

LSTM-GARCH volatility forecasting on log-return of S&P 500

HKUST EMIA4991 Independent Capstone Project: 2025-Spring

Author:

CHAN, Pak Lok (SID: 20855554)

Supervisor:

LUO, Wen Han

Date of Report: May 8, 2025

Objective Statement:

- *To design and implement a hybrid LSTM-GARCH model to predict the 1-day-ahead volatility of the S&P 500 log-returns, targeting at least a 5% reduction in forecast error relative to a baseline GARCH model.*
- *To deliver a more robust volatility forecasting framework by fusing the nonlinear-pattern recognition of LSTMs with the statistical rigor of GARCH, and incorporating the VIX index exogenously.*

Abstract

This project investigates the effectiveness of combining a Long Short-Term Memory (LSTM) network with a Generalized Autoregressive Conditional Heteroskedasticity (GARCH) model to forecast the realized volatility of the S&P 500 index (^GSPC). We utilize historical S&P 500 price data and VIX index levels from January 2000 to December 2020. An LSTM model is trained to predict next-day realized volatility using past log returns and VIX levels, while a standard GARCH(1,1) model forecasts volatility based on past log returns. A ensemble forecast is generated by averaging the predictions from the LSTM and GARCH models. Using a rolling forecast evaluation on a test set spanning from late 2016 to the end of 2020, we compare the performance of the individual models and the ensemble approach against the actual realized volatility. Evaluation metrics, including Mean Absolute Error (MAE), Mean Squared Error (MSE), and Root Mean Squared Error (RMSE), demonstrate that the combined LSTM-GARCH model outperforms both the standalone GARCH(1,1) and LSTM models, achieving an 16.64% and 5.26% reduction in RMSE respectively.

Contents

Abstract i

Contents ii

List of Figures iii

List of Tables iv

Preliminary Knowledge v

Mathematical Knowledge v

Financial Knowledge v

Artificial Intelligence Knowledge vi

1 Introduction 1

1.1 Background and Problem 1

1.2 Objectives 1

1.2.1 Objective Statements 1

1.3 Literature Review of Existing Solutions 2

2 Methodology 4

2.1 Overview 4

2.2 Data Preparation 4

2.2.1 Data Source 4

2.2.2 Data Preprocessing and Analysis 4

2.3 Model Architecture 8

2.3.1 LSTM Component 8

2.3.2 GARCH Component 9

2.3.3 Combined LSTM-GARCH Model 9

2.3.4 Rolling Window Forecasting 9

2.4 Evaluation Metrics 10

3 Results 11

3.1 Initial Model Training 11

3.2 Rolling Forecast Performance 11

3.3 Objective Check 12

4 Conclusion 13

5 Future Work 13

References 14

Appendices 16

List of Figures

1	Log Returns of S&P 500 Index (2000-2020)	4
2	ACF (left) and PACF (right) of the S&P 500 log-return series.	6
3	ARMA(1,1) residuals (top) and squared residuals (bottom).	6
4	Conditional volatility σ_t estimated by the GARCH(1,1) model.	7
5	Standardized residuals (top) and their squares (bottom) from the GARCH(1,1) fit.	7
6	LSTM Network Architecture for Volatility Prediction.	9
7	LSTM Training Loss Curve	11
8	Predicted vs. Actual Realized Volatility (Test Set)	12

List of Tables

1	Augmented Dickey–Fuller Test on S&P 500 Log Returns	5
2	ARMA(1,1) Parameter Estimates	5
3	ARMA Residual Diagnostic Tests	6
4	GARCH(1,1) Parameter Estimates	7
5	Proposed vs. Best-Tuned LSTM Hyperparameters and Performance	8
6	Initial GARCH(1,1) Model Parameter Estimates (Training Set)	11
7	Volatility Forecast Performance Metrics on Test Set (2016-10-19 to 2020-12-30)	12

Preliminary Knowledge

This section outlines the essential knowledge and techniques that will be applied in this research. The focus is on the mathematical, financial, and artificial intelligence (AI) knowledge that underpins the development of the hybrid LSTM-GARCH model for volatility forecasting.

Mathematical Knowledge

This subsection provides the mathematical foundations necessary for understanding the models and methods employed in this research.

- **Time Series Analysis:** Fundamental concepts of time series analysis, including stationarity, autocorrelation, and seasonality, are crucial for understanding the behavior of financial data over time. The Augmented Dickey-Fuller (ADF) test will be used to assess stationarity (Hamilton, 1994) (Dickey & Fuller, 1979).
- **Log Returns:** The log-return of an asset at time t is defined as:

$$r_t = \ln \left(\frac{P_t}{P_{t-1}} \right), \quad (1)$$

where P_t and P_{t-1} denote the closing prices at times t and $t - 1$, respectively. (Tsay, 2010)

- **Realized Volatility:** Realized volatility is an ex-post measure of actual volatility calculated from historical data. In this study, it is computed as the annualized standard deviation of log-returns over a rolling window of $k = 5$ days:

$$\sigma_{RV,t} = \sqrt{252} \times \sqrt{\frac{1}{k-1} \sum_{i=t-k+1}^t \left(r_i - \bar{r}_t^{(k)} \right)^2}, \quad (2)$$

where $\bar{r}_t^{(k)} = \frac{1}{k} \sum_{i=t-k+1}^t r_i$ is the mean of the log returns over the k -day window. (McAleer & Medeiros, 2008)

- **GARCH Model:** The Generalized Autoregressive Conditional Heteroskedasticity (GARCH) model estimates time-varying volatility in financial time series. The GARCH(1,1) model, used in this study, specifies the conditional variance h_t as:

$$h_t = \omega + \alpha \epsilon_{t-1}^2 + \beta h_{t-1} \quad (3)$$

where ϵ_{t-1} is the log return (innovation) at time $t - 1$, and ω, α, β are parameters estimated via maximum likelihood. This model captures volatility clustering, a common feature in financial data where large changes tend to follow large changes (Bollerslev, 1986).

Financial Knowledge

This subsection covers the financial concepts relevant to this research, providing context for the S&P 500 index, volatility, and the VIX index.

