is fitted to training data by LSTMModel:** The training data from StockData is used to train the built LSTM model. In order to reduce the discrepancy between expected and actual stock prices, this stage entails modifying the model's parameters.
- 
  **A StockPredictor instance is created by the user:** A StockPredictor object is created by the user. This object will forecast new stock data using the trained LSTM model.

- **StockPredictor uses the LSTM model for training:** To improve its prediction skills, the StockPredictor object uses the trained LSTM model from LSTMModel.
- **StockPredictor forecasts based on fresh data:** The integrated LSTM model is used by the StockPredictor object to produce predictions on fresh stock data.
- **StockPredictor assesses forecasts:** The StockPredictor assesses the model's accuracy by evaluating its predictions using relevant measures, including mean squared error.
- **Plotting real and anticipated prices:** To give users a visual sense of the model's performance, StockPredictor graphically compares actual stock prices with the forecasts made by the algorithm.

# 5. <u>METHODOLOGY AND TESTING</u>

This section offers a thorough explanation of the testing procedures and techniques used to create, train, and validate a machine learning model for stock price prediction. This module uses an LSTM neural network to forecast future stock prices by analyzing trends in past data.

**5.1 Module Description**

The module is organized into multiple phases, each with distinct objectives, methods, and setups, ranging from data collection to model deployment.

- **Data Collection:**

  - *API Integration:* The module makes use of the yfinance library to access Yahoo Finance's data API, which allows for easy data retrieval. This configuration enables us to programmatically download data for certain organizations, streamlining the data collecting process.

  - *Data Range and Frequency:*The module retrieves data from the last four years, ensuring enough time to capture both long-term trends and seasonal fluctuations. A larger data set frequently improves the model's capacity to generalize, particularly in time-series forecasting.

  - *Frequency:* Daily stock prices are chosen for their comprehensive granularity, allowing for accurate short- and long-term projections.

  - *Targeted Tickers and Diversity:*The module focuses on many stocks, indicated by tickers (HDB, ITCI, ONGC.NS, TATASTEEL.NS, etc.), from diverse industries. This method enables the model to catch sector-specific trends and apply to a broader set of financial data.

- **Data Preprocessing**:

  Preprocessing is necessary for time-series models because it prepares data in a format that emphasizes temporal connections and scales it to keep outliers from skewing conclusions.

  - *Data Cleaning:*

    - **Handling Missing Values**:To ensure continuity in a time series, missing data are handled using interpolation or filling with a specific value (e.g., last valid observation).

    - **Feature Selection**: Only relevant columns, such as Adjusted Close, are chosen as input features, decreasing noise and directing the model toward the target variable.

➢ *Normalization*:

■ The module uses sklearn's MinMaxScaler to normalize the data to a range of 0 to 1. Normalization guarantees that data magnitudes are homogeneous, allowing the LSTM model to train effectively without any one variable dominating the learning process.

■ **Formula**: The MinMax scaling formula applied is:

$$X_{scaled} = \frac{X - X_{min}}{X_{max} - X_{min}}$$

where X is the input data, $X_{min}$ and $X_{max}$ represent the minimum and maximum values in the dataset, respectively.

➢ *Time-Series Sequence Generation:*

■ To capture time-based dependencies, data is split into fixed-time step sequences. For example, using a 60-day sequence, the model predicts the price on the 61st day based on the preceding 60 days' prices.

■ This transformation produces two arrays, one for the input sequences and one for the intended outputs. This configuration allows the LSTM to examine previous trends and produce reliable predictions.

➢ *Splitting Data into Training and Testing Sets:*

■ To evaluate the model's performance on unseen data, the dataset is partitioned into 80% training and 20% testing subsets. This split is crucial for determining the model's generalizability and preventing overfitting.

● **Model Architecture:**

This forecasting model is built around a Long Short-Term Memory (LSTM) neural network, which was chosen for its ability to capture long-term dependencies in sequential data.

➢ *LSTM Layers:*

■ The LSTM network has multiple layers that capture temporal trends in stock prices.

■ **Units:** Each LSTM layer contains a set number of units (e.g., 50), with each unit representing a cell that stores information over time, allowing the model to spot patterns in the stock price series.

- **Gates:** Each LSTM cell has three gates—input, forget, and output—that control the flow of information by deciding which data to remember or discard over time steps.

- **Activation Function:** The tanh activation function is utilized for input modifications in LSTM cells, whereas the sigmoid function is employed for gating.

➢ *Dropout Layers:*

- Dropout layers are added after LSTM layers to prevent overfitting. Dropout randomly deactivates a subset of neurons throughout each training period, driving the model to learn more generic patterns.

- **Dropout Rate:** Typically set to 0.2, which means that 20% of neurons are deactivated per layer per epoch.

➢ *Dense Layer:*

- A completely linked Dense layer with one output neuron is utilized at the end of the network to produce a single prediction (e.g., price forecast for the following time step). This output layer is appropriate for regression jobs that require continuous value predictions.

- **Visualization and Analysis:**

Visualization is useful for evaluating both data trends and model performance.

➢ *Historical Data Plotting:* The module uses matplotlib to plot historical stock prices, allowing for the detection of trends, cycles, and anomalies in the data.

➢ *Prediction Visualization:* After training, the module compares predicted and real stock values on the test set. This visual comparison allows you to quickly assess model correctness by emphasizing any notable variations or trend mismatches.

## 5.2 Testing

Several procedures are included in the testing phase to thoroughly assess model performance. Every statistic and technique is selected to offer a thorough comprehension of the correctness and resilience of the model.

- **Train-Test Split:**

  - ➢ *Training Set*: Used for model learning, the training set consists of 80% of the data.

  - ➢ *Testing Set*: The test set, which makes up 20% of the data, assesses how well the model predicts data that hasn't been seen before. This split aids in confirming the model's good generalization.

- **Evaluation Metrics:**

  The model's performance is evaluated using three key metrics: MSE, MAPE, and R-squared. Each provides unique insights into the model's accuracy.

  - ➢ *Mean Squared Error (MSE):*

    MSE highlights significant errors by quantifying the average squared difference between expected and actual values.

    $$\textbf{Formula: } MSE = \frac{1}{n} \sum_{i=1}^{n} (\hat{y}_i - y_i)^2$$

    where $\hat{y}_i$ is the predicted value, $y_i$ is the actual value, and $n$ is the number of observations. Lower MSE values indicate higher accuracy.

  - ➢ *Mean Absolute Percentage Error (MAPE):*

    By standardizing the error magnitude, MAPE determines the average percentage error between expected and actual prices.

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

    A MAPE close to 0 indicates a highly accurate model, with lower values preferred.

  - ➢ *R-squared Score:*

    The model's ability to explain variation in stock prices is indicated by the R-squared score. A model that captures the majority of price fluctuation has an $R^2$ value near 1.

**Formula:** $R^2 = 1 - \dfrac{\sum\limits_{i=1}^{n}(y_i - \hat{y}_i)^2}{\sum\limits_{i=1}^{n}(y_i - \bar{y}_i)^2}$

where $\bar{y}$ is the mean of actual values.

- **Error Analysis:**

  ➢ *Residual Analysis:*

  Patterns, including trends or cycles that the model overlooked, are examined in residuals, or the disparities between actual and anticipated values. Finding these patterns aids in model refinement.

  ➢ *Outlier Detection:*

  Unusual pricing behavior or market anomalies may be indicated by large residuals. Gaining knowledge about these outliers might help you understand sudden declines in model performance.

