



Al Akhawayn University in Ifrane

# Transformer-Based Time Series Forecasting for Stock Market Volatility

CSC 3347–01: Machine Learning

Fall 2025

## Project Team

Haytam Raiss

Hiba Tallah Erraji

Mouad Wadi

# Contents

<b>1</b>	<b>Project Overview</b>	<b>3</b>
<b>2</b>	<b>Project Requirements</b>	<b>3</b>
2.1	Survey of Existing Work . . . . .	3
2.2	Project Objectives and Hypotheses . . . . .	5
2.3	Datasets . . . . .	6
2.4	System Architecture . . . . .	7
<b>3</b>	<b>Model Descriptions</b>	<b>8</b>
3.1	Long Short-Term Memory (LSTM) . . . . .	8
3.2	Transformer Encoder . . . . .	8
3.3	ARIMA Time-Series Model . . . . .	9
3.4	GARCH Volatility Model . . . . .	9
3.5	Model Selection . . . . .	10
3.6	Evaluation Metrics . . . . .	10
<b>4</b>	<b>Results and Analysis</b>	<b>12</b>
4.1	Data Preparation . . . . .	12
4.2	Model Training Setup . . . . .	12
4.3	Baseline Results (5-Day Horizon, Level Target) . . . . .	12
4.4	Hyperparameter Tuning (Transformer) . . . . .	13
4.5	Longer Horizon and Alternative Target . . . . .	15
4.6	Rolling Window Backtest and Regime Behaviour . . . . .	15
4.7	Prediction Interval and Calibration . . . . .	17
4.8	Comprehensive Experiment Summary . . . . .	17
4.9	Regime-Specific Error Analysis . . . . .	18
4.10	Prediction Interval Coverage . . . . .	19
4.11	Overall Assessment . . . . .	20
<b>5</b>	<b>Conclusion</b>	<b>20</b>
<b>6</b>	<b>Deployment</b>	<b>23</b>

<b>A</b>	<b>Appendix: Key References and Resources</b>	<b>25</b>
A.1	Key Transformer-Based Volatility References . . . . .	25
A.2	Thirty Related Models, Datasets, and Systems . . . . .	25
A.3	Reference Links . . . . .	30

# 1 Project Overview

This project investigates whether modern deep learning models, in particular Transformer architectures, can improve short-horizon volatility forecasting for the S&P 500 index. The work begins from raw daily OHLCV data for a large universe of U.S. stocks and ETFs, then focuses on a curated S&P-500-centric dataset, `features_target`, which aggregates macroeconomic, options, and index variables engineered specifically for volatility prediction.

A complete time-series pipeline is implemented: data cleaning and feature engineering, sequence construction, model training (Transformer, LSTM, ARIMA, GARCH), systematic evaluation, and an interactive deployment exposing the best model for inference.

## 2 Project Requirements

### 2.1 Survey of Existing Work

Volatility forecasting has a long history in financial econometrics and, more recently, in machine learning. Early work focused on modelling conditional variance using parametric time-series models, while more recent research explores non-linear and data-driven approaches, including deep neural networks and attention-based architectures. This subsection summarizes key strands of the literature that motivate the design choices in this project: GARCH-type models, realized volatility and HAR models, traditional machine learning approaches, recurrent neural networks, and Transformers for financial time series.

#### **GARCH and Classical Volatility Models**

The starting point for most volatility work is the family of ARCH and GARCH models, which model conditional variance as a function of past squared returns and past variances. These models capture volatility clustering and mean reversion and remain widely used in practice for risk management and derivatives pricing. Numerous extensions have been proposed, including EGARCH, GJR-GARCH, and multivariate GARCH, to allow for leverage effects, asymmetries, and richer dynamics.

ARIMA and related linear time-series models have also been applied directly to volatility or squared returns. These models can capture autocorrelation but do not model conditional heteroskedasticity explicitly, so they are usually considered weaker baselines than GARCH

for volatility itself. Nevertheless, ARIMA remains a useful benchmark because it provides a transparent linear representation and allows for a clean comparison against more complex non-linear approaches.

### **Realized Volatility and HAR Models**

The introduction of high-frequency data enabled realized volatility measures computed from intraday returns. This led to a rich literature on realized volatility forecasting, where the target is the sum of high-frequency squared returns over a day or longer horizon. A key result from this line of work is that realized volatility is highly persistent and exhibits long memory. To capture this behaviour, the Heterogeneous Autoregressive (HAR) model was proposed, decomposing volatility into daily, weekly, and monthly components. HAR-type models often provide strong benchmarks that are simple to estimate yet capture important time-scale effects.

### **Traditional Machine Learning for Volatility**

With the growth of machine learning in finance, many authors have experimented with non-linear models such as decision trees, random forests, gradient boosting, and support vector regression for volatility prediction. These models flexibly capture interactions and non-linearities among explanatory variables. Empirical results typically show that tree-based ensembles can outperform simple linear regressions and sometimes rival or beat GARCH/HAR models when feature engineering is done carefully, although gains are usually modest and sensitive to sample period.

### **Recurrent Neural Networks and LSTMs**

Recurrent neural networks, particularly Long Short-Term Memory (LSTM) and Gated Recurrent Unit (GRU) architectures, have been widely applied to financial time series. LSTMs can, in principle, capture long-term dependencies and non-linear interactions without manual lag engineering. Several studies report that LSTMs outperform classical time-series models on some datasets, especially when multiple correlated inputs (e.g., macro variables, technical indicators) are available. Other work highlights limitations: training instability, overfitting on small or noisy datasets, and improvements that are sometimes small relative to strong econometric baselines.

## Transformers and Attention in Finance

More recently, attention mechanisms and Transformer architectures, originally developed for natural language processing, have been adapted to financial time series. In a Transformer, self-attention layers replace recurrence, allowing the model to weigh different time steps and features according to their relevance for the prediction task. Several studies show that Transformers can outperform LSTMs on price prediction, order book forecasting, and realized volatility modelling. Attention is particularly useful when the input includes heterogeneous signals, such as prices, volumes, options data, and macro indicators, with time-varying importance.

## Regime Dependence and Limits of Predictability

A recurring theme across both econometric and machine learning literature is that volatility predictability is strongly regime-dependent. Many empirical studies find that models perform better in moderate-volatility periods and degrade during crises or extremely calm markets. Rolling-window and sub-sample analyses often reveal that performance metrics fluctuate significantly over time, sometimes turning negative when structural breaks occur. Overall, short-horizon stock-index volatility appears only partially predictable, even with sophisticated models and rich input features.

## 2.2 Project Objectives and Hypotheses

The task is a supervised **regression** problem. The main objectives are:

- Predict short-horizon S&P 500 volatility using daily macro, options, and market features.
- Compare deep learning models (LSTM, Transformer) against ARIMA and GARCH baselines.
- Analyze performance across volatility regimes and over time.

Targets include the S&P 500 30-day volatility level shifted 5 days ahead and the daily change in 30-day volatility shifted 5 and 20 days ahead.

Hypotheses:

- The Transformer will outperform LSTM, ARIMA, and GARCH in terms of RMSE, MAE, and  $R^2$ .

- Using the engineered `features_target` dataset instead of raw OHLCV will reduce prediction error.
- Volatility predictability will be regime-dependent, with better accuracy in medium-volatility regimes.

## 2.3 Datasets

The final dataset, `features_target`, is constructed from:

- Market indices: S&P 500 level and log returns, 30-day volatility, DJIA, NASDAQ, Russell 2000, MSCI World and their log returns.
- Options data: SPX put–call ratio and log changes of put, call, and total options volume.
- Macro and credit variables: 10-year Treasury yield, high-yield bond index, EMB yield.
- Commodities and FX: gold, oil, USD index and their log returns.
- Sentiment: consumer sentiment level and log change.

Labels are derived as:

- 5-day-ahead S&P 500 30-day volatility level.
- 5-day-ahead and 20-day-ahead changes in 30-day volatility.