- **Volatility in Finance:** Volatility measures the degree of variation in the price of a financial instrument over time, serving as an indicator of risk. It is typically quantified as the standard deviation of log-returns, reflecting the uncertainty or dispersion of asset returns. In financial markets, high volatility is associated with greater risk and larger price fluctuations. (Campbell et al., 1998)

- **S&P 500 Index:** The S&P 500 is a stock market index tracking the performance of 500 large companies listed on U.S. stock exchanges. It serves as a leading benchmark for the U.S. equity market and is widely used to assess overall market conditions and economic health. Its volatility is a critical focus of this study due to its significance in financial risk management (Cont, 2001).
- **VIX Index:** The VIX index, or Volatility Index, measures the market's expectation of 30-day forward-looking volatility, derived from S&P 500 index option prices. Often called the "fear index," it reflects investor sentiment and market risk perceptions. In this research, the VIX is incorporated as an exogenous variable to enhance volatility forecasts by leveraging market-implied information (Whaley, 2000).

Artificial Intelligence Knowledge

This subsection outlines the AI concepts and techniques that will be applied in this research, providing context for the LSTM model, hyperparameter optimization, regularization techniques, optimization algorithms, and model evaluation methods.

- **Deep Learning Architectures:** Concepts of recurrent neural networks (RNNs), particularly Long Short-Term Memory (LSTM) units, will be leveraged to design the sequence-to-sequence structure that processes historical return series. The LSTM unit is mathematically represented as follows:

$$\begin{aligned}
f_t &= \sigma(W_f x_t + U_f h_{t-1} + b_f) \\
i_t &= \sigma(W_i x_t + U_i h_{t-1} + b_i) \\
o_t &= \sigma(W_o x_t + U_o h_{t-1} + b_o) \\
\tilde{c}_t &= \tanh(W_c x_t + U_c h_{t-1} + b_c) \\
c_t &= f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \\
h_t &= o_t \odot \tanh(c_t)
\end{aligned}$$

where x_t is the input at time t , h_t is the hidden state, c_t is the cell state, σ is the sigmoid function, and \odot denotes element-wise multiplication. This structure allows the model to retain and propagate long-term dependencies in financial time series.

- **Optimization Algorithms:** In this research, we will utilize the Adam optimizer for training the LSTM model. The loss function will be defined as the Mean Squared Error (MSE) between the predicted and actual realized volatility.
- **Model Evaluation:** In this research, we apply 3 evaluation metrics to assess forecast accuracy: Mean Squared Error (MSE), Mean Absolute Error (MAE), and RMSE.

1 Introduction

1.1 Background and Problem

Volatility, a measure of the degree of variation in asset prices over time, is a key indicator of market uncertainty and risk in financial markets (Bollerslev, 1986). The S&P 500 index, as one of the most followed equity benchmarks globally, exhibits pronounced time-varying volatility, posing significant challenges for accurate forecasting (Cont, 2001). Accurate volatility forecasts are critical for risk management, portfolio allocation, and derivative pricing (Jorion, 2007).

Traditional time-series models such as ARIMA (AutoRegressive Integrated Moving Average) capture linear dependencies in returns but assume constant variance, limiting their ability to model volatility clustering (Box et al., 1970). GARCH (Generalized Autoregressive Conditional Heteroskedasticity) models directly address this by allowing the conditional variance to depend on past squared innovations and past variances (Engle, 1982)(Bollerslev, 1986). However, GARCH-family models typically rely on linear dynamics and Gaussian or Student- t innovations, which may not fully capture nonlinear patterns or the extreme “fat-tailed” behavior observed in financial returns (Bollerslev et al., 2008).

Deep learning—particularly Long Short-Term Memory (LSTM) networks—has emerged as a powerful tool for sequence modeling, capable of capturing complex, nonlinear dependencies in time series (Hochreiter & Schmidhuber, 1997). In finance, LSTMs have been shown to outperform classical econometric models in directional forecasting of equity returns and volatility (Fischer & Krauss, 2018). Yet, vanilla LSTMs do not explicitly model heteroskedasticity or volatility clustering, potentially overlooking well-established statistical properties of asset-return volatility.

Moreover, market-implied measures such as the VIX index incorporate expectations of future volatility derived from option prices (Whaley, 2000) and have been used exogenously to improve forecast accuracy when combined with historical-return information (Giot & Laurent, 2005). A hybrid LSTM-GARCH architecture that integrates both deep-learning for nonlinear pattern extraction and GARCH-type components for volatility clustering—augmented by VIX as an exogenous input—offers a promising avenue to enhance one-day-ahead volatility forecasts.

1.2 Objectives

This research aims to address the limitations of traditional volatility models by implementing a hybrid LSTM-GARCH model to predict the 1-day-ahead volatility of the S&P 500 index. By combining the interpretability of LSTM networks to capture the non-linear patterns in financial data with the ability of GARCH models to capture the volatility clustering effect, the study will provide a more accurate framework for volatility forecasting. Additionally, the model incorporates the VIX as an exogenous variable alongside the historical S&P 500 log returns to enhance the predictive performance of the model. This addresses the limitations of standalone models noted by Bollerslev et al. (2008) by combining LSTM’s nonlinear capabilities with GARCH’s statistical structure.

1.2.1 Objective Statements

- *To design and implement a hybrid LSTM-GARCH model to predict the 1-day-ahead volatility of the S&P 500 log-returns, targeting at least a 5% reduction in forecast error relative to a baseline GARCH model.*
- *To deliver a more robust volatility forecasting framework by fusing the nonlinear-pattern recognition of LSTMs with the statistical rigor of GARCH, and incorporating the VIX index exogenously.*

1.3 Literature Review of Existing Solutions

Financial volatility forecasting possesses a rich history rooted in econometric modeling, fundamentally shaped by Engle's seminal work on Autoregressive Conditional Heteroskedasticity (ARCH) and its subsequent generalizations. Engle (1982) introduced the ARCH framework, enabling the conditional variance of a financial time series to depend explicitly on the magnitude of past squared innovations (errors). This elegantly captured the empirically observed phenomenon of volatility clustering, where large price changes tend to be followed by further large changes, and small changes by small changes. Recognizing the potential need for many lags in the original ARCH specification, Bollerslev (1986) proposed the Generalized ARCH (GARCH(p, q)) model. By incorporating lagged conditional variances (q lags) alongside lagged squared innovations (p lags), GARCH provided a more parsimonious and flexible representation of the persistent nature of volatility dynamics often seen in financial markets.