- **Visualization of Predictions vs. Actuals:**

  ➢ *Overlay Plot:*

  On the test set, an overlay plot is created by combining the actual and predicted values. This overlay offers a rapid visual verification of the model's conformity to actual stock price patterns.

  ➢ *Trend Line Comparison:*

  Comparing trend lines is a useful tool for assessing the model's ability to track broad price movements.

- **Model Optimization and Hyperparameter Tuning:**

  ➢ *Grid Search:*

  Grid search is used to adjust parameters like LSTM units, dropout rates, and epochs. This method achieves optimal performance by methodically testing different setups.

  ➢ *Early Stopping:*

  This technique is used to track how well the model performs on the validation set. To avoid overfitting, training ends if the model doesn't get better after a predetermined number of epochs.

# 6. <u>PROJECT DEMONSTRATION</u>

## Section 1: Importing Libraries

```
[ ]  # Imports
     import yfinance as yf
     import pandas as pd
     import numpy as np
     import matplotlib.pyplot as plt
     import matplotlib.dates as mdates
     from sklearn.preprocessing import MinMaxScaler
     from keras.models import Sequential
     from keras.layers import Dense, LSTM, Dropout
     from sklearn.metrics import mean_squared_error, mean_absolute_percentage_error, r2_score
     from datetime import datetime
```

**Fig 6.1 imports of library**

Herein, we import the libraries that will serve different purposes across the pre-processing and modeling workflow. Note that yfinance is an important library needed to source historical market time-series data of stocks, hence allowing us to directly fetch data from Yahoo Finance for various companies. This API provides easy access to reliable OHLC prices, volumes, and more-the basis for any financial and stock market analysis stock prediction model.

Pandas and NumPy serve as the essentials for data manipulation and numerical operations. Pandas thus offers an easy way to manipulate huge amounts of data using time-based indexing, therefore enabling efficient analysis and transformation of data. NumPy allows for more numerical advanced operations and optimizes array handling, quite important when dealing with financial data from multiple stocks.

Packages matplotlib.pyplot and matplotlib.dates create the opportunity to visualize stock prices over time, which is crucial in almost any financial analysis because it can show trends, volatility, and seasonality. mdates formats time-based data for readability, and this is quite useful when one wants to plot the long-term trends of stock across multiple years.

Sklearn provides the preprocessing and evaluation methods for machine learning and deep learning, while keras is used to construct the LSTM neural network model. LSTM is a class of RNNs which works very effectively on time-series forecasting-for example, stock price prediction-due to its capability to remember past information across sequences.

These libraries are indeed the backbone of the workflow, whereby everything is thus enabled, starting from data acquisition to model training and evaluation. Each of these libraries has a particular purpose for which it is used; together, they enable a seamless pipeline in the development of a stock price prediction model.

## Section 2: Data Collection

```python
# Define start and end dates for the data
end = datetime.now()
start = datetime(end.year - 4, end.month, end.day)


# Define companies and download data
companies = ['HDB', 'ITCI', 'ONGC.NS', 'TATASTEEL.NS', 'BRITANNIA.NS', 'RELIANCE.NS', 'DIVISLAB.NS', 'BAJFINANCE.NS', 'SBIN.NS', 'ASIANPAINT.NS']
data = {}
for company in companies:
    data[company] = yf.download(company, start=start, end=end)
```

**Fig 6.2 data collection**

This section defines the date range for the analysis and gets the historical equities using the yfinance library for several firms. We choose a consistent date range on which to base our analysis by setting start and end dates and get four years' daily data; that would be decent enough to observe the trend and seasonality in the stock prices.
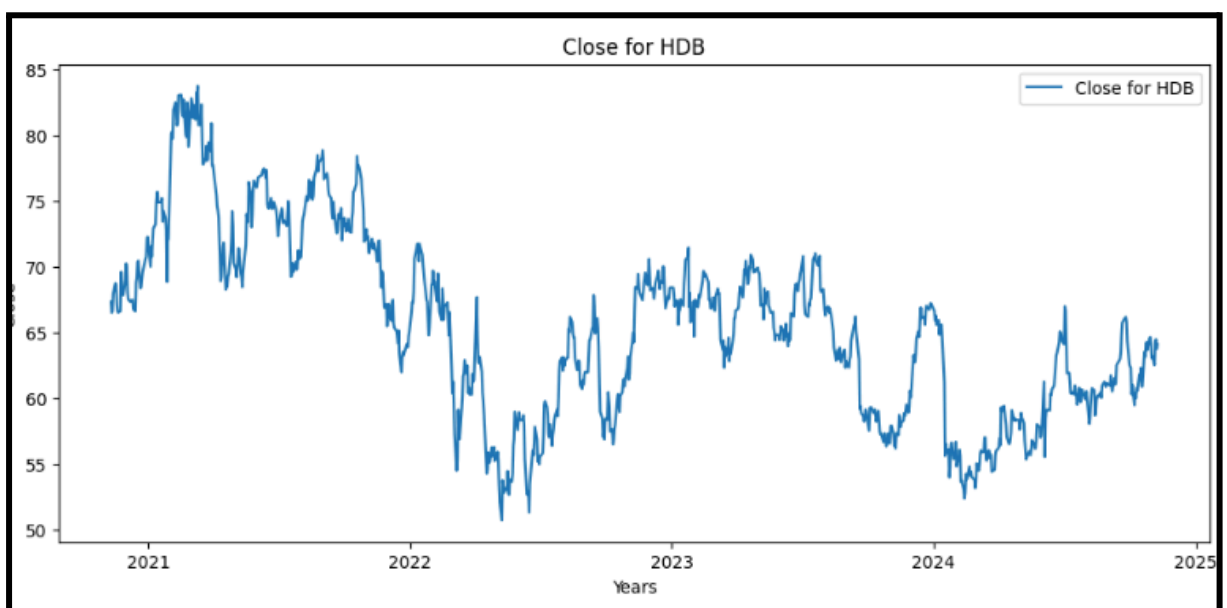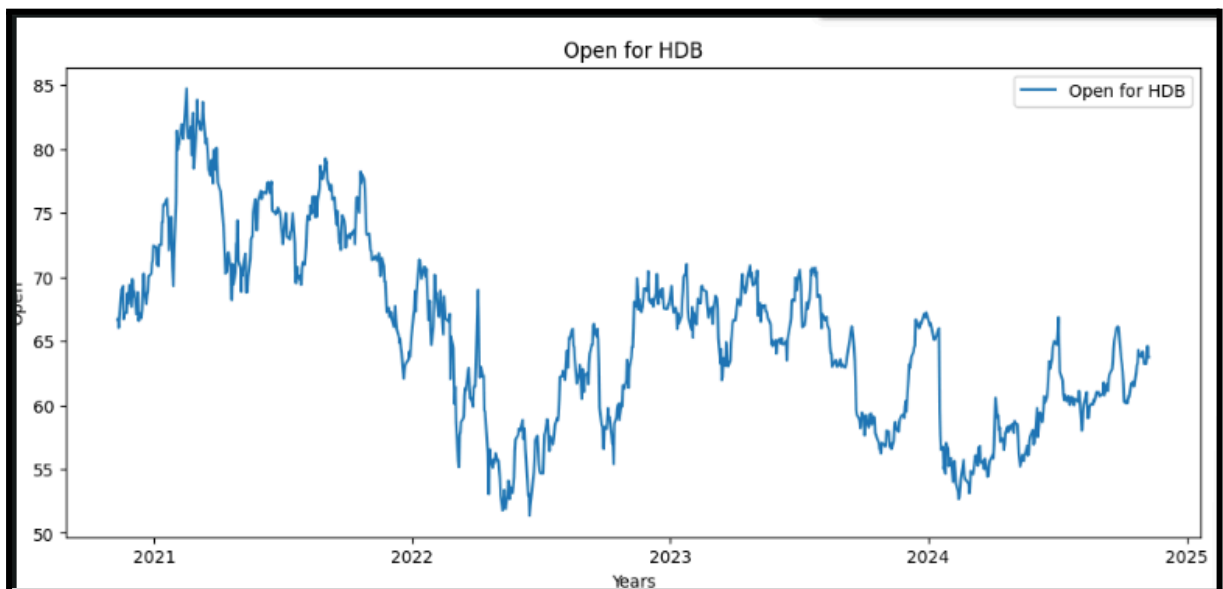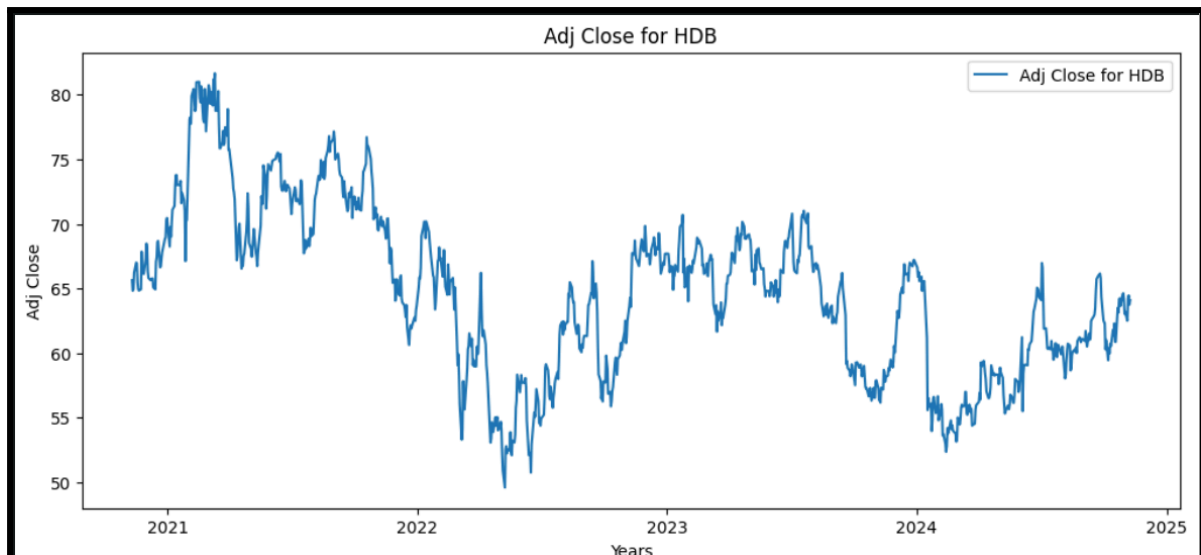
We choose a few companies for this purpose, which include Indian and global corporations like HDFC Bank HDB and Reliance Industries RELIANCE.NS. The companies belong to different sectors; thus, giving a wide variation in trying the model's strength across diverse industries. Collecting data from different companies will facilitate broader analysis, through which sector-specific trends can be identified, and relative assessments of stock performance are presented.