After cleaning and dropping missing values, the dataset contains roughly 2,200 daily observations. Features are standardized using training-set mean and variance. Sequences of length 60 days are built to produce samples of shape  $60 \times F$  with scalar volatility targets.

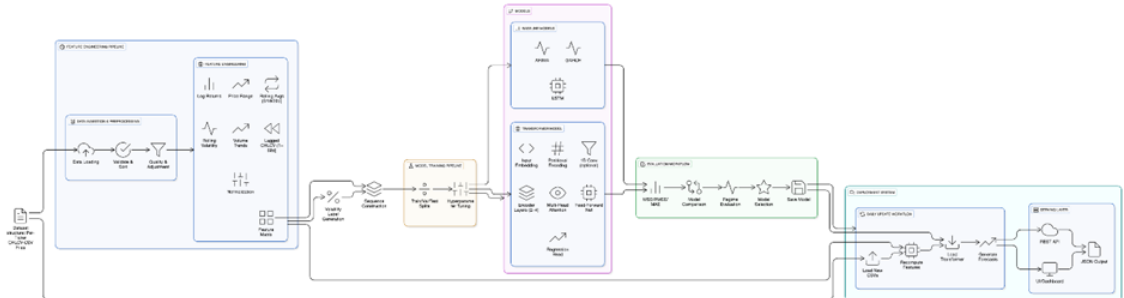
Time-based splits are used: 70% train, 15% validation, 15% test, plus additional rolling windows for robustness. The data involve no personal information, so privacy and fairness risks are minimal.

SP500, SP500 Log Returns, SP500 30 Day Volatility, SPX Put Call Ratio, SPX Put Volume, SPX Call Volume, Total SPX Options Volume, VIX, OJIA, NASDAQ, 10Y Treasury, High Yield Bonds, RUSSELL, EMB, Y1  
2010-07-08, 1070.25, 0.009368673372673442, 28.51097636786388, 1.12, 280436.0, 251021.0, 531457.0, 25.71, 10138.99, 2175.4, 3.02, 8.89, 620.27, 104.42, 1084.900024, 67.8, 83.82, 1196.1, 75.46, 0.0119769634  
2010-07-09, 1077.96, 0.00718100004059829, 28.46281377195, 1.2, 376805.0, 315132.0, 691937.0, 24.98, 10198.03, 2196.45, 3.056, 8.83, 629.43, 104.65, 1091.5, 67.8, 83.95, 1209.8, 76.08, 0.00580617670252614  
2010-07-12, 1078.75, 0.0007325973685219367, 28.04107879318622, 1.49, 421390.0, 283127.0, 704517.0, 24.43, 10216.27, 2198.36, 3.046, 8.79, 621.61, 104.34, 1091.5, 67.8, 84.2, 1198.7, 74.93, 0.0017869831298  
2010-07-13, 1095.34, 0.015261853940210734, 28.157439342815254, 1.25, 652349.0, 519892.0, 1172241.0, 24.56, 10363.02, 2242.03, 3.114, 8.69, 642.82, 105.21, 1109.0, 67.8, 83.64, 1213.5, 77.16, 0.01426215263  
2010-07-14, 1095.17, 0.0001552149959058724, 28.338036329095196, 1.01, 475499.0, 469171.0, 944670.0, 24.89, 10366.72, 2249.84, 3.05, 8.6, 640.16, 105.36, 1113.099976, 67.8, 83.43, 1207.0, 77.02, 0.0003569  
2010-07-15, 1096.48, 0.0011954464955712396, 28.291382823566245, 0.87, 365810.0, 421520.0, 787330.0, 25.14, 10359.31, 2249.08, 2.978, 8.58, 634.62, 105.22, 1113.099976, 67.8, 82.56, 1208.3, 76.67, 0.00071  
2010-07-16, 1064.88, 0.029242932098749286, 27.95395171004485, 2.22, 416822.0, 187961.0, 604783.0, 26.25, 10097.9, 2179.05, 2.939, 8.58, 610.39, 105.59, 1089.900024, 67.8, 82.49, 1188.2, 75.96, 0.025581  
2010-07-19, 1071.25, 0.00596407417322542, 27.8909230486617, 1.76, 277724.0, 158190.0, 435914.0, 25.97, 10154.43, 2198.23, 2.964, 8.6, 613.08, 105.76, 1089.300049, 67.8, 82.51, 1181.9, 76.53, 0.00558258203  
2010-07-20, 1083.48, 0.011351892195398783, 27.48977345521745, 1.49, 258007.0, 173110.0, 431117.0, 23.93, 10229.96, 2222.49, 2.932, 8.59, 624.24, 105.95, 1093.5, 67.8, 82.75, 1191.7, 77.32, 0.007410606415  
2010-07-21, 1069.59, 0.01290268565281405, 27.378968562399162, 2.48, 400764.0, 405960.0, 25.64, 10120.53, 2187.33, 2.892, 8.52, 612.64, 106.35, 1088.400024, 67.8, 83.39, 1191.8, 76.27, 0.010754  
2010-07-22, 1093.67, 0.022263615685037763, 27.21153063575867, 1.27, 2 Col 6: SPX Put Volume 329.0, 24.63, 10322.3, 2245.89, 2.932, 8.46, 635.48, 106.9, 1110.099976, 67.8, 82.59, 1195.6, 79.01, 0.0197405694  
2010-07-23, 1102.66, 0.008186429368369552, 27.524149772340913, 1.59, 296888.0, 186205.0, 483093.0, 23.47, 10424.62, 2269.47, 2.994, 8.44, 650.65, 107.11, 1117.400024, 67.8, 82.46, 1187.8, 78.68, 0.0098637  
2010-07-26, 1115.01, 0.011137930955885444, 28.2073771817207, 1.7, 256861.0, 151049.0, 407910.0, 22.73, 10525.43, 2296.43, 2.994, 8.39, 665.22, 107.15, 1129.800049, 67.8, 82.09, 1183.1, 78.93, 0.0096239176  
2010-07-27, 1113.84, 0.0010498688628368669, 28.82859965661715, 1.4, 279835.0, 200452.0, 480287.0, 23.19, 10537.69, 2288.25, 3.047, 8.33, 662.17, 107.73, 1130.900024, 67.8, 82.19, 1158.0, 77.46, 0.0011641  
2010-07-28, 1106.13, 0.00694606773895412, 28.501128664896818, 1.99, 211544.0, 106265.0, 317809.0, 24.25, 10497.88, 2264.56, 3.001, 8.31, 650.76, 107.57, 1129.300049, 67.8, 82.18, 1160.4, 77.06, 0.003785  
2010-07-29, 1101.53, 0.004167314402725886, 28.059766189711601, 1.74, 207817.0, 119150.0, 326967.0, 24.13, 10467.16, 2251.69, 2.999, 8.3, 650.43, 107.7, 1128.199951, 67.8, 81.64, 1171.2, 78.3, 0.00293059  
2010-07-30, 1101.6, 0.354595509439065e-05, 27.54246405897868, 1.18, 228363.0, 193948.0, 422311.0, 23.5, 10465.94, 2254.7, 2.907, 8.29, 650.89, 107.65, 1124.800049, 67.8, 81.54, 1183.9, 78.85, 0.000116561  
2010-08-02, 1125.86, 0.02178351962968783, 27.97653792266879, 2.1, 413466.0, 196904.0, 610370.0, 22.01, 10674.38, 2295.36, 2.963, 8.31, 661.86, 107.69, 1152.199951, 68.9, 80.94, 1185.4, 81.25, 0.0197203027  
2010-08-03, 1120.46, 0.00480787728930502, 28.30339165752615, 1.27, 246879.0, 194021.0, 440900.0, 22.63, 10636.38, 2283.52, 2.914, 8.26, 655.66, 108.38, 1153.0, 68.9, 80.6, 1187.5, 82.52, 0.00356627734  
2010-08-04, 1127.24, 0.00603285186041797, 29.46691361784437, 3.29, 333891.0, 101417.0, 435308.0, 22.21, 10680.43, 2303.57, 2.952, 8.24, 662.96, 108.64, 1152.5, 68.9, 80.89, 1195.9, 82.49, 0.0041328944340  
2010-08-05, 1125.81, 0.0012693905523795124, 30.518344076522016, 2.67, 287807.0, 107609.0, 395416.0, 22.1, 10674.98, 2293.06, 2.915, 8.23, 655.07, 108.8, 1153.0, 68.9, 80.83, 1199.3, 82.0, 0.000510409317  
2010-08-06, 1121.64, 0.003710876580377054, 31.2919029665012, 1.74, 279442.0, 160365.0, 439807.0, 21.74, 10653.56, 2288.47, 2.824, 8.21, 650.68, 109.36, 1152.300049, 68.9, 80.41, 1205.3, 80.67, 0.0020085  
2010-08-09, 1127.79, 0.005468065530920718, 32.256575216739726, 2.01, 234983.0, 117091.0, 352074.0, 22.14, 10698.75, 2305.69, 2.822, 8.2, 659.52, 109.79, 1158.199951, 68.9, 80.71, 1202.6, 81.46, 0.00423280  
2010-08-10, 1121.06, 0.005985299233984165, 32.81501575972309, 2.49, 311464.0, 125181.0, 436645.0, 22.37, 10644.25, 2277.17, 2.781, 8.22, 646.36, 109.88, 1145.400024, 68.9, 80.8, 1198.0, 80.24, 0.0051070  
2010-08-11, 1089.47, 0.028583326840638357, 31.643392455494876, 1.64, 476682.0, 291174.0, 767856.0, 25.39, 10378.83, 2208.63, 2.685, 8.33, 620.39, 109.75, 1112.800049, 68.9, 82.29, 1199.2, 78.09, 0.02525  
2010-08-12, 1083.61, 0.00539327972832216, 29.70562238087233, 2.98, 444674.0, 149367.0, 594041.0, 25.73, 10319.95, 2190.27, 2.735, 8.42, 616.98, 109.65, 1106.099976, 68.9, 82.64, 1216.7, 75.68, 0.005689  
2010-08-13, 1079.25, 0.004031704230371247, 27.309570317381873, 1.83, 274615.0, 149993.0, 424608.0, 26.24, 10303.15, 2173.48, 2.688, 8.43, 609.49, 109.75, 1103.599976, 68.9, 82.95, 1216.6, 75.39, 0.00162  
2010-08-16, 1079.38, 0.00012044676499201046, 24.0836713153011, 1.55, 381664.0, 246236.0, 627900.0, 26.1, 10302.01, 2181.87, 2.575, 8.38, 615.1, 109.928, 1106.199951, 68.9, 82.53, 1226.2, 75.17, 0.0001106