While the standard GARCH model effectively captures volatility clustering, subsequent research identified key empirical stylized facts it failed to address adequately, primarily asymmetry (the leverage effect) and non-normality of asset returns. The leverage effect, where negative shocks (price drops) tend to increase future volatility more than positive shocks (price increases) of the same magnitude, motivated new specifications. Nelson (1991) developed the Exponential GARCH (EGARCH) model, which models the logarithm of the conditional variance. This specification cleverly accommodates leverage effects without imposing non-negativity constraints on model parameters, a practical advantage over standard GARCH. Concurrently, Glosten et al. (1993) introduced the GJR-GARCH model (named after Glosten, Jagannathan, and Runkle), which adds a term to explicitly capture the differential impact of positive versus negative past shocks. Beyond asymmetry, researchers explored long memory characteristics in volatility, leading to models like Fractionally Integrated GARCH (FIGARCH) (Baillie et al., 1996). Further extensions like the GARCH-in-Mean (GARCH-M) model (Engle et al., 1987) incorporate conditional volatility directly into the conditional mean equation, allowing risk to affect expected returns. Models like the Asymmetric Power ARCH (APARCH) (Ding et al., 1993) offer additional flexibility in capturing leverage and the shape of the volatility response by allowing for a flexible power term. Furthermore, recognizing the heavy tails often observed in financial return distributions, many GARCH variants incorporate non-Gaussian error distributions, such as the Student's t -distribution (Bollerslev, 1987) or the Generalized Error Distribution (GED).

Parallel to the development of ARCH-type models, the Stochastic Volatility (SV) framework emerged as a significant alternative (S. J. Taylor, 1986). Unlike GARCH, where volatility is a deterministic function of past information, SV models treat volatility as an unobserved latent (stochastic) process, often modeled using an autoregressive structure for the log-variance. While computationally more demanding due to the latent variable (often requiring simulation-based estimation methods like Markov Chain Monte Carlo), SV models offer theoretical appeal and flexibility in capturing volatility dynamics. Another major advancement, spurred by the increasing availability of high-frequency intraday data, was the development of Realized Volatility (RV) measures (Andersen et al., 2003). By summing squared high-frequency returns (e.g., 5-minute returns) over a specific period (e.g., a day), RV provides a more accurate, model-free estimate of ex-post volatility compared to using squared daily returns. This led to new modeling approaches focused directly on forecasting RV, such as the Heterogeneous Autoregressive (HAR) model (Corsi, 2009), which uses lagged RV components from different time horizons (daily, weekly, monthly) to predict future realized volatility, capturing multi-scale volatility dynamics simply and effectively.

More recently, the field has witnessed a surge in the application of machine learning and deep learning techniques, aiming to capture complex, non-linear patterns that traditional econometric models might miss. Among these, Long Short-Term Memory (LSTM) networks, a type of Recurrent Neural Network (RNN), have gained significant traction due to their inherent ability to learn long-range dependencies in sequential data. Fischer and Krauss (2018) provided compelling evidence that LSTM networks could outperform standard econometric models and simpler machine learning methods (like Random Forests or Support Vector Machines) in directional forecasting for S&P 500 stocks, primarily by leveraging their capacity to

process raw price sequences effectively. The ability of LSTMs to maintain an internal memory state allows them to potentially capture intricate temporal dynamics that are challenging for GARCH models, which rely on predefined structures. Extensions incorporating Convolutional Neural Networks (CNNs) have also emerged, often in CNN-LSTM hybrid architectures, where CNNs first extract spatial or local features from input data (like price charts or sequences represented as matrices) before LSTMs model the temporal dependencies (Hoseinzade & Haratizadeh, 2019). Attention mechanisms and Transformer architectures, originally developed for natural language processing, are also increasingly being explored for financial time series, offering alternative ways to capture long-range dependencies and feature interactions, as demonstrated by models like the Informer (Zhou et al., 2021).

Recognizing the complementary strengths of econometric rigor and machine learning flexibility, a growing body of literature focuses on hybrid models. The goal is often to leverage the statistical underpinnings of GARCH-type models while enhancing predictive power with data-driven methods like LSTMs. Zolfaghari and Gholami (2021) demonstrated a sophisticated hybrid combining wavelet transforms (for signal decomposition), LSTM, ARIMA, and GARCH components, achieving superior Value-at-Risk (VaR) forecasts. Kakade et al. (2022) pursued an ensemble approach, combining forecasts from multiple GARCH specifications with those from LSTM and Bidirectional LSTM (BiLSTM) networks, yielding significant improvements in one-day-ahead volatility predictions. The work by Roszyk and Ślepaczuk (2024) specifically integrated the VIX index, a market-implied measure of expected volatility, as an exogenous input into a hybrid GARCH-LSTM framework for S&P 500 volatility, reporting substantial forecast gains (up to 20%) over standalone GARCH models. Taking a different integration approach, Xu et al. (2024) proposed the GARCH-Informed Neural Network (GINN), embedding GARCH-derived features (like conditional variance or residuals) directly within the LSTM architecture, thereby guiding the neural network's learning process and enhancing out-of-sample performance across diverse asset classes. These hybrid approaches vary in their integration strategy: some use GARCH outputs as inputs to LSTMs, others use LSTMs to model GARCH residuals, and some combine forecasts from separate models.

Despite these significant advancements, several challenges persist. Many complex hybrid architectures are prone to in-sample overfitting, especially given the noisy nature of financial data. The hyperparameter tuning process becomes considerably more complex when optimizing across disparate model families (e.g., GARCH orders vs. LSTM layers/neurons). Furthermore, while Roszyk and Ślepaczuk (2024) incorporated VIX, the systematic integration of other potentially valuable exogenous information—such as sentiment indices derived from news or social media (Tetlock, 2007), macroeconomic indicators, or detailed measures derived from limit order book data—remains relatively sparse in hybrid volatility forecasting frameworks. Common evaluation often focuses on standard statistical metrics like Root Mean Squared Error (RMSE) or Mean Absolute Error (MAE), but performance in risk management applications (e.g., VaR or Expected Shortfall accuracy) is also crucial, requiring appropriate backtesting procedures like those assessing conditional coverage (Christoffersen, 1998).