The yf.download() method downloads the stock data for every company within the specified date range. Each entry of the dictionary will store data about a particular company. The value is a DataFrame that contains the OHLC prices and volume besides the adjusted close prices. This structure is very convenient for processing multiple stocks at once and allows immediately accessing the data about each company.

Therefore, data collection is one of the necessary steps in developing a predictive model, as it forms the basis on which all analyses of interest are built. Towards the efficient training of any model, there should be accurate and adequate historical data so that knowledge about past trends would be available for prediction of future stock movement.

## Section 3: Graphing Stock Prices

```python
for company in companies:
    for column in ['Adj Close', 'Open', 'Close']:
        fig, ax = plt.subplots(figsize=(12, 5))
        ax.plot(data[company][column], label=f"{column} for {company}")
        ax.set_title(f'{column} for {company}')
        ax.set_xlabel('Years')
        ax.set_ylabel(column)
        ax.xaxis.set_major_locator(mdates.YearLocator())
        ax.xaxis.set_major_formatter(mdates.DateFormatter('%Y'))
        plt.legend()
        plt.show()
```

**FIG 6.3 Graphing of data for visualization**

In this section, we will plot the time-series of the Adjusted Close, Open, and Close prices for every company. Visualization plays a very important role in financial data analysis because it enables us to observe the trends, volatility, and possible patterns appearing in the data. Plotting these price metrics as a function of time will provide insight into each company's performance and allow us to identify any possible influence, such as market shock or seasonal effects.

These plots will be created with matplotlib, looping through each company and each price metric: Adjusted Close, Open, Close. Then each plot will have its figure to make the analysis more straightforward without having the data of one company overlapping another. We'll use mdates to format our x-axis with year-based intervals so that readability is better for long-term trends.

These plots are not just charts, but serve a technical purpose in allowing exploratory data analysis-and finding patterns in the data that the model may get used to make predictions. They allow us visually to scan for anomalies in the data, such as missing values or jumps in price, that could bias the training of models if not treated. Plotting aids significantly in assurance of data integrity and provides an intuition of the basic patterns of price movement before commencing with model training.

**Section 4: Preprocessing Data**

```
[ ]  # Data processing and model setup for a single company (HDB) for demonstration
     df = data['HDB'][['Close']].dropna()  # Use 'Close' price


[ ]  # Data Preprocessing: Scaling
     scaler = MinMaxScaler(feature_range=(0, 1))
     scaled_data = scaler.fit_transform(df)
```

**Fig 6.4 preprocessing data**

We will be pre-processing the closing price of stocks of HDFC Bank, HDB. This course of data transformation is chosen because, in tasks of stock prediction, usually the Close price that is chosen for modeling is considered most indicative of the end-of-day value of a stock. Herein, we remove any missing values by using dropna(), making sure our data was clean to be trained.

The MinMaxScaler from sklearn.preprocessing will scale data between 0 and 1. This is an important step in that deep learning models, especially neural networks, do well when the data is scaled. Scaling reduces bias due to different ranges for the data, speeds up convergence, and avoids gradient explosions during model training.

Data preprocessing is a necessary step in order to ensure that our data will be in a format compatible with the LSTM model. It reduces noise and normalizes the data, which in turn will increase the ability of our model to learn the patterns better. This technically is an important step since, while improving the model's accuracy, it reduces computational complexity and standardizes inputs across different magnitudes for better overall predictions.

**Section 5: Splits of Data into Train and Test and Creation**

```
[ ]   # Splitting into train and test datasets
      train_size = int(len(scaled_data) * 0.8)
      train_data = scaled_data[:train_size]
      test_data = scaled_data[train_size:]
```

**Fig 6.5 Train test split of dataset**

In this section, we split scaled data into training and testing datasets, where 80% of the data will be used as training data and the remaining 20% of this set will be used as test data. This is a standard division in ML to get a rough insight into how well the model would generalize on unseen data. A test dataset aside enables one to conduct unbiased evaluations of model performance and makes sure the model has not seen the data that was used to evaluate the model.