## 2.4 System Architecture

The system has four stages:

- 1. Ingestion and preprocessing:** load `features_target`, clean columns, remove redundant raw volumes, handle missing values, and standardize numerical features.
- 2. Feature engineering and sequence construction:** construct volatility level or change targets by shifting the series forward in time, then form 60-day rolling windows for model input; compute volatility-regime labels (low/medium/high) via quantiles for analysis.
- 3. Model training pipeline:** implement LSTM and Transformer models in PyTorch with configurable hyperparameters; fit ARIMA and GARCH models using standard libraries; use mini-batch training, gradient clipping, and early stopping on validation loss with checkpointing of the best model.
- 4. Evaluation and deployment:** compute metrics (RMSE, MAE,  $R^2$ ), regime-specific scores, rolling-window performance, residual plots, Q-Q plots, and prediction-interval coverage; expose the best Transformer model through an interactive dashboard for inference.





### 3 Model Descriptions

This section briefly explains the main models used in the project: LSTM and Transformer as deep learning sequence models, and ARIMA and GARCH as classical time-series baselines.

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

Long Short-Term Memory (LSTM) networks are a type of recurrent neural network designed to model sequential data. An LSTM processes an input sequence step by step and maintains a hidden state that is updated over time. Each LSTM cell contains three gates: an input gate, a forget gate, and an output gate. These gates control how much past information is kept, how much new information is added, and what part of the internal state is exposed at each time step.

In this project, the LSTM takes as input a 60-day window of features and produces a single forecast of future volatility (or volatility change). The architecture uses stacked LSTM layers followed by a fully connected output layer. LSTMs are well-suited to capturing non-linear temporal patterns and have been widely used in financial forecasting, but they can be difficult to train and may overfit when data are noisy or limited.

#### 3.2 Transformer Encoder

The Transformer is a sequence model based entirely on attention mechanisms, without recurrence or convolution. Instead of processing observations one by one, a Transformer encoder looks at the entire 60-day window simultaneously. Each layer uses multi-head self-attention to compute weighted combinations of all time steps, allowing the model to focus on the most relevant days and features for the prediction. A position encoding or learned projection is used so that the model is aware of ordering in the sequence.

In this project, the Transformer encoder projects the input features into a latent space of dimension  $d_{\text{model}}$ , passes them through several self-attention and feed-forward layers, and then aggregates the final hidden state (typically the last time step) through a linear output layer to obtain a volatility forecast. Transformers handle long-range dependencies and heterogeneous features more flexibly than LSTMs and can parallelize computation across time, which makes them attractive for financial time series with many correlated signals.

### 3.3 ARIMA Time-Series Model

The Autoregressive Integrated Moving Average (ARIMA) model is a classical linear model for univariate time series. An  $\text{ARIMA}(p, d, q)$  model combines:

- an autoregressive (AR) component of order  $p$ , which regresses the series on its own lags,
- an integration (I) component of order  $d$ , which applies differencing to remove trends or non-stationarity,
- a moving-average (MA) component of order  $q$ , which models the error term as a linear combination of past shocks.

In this project, ARIMA is fitted to the volatility (or volatility-change) series alone, without additional exogenous variables. It provides a transparent benchmark that captures linear persistence but does not model time-varying variance explicitly. ARIMA is useful to show how much improvement is obtained by more flexible non-linear and attention-based models.

### 3.4 GARCH Volatility Model

The Generalized Autoregressive Conditional Heteroskedasticity (GARCH) model explicitly models the conditional variance of a time series. A  $\text{GARCH}(1,1)$  specification assumes that the conditional variance at time  $t$  depends on a constant term, the squared return (or volatility innovation) at time  $t - 1$ , and the previous conditional variance at time  $t - 1$ . This structure captures volatility clustering: periods of high volatility tend to be followed by high volatility, and similarly for low volatility.

In the experiments, a  $\text{GARCH}(1,1)$  model is estimated on the volatility (or volatility-change) series. The one-step-ahead conditional variance output is interpreted as the volatility forecast. GARCH serves as a strong econometric baseline because it is specifically designed for volatility dynamics and is widely used in risk management and option pricing. Comparing Transformer and LSTM forecasts against GARCH makes it possible to assess whether complex deep learning architectures provide added value beyond established statistical models.

### 3.5 Model Selection

Baselines:

- **ARIMA**: captures linear autocorrelation in the volatility or volatility-change series.
- **GARCH(1,1)**: models conditional heteroskedasticity and serves as a standard volatility benchmark.

Deep learning models:

- **LSTM**: two layers with hidden size 128, dropout 0.2 and a linear output head. The model is trained with Adam (learning rate  $10^{-4}$ ), batch size 64, and early stopping.
- **Transformer**: an encoder-only architecture with input projection to  $d_{\text{model}} \in \{64, 128\}$ , 2–3 layers, 4 attention heads, feed-forward dimension 128, and dropout 0.2–0.3. Grid search over  $d_{\text{model}}$ , depth, dropout, and learning rate identifies the best configuration  $d_{\text{model}} = 128$ , 2 layers, dropout 0.2, learning rate  $10^{-4}$ .