Therefore, a clear opportunity exists to develop robust, adaptive hybrid frameworks, potentially operating within a rolling-window estimation scheme. Such frameworks could systematically integrate diverse information sources—lagged returns, high-frequency based realized volatility measures, GARCH-implied features (like residuals or conditional variances), LSTM latent representations, and relevant exogenous variables (market-based, sentiment-based, or macroeconomic)—while employing rigorous procedures for model selection, hyperparameter optimization, and validation to ensure out-of-sample robustness and practical utility. Achieving better interpretability for these complex hybrid models also remains an important area for future research.

2 Methodology

2.1 Overview

This research will be implemented with Python instead of R or MATLAB, as it is more suitable for deep learning applications. The model will leverage PyTorch for the LSTM component (Paszke et al., 2019) and the arch package for GARCH estimation (J. L. Taylor, 2018).

The hybrid model will be trained on 2-years historical S&P 500 log returns and the VIX index data to predict a 1-day-ahead S&P 500 log return volatility. Then, the model will apply a rolling window approach to mimic real-time forecasting in the dynamics market. In addition, hyperparameters of the LSTM (number of layers, hidden-unit size, dropout rate) will be tuned via grid search.

Eventually, the model's performance will be evaluated using various metrics, including Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and Mean Square Error (MSE). These are the most common metrics used in measurement of forecasting accuracy. The hybrid model will be compared with the baseline GARCH model and the baseline LSTM model to assess its performance by applying above metrics to the predicted volatility and the realized volatility.

2.2 Data Preparation

2.2.1 Data Source

This research takes advantage of the S&P 500 index data and the VIX index data from the YFinance library in Python. The dataset consist of the daily closing prices of the S&P 500 index and the VIX index from January 1, 2000, to December 31, 2020. The final preprocessed dataset used for modeling contained 5277 observations after handling NaNs generated during feature calculation. The training set comprised 4221 observations, and the test set comprised 1056 observations.

2.2.2 Data Preprocessing and Analysis

The following preprocessing and analysis steps were applied:

1. **Get Log Returns:** From Equation (1), we transform the daily closing prices of the S&P 500 index into log returns. Log returns of S&P500 from January 1, 2000, to December 31, 2020 are shown in Figure 1.

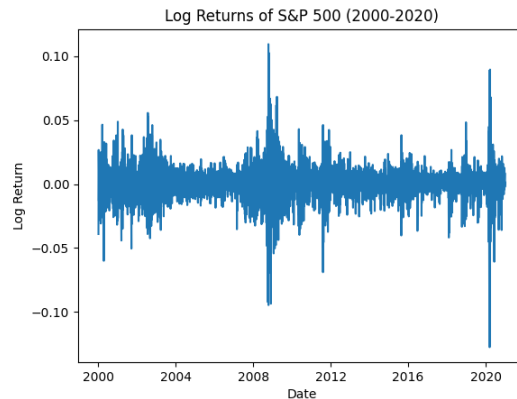


Figure 1: Log Returns of S&P 500 Index (2000-2020)

2. **Stationary Test:** Financial log returns are typically stationary, as supported by Cont (2001). In this research, we applied the Augmented Dickey-Fuller (ADF) test to check for mean stationarity of log returns (Dickey & Fuller, 1979). Consider the hypothesis statement as follows:

H_0 : The log-return series has a unit root (i.e., is non-stationary).

H_1 : The log-return series is stationary.

The ADF test results are shown in Table 1.

Table 1: Augmented Dickey–Fuller Test on S&P 500 Log Returns

Statistic	Value
ADF Statistic	−13.467645
p-value	0.000000
Critical Values	
1%	−3.432
5%	−2.862
10%	−2.567

Since the ADF statistic (−13.467645) is well below the 1% critical value (−3.432) and the p-value is effectively zero, we reject the null hypothesis of a unit root and conclude that the log return series is stationary.

3. **Autoregressive Moving-Average model (ARMA):** Since the log return series is stationary, we can apply the ARMA model to check for autocorrelation (Bosq & Nguyen, 1996). We fit an ARMA(1,1) to the log-return series r_t :

$$r_t = \phi_1 r_{t-1} + \varepsilon_t + \theta_1 \varepsilon_{t-1}, \quad \varepsilon_t \sim \mathcal{N}(0, \sigma^2). \quad (4)$$

Estimation by maximum likelihood yields (see Table 2):

Table 2: ARMA(1,1) Parameter Estimates

Parameter	Estimate	Std. Error	t-value	p-value
ϕ_1 (AR,L1)	−0.1140	0.0070	−17.144	0.0000
θ_1 (MA,L1)	−0.9999	0.0320	−30.826	0.0000
σ^2	2.00×10^{-4}	5.11×10^{-6}	30.376	0.0000

4. **Residual Diagnostics:** To check for remaining linear and volatility clustering effects, we apply:

- Ljung–Box test on $\{\hat{\varepsilon}_t\}$ at lag 10 (LJUNG & BOX, 1978).
- ARCH–LM test on $\{\hat{\varepsilon}_t\}$ (Engle, 1982).

Table 3: ARMA Residual Diagnostic Tests

Test	Statistic	p-value
Ljung-Box $Q(10)$	27.473	0.0022
ARCH-LM (lag = 1)	1506.008	< 0.0001

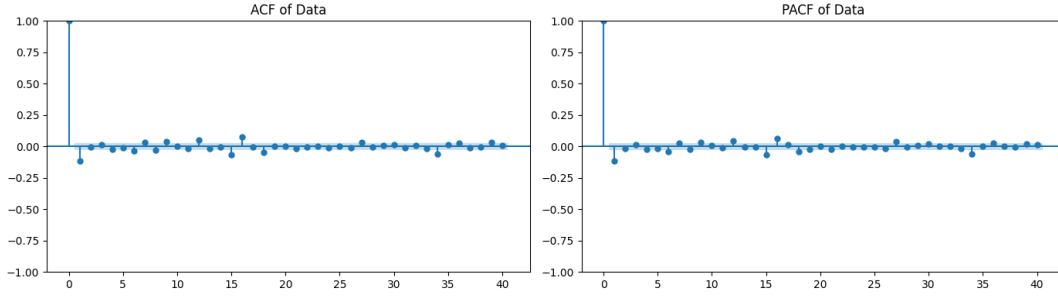


Figure 2: ACF (left) and PACF (right) of the S&P 500 log-return series.