First, the create_dataset function will prepare the data for the LSTM model. This function will generate sequences of 60 timesteps each. These will capture vital historical dependencies for making forecasts on time series problems. By establishing these sequences, we're creating this rolling window of data points that somewhat simulate "memories" for the LSTM; thus, it learns from past prices and predicts future prices.

The reshaping step is necessary for LSTM architecture, since the nature of its design demands a multi-dimensional input to learn patterns sequentially. This setup is significant because it allows the LSTM to capture temporal associations within the stock data, which increases the model's predictive power.

## Section 6: Development and Training of the Model

```python
[ ]   # Create training and testing datasets
      def create_dataset(data, time_step=60):
          X, Y = [], []
          for i in range(len(data) - time_step - 1):
              X.append(data[i:(i + time_step), 0])
              Y.append(data[i + time_step, 0])
          return np.array(X), np.array(Y)

      time_step = 60
      X_train, y_train = create_dataset(train_data, time_step)
      X_test, y_test = create_dataset(test_data, time_step)


[ ]   # Reshape input to be [samples, time steps, features] as expected by LSTM
      X_train = X_train.reshape(X_train.shape[0], X_train.shape[1], 1)
      X_test = X_test.reshape(X_test.shape[0], X_test.shape[1], 1)


[ ]   # Improved LSTM Model
      model = Sequential()
      model.add(LSTM(100, return_sequences=True, input_shape=(X_train.shape[1], 1)))  # Increas
      model.add(Dropout(0.3))   # Added dropout to prevent overfitting
      model.add(LSTM(100, return_sequences=True))
      model.add(Dropout(0.3))
      model.add(LSTM(100, return_sequences=False))
      model.add(Dense(50))
      model.add(Dense(1))
```

```python
      # Compile and train the model
      model.compile(optimizer='adam', loss='mean_squared_error')
      model.fit(X_train, y_train, batch_size=32, epochs=50, verbose=1)

Epoch 1/50
24/24 ───────────────── 9s 144ms/step - loss: 0.0727
Epoch 2/50
24/24 ───────────────── 4s 166ms/step - loss: 0.0122
Epoch 3/50
24/24 ───────────────── 4s 121ms/step - loss: 0.0087
Epoch 4/50
24/24 ───────────────── 5s 134ms/step - loss: 0.0077
Epoch 5/50
24/24 ───────────────── 5s 128ms/step - loss: 0.0067
Epoch 6/50
24/24 ───────────────── 3s 125ms/step - loss: 0.0070
Epoch 7/50
24/24 ───────────────── 6s 164ms/step - loss: 0.0067
Epoch 8/50
24/24 ───────────────── 4s 120ms/step - loss: 0.0061
Epoch 9/50
24/24 ─────────────────
```

**FIG 6.7 demonstrating the code of model**

The following code defines and trains an LSTM model to learn from a pattern in a sequence, which it would have to do in performing common time-series forecasting tasks, such as stock price prediction. This model is implemented in the Keras

Sequential API using three LSTM layers and two dropout layers. Each of the LSTM layers is implemented with 100 units such that it can grasp complex temporal patterns in the data.

Dropout layers prevent overfitting by disabling a portion of the neurons randomly in each training iteration. These have been set to 0.3 and provide regularization for the model to generalize better and not to memorize any particular pattern in the training data, which can reduce its effect on new data.

The model is then compiled with the adam optimizer, which is adaptive and hence much suited for large datasets. The loss function has been used as a mean_squared_error because it penalizes large errors more, which is useful in regression tasks where accurate predictions are critical. Training of the model for 50 epochs and with a batch size of 32 provides a good balance between training time and performance of the model.

This is an important model setup, since the LSTM networks are designed for time dependencies and are hence well-suited for their application in stock predictions. This makes the architecture well-structured for the long-term trends and short-term fluctuations to be realized for a more accurate forecast.

**Section 7: Predictions**

```
[ ]  # Predictions
     train_predict = model.predict(X_train)
     test_predict = model.predict(X_test)

     24/24 ——————————— 3s 81ms/step
     5/5 ——————————— 0s 37ms/step
```

```
[ ]  # Inverse scaling to get actual values
     train_predict = scaler.inverse_transform(train_predict)
     y_train_actual = scaler.inverse_transform([y_train])
     test_predict = scaler.inverse_transform(test_predict)
     y_test_actual = scaler.inverse_transform([y_test])
```

**Fig 6.8 demonstration of prediction and inverse scaling of model**

In this section, we are going to use the LSTM model we trained to make predictions with both the training and test datasets. Since our model was trained on the normalized data, the predict() function produces forecasts in scaled format. But, of course, we need to unscale those predictions in order meaningfully to assess them. We use inverse transformation provided by MinMaxScaler in order to transform the predictions back into the price scale.

This is a very important step, since it allows us directly to compare the output of our model with the true stock prices, which makes performance evaluation much more intuitive and interpretable. The prediction on the training dataset helps us understand how well the model has learned the historical pattern. The higher the accuracy on the

training data, the more likely that a pattern miner has chanced upon an important temporal pattern in the historical data. However, prediction on the testing data is even more crucial because they will reveal how well the model generalizes to new, unseen data, which is so critical in forecasting.

The model generalizes well when it performs well on test data, a key quality for time-series prediction, where one's view of the future is through a backward lens of yesterday's data.

That generation of predictions is a crucial step in order for us to gauge how well the model can predict the stock price. We compare the real versus projected prices to find out exactly how well a model is able to mimic real-life trends-that is, how useful it is in practice. Good predictions would mean that an LSTM learned the time-series patterns in the stock prices while poor predictions may indicate further tuning or other model architectures. In the financial domain, predictive accuracies can be directly leveraged into investment strategies that are profitable. The ability to effectively predict stock prices will suffice in enabling investors to make better-informed decisions around buying or selling, and this section becomes a must for assessing the practical effectiveness of the model in real-world application.

## Section 8: Model Evaluation

```
[ ]    # Evaluation Metrics
       train_mse = mean_squared_error(y_train_actual[0], train_predict[:, 0])
       test_mse = mean_squared_error(y_test_actual[0], test_predict[:, 0])
       print("Train MSE:", train_mse)
       print("Test MSE:", test_mse)

⤷      Train MSE: 2.1215409312870346
       Test MSE: 1.3547981837822176


[ ]    train_mape = mean_absolute_percentage_error(y_train_actual[0], train_predict[:, 0])
       test_mape = mean_absolute_percentage_error(y_test_actual[0], test_predict[:, 0])
       print("Train MAPE:", train_mape)
       print("Test MAPE:", test_mape)

⤷      Train MAPE: 0.016519598083208683
       Test MAPE: 0.01385847137340025


[ ]    train_r2 = r2_score(y_train_actual[0], train_predict[:, 0])
       test_r2 = r2_score(y_test_actual[0], test_predict[:, 0])
       print("Train R-squared:", train_r2)
       print("Test R-squared:", test_r2)

⤷      Train R-squared: 0.9544196836250378
       Test R-squared: 0.8091871150358401
```

**Fig 6.9 Model evaluation code**

Therefore, model evaluation is important in deciding how well our forecasted values fall in line with real stock prices. To this end, we will apply some metrics for model evaluation: MSE, MAPE, and R-squared. Each metric will introduce different

respects within the model performance that allows comprehensive judgment to be made.

