# Cracking the Market Code with AI-Driven Stock Price Prediction using Time Series Analysis

## PHASE -3

*STUDENT NAME: S AKASH*

*REGISTER  NUMBER:620523243006*

*INSTITUTION: CMS COLLEGE OF ENGINEERING*

*DEPARTMENT: ARTIFICIAL INTELLIGENCE AND DATA SCIENCE*

*DATE OF SUBMISSION*

*GITHUB REPOSITORY LINK:*

## Problem Statement:

Predicting stock prices accurately has been a long-standing challenge in the financial domain. The stock market is a complex and dynamic system influenced by a multitude of interconnected factors, including economic indicators, company performance, investor sentiment, global events, and even random noise. Traditional statistical and econometric models often struggle to capture these intricate relationships and non-linear patterns, leading to limited predictive power. The inherent volatility and unpredictability of the market result in significant financial risks for investors and hinder effective investment strategies. Therefore, there is a need for more sophisticated and adaptive approaches capable of learning complex patterns from historical data to provide more reliable stock price predictions.

## Abstract:

This project aims to develop an AI-driven system for predicting stock prices using time series analysis techniques. By leveraging historical stock market data and

potentially incorporating other relevant features like technical indicators, fundamental data, and sentiment analysis, we will explore and implement various machine learning and deep learning models. The project will involve data collection, preprocessing, exploratory data analysis, feature engineering, model building (including models like ARIMA, LSTMs, and Transformers), rigorous model evaluation, and a discussion on potential deployment strategies. The goal is to create a robust and informative system that can provide insights into future stock price movements and contribute to more informed investment decisions.

## System Requirements:

- ➔ **Software Requirements:**
  - ◆ Operating System: Windows, macOS, or Linux
  - ◆ Programming Language: Python (version 3.x)
  - ◆ Libraries:
    - ● Pandas
    - ● NumPy
    - ● Matplotlib
    - ● Seaborn
    - ● Statsmodels
    - ● Scikit-learn
    - ● TensorFlow (with Keras) or PyTorch
    - ● Prophet (optional)
    - ● Requests (for API interaction)
  - ◆ Development Environment: Jupyter Notebooks, VS Code, PyCharm
  - ◆ Version Control: Git
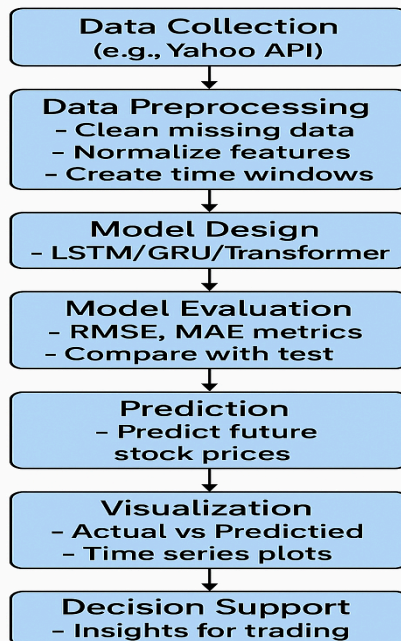- ➔ **Hardware Requirements:**
  - ◆ Processor: Intel Core i5 or equivalent (or higher recommended for complex models)
  - ◆ RAM: 8 GB or higher (16 GB or more recommended for deep learning)
  - ◆ Storage: Sufficient disk space to store datasets and model files
  - ◆ Internet Connectivity: For accessing data sources and online resources
  - ◆ GPU (Recommended for faster training of deep learning models): NVIDIA GPU with CUDA support
  - ◆

## Objectives:

- ➔ Develop an AI-driven model for stock price prediction using time series analysis.
- ➔ Evaluate the effectiveness of various AI algorithms for stock price forecasting.
- ➔ Identify and analyze relevant features beyond historical price data to improve prediction accuracy.
- ➔ Establish a robust methodology for data handling, model training, validation, and evaluation.
- ➔ Provide insights into factors influencing stock price movements based on model analysis.
- ➔ Potentially develop a prototype for real-time stock price prediction.

## Flow Chart of the Project Workflow:

Cracking the Market Code with
AI-Driven Stock Price Prediction
Using Time Series Analysis

**Data Collection**
(e.g., Yahoo API)

↓

**Data Preprocessing**
– Clean missing data
– Normalize features
– Create time windows

↓

**Model Design**
– LSTM/GRU/Transformer

↓

**Model Evaluation**
– RMSE, MAE metrics
– Compare with test

↓

**Prediction**
– Predict future
stock prices

↓

**Visualization**
– Actual vs Predictied
– Time series plots

↓

**Decision Support**
– Insights for trading

## Dataset Description:

The primary dataset will consist of historical time series data for the selected stocks. This will include:

➜ **Ticker Symbol:** Unique identifier for each stock.
➜ **Date:** Timestamp of the data point.
➜ **Open Price:** The price at which the stock started trading on a particular day.
➜ **High Price:** The highest price reached by the stock during the trading day.
➜ **Low Price:** The lowest price reached by the stock during the trading day.
➜ **Close Price:** The price at which the stock ended trading on a particular day (our primary target variable).
➜ **Adjusted Close Price:** The closing price adjusted for dividends and stock splits, providing a more accurate reflection of long-term returns.
➜ **Volume:** The number of shares traded during the day.