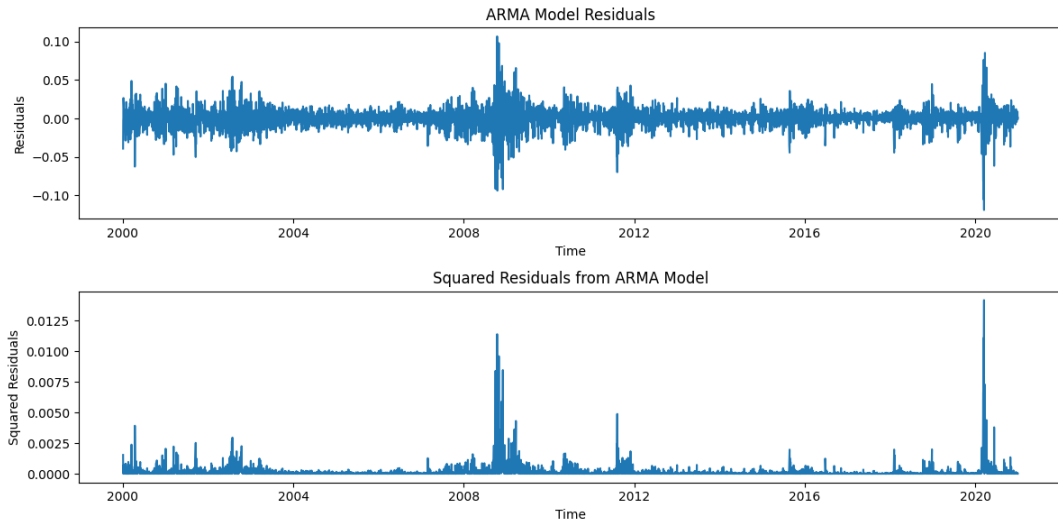


Figure 3: ARMA(1,1) residuals (top) and squared residuals (bottom).

- Both the AR(1) coefficient $\hat{\phi}_1 = -0.1140$ and MA(1) coefficient $\hat{\theta}_1 = -0.9999$ are highly significant (p-values < 10^{-4}), indicating a clear short-memory structure in returns.
- The innovation variance $\hat{\sigma}^2 = 2.00 \times 10^{-4}$ is small, consistent with typical daily return volatility around 1–2%.
- However, the Ljung–Box test at lag 10 yields $Q = 27.47$ with $p = 0.0022$, and the ARCH–LM test is highly significant ($p < 10^{-4}$), implying remaining volatility clustering in the ARMA residuals.
- These diagnostics motivate a GARCH specification to capture the time-varying conditional variance.

5. **GARCH Model:** We fit a GARCH(1,1) model to the ARMA residuals $\{\hat{\epsilon}_t\}$ to capture the time-varying conditional variance. The results are shown in Table 4, Figure 4 and Figure 5.

Table 4: GARCH(1,1) Parameter Estimates

Parameter	Estimate	Std. Error	t-value
ω	3.1519×10^{-6}	3.819×10^{-11}	82,530.0
α_1	0.1000	0.000126	794.006
β_1	0.8800	0.002650	332.039
$\alpha_1 + \beta_1$	0.9800		

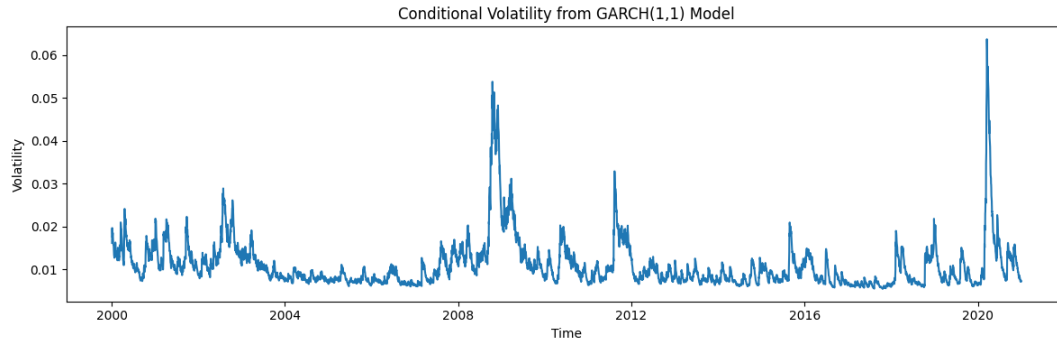


Figure 4: Conditional volatility σ_t estimated by the GARCH(1,1) model.

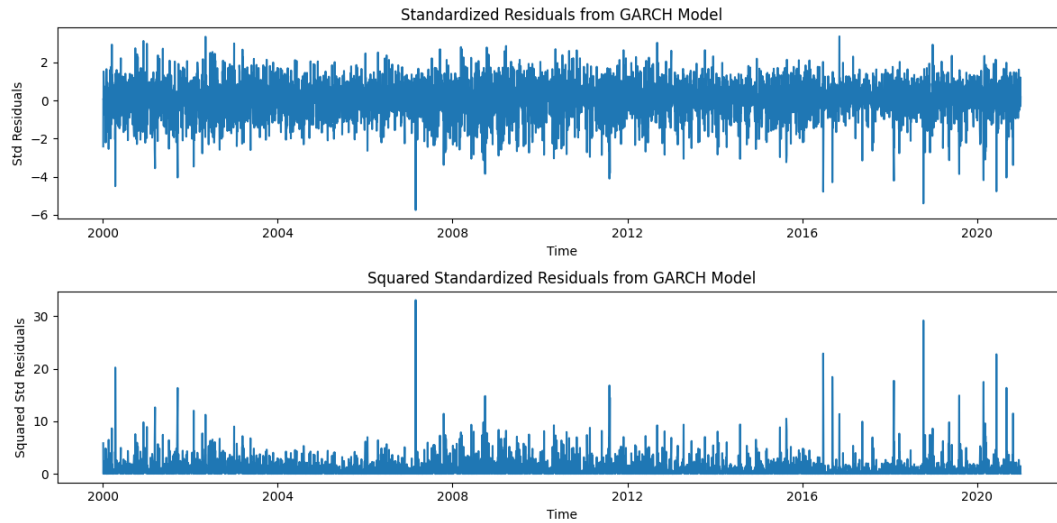


Figure 5: Standardized residuals (top) and their squares (bottom) from the GARCH(1,1) fit.