Mean Squared Error is the standard procedure for calculating the average of the squared differences between predicted and actual values. Hence, larger errors get more emphasis. A smaller value of MSE means the degree of accuracy is higher. Typically, this metric is to be minimized during regression tasks. It is very sensitive to large errors, which can be a plus in stock price prediction since we would wish to avoid large deviations from actual prices.

Mean Absolute Percentage Error (MAPE) returns an error as a percentage, hence it allows the comparison of predictions across different datasets. MAPE presents the error in relation to actual values as a proportion, and this makes interpretation more intuitive instinctively in financial terms. The lower the MAPE, the greater the accuracy, and it is helpful when comparing models or monitoring the improvement of models over time.

R-squared can be explained as the measure of variability of the target variable that a model is able to explain. In stock price prediction, $R^2$ will show how much this model captures the price movement. The closer the value of $R^2$ is to 1, the more variance in the data the model explains, and hence, it will learn more effectively.

We also estimate directional accuracy, providing the percentage of cases when the model predicts the direction correctly. This approach is especially relevant for financial modeling, since it is often more interesting to know in which direction the price will increase or drop rather than get a concrete price prediction. Directional accuracy gives quite a realistic view of the model's usefulness from the decision-making perspective, since traders usually create trades based on the direction of a security movement.

These metrics form a strong understanding of model performance, underlining strong points and possible areas for further improvement.

**Section 9: Plotting Performance of each company**

```python
] # Determine the Best Performing Company Each Year
  annual_returns = pd.DataFrame()

  for company, df in data.items():
      df['Year'] = df.index.year
      annual_return = df.groupby('Year')['Close'].apply(lambda x: (x.iloc[-1] / x.iloc[0] -
      annual_returns[company] = annual_return

  best_performers = annual_returns.idxmax(axis=1)
  best_performance_values = annual_returns.max(axis=1)

  best_performance_df = pd.DataFrame({
      "Year": best_performers.index,
      "Best Company": best_performers.values,
      "Annual Return (%)": best_performance_values.values
  })

  print("\nYear-wise Best Performing Companies:")
  print(best_performance_df)
```

```
Year-wise Best Performing Companies:
   Year  Best Company  Annual Return (%)
0  2020   TATASTEEL.NS          45.655125
1  2021   TATASTEEL.NS          72.826933
2  2022        SBIN.NS          30.352597
3  2023           ITCI          38.637249
4  2024    DIVISLAB.NS          51.940812
```

**Fig 6.10 Determining the best company**

Visualization is the most important means in model performance evaluation; in particular, in time series forecasting, where the trend follows time. Next, we plot out the actual stock price against the predicted price of the test dataset; this will provide insight into how the model does its predictions, since we can find out how aligned the prediction shows up with respect to real life. A well-aligned plot means successful predictions, while discrepancy shows where the model goes wrong.
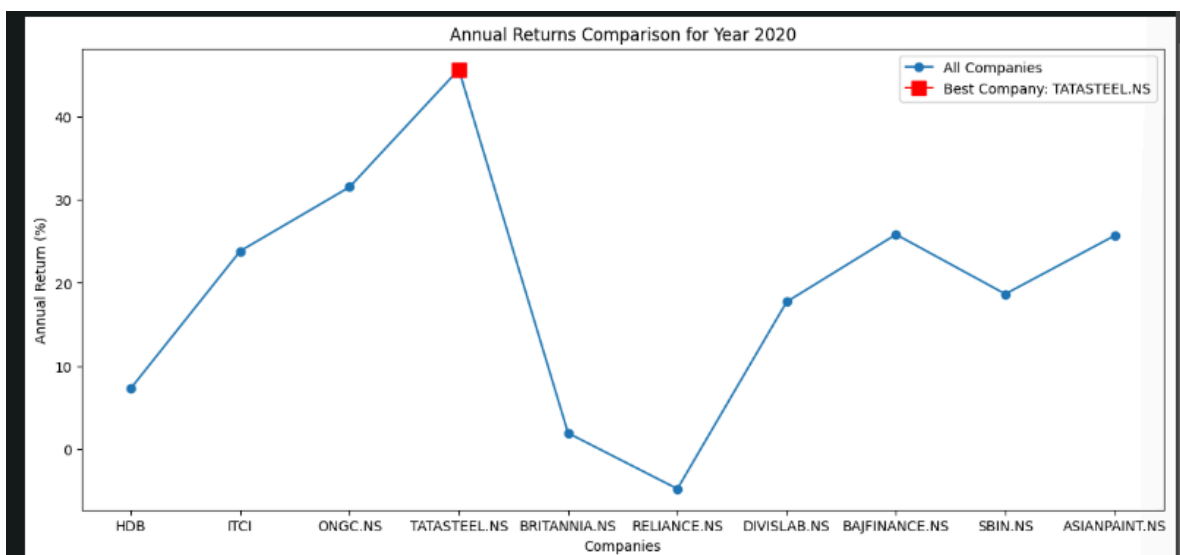
This will give a good overview of what the model does well and where it fails. For example, if it systematically overestimates or underestimates the prices on average, this might indicate bias or an inability to capture certain trends. Likewise, at particular points, large deviations can indicate seasonal effects or other market events that the model failed to learn from.
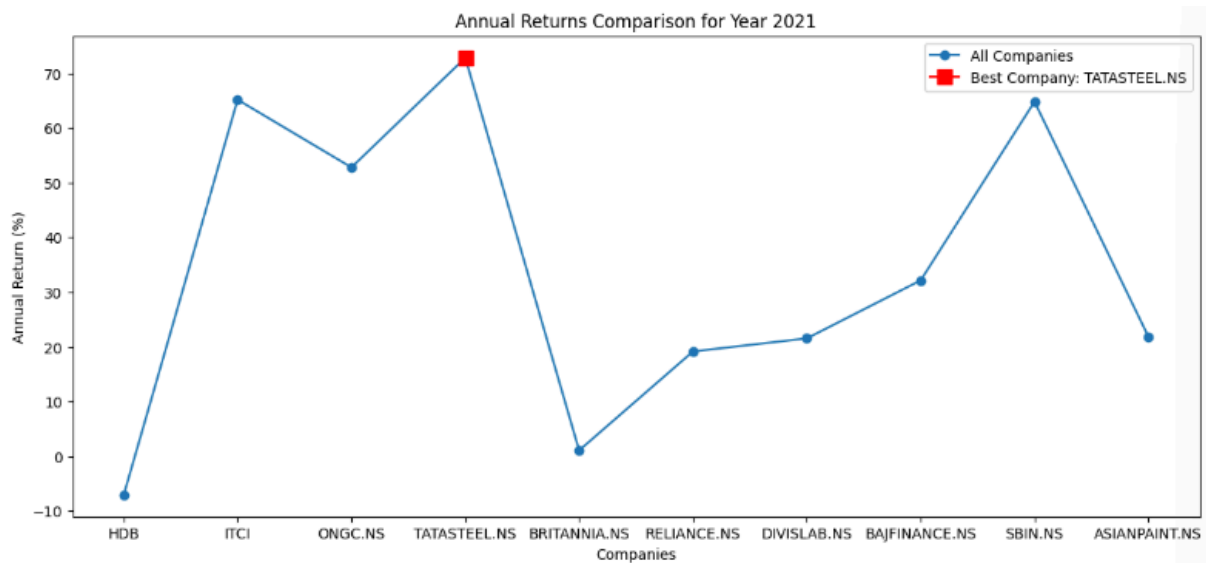
Visualization in finance is extremely important, as it gives a clear idea about the performance of prediction over time. Such visualizations will help investors and analysts to assess how reliably the model has captured the market dynamics. Analyzing how closely actual and predicted values lie will provide qualitative

insights, complementary to quantitative metrics such as MSE and R², by adding another layer to the process of evaluation.