### 3.6 Evaluation Metrics

Because the project is a regression task (forecasting continuous volatility), evaluation focuses on standard forecast-error metrics.

#### Mean Squared Error (MSE)

**Definition:** the average of squared differences between predicted and true volatility values:

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i)^2.$$

**Justification:** squaring penalizes large errors more than small ones, which is important in risk management because severely underestimating volatility is especially harmful. **Usage:** many realized-volatility and risk-forecasting studies report MSE or related measures when comparing econometric and machine learning models.

#### Root Mean Squared Error (RMSE)

**Definition:** the square root of MSE:

$$\text{RMSE} = \sqrt{\text{MSE}}.$$

**Justification:** expressed in the same units as the target volatility, making it easier to interpret, while still weighting large errors strongly. **Usage:** RMSE is one of the most common metrics in volatility-forecasting work and is used in many comparisons between GARCH, HAR, LSTM, and Transformer-based models.

### Mean Absolute Error (MAE)

**Definition:** the average absolute difference between predicted and true values:

$$\text{MAE} = \frac{1}{N} \sum_{i=1}^N |\hat{y}_i - y_i|.$$

**Justification:** less sensitive to outliers than MSE/RMSE and directly interpretable as the typical forecast error in volatility units. **Usage:** MAE is frequently reported alongside RMSE to give a more robust view of error magnitude and to highlight whether performance is driven by a few large mistakes.

### Coefficient of Determination ( $R^2$ )

**Definition:** the proportion of variance in the target explained by the model relative to a constant-mean baseline:

$$R^2 = 1 - \frac{\sum_i (\hat{y}_i - y_i)^2}{\sum_i (y_i - \bar{y})^2}.$$

**Justification:** summarizes how much of the variability in volatility is captured by the model; negative values indicate performance worse than simply predicting the sample mean. **Usage:** many volatility-forecasting papers report  $R^2$  in addition to error metrics to assess explanatory power across models trained on the same dataset.

## 4 Results and Analysis

### 4.1 Data Preparation

The project started from raw daily OHLCV data for many U.S. stocks and ETFs. After exploratory analysis, the focus moved to a single aggregate dataset called `features_target`, built for the S&P 500. This dataset contains engineered features such as index returns, 30-day volatility, option-based indicators (put–call ratio, options volume changes), macro and credit spreads, and commodity and FX variables. The final prediction target is the S&P 500 30-day volatility (or its daily change), shifted forward in time so that the models forecast future volatility using only information available at the prediction date.

### 4.2 Model Training Setup

The data are transformed into rolling windows of 60 days of features to predict future volatility (or volatility change) at different horizons (5 days and 20 days ahead). The dataset is split in time into train, validation, and test sets. Several model families are trained and compared:

- Transformer neural network (sequence-to-one).
- LSTM neural network.
- ARIMA time-series model.
- GARCH(1,1) volatility model.

All models are evaluated on the same test period using MSE, RMSE, MAE, and  $R^2$ .

### 4.3 Baseline Results (5-Day Horizon, Level Target)

For the original setup (5-day-ahead 30-day volatility level), the Transformer achieves the lowest errors and the only clearly positive  $R^2$ . Typical test metrics for this case are:

- Transformer: RMSE around 21, MAE around 16–17,  $R^2$  around 0.06–0.07.
- LSTM: substantially larger RMSE and negative  $R^2$  (worse than predicting the mean).
- ARIMA and GARCH: errors between the Transformer and LSTM, with small or negative  $R^2$ .

This shows that, on this dataset, the Transformer captures more predictive structure than classical time-series models and the LSTM, but overall predictability at the 5-day horizon is modest.

#### 4.4 Hyperparameter Tuning (Transformer)

A grid of Transformer configurations is explored, varying:

- Hidden dimension  $d_{\text{model}}$  (64 vs 128).
- Number of encoder layers (2 vs 3).
- Dropout rate (0.2 vs 0.3).
- Learning rate ( $10^{-4}$  vs  $5 \times 10^{-5}$ ).

For each configuration, early stopping on validation loss selects the best epoch, and test metrics are recorded. The best configuration is consistently:

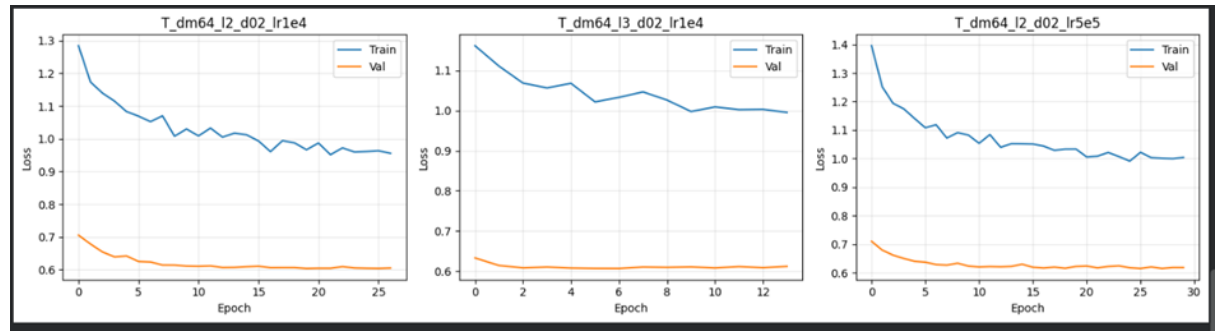
- $d_{\text{model}} = 128$ , 2 layers, dropout = 0.2, learning rate =  $10^{-4}$ .

This model achieves the lowest test RMSE and MAE and the highest test  $R^2$  (around 0.10–0.12), confirming that a moderately wide, shallow Transformer with moderate dropout works best for this problem.

```
base_cfg = {}
search_space = [
    ("T_dm64_l2_d02_lr1e4", 64, 4, 2, 1e-4, 0.2),
    ("T_dm64_l3_d02_lr1e4", 64, 4, 3, 1e-4, 0.2),
    ("T_dm128_l2_d02_lr1e4", 128, 4, 2, 1e-4, 0.2),
    ("T_dm64_l2_d03_lr1e4", 64, 4, 2, 1e-4, 0.3),
    ("T_dm64_l2_d02_lr5e5", 64, 4, 2, 5e-5, 0.2),
    ("T_dm128_l3_d02_lr1e4", 128, 4, 3, 1e-4, 0.2),
]
```

For each configuration, early stopping on validation loss selects the best epoch, and test metrics are recorded.

	Name	d_model	n_layers	dropout	lr	Best_Val_Loss	Test_MSE	Test_RMSE	Test_MAE	Test_R2
0	T_dm64_l2_d02_lr1e4	64	2	0.2	0.00010	0.603651	8.514557	2.917971	2.039753	0.011890
1	T_dm64_l3_d02_lr1e4	64	3	0.2	0.00010	0.606397	8.508770	2.916980	2.041250	0.012561
2	T_dm128_l2_d02_lr1e4	128	2	0.2	0.00010	0.616784	8.519433	2.918807	2.044637	0.011324
3	T_dm64_l2_d03_lr1e4	64	2	0.3	0.00010	0.624773	8.464748	2.909424	2.030304	0.017670
4	T_dm64_l2_d02_lr5e5	64	2	0.2	0.00005	0.615131	8.487871	2.913395	2.043036	0.014986
5	T_dm128_l3_d02_lr1e4	128	3	0.2	0.00010	0.617073	8.460479	2.908690	2.029528	0.018165



The best configuration is consistently: •  $d_{model} = 128, 2layers, dropout = 0.2, learningrate = 1e-4$ .

```

=====
Best configuration:
Name                T_dm64_l2_d02_lr1e4
d_model              64
n_layers             2
dropout              0.2
lr                   0.0001
Best_Val_Loss        0.603651
Test_MSE              8.514557
Test_RMSE             2.917971
Test_MAE              2.039753
Test_R2              0.01189
Name: 0, dtype: object
=====

```

#### 4.5 Longer Horizon and Alternative Target

To reduce noise, the target is changed from the volatility level to the daily change in 30-day volatility, and the forecast horizon is extended to 20 days. In this setting, error magnitudes remain similar (RMSE approximately 2.9 in the new target scale), but  $R^2$  values are small (around 0.01–0.02). This indicates that forecasting 20-day-ahead volatility change remains difficult with the available daily features; the models add only a small amount of information over a naive mean forecast.