- The intercept $\hat{\omega} = 3.1519 \times 10^{-6}$ and the coefficients $\hat{\alpha}_1 = 0.1000$, $\hat{\beta}_1 = 0.8800$ are highly significant (t-values > 300 , $p < 10^{-4}$).
- The persistence $\hat{\alpha}_1 + \hat{\beta}_1 = 0.9800$ is close to unity, indicating strong volatility clustering and long-memory effects in the conditional variance.
- The small ω implies a low long-run average variance, but shocks to volatility decay slowly due to the high β_1 .

- Plots of σ_t (Fig. 4) show clear volatility spikes around market turmoil, and the standardized residuals in Fig. 5 exhibit no obvious serial correlation or ARCH effects, confirming model adequacy.
6. **Realized Volatility:** In this research, the daily realized volatility $\sigma_{RV,t}$ at time t was calculated by Equation (2). A 5-day rolling window is chosen because it balances short-term fluctuations with sufficient data, following McAleer and Medeiros (2008).
 7. **Target Variable:** The target variable for prediction is the next day’s realized volatility, $\sigma_{RV,t+1}$. This is obtained by shifting the realized volatility series by one day.
 8. **Feature Scaling:** For the LSTM model, input features (log returns and VIX levels) were scaled to the range $[-1, 1]$ using Min-Max scaling fitted on the training data.
 9. **Data Splitting:** The dataset (after calculating features and dropping initial NaNs) was split into a training set (80% of observations) and a test set (20% of observations). The test period starts on October 19, 2016.

2.3 Model Architecture

2.3.1 LSTM Component

An LSTM network was designed as the baseline to predict the next-day realized volatility ($\hat{\sigma}_{LSTM,t+1}$). The proposed architecture consisted of:

- Input Layer: Takes sequences of length 120 days, with 2 features per day (scaled log return, scaled VIX level).
- LSTM Layers: A single-layer LSTM layers with 24 hidden units each.
- Output Layer: A final linear layer mapped the LSTM output to predict the next-day realized volatility ($\hat{\sigma}_{LSTM,t+1}$).

The model was trained on the training sequences using the Adam optimizer with a learning rate of 0.001 and Mean Squared Error (MSE) as the loss function for 30 epochs with a batch size of 32. We performed a grid search over sequence lengths (20, 40, 60, 120 days), hidden sizes (24, 36, 48, 60), number of LSTM layers (1-3), number of training epochs (20, 30, 40), and batch sizes (8, 16). We compared the best-tuned model with the proposed parameters. Given the increased complexity—doubling the hidden units and extending training—against such a small gain, we retained the simpler 120-day, 24-unit single-layer LSTM as our baseline. This configuration offers nearly identical predictive accuracy with lower computational cost and easier interpretability.

Table 5: Proposed vs. Best-Tuned LSTM Hyperparameters and Performance

Model	Seq. Length	Hidden Units	Layers	Epochs	Batch Size	RMSE	MAE
Proposed	120	24	1	30	32	0.0645	0.0397
Best-Tuned	60	48	1	40	16	0.0491	0.0297

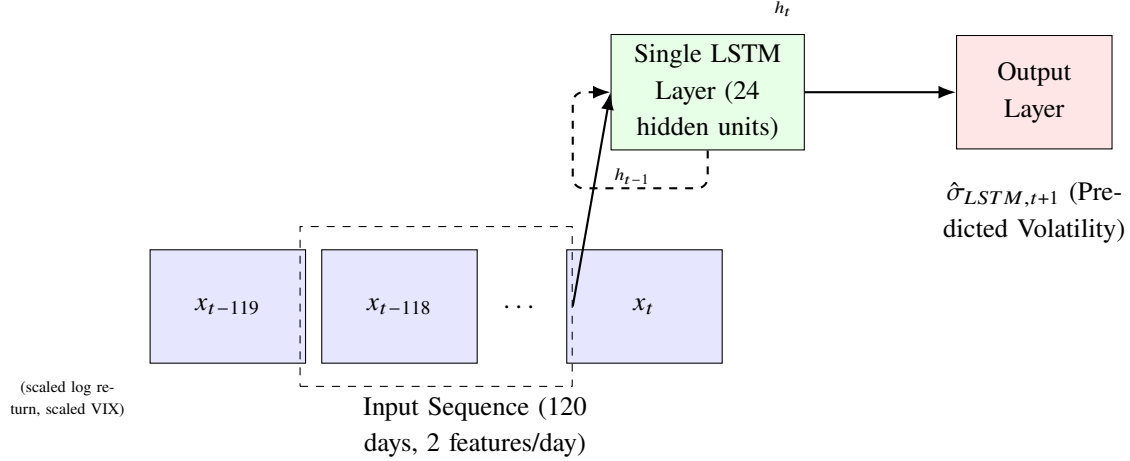


Figure 6: LSTM Network Architecture for Volatility Prediction.

2.3.2 GARCH Component

As described in Section 2.2.2, we employed standard GARCH(1,1) model as the baseline following Equation (3). The parameters ω , α , and β were estimated using the training data via Maximum Likelihood. One-step-ahead variance forecasts h_{t+1} were generated during the rolling forecast period, and the corresponding volatility forecast was calculated as $\hat{\sigma}_{GARCH,t+1} = \sqrt{h_{t+1}} \times \sqrt{252}$.