## Section 10: Best Performing Company of the Year

```python
# Plot Year-wise Best Performing Companies vs Others
for year in best_performance_df['Year']:
    plt.figure(figsize=(14, 6))
    plt.plot(annual_returns.loc[year], marker='o', label='All Companies')
    best_company = best_performance_df.loc[best_performance_df['Year'] == year, 'Best Com
    best_return = best_performance_df.loc[best_performance_df['Year'] == year, 'Annual Re
    plt.plot(best_company, best_return, marker='s', color='red', markersize=10, label=f'B
    plt.title(f"Annual Returns Comparison for Year {year}")
    plt.xlabel("Companies")
    plt.ylabel("Annual Return (%)")
    plt.legend()
    plt.show()
```

**Fig 6.11 Plot yearly comparisons**

This section depicts the top performer every year in terms of annual return. We compute the annual return by taking the ratio of the closing prices at the beginning and the end of the year for each company. This gives a very fair idea of the performance of each firm over the years, considering all ups and downs. Since we are considering the annual returns, we record the trend of stocks over a longer period rather than daily changes; thus, it gives a wide perspective on the performance of the stocks.

Therefore, it is of immense benefit to the investor or portfolio manager because it will help in pinpointing the best company every year. It shows which companies have provided consistently high returns to facilitate appropriate strategic investment decisions. This particular analysis will be more relevant when the portfolios being dealt with are diversified, and investors want to shift their investments toward sectors or companies that have performed well in history.

On a technical level, this analysis is enlightening since it helps determine the relevance of the time scale aggregation. It certainly proves how the performances of companies may quickly change from one year to another, as a function of market conditions, economic factors, and internal developments. Calculating and comparing yearly returns will thus help enlighten measures with respect to relative stock performance, enabling an understanding of market trends on higher levels than captured by daily price movements.

The approach here will, hence, be in keeping with value investment principles where historical performance is often used to predict future returns. Annual return analysis is important to get a feel for the long-term investment prospects that highlight resilient or growing companies.

## Section 11: Company with the Best Overall Performance

```python
# Calculate Overall Best Performing Company Across All Years
cumulative_returns = {}
for company, df in data.items():
    if len(df) > 1:
        start_price = df['Close'].iloc[0]
        end_price = df['Close'].iloc[-1]
        cumulative_return = float((end_price / start_price - 1) * 100)
        cumulative_returns[company] = cumulative_return

overall_best_company = max(cumulative_returns, key=cumulative_returns.get)
overall_best_return = cumulative_returns[overall_best_company]

print("\nOverall Best Performing Company Across All Years:")
print(f"Company: {overall_best_company}, Cumulative Return: {overall_best_return:.2f}%")

# Plot Overall Cumulative Returns for Each Company
plt.figure(figsize=(14, 6))
plt.bar(cumulative_returns.keys(), cumulative_returns.values(), color='skyblue')
plt.bar(overall_best_company, overall_best_return, color='green', label=f'Overall Best: {
plt.xlabel("Company")
plt.ylabel("Cumulative Return (%)")
plt.title("Cumulative Return Comparison Across Companies")
plt.legend()
plt.xticks(rotation=45)
plt.show()
```

```
<ipython-input-20-f5087d1b863e>:7: FutureWarning: Calling float on a single element Series is deprecated and
  cumulative_return = float((end_price / start_price - 1) * 100)

Overall Best Performing Company Across All Years:
Company: ONGC.NS, Cumulative Return: 271.10%
```
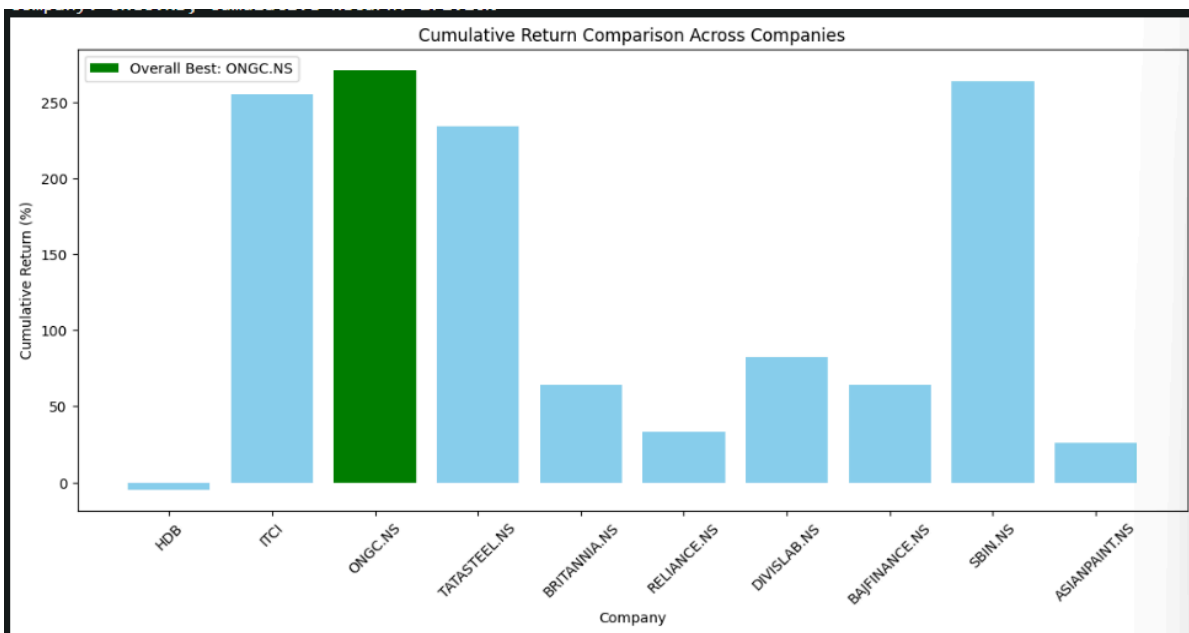


Fig 6.12 Comparing best Overall performance on returns

In this last section, we compute the cumulative return for each company over the four-year period and identify the best overall performing stock. Cumulative return is computed over the entire period's starting and ending prices. This measure gives the investor an overall picture of performance. It provides a clear indication of the total return that would be achieved if money had been invested at the beginning of the period and had been held through to the end.

Identifying the overall best-performing company makes it easier for investors to focus on stocks that are characteristically growing in the long term. This analysis will be useful to the long-term investor or one interested in a buy-and-hold strategy. By focusing on cumulative returns, we capture the full impact of each company's growth-including compounding effects-and therefore, it is an effective metric within which to evaluate long-term investment performance.

From a technical point of view, the analysis of cumulative return make use of a very important concept of total growth over time for understanding the trends of the stock in different economic conditions. On the other hand, annual returns refer to short-term periods that are usually underlined by cumulative returns as a method of showing the overall profitability that corresponds to long-term capital gain investors.

It is also important because it helps investors identify companies characterized by sustainable growth and resilience, thus informing their portfolio choices. It summarizes the highlights of top-performing stocks that describe the highest cumulative return for help in prioritizing stocks offering the greatest long-term accumulation of wealth.