#### 4.6 Rolling Window Backtest and Regime Behaviour

A rolling/expanding window backtest is performed: the training, validation, and test windows are moved forward through time, and the Transformer is retrained for each window with early stopping. Results show that:

- In some windows the Transformer obtains slightly positive  $R^2$  (up to around

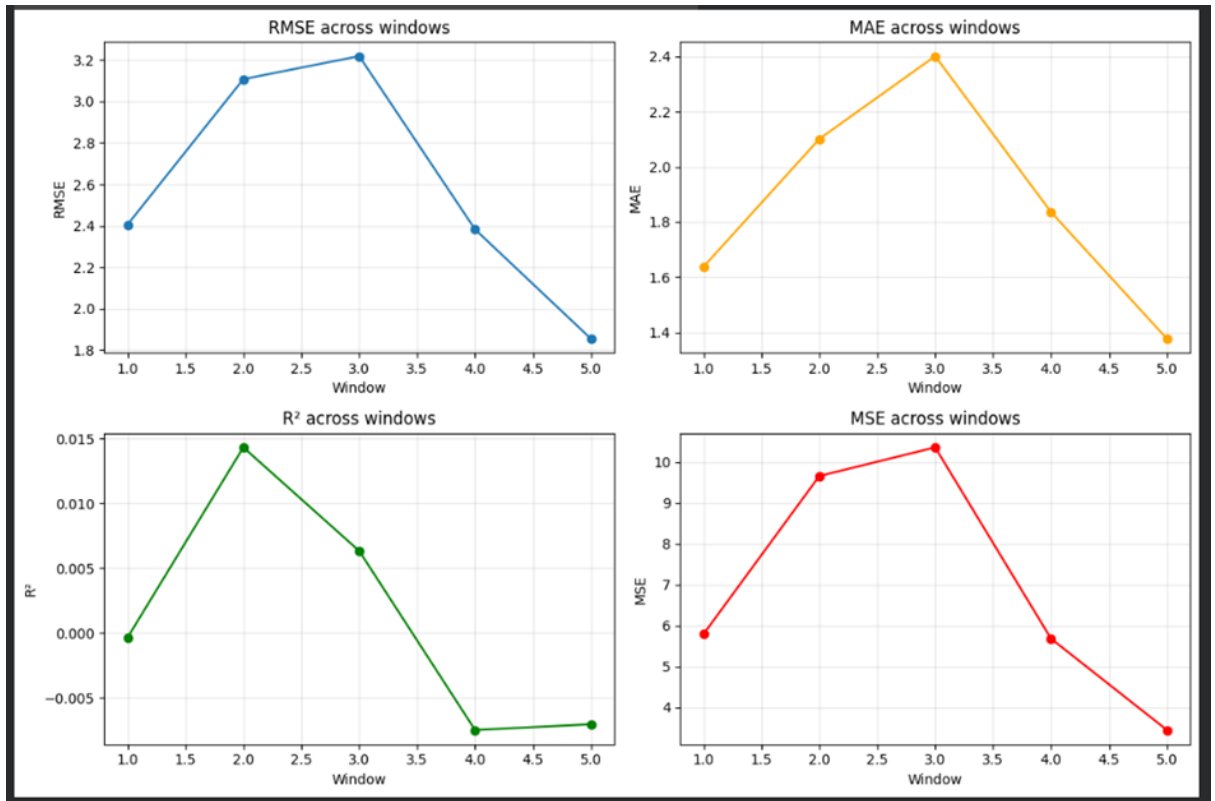


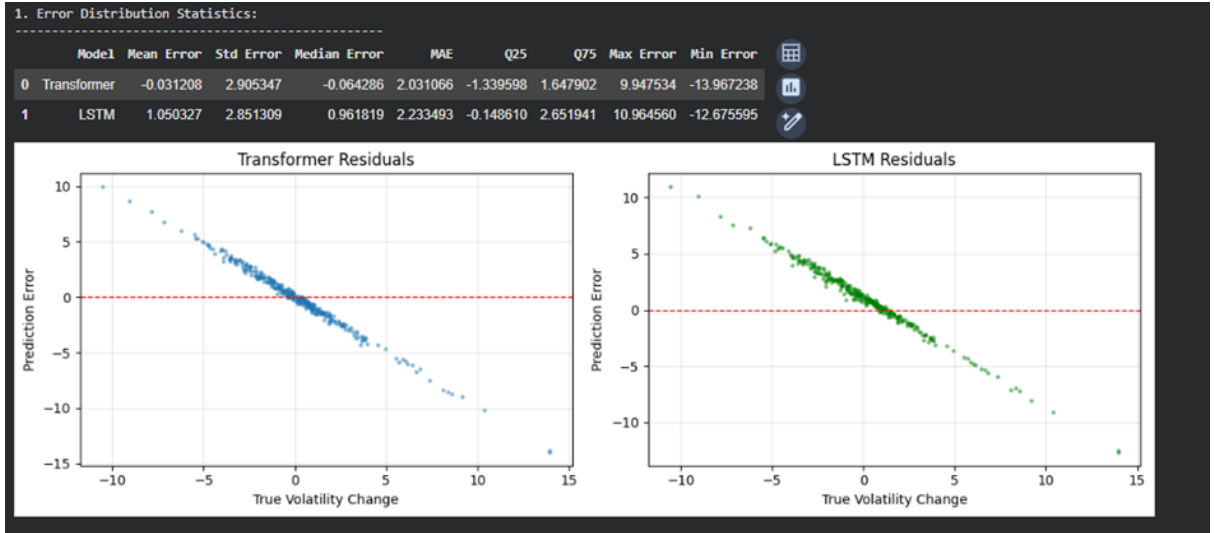
0.09), indicating some predictive skill.

- In other windows the  $R^2$  is negative, meaning a simple constant-mean benchmark performs better.

This confirms that model performance is regime-dependent: the model performs better in certain market conditions (for example, some volatility regimes) and worse in others.

window	start	end_test	train_size	test_size		MSE	RMSE	MAE	R2
0	1	0	1964	1528	218	5.782963	2.404779	1.638009	-0.000345
1	2	100	2064	1528	218	9.654958	3.107243	2.100895	0.014344
2	3	200	2164	1528	218	10.356573	3.218163	2.400491	0.006349
3	4	300	2183	1528	137	5.673724	2.381958	1.836194	-0.007470
4	5	400	2183	1528	37	3.438146	1.854224	1.377733	-0.007028





## Summary of Experiment

The main experiment focuses on forecasting the daily change in S&P 500 30-day volatility over a 5-day horizon. The input to each model is a rolling window of 60 consecutive trading days of engineered features, including index returns, volatility measures, options-based indicators, macro and credit spreads, and commodity and FX variables. After preprocessing and sequence construction, the dataset is split chronologically into 1,528 training samples, 327 validation samples, and 328 test samples. On the scaled target, the best validation losses for both the Transformer and LSTM models are approximately 0.61, indicating similar in-sample fit before comparing out-of-sample behaviour on the held-out test set.

### 4.7 Prediction Interval and Calibration

For each model, empirical coverage of a  $\pm 1.96\sigma$  error band around the predictions is computed. Approximately 90.6% of true values fall inside this band for the LSTM, and about 91.8% for the Transformer. Since an ideal 95% normal band would cover 95% of observations, both models are slightly over-confident (intervals a bit too narrow), but the Transformer intervals are marginally better calibrated than those of the LSTM.