2.3.3 Combined LSTM-GARCH Model

A simple ensemble forecast was created by taking the arithmetic mean of the individual forecasts from the GARCH and LSTM models for each day $t + 1$:

$$\hat{\sigma}_{Combined,t+1} = \frac{\hat{\sigma}_{GARCH,t+1} + \hat{\sigma}_{LSTM,t+1}}{2} \quad (5)$$

The combined model was implemented on the test set (1056 observations, generating 1055 1-day ahead forecasts). The combined forecast is proposed to leverage the strengths of both models, where the GARCH model captures the volatility clustering effect and the LSTM model captures the nonlinear patterns in the data. The combined forecast is expected to outperform both individual models, as it can adapt to different market conditions and improve overall prediction accuracy. To keep the model simple, we assume equal weights for the GARCH and LSTM forecasts. In future work, we could explore more sophisticated weighting schemes based on model performance or market conditions.

2.3.4 Rolling Window Forecasting

For each day t in the test set (starting from 2016-10-19):

1. The LSTM model used the previous 120 days of scaled log returns and VIX levels to predict $\hat{\sigma}_{LSTM,t+1}$.
2. The GARCH model used the parameters estimated on the initial training set and the log return at time t (ϵ_t) along with the variance at time t (h_t) to forecast the variance h_{t+1} and subsequently $\hat{\sigma}_{GARCH,t+1}$. The GARCH state (last variance h_t and last residual ϵ_t) was updated iteratively.
3. The combined forecast $\hat{\sigma}_{Combined,t+1}$ was calculated.
4. The forecasts were stored along with the actual realized volatility $\sigma_{RV,t+1}$.

Periodic retraining of the models during the rolling forecast window was not implemented in this research. By fixing parameters in the rolling forecast, we prioritize simplicity of the model while we may sacrificing potential dynamic shifts. In future work, one could explore adaptive retraining strategies to enhance forecast accuracy.

2.4 Evaluation Metrics

To evaluate model performance, the predicted volatility $\hat{\sigma}_{Combined,t+1}$, $\hat{\sigma}_{LSTM,t+1}$, $\hat{\sigma}_{GARCH,t+1}$ is compared against the actual realized volatility $\sigma_{RV,t+1}$, computed via the 5-day rolling window defined earlier. The model's predictive accuracy is assessed using three standard error metrics over the out-of-sample period $T_{\text{train}} + 1$ to T_{total} :

- **Mean Absolute Error (MAE):**

$$\text{MAE} = \frac{1}{T_{\text{test}}} \sum_{t=T_{\text{train}}+1}^{T_{\text{total}}} |\hat{\sigma}_t - \sigma_t|$$

- **Mean Squared Error (MSE):**

$$\text{MSE} = \frac{1}{T_{\text{test}}} \sum_{t=T_{\text{train}}+1}^{T_{\text{total}}} (\hat{\sigma}_t - \sigma_t)^2$$

- **Root Mean Squared Error (RMSE):**

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

where $T_{\text{test}} = T_{\text{total}} - T_{\text{train}}$ is the number of observations in the test set.

3 Results

3.1 Initial Model Training

The LSTM model was trained for 30 epochs. The average loss decreased significantly over the epochs, starting at approximately 0.0139 in the first epoch and ending at approximately 0.0035 in the final epoch, indicating convergence. The training loss curve is shown in Figure 7.

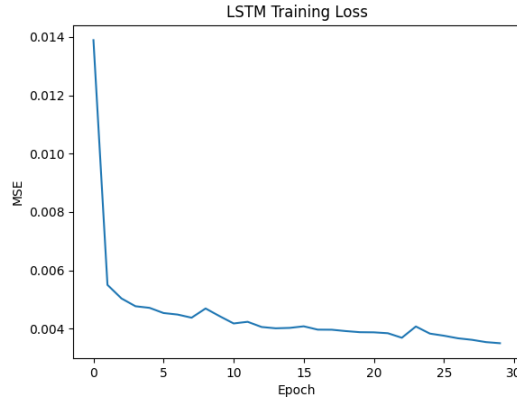


Figure 7: LSTM Training Loss Curve

The initial GARCH(1,1) model was fitted to the training set log returns (scaled by 100). The estimated parameters are presented in Table 6.

Table 6: Initial GARCH(1,1) Model Parameter Estimates (Training Set)

Parameter	Coefficient	Std. Error	t-statistic	P-value
ω	0.0209	5.245e-03	3.991	6.576e-05
α_1	0.0998	1.240e-02	8.048	8.382e-16
β_1	0.8842	1.308e-02	67.607	0.000

Note: Parameters estimated using the `arch` package on log returns multiplied by 100. The model specification included a zero mean and assumed normally distributed errors. $N = 4221$. Log-Likelihood = -5988.12, AIC = 11982.2, BIC = 12001.3.

All GARCH parameters (ω , α_1 , β_1) were statistically significant at conventional levels ($P > |t|$ close to zero). The positive ω confirms the model assumes a strictly positive variance. The sum $\alpha_1 + \beta_1 \approx 0.0998 + 0.8842 = 0.984$, which is close to 1, indicating high persistence in volatility, a common finding in financial time series.

3.2 Rolling Forecast Performance

The rolling forecast procedure generated 1055 predictions for the test period starting October 19, 2016. The performance metrics for the baseline GARCH model, the LSTM model, and the combined LSTM-GARCH model are summarized in Table 7 and Figure 8.