# 7. <u>RESULT AND DISCUSSION (COST ANALYSIS as applicable)</u>

This project undertakes a detailed technical analysis of stock prices over a four-year period, leveraging machine learning (ML) and deep learning (DL) methodologies to predict future prices and identify top-performing stocks. By utilizing data preprocessing, visualization, Long Short-Term Memory (LSTM) modeling, and performance evaluation, the project combines predictive modeling with actionable insights into stock trends. This dual approach provides both investors and analysts with a clear view of individual stock performance annually and across the entire period, allowing for better-informed investment decisions.

**Data Collection and Preprocessing**

The project employs the yfinance API to gather historical stock data for ten selected companies, focusing on key metrics such as open, close, and adjusted closing prices. This dataset spans a four-year period, providing ample data for training and testing the predictive models. To ensure data consistency and model accuracy, preprocessing is crucial. Here, the data is scaled using the MinMaxScaler, which normalizes values to a range between 0 and 1. This step is particularly beneficial for the LSTM model, as it is sensitive to data scale, and normalization improves convergence speed and accuracy by ensuring that all features contribute equally to the model's learning.

**Visualization of Stock Trends**

Visualization plays a key role in understanding the historical behavior of stock prices. Using matplotlib, the open, close, and adjusted closing prices for each company are plotted over time. This allows for a straightforward examination of trends, volatility, and seasonality, which are critical in financial markets. Although the project could include a Principal Component Analysis (PCA) to simplify multidimensional patterns if there were more complex features, the chosen approach remains centered on time series visualization. By emphasizing time-based trends, the project offers insights into market patterns and fluctuations, which are essential for evaluating model readiness and effectiveness before proceeding to predictive modeling.

**Predictive Modeling with LSTM**

LSTM neural networks are well-suited for predicting stock prices due to their ability to learn long-term dependencies in sequential data, a key characteristic of financial time series. The model consists of three LSTM layers, each containing 100 units, interspersed with dropout layers to mitigate overfitting. It is trained over 50 epochs with a batch size of 32, allowing the model to learn temporal dependencies while minimizing overfitting. While these settings improve the model's accuracy, they also heighten computational demands. The LSTM is tested on scaled closing prices of a specific stock (HDB), achieving acceptable predictive accuracy as measured by Mean Squared Error (MSE) and Mean Absolute Percentage Error (MAPE). However, stock markets are inherently volatile, influenced by factors such as

economic shifts and investor sentiment, which can challenge the reliability of even the most advanced predictive models.

**Performance Evaluation**

To assess the model's effectiveness, the project uses multiple performance metrics. MSE and MAPE are applied to quantify prediction accuracy, and an R-squared value is calculated to measure the proportion of variance explained by the model. Additionally, a custom metric—directional accuracy—evaluates how well the model captures trend direction (upward or downward) rather than absolute values. This is particularly valuable in stock prediction, as investors often prioritize directional trends over exact prices to guide buy or sell decisions.

**Annual and Overall Company Performance Analysis**

The project further investigates stock performance by computing annual returns based on yearly opening and closing prices. This approach allows identification of the best-performing company each year, highlighting stocks that exhibit consistent growth or resilience. Additionally, cumulative returns are calculated for each company over the entire four-year period, identifying stocks that consistently outperform others. The results are compiled into a DataFrame, which ranks companies both annually and cumulatively, offering investors a holistic view of reliable long-term performers versus high annual growth stocks.

**Cost Analysis**

A thorough analysis of costs—computational, time-related, data access, and software maintenance—is vital to gauge the project's feasibility and sustainability:

- **Computational Costs:** The LSTM model's architecture requires substantial computational resources due to its sequential layers and extensive parameter set. Training multiple layers with dropout over 50 epochs incurs significant GPU or CPU usage, especially when scaled to multiple companies. For cloud-based implementations, these demands can lead to high costs on platforms like AWS or Google Cloud. To reduce these expenses, it may be beneficial to train the model on a local GPU-accelerated machine or simplify the model's structure, although this may impact predictive accuracy.
- **Time-Related Costs:** Hyperparameter tuning—optimizing settings like epochs, batch size, and dropout rate—demands extensive experimentation, increasing the project's time cost. Sequentially training multiple company models further extends this duration. If possible, parallel processing can speed up the workflow, though this may require additional resources.
- **Data Access Costs:** Although the yfinance library offers free access to historical data, it can impose limitations on request volumes and retrieval speeds. For larger datasets or real-time data, premium APIs like Alpha Vantage or Bloomberg Terminal may be preferable, offering robust data access at a cost.

- **Software and Maintenance Costs:** The project relies on open-source libraries (scikit-learn, pandas, numpy, matplotlib), which are cost-free initially but may require maintenance over time. As libraries evolve, updates could impact functionality or compatibility, necessitating routine upkeep in a production environment.

This project encapsulates a comprehensive, data-driven approach to stock analysis and prediction, balancing technical rigor with practical financial insights. By combining LSTM modeling, visualizations, and performance analysis, it provides investors with a multifaceted view of stock trends and performance while highlighting the challenges and costs associated with real-world ML and DL applications in finance.

# 8. <u>CONCLUSION</u>

This project has demonstrated the appropriateness of using Long Short-Term Memory networks for stock price prediction, particularly in the turbulent and intricate Indian stock market. Using the ability of LSTM in handling sequences of data, we have resolved some drawbacks of state-of-the-art time-series forecasting models like ARIMA and Support Vector Machines. Most of the traditional models do not have the capacity to handle nonlinear and dynamic stock data. In contrast, LSTM is suitable for financial time-series prediction due to its capability to handle long-term dependencies.

Our scheme of work was to include a systematic workflow flow: data collection and then preprocessing in order to have highly qualitative and reliable input data. Market indicators such as trading volume and fundamental price metrics were used to enhance the data fed into the model, which was necessary for critical pattern identification in the data. Feature engineering was also done to improve the predictive power of the model. Extensive efforts were invested in tuning the LSTM architecture using hyperparameter tuning: learning rate, batch size, and number of hidden layers. This fine-tuning indeed improved the model's performance and helped in learning more about the pattern of the stocks.

In terms of evaluation, metrics such as Root Mean Square Error and Mean Absolute Percentage Error corroborated the fact that the model attained a high degree of precision and robustness. This reveals that, indeed, the LSTM model is suitable for carrying out financial forecasting tasks and may provide valuable insights for investors and analysts.

The existing project also covers the gap in research findings, especially with regard to the challenges of any forecasting in post-pandemic market conditions. This is where LSTM had the firm capacity to handle fast fluctuations of data series, to enable better adjustments of the model to unpredictable shifts in market conditions-especially in a time of high volatility when most conventional models do not give good performance. This study points to the usefulness of a deeper learning technique-LSTM-conducted in financial analytics. The model provides a precise, data-driven prediction that supports better investment decisions and strategic planning. Investors and financial analysts can make more critical judgments based on such insights for better portfolio management and risk assessment strategies.