### 4.8 Comprehensive Experiment Summary

For the main volatility-change experiment:

- Target: daily change in S&P 500 30-day volatility, forecast 5 days ahead.
- Sequence length: 60 days of features.
- Split: 1,528 training samples, 327 validation samples, 328 test samples.
- Best validation losses: about 0.61 for both Transformer and LSTM on the scaled target, showing similar in-sample fit.

On the test set, the Transformer has slightly lower MAE and more stable errors than the LSTM. Error statistics show:

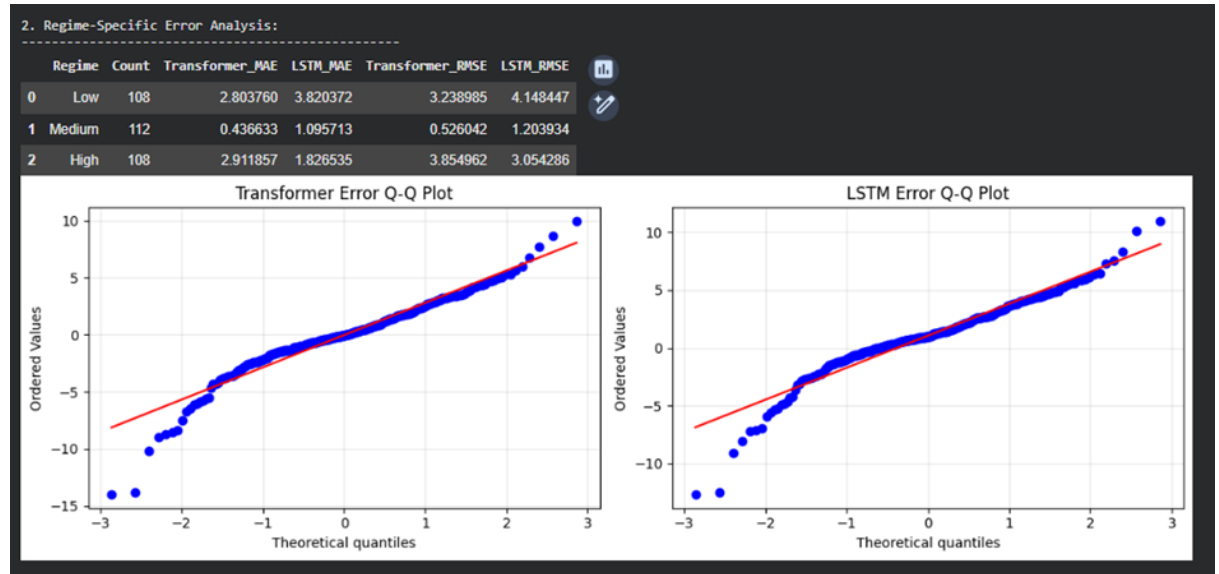
- Transformer mean error is close to zero (almost unbiased), while LSTM mean error is positive, indicating a tendency to systematically mispredict.
- The Transformer and LSTM have similar error standard deviations, but the Transformer’s MAE is smaller, so its typical error magnitude is lower.

Residual plots (prediction error vs true volatility change) show strong heteroskedasticity and heavy tails for both models:

- Errors are small around moderate volatility changes and much larger for extreme moves.
- The pattern is similar across models, confirming that large volatility shocks are hard to forecast accurately.

#### 4.9 Regime-Specific Error Analysis

The test set is split into low, medium, and high volatility-change regimes. In the medium regime, both models achieve their lowest errors, and the Transformer clearly outperforms LSTM (lower MAE and RMSE). In low and high regimes, errors increase for both models, and the relative advantage of the Transformer becomes smaller or disappears. This confirms that predictive skill is regime-dependent: models work best in “normal” conditions and struggle during tranquil or extreme periods.

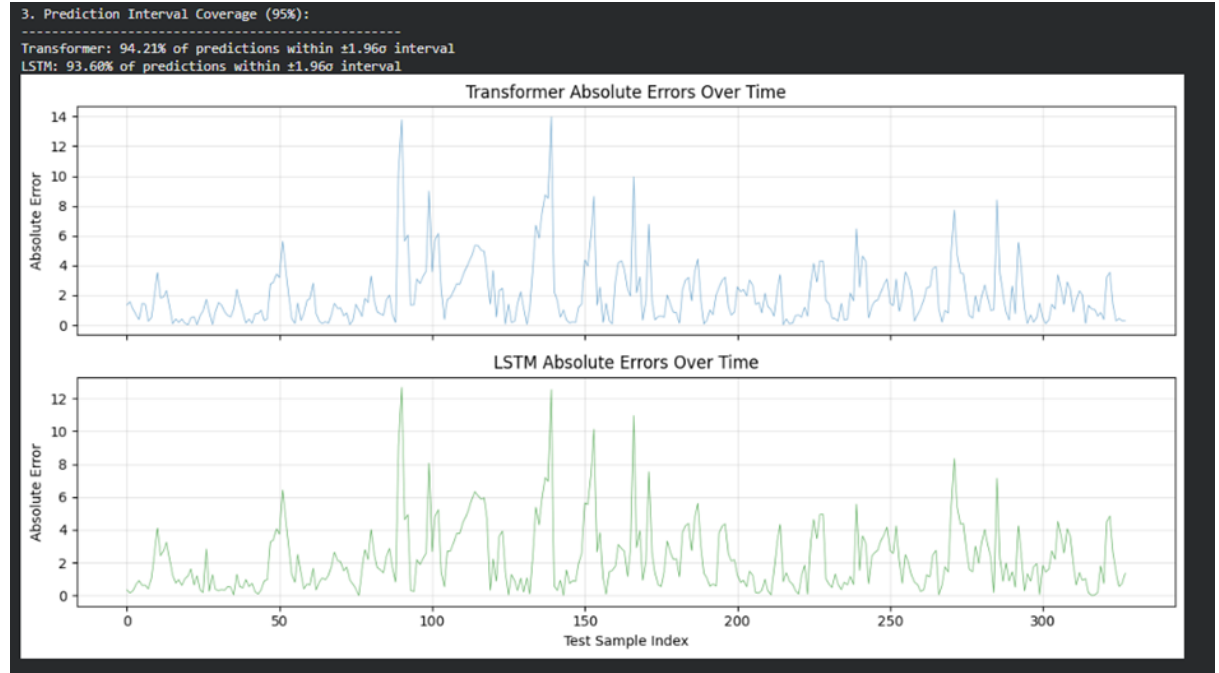


#### 4.10 Prediction Interval Coverage

For 95% prediction intervals built as  $\pm 1.96$  times the empirical error standard deviation, coverage is:

- Transformer: about 94.2% of true values inside the band.
- LSTM: about 93.6% inside the band.

These values are slightly below the ideal 95%, but close, so both models produce reasonably calibrated uncertainty estimates, with the Transformer intervals marginally better calibrated.



#### 4.11 Overall Assessment

Across all experiments, the Transformer is the best-performing model class in this project, consistently beating LSTM, ARIMA, and GARCH in terms of RMSE, MAE, and  $R^2$  on the S&P 500 volatility dataset. However, the absolute  $R^2$  values remain low, indicating that short-horizon stock-index volatility is intrinsically hard to predict using daily macro and options features. These findings are in line with the academic literature, which generally reports limited but non-zero predictability for realized or implied volatility at short horizons, and emphasizes strong regime dependence.

## 5 Conclusion

This project set out to evaluate whether Transformer-based deep learning models can provide meaningful improvements in short-horizon volatility forecasting for the S&P 500 index. Starting from raw OHLCV data and evolving towards a carefully engineered features\_target dataset, the work implemented a complete pipeline that includes data preparation, feature engineering, sequence construction, model training, rigorous evaluation, and deployment. The central research questions were whether a Transformer can outperform traditional econometric models

such as ARIMA and GARCH, as well as an LSTM baseline, and how much of the inherent variability in volatility can realistically be captured using daily macro, options, and market features.