Depending on the project scope, we might incorporate additional datasets containing:

➜ **Fundamental Data:** Company financial statements (e.g., income statements, balance sheets, cash flow statements) reported quarterly or annually.
➜ **Economic Indicators:** Macroeconomic data released periodically (e.g., inflation rates, interest rates, GDP growth).
➜ **Sentiment Data:** Scores derived from news articles, social media, or financial reports, reflecting market sentiment towards specific stocks or the overall market.
➜ The data will be collected from reliable financial data providers, ensuring a

sufficient historical depth for effective model training. The frequency of the data will likely be daily, but we might explore other granularities (e.g., hourly) depending on the chosen models and computational resources.

## Data Preprocessing:

1. This stage will involve the following steps:
2. **Data Collection:** Fetching the datasets from the identified sources.
3. **Data Cleaning:**
   a. Handling missing values (imputation or removal).
   b. Identifying and treating outliers using statistical methods or domain knowledge.
   c. Ensuring data consistency and format uniformity.
4. **Data Transformation:**
   a. Scaling or normalization of numerical features (e.g., Min-Max scaling, Standardization) to improve model convergence.
   b. Converting categorical features (if any, related to events or news) into numerical representations (e.g., one-hot encoding).
5. **Time Series Specific Preprocessing:**
   a. Ensuring the data is ordered chronologically.
   b. Checking for stationarity and applying transformations (e.g., differencing) if necessary for certain models like ARIMA.
   c. Creating lagged features of the target variable and other relevant features.
6. **Data Splitting:** Dividing the data into training, validation, and testing sets to train the models, tune hyperparameters, and evaluate final performance on unseen data.

**Exploratory Data Analysis (EDA):**

EDA will be performed to gain insights into the characteristics of the data and identify potential patterns. This will involve:
1. **Visualizations:**
   ○ Time series plots of stock prices and volume.
   ○ Distribution plots (histograms, box plots) for individual features.
   ○ Correlation matrices to understand the relationships between different variables.
   ○ Decomposition of time series data into trend, seasonality, and residual components.
   ○ Scatter plots to visualize the relationship between features and the target variable.
2. **Statistical Analysis:**
   ○ Calculating descriptive statistics (mean, median, standard deviation, etc.).
   ○ Performing stationarity tests (e.g., Augmented Dickey-Fuller test).
   ○ Analyzing autocorrelation and partial autocorrelation functions (ACF and PACF) for time series data.
   ○ Identifying potential seasonality patterns.

## Feature Engineering:

This crucial step involves creating new features from the existing data that might improve the predictive power of the models. Examples include:
1. **Technical Indicators:** Calculating various technical indicators (as listed

previously: SMA, EMA, RSI, MACD, Bollinger Bands, etc.) based on historical price and volume data.

2. **Lagged Variables:** Creating past values of the target variable and other relevant features as input to the models. The number of lag periods will be determined through analysis and experimentation.
3. **Rolling Statistics:** Calculating rolling mean, standard deviation, and other statistical measures over different time windows.
4. **Volatility Measures:** Calculating historical volatility based on price fluctuations.
5. **Derived Features:** Creating ratios or combinations of existing features that might capture relevant relationships.
6. **External Feature Integration (if scope allows):** Incorporating sentiment scores, economic indicators, or global market indices as additional features.

# Model Building:

We will explore and implement several time series forecasting models, including:
- ➔ **Statistical Models:**
    - ◆ **ARIMA (Autoregressive Integrated Moving Average):** A traditional linear model that captures autocorrelations in the data.
    - ◆ **SARIMA (Seasonal ARIMA):** An extension of ARIMA to handle seasonality.
- ➔ **Machine Learning Models:**
    - ◆ **Regression Models (e.g., Linear Regression, Support Vector Regression, Random Forest, Gradient Boosting):** These models can be adapted for time series forecasting by using lagged features.
- ➔ **Deep Learning Models:**
    - ◆ **Recurrent Neural Networks (RNNs):** Specifically, Long Short-Term Memory (LSTM) and Gated Recurrent Units (GRUs) which are well-suited for capturing sequential dependencies in time series data.
    - ◆ **Transformer Networks:** Models that utilize attention mechanisms to capture long-range dependencies and have shown promising results in sequence modeling.
    - ◆ **Hybrid Models:** Combining different model architectures or statistical and machine learning approaches to leverage their respective strengths.