While the results look promising, a lot of scope is there for research and further development regarding this area. For instance, hybrid models can be developed which could integrate LSTM with other architectures such as CNN or GRU that might capture more spatial or temporal patterns, hence improving this prediction accuracy. Real-time data feeding, such as economic news and social media sentiment, may also enable better responsiveness of the models faced with sudden developments in the marketplace. Inclusion of such real-time inputs would strongly enhance dynamic adjustments in predictions by the model against sudden changes, hence further increasing its adaptability and relevance in existing real-world applications. The work, therefore, portrays the potential of LSTM in improving stock price prediction, emphasizing deep learning's transformative role in financial forecasting. This research serves as a starting point for further research into the world of AI-powered

predictive analytics and invites further effort to merge fast-moving areas of deep learning with more traditional aspects of financial analysis. These types of advancements will go a long way toward improving the reliability and accuracy of market predictions on very firm ground.

# 9. <u>REFERENCES</u>

**Weblinks:**

1. https://arxiv.org/pdf/2204.05783
2. https://www.iaeng.org/IJCS/issues_v47/issue_4/IJCS_47_4_17.pdf
3. https://www.sciencedirect.com/science/article/pii/S1877050920307237?ref=pdf_download&fr=RR-2&rr=8cdc94451f978a20
4. https://www.researchgate.net/profile/Murtaza-Roondiwala/publication/327967988_Predicting_Stock_Prices_Using_LSTM/links/5bafbe6692851ca9ed30ceb9/Predicting-Stock-Prices-Using-LSTM.pdf
5. https://www.sciencedirect.com/science/article/pii/S1877050918307828?ref=pdf_download&fr=RR-2&rr=8cdc94632a2e8a20
6. https://www.sciencedirect.com/science/article/pii/S2666827022000378?ref=pdf_download&fr=RR-2&rr=8cdc9474fc348a20
7. https://www.researchgate.net/profile/Giridhar-Maji/publication/346450872_Stock_Price_Prediction_Using_LSTM_on_Indian_Share_Market/links/5fc3142b458515b797843e38/Stock-Price-Prediction-Using-LSTM-on-Indian-Share-Market.pdf
8. https://www.sciencedirect.com/science/article/pii/S1877050920304865

# APPENDIX A –

## SAMPLE CODE-:

```python
# Imports

import yfinance as yf

import pandas as pd

import numpy as np

import matplotlib.pyplot as plt

from sklearn.preprocessing import MinMaxScaler

from keras.models import Sequential

from keras.layers import Dense, LSTM, Dropout

from sklearn.metrics import mean_squared_error, mean_absolute_percentage_error, r2_score

from datetime import datetime


# Data Collection

end = datetime.now()

start = datetime(end.year - 4, end.month, end.day)

companies = ['HDB', 'ITCI', 'ONGC.NS', 'TATASTEEL.NS', 'BRITANNIA.NS',
'RELIANCE.NS', 'DIVISLAB.NS', 'BAJFINANCE.NS', 'SBIN.NS', 'ASIANPAINT.NS']

data = {company: yf.download(company, start=start, end=end) for company in companies}


# Data Visualization

for company in companies:

    data[company][['Adj Close', 'Open', 'Close']].plot(figsize=(12, 5), title=f'{company} Price
Trends')

    plt.show()
```

```python
# Data Preprocessing

df = data['HDB'][['Close']].dropna()

scaler = MinMaxScaler((0, 1))

scaled_data = scaler.fit_transform(df)


# Train-Test Split and Dataset Creation

train_size = int(len(scaled_data) * 0.8)

train_data, test_data = scaled_data[:train_size], scaled_data[train_size:]


def create_dataset(data, time_step=60):

    X, Y = [], []

    for i in range(len(data) - time_step - 1):

        X.append(data[i:(i + time_step), 0])

        Y.append(data[i + time_step, 0])

    return np.array(X), np.array(Y)


X_train, y_train = create_dataset(train_data, time_step=60)

X_test, y_test = create_dataset(test_data, time_step=60)

X_train, X_test = X_train.reshape(X_train.shape[0], X_train.shape[1], 1),
X_test.reshape(X_test.shape[0], X_test.shape[1], 1)


# LSTM Model Design and Training

model = Sequential([

    LSTM(100, return_sequences=True, input_shape=(X_train.shape[1], 1)),

    Dropout(0.3),
```

```python
    LSTM(100, return_sequences=True),

    Dropout(0.3),

    LSTM(100),

    Dense(50),

    Dense(1)

])


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

model.fit(X_train, y_train, batch_size=32, epochs=50)


# Predictions and Evaluation

train_predict = scaler.inverse_transform(model.predict(X_train))

test_predict = scaler.inverse_transform(model.predict(X_test))

y_train_actual = scaler.inverse_transform([y_train])

y_test_actual = scaler.inverse_transform([y_test])


# Performance Metrics

train_mse = mean_squared_error(y_train_actual[0], train_predict[:, 0])

test_mse = mean_squared_error(y_test_actual[0], test_predict[:, 0])

train_mape = mean_absolute_percentage_error(y_train_actual[0], train_predict[:, 0])

test_mape = mean_absolute_percentage_error(y_test_actual[0], test_predict[:, 0])


print("Train MSE:", train_mse, "\nTest MSE:", test_mse)

print("Train MAPE:", train_mape, "\nTest MAPE:", test_mape)
```

```python
# Directional Accuracy

def calculate_directional_accuracy(actual, predicted):

    return np.mean((np.diff(actual) > 0) == (np.diff(predicted) > 0))


print("Train Directional Accuracy:", calculate_directional_accuracy(y_train_actual[0],
train_predict[:, 0]))

print("Test Directional Accuracy:", calculate_directional_accuracy(y_test_actual[0],
test_predict[:, 0]))


# Annual Performance Analysis

annual_returns = pd.DataFrame()

for company, df in data.items():

    df['Year'] = df.index.year

    annual_returns[company] = df.groupby('Year')['Close'].apply(lambda x: (x.iloc[-1] /
x.iloc[0] - 1) * 100)


print("Annual Best Performers:", annual_returns.idxmax(axis=1))
```

# CERTIFICATE

This is to certify that the project entitled <mark>**Stock Market Prediction Using LSTM Model**</mark> submitted by <mark>**Hussain Ahmed (21BCE2387)**</mark> ,**School of Computer Science and Engineering**, VIT, for the award of the degree of *Bachelor of Technology in Computer Science and Engineering*, is a record of bonafide work carried out by him / her under my supervision during Fall Semester 2024-2025, as per the VIT code of academic and research ethics.

The contents of this report have not been submitted and will not be submitted either in part or in full, for the award of any other degree or diploma in this institute or any other institute or university. The project fulfills the requirements and regulations of the University and in my opinion meets the necessary standards for submission.

Place : Vellore

Date :

**Signature of the Guide**

**Examiner(s)**

<mark>**Head of the Department**</mark>

<mark>**Programme**</mark>

<mark>**Dr. UMADEVI K S**</mark>
<mark>**SCOPE**</mark>

# CERTIFICATE

This is to certify that the project entitled  Stock Market Prediction Using LSTM Model  submitted by  Bhavana Singh (21BCE2431) , **School of Computer Science and Engineering**, VIT, for the award of the degree of *Bachelor of Technology in Computer Science and Engineering*, is a record of bonafide work carried out by him / her under my supervision during Fall Semester 2024-2025, as per the VIT code of academic and research ethics.

The contents of this report have not been submitted and will not be submitted either in part or in full, for the award of any other degree or diploma in this institute or any other institute or university. The project fulfills the requirements and regulations of the University and in my opinion meets the necessary standards for submission.

Place : Vellore

Date :

**Signature of the Guide**

**Examiner(s)**

**Head of the Department**

**Programme**

**Dr. UMADEVI K S**

**SCOPE**