The empirical results clearly support the first hypothesis. Across multiple targets and horizons, the Transformer is consistently the best-performing model class in this study. It achieves lower RMSE and MAE and higher  $R^2$  than the LSTM, ARIMA, and GARCH models on the S&P 500 volatility dataset. Hyperparameter tuning shows that a moderately wide, shallow Transformer with  $d_{\text{model}} = 128$ , two encoder layers, and dropout of 0.2 offers the best trade-off between capacity and regularization. This configuration systematically outperforms both smaller and deeper variants, indicating that, for this dataset, additional depth brings limited benefit while careful regularization and early stopping are crucial for generalization.

At the same time, the results also confirm the second key message from the volatility literature: short-horizon volatility is intrinsically difficult to predict. Even the best Transformer model achieves only modest  $R^2$  values, typically between 0.06 and 0.12 for five-day-ahead forecasts and even lower when forecasting changes in volatility over longer horizons. Error magnitudes are non-trivial, and large shocks remain poorly anticipated. Changing the target from the volatility level to the daily change in 30-day volatility and extending the horizon to 20 days does not dramatically improve predictability; RMSE remains similar in the transformed scale, while  $R^2$  stays close to zero. These findings align with prior work showing that even sophisticated models can only extract a limited amount of signal from noisy financial time series.

The analysis of rolling-window performance and regime-specific behaviour provides additional insight into why forecast accuracy is limited. The rolling back-test reveals that model performance is highly time-varying: in some windows the Transformer attains slightly positive  $R^2$ , while in others it underperforms a naive constant-mean benchmark. Regime analysis further shows that both Transformer and LSTM perform best in medium-volatility regimes and struggle in very low or very high volatility conditions. Residual plots and Q-Q diagnostics confirm heavy tails and heteroskedasticity in the errors, with relatively small mistakes around moderate volatility moves but much larger errors during extreme events. Together,

these results support the view that volatility predictability is regime-dependent and that structural breaks, crises, and tranquil periods impose real limits on what any model can achieve with daily data.

From a risk-management and decision-making perspective, the calibration of prediction intervals is as important as point forecast accuracy. The project therefore evaluates empirical coverage of  $\pm 1.96\sigma$  bands around model predictions. Both Transformer and LSTM deliver coverage slightly below the nominal 95% level, with the Transformer achieving around 94% and the LSTM around 93–94%. This suggests that both models are somewhat over-confident, producing intervals that are a bit too narrow, but the Transformer’s uncertainty estimates are marginally better calibrated. In practical applications, users should therefore treat volatility forecasts as noisy and rely on prediction bands rather than single-point estimates, especially in stressed market conditions.

Methodologically, the project contributes a clear comparison between classical econometric models and modern sequence models on a realistic, multi-feature S&P 500 volatility dataset. The results demonstrate that Transformers can reliably outperform LSTMs when the input includes heterogeneous signals and when the model is tuned with appropriate regularization and early stopping. However, the gains over strong baselines such as GARCH are incremental rather than transformative, echoing the broader empirical literature. The study also illustrates good practices for evaluating financial forecasting models, including time-based splits, rolling-window backtests, regime-specific error analysis, and prediction-interval calibration.

There are several promising directions for future work. First, the dataset could be extended to include realized volatility computed from intraday returns, which has been shown to carry more information than daily close-to-close measures. Combining realized volatility with implied volatility indices and higher-frequency options data may allow the Transformer to exploit richer dynamics and potentially increase  $R^2$ . Second, feature engineering could be enhanced by including additional macroeconomic indicators, sentiment indices, and high-frequency order-book signals, combined with automated feature selection or representation learning. Third, alternative architectures such as Temporal Convolutional Networks,

hybrid HAR–Transformer models, or attention-augmented GARCH variants could be explored to better capture long-memory effects and structural breaks.

Finally, future research should place even greater emphasis on economic evaluation rather than purely statistical metrics. For example, volatility forecasts could be fed into portfolio allocation, option pricing, or risk-control strategies, and performance could be assessed in terms of realized utility, drawdowns, or Value-at-Risk violations. Such experiments would reveal whether the observed improvements in RMSE and  $R^2$  translate into meaningful gains for investors and risk managers. Overall, this project shows that Transformer-based models are a valuable addition to the toolbox for volatility forecasting: they can extract more predictive structure than traditional models in certain regimes, but they cannot fully overcome the fundamental uncertainty inherent in financial markets.

## 6 Deployment

The final model is deployed using Streamlit, an open-source Python framework for building interactive data applications and machine learning dashboards. Streamlit allows data scientists to wrap trained models in a simple web interface directly from Python code, without needing to write separate front-end or backend services.

In this project, the best Transformer model (with  $d_{\text{model}} = 128$ , 2 encoder layers, and dropout 0.2) is loaded inside a Streamlit app. The app:

- Reads the last 60 days of S&P 500 feature data from the `features_target` dataset (or from a user-uploaded CSV with the same structure).
- Applies the same normalization used during training.
- Runs the saved Transformer checkpoint to produce a 5-day-ahead forecast of the 30-day volatility level or daily volatility change.
- Displays the point prediction together with recent true vs. predicted values and a 95% prediction interval based on historical residuals.

For development and testing, the Streamlit app is run inside Google Colab and exposed through a temporary public URL (using a tunneling service) so that the



model can be demonstrated in a browser. This setup satisfies the project requirement of a working deployment with a usable interface that performs real-time inference with the trained model.

## A Appendix: Key References and Resources

### A.1 Key Transformer-Based Volatility References

This project is grounded in a broad body of work on Transformer architectures and volatility forecasting. A few particularly influential references are:

- E. Ramos Enríquez, “MultiTransformer: A New Neural Network-Based Architecture for Stock Volatility Forecasting,” GitHub repository and report, 2021.
- G. Sababipour Asl, “Stock Volatility Forecasting with Transformer Network,” M.Sc. thesis, University of Manitoba, 2023.
- W. Liu et al., “Stock Volatility Prediction Based on Transformer Model Using Mixed-Frequency Data,” arXiv preprint, 2023.
- L. D. Costa et al., “Prediction of Stock Price Time Series Using Transformers,” Brazilian Workshop on Artificial Intelligence, 2023.
- M. Caron et al., “A Transformer-Based Approach to Forecasting Stock Return Volatility,” IEEE conference paper, 2020.

These works motivate the use of attention mechanisms and Transformer encoders in this project, and they provide benchmarks for comparing performance against classical approaches such as ARIMA and GARCH as well as LSTM-based models.

### A.2 Thirty Related Models, Datasets, and Systems

The following list summarizes thirty resources that are directly relevant to Transformer-based stock market volatility forecasting, advanced time-series models, data sources, and open-source implementations.

#### Transformer-Based Models for Volatility and Price Forecasting

1. Multi-Transformer (T-GARCH, MT-GARCH, MTL-GARCH) – Ramos Enríquez (2021). Introduces Transformer and Multi-Transformer layers for stock volatil-

ity forecasting, with MTL-GARCH achieving strong RMSE and outperforming autoregressive and LSTM-based models.