For each model, we will:
- ➔ Define the model architecture and hyperparameters.
- ➔ Train the model using the training data.
- ➔ Tune the hyperparameters using the validation set to optimize performance.
- ➔

# Model Evaluation:

The performance of the trained models will be evaluated using the following metrics on the unseen test data:
- ➔ **Mean Squared Error (MSE):** Average of the squared differences between predicted and actual values. $\qquad MSE = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$
- ➔ **Root Mean Squared Error (RMSE):** Square root of MSE, providing an error in the

same units as the target variable. $RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2}$

- ➔ **Mean Absolute Error (MAE):** Average of the absolute differences between predicted and actual values. $MAE = \frac{1}{n} \sum_{i=1}^{n} |y_i - \hat{y}_i|$
- ➔ **Directional Accuracy:** Percentage of times the model correctly predicts the direction of price movement (up or down). $\text{Directional Accuracy} = \frac{\text{Number of correct direction predictions}}{\text{Total number of predictions}} \times 100\%$
- ➔ **R-squared (Coefficient of Determination) (for some models):** Measures the proportion of the variance in the dependent variable that is predictable from the independent variables.

# Deployment:

- ➔ **Developing a Web Application:** Creating a user-friendly interface where users can input stock tickers and get price predictions.
- ➔ **Building an API:** Exposing the model's prediction capabilities through an API that can be integrated into other applications or trading platforms.
- ➔ **Creating a Desktop Application:** A standalone application that provides stock price forecasts.
- ➔ **Integration with Trading Bots (Future advanced scope):** Incorporating the model's predictions into automated trading strategies (with significant risk considerations and regulatory compliance).

# Source Code:

```
import yfinance as yf
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.preprocessing import MinMaxScaler
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import LSTM, Dense, Dropout
from sklearn.metrics import mean_squared_error

# Step 1: Data Collection
ticker = 'AAPL'
data = yf.download(ticker, start='2015-01-01', end='2024-12-31')
df = data[['Close']]

# Step 2: Preprocessing
scaler = MinMaxScaler(feature_range=(0, 1))
scaled_data = scaler.fit_transform(df)

# Step 3: Creating time-series sequences
def create_dataset(data, time_step=60):
    X, y = [], []
```

```python
    for i in range(time_step, len(data)):
        X.append(data[i-time_step:i, 0])
        y.append(data[i, 0])
    return np.array(X), np.array(y)

time_step = 60
X, y = create_dataset(scaled_data, time_step)
X = X.reshape(X.shape[0], X.shape[1], 1)

# Step 4: Model Building
model = Sequential()
model.add(LSTM(units=50, return_sequences=True, input_shape=(X.shape[1], 1)))
model.add(Dropout(0.2))
model.add(LSTM(units=50))
model.add(Dropout(0.2))
model.add(Dense(1))

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

# Step 5: Training
model.fit(X, y, epochs=20, batch_size=32)

# Step 6: Prediction
predicted = model.predict(X)
predicted_prices = scaler.inverse_transform(predicted)
real_prices = scaler.inverse_transform(y.reshape(-1, 1))

# Step 7: Evaluation
rmse = np.sqrt(mean_squared_error(real_prices, predicted_prices))
print(f'RMSE: {rmse:.2f}')

# Step 8: Visualization
plt.figure(figsize=(12,6))
plt.plot(real_prices, label='Actual Price')
plt.plot(predicted_prices, label='Predicted Price')
plt.title(f'{ticker} Stock Price Prediction')
plt.xlabel('Time')
plt.ylabel('Price')
plt.legend()
plt.show()
```

## Future Scope:

1. **Incorporating Real-time Data:** Integrating the model with real-time stock market data feeds for more up-to-date predictions.
2. **Sentiment Analysis Integration:** Developing and incorporating sophisticated sentiment analysis techniques on news and social media data.
3. **Fundamental Analysis Integration:** Incorporating and analyzing fundamental company data to improve long-term prediction accuracy.
4. **Risk Management Integration:** Developing methods to quantify the uncertainty

associated with the predictions and provide risk assessments.

5. **Portfolio Optimization:** Extending the model to predict prices for multiple stocks and develop portfolio optimization strategies.
6. **Explainable AI (XAI):** Implementing techniques to understand and interpret the model's predictions, making them more transparent and trustworthy.
7. **High-Frequency Trading Analysis (Advanced):** Exploring the applicability of AI techniques to high-frequency trading data (requires significant computational resources and expertise).
8. **Cross-Market Analysis:** Extending the model to analyze and predict stock prices across different geographical markets.

## Team Members and Roles:

- **[Team Member 1 Name]:** Project Lead, responsible for overall project management, planning, and coordination.
- **[Team Member 2 Name]:** Data Scientist, responsible for data collection, preprocessing, EDA, and feature engineering.
- **[Team Member 3 Name]:** Machine Learning Engineer, responsible for model building, training, evaluation, and deployment.
- **[Team Member 4 Name]:** Research Analyst, responsible for literature review, identifying relevant features, and analyzing results.