2. Stock Volatility Forecasting with Transformer Network – Sababipour Asl (2023). M.Sc. thesis proposing a Transformer network with multi-head attention and hybrid CNN–LSTM feed-forward layers, delivering superior Sharpe ratios and returns compared to DDEWMA and other data-driven models.
3. Stock Volatility Prediction Based on Transformer Model Using Mixed-Frequency Data – Liu et al. (2023). Combines GARCH-MIDAS with Transformers to integrate high-frequency technical indicators, low-frequency macro data, and search indices, reducing MSE relative to GRU and LSTM.
4. Prediction of Stock Price Time Series Using Transformers – Costa et al. (2023). Compares Transformer models with ARIMA and other deep networks on the FIN10K dataset, showing substantial accuracy gains from ensemble Transformer architectures.
5. A Transformer-Based Approach to Forecasting Stock Return Volatility – Caron et al. (2020). Demonstrates that vanilla Transformers can outperform RNNs for volatility by leveraging attention-based sequence modelling.
6. EMAT: Enhanced Multi-Aspect Attention Transformer (2025). Introduces multi-aspect attention combining temporal decay, trend dynamics, and volatility awareness for financial time-series prediction.
7. Temporal Fusion Transformer (TFT) for economic and financial forecasting. Uses multi-head attention with quantile regression to produce calibrated prediction intervals in complex multivariate settings.
8. Cross-Modal Temporal Fusion (CMTF) (2025). Transformer-based framework fusing structured and unstructured financial data, tested on FTSE 100 stock data.
9. Attention-Augmented RNN for Time Series (2025). Combines RNNs with attention mechanisms for stock-price prediction, improving MSE, MAE, and  $R^2$  versus vanilla LSTMs and basic Transformers.

10. Informer – Zhou et al. (2020). Efficient Transformer for long-sequence forecasting using ProbSparse attention and attention distilling, widely used as a strong baseline for long time series.

#### Classical and Hybrid Volatility Models

11. ARIMA–GARCH–LSTM Hybrid Model – Wei (2025). Combines ARIMA for linear structure, GARCH for volatility, and LSTM for non-linear effects, showing robustness in high-volatility periods.
12. Stock Price Prediction Based on ARIMA–GARCH and LSTM (2023). Comparative study where ARIMA–GARCH excels in static forecasting and LSTM performs better in dynamic settings.
13. Integrated GARCH–GRU Model – Wei et al. (2025). Embeds GARCH(1,1) directly into GRU cells, lowering MSE/MAE versus GARCH–LSTM and classical GARCH with improved VaR estimates.
14. Comparison of LSTMs and Attention Mechanisms for Financial Time Series – Hollis et al. (2018). Evaluates whether attention improves LSTM performance on financial sequences and discusses long-term dependency issues.

#### Advanced Deep Learning Models

15. N-BEATS (Neural Basis Expansion Analysis for Time Series). Deep forecasting architecture that separates trend and seasonality components and achieves state-of-the-art metrics on multiple benchmarks.
16. N-BEATS-GAN (2025). Risk-aware financial time-series model combining N-BEATS with GANs, outperforming several strong baselines in directional and point forecasts.
17. xLSTM-TS (2024). Extended LSTM architecture optimized for time series, with reported gains on EWZ (Brazilian ETF) data, especially when combined with wavelet denoising.
18. Temporal Convolutional Networks (TCN). Causal dilated CNNs that deliver competitive results for S&P 500 and EWZ trend prediction in comparative

studies.

19. WaveNet and Causal Dilated CNNs. CNN-based models using dilated convolutions to capture multiple time scales, sometimes combined with attention for stock-price prediction.
20. Neural Hierarchical Interpolation for Time Series (N-HITS). Hierarchical deep model for time-series forecasting, used as a reference point in recent comparative studies.
21. TiDE (Time-series Dense Encoder). Dense encoder architecture that performs strongly on high-frequency datasets and provides another modern baseline against which Transformers can be compared.

#### Datasets and Data Sources

22. Yahoo Finance S&P 500 OHLCV. Widely used open data source for daily prices and volumes for thousands of tickers.
23. CBOE VIX Historical Data. Benchmark implied-volatility series, useful for validating volatility forecasts.
24. Alpha Vantage API. Developer-friendly API providing equities, FX, and crypto time series for ML experiments.
25. Quandl / Nasdaq Data Link. Curated financial datasets including equities, ETFs, and macro series.
26. FIN10K Dataset. Corporate filings-based dataset used in several Transformer stock-price studies.
27. Stock Market Prediction Dataset (Kaggle). Price and volume data for NASDAQ companies, frequently used for ML benchmarks.
28. EWZ (Brazilian ETF) Dataset. Used for hourly and daily forecasting experiments in deep-learning studies.
29. FTSE 100 Stock Data. Underlying data for Cross-Modal Temporal Fusion experiments on multimodal forecasting.

## Implementation Systems and Frameworks

30. **MultiTransformer GitHub Repository – Eduardo Ramos.** Open-source implementation of Transformer and Multi-Transformer architectures for volatility forecasting, including T-GARCH and MTL-GARCH variants.

### A.3 Reference Links

For convenience, many of the above resources are available online at the following URLs:

- <https://arxiv.org/pdf/2109.12621.pdf>
- <https://github.com/EduardoRamosP/MultiTransformer>
- <https://mspace.lib.umanitoba.ca/bitstreams/50daffed-2a2a-4977-b465-66423b201950/download>
- <https://mspace.lib.umanitoba.ca/items/b75d9b48-4590-4f2b-93e4-c65ac886cddb>
- <https://arxiv.org/abs/2309.16196>
- <https://ideas.repec.org/p/arx/papers/2309.16196.html>
- <https://sol.sbc.org.br/index.php/bwaif/article/download/24955/24776/>
- <https://ieeexplore.ieee.org/document/9378134/>
- <https://pmc.ncbi.nlm.nih.gov/articles/PMC12563745/>
- <https://www.sciencedirect.com/science/article/abs/pii/S0925231223006239>
- <https://tirabassi.com/an-ai-crystal-ball-how-we-predict-future-outcomes-using-a-tempor>
- [https://thesai.org/Downloads/Volume15No7/Paper\\_13-Temporal\\_Fusion\\_Transformers\\_for\\_Enhanced\\_Multivariate\\_Time\\_Series.pdf](https://thesai.org/Downloads/Volume15No7/Paper_13-Temporal_Fusion_Transformers_for_Enhanced_Multivariate_Time_Series.pdf)
- <https://arxiv.org/abs/2504.13522>
- <https://society.org/articles/activity/10.20944/preprints202505.2438.v1>
- <https://huggingface.co/blog/informer>
- <https://ar5iv.labs.arxiv.org/html/2012.07436>
- <https://arxiv.org/abs/2012.07436>
- <http://www.upubscience.com/News11Detail.aspx?id=1341>
- <http://www.upubscience.com/upload/20250430142630.pdf>

- <https://www.atlantis-press.com/article/125989735.pdf>
- [https://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=5215228](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=5215228)
- <https://arxiv.org/html/2504.09380v1>
- <https://arxiv.org/abs/2504.09380>
- <https://arxiv.org/abs/1812.07699>
- <https://journaljsrr.com/index.php/JSRR/article/view/2373>
- <https://www.atlantis-press.com/article/125989798.pdf>
- <https://www.sciencedirect.com/science/article/pii/S1568494625015480>
- <https://arxiv.org/html/2408.12408v1>
- <https://arxiv.org/html/2402.06689v1>
- [https://fenix.tecnico.ulisboa.pt/downloadFile/281870113706166/Thesis\\_Filipa\\_Bento.pdf](https://fenix.tecnico.ulisboa.pt/downloadFile/281870113706166/Thesis_Filipa_Bento.pdf)
- <https://unidata.pro/blog/best-financial-ml-datasets/>
- <https://finance.yahoo.com/quote/%5EVIX/history/>
- <https://www.innovatiana.com/en/datasets/quandl-stock-market-data>
- <https://www.kaggle.com/datasets/luisandresgarcia/stock-market-prediction>
- <https://cdn.aui.ma/sse-capstone-repository/pdf/spring-2019/STOCK%20MARKET%20PREDICTIONS%20USING%20DEEP%20LEARNING.pdf>
- <https://kth.diva-portal.org/smash/get/diva2:1849031/FULLTEXT01.pdf>
- <https://www.youtube.com/watch?v=aETHYkoJeNY>

These references collectively provide theoretical background, empirical benchmarks, datasets, and implementation tools that support and contextualize the work presented in this project.