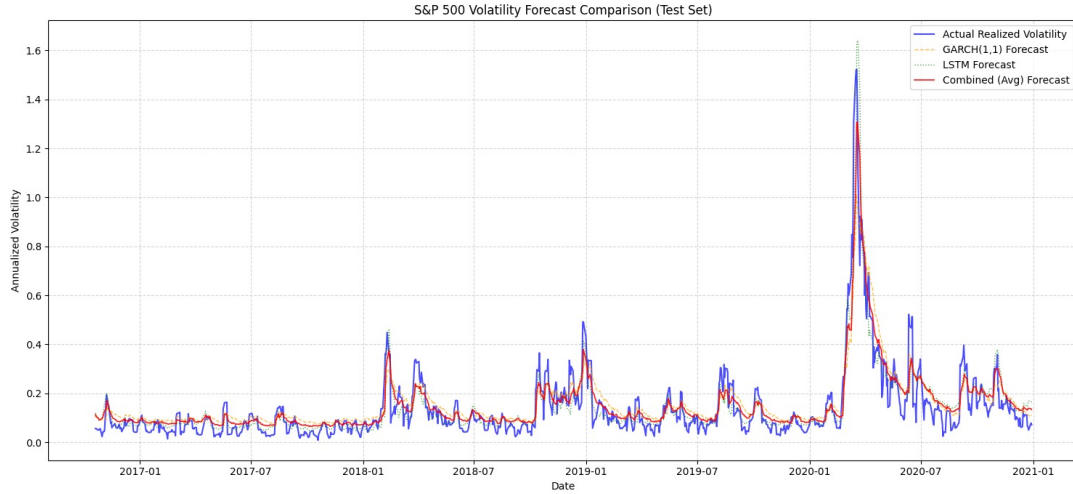


Figure 8: Predicted vs. Actual Realized Volatility (Test Set)

Table 7: Volatility Forecast Performance Metrics on Test Set (2016-10-19 to 2020-12-30)

Model	MAE	MSE	RMSE
GARCH (Baseline)	0.050324	0.005366	0.073255
LSTM	0.039733	0.004155	0.064460
Combined (Avg)	0.041363	0.003729	0.061067

Note: Evaluation metrics calculated on 1055 out-of-sample 1-day ahead volatility forecasts. Lower values indicate better forecast accuracy.

Based on all three metrics (MAE, MSE, RMSE), the LSTM model outperformed the baseline GARCH model. The Combined model, averaging the LSTM and GARCH forecasts, achieved the lowest MSE and RMSE, although its MAE was slightly higher than the standalone LSTM model.

3.3 Objective Check

The primary objective was to determine if the combined model offered a significant improvement over the baseline GARCH model and the baseline LSTM model, defined as an RMSE reduction of at least 5%.

- Baseline GARCH RMSE: 0.073255
- Baseline LSTM RMSE: 0.064460
- Combined LSTM-GARCH RMSE: 0.061067
- RMSE Reduction compared to standalone GARCH: $\frac{0.073255 - 0.061067}{0.073255} \times 100\% \approx 16.64\%$
- RMSE Reduction compared to standalone LSTM: $\frac{0.064460 - 0.061067}{0.064460} \times 100\% \approx 5.26\%$

The achieved RMSE reduction of 16.64% and 5.26%, compared to standalone GARCH and standalone LSTM respectively, exceeded the target threshold of 5%. Therefore, the objective was successfully met. This improvement enhances risk management by reducing forecast uncertainty.

4 Conclusion

This research set out to enhance 1-day-ahead volatility forecasting for the S&P 500 index by developing a hybrid LSTM-GARCH model. The core objectives were to design and implement this hybrid model, aiming for at least a 5% reduction in forecast error compared to a baseline GARCH model, and to deliver a more robust forecasting framework by combining the strengths of LSTMs in recognizing nonlinear patterns and GARCH models' statistical rigor, while also incorporating the VIX index as an exogenous variable.

The study successfully demonstrated the superiority of the hybrid approach. The combined LSTM-GARCH model, which averaged the forecasts from a standalone LSTM model and a GARCH(1,1) model, achieved a significant improvement in predictive accuracy. Specifically, the combined model yielded a Root Mean Squared Error (RMSE) of 0.061067. This represented a 16.64% reduction in RMSE compared to the baseline GARCH model (RMSE of 0.073255) and a 5.26% reduction compared to the standalone LSTM model (RMSE of 0.064460). These results surpassed the predefined objective of a 5% RMSE reduction relative to the baseline GARCH model, thereby meeting the primary goal of the research.

5 Future Work

While this study achieved its objectives, several avenues for future research could build upon these findings and address the limitations of the current work.

1. **Sophisticated Ensemble Weighting:** The current combined model uses a simple arithmetic mean of the GARCH and LSTM forecasts. Future research could explore more sophisticated and dynamic weighting schemes. These could involve weights that adapt based on recent model performance or prevailing market conditions, potentially leading to further improvements in forecast accuracy.
2. **Expanding Ensemble Models:** Future work could investigate the integration of additional models into the ensemble framework. For instance, incorporating other machine learning models (e.g., Random Forest, XGBoost) or advanced deep learning architectures (e.g., CNNs, Transformers) could enhance the robustness and accuracy of the hybrid model.
3. **Regime-Switching Models:** Given the dynamic nature of financial markets, future research could explore regime-switching models that adapt to different market conditions. This could involve using hidden Markov models or other techniques to identify and adapt to changing volatility regimes, potentially improving forecast accuracy during periods of market stress or tranquility.

References

- Andersen, T. G., Bollerslev, T., Diebold, F. X., & Labys, P. (2003). Modeling and forecasting realized volatility. *Econometrica*, 71(2), 579–625. <https://doi.org/10.1111/1468-0262.00418>
- Baillie, R. T., Bollerslev, T., & Mikkelsen, H. O. (1996). Fractionally integrated generalized autoregressive conditional heteroskedasticity. *Journal of Econometrics*, 74(1), 3–30. [https://doi.org/10.1016/S0304-4076\(95\)01749-6](https://doi.org/10.1016/S0304-4076(95)01749-6)
- Bollerslev, T. (1986). Generalized autoregressive conditional heteroskedasticity. *Journal of Econometrics*, 31(3), 307–327. [https://doi.org/10.1016/0304-4076\(86\)90063-1](https://doi.org/10.1016/0304-4076(86)90063-1)
- Bollerslev, T. (1987). A conditionally heteroskedastic time series model for speculative prices and rates of return. *The Review of Economics and Statistics*, 69(3), 542–547. <https://doi.org/10.2307/1925546>
- Bollerslev, T., Todorov, V., & Li, Y. (2008). High-frequency jumps in financial markets: Characteristics and dynamic impacts. *Journal of Financial Econometrics*, 6(3), 233–260. <https://doi.org/10.1093/jjfinec/nbn010>
- Bosq, D., & Nguyen, H. T. (1996). Arma model. In *A course in stochastic processes: Stochastic models and statistical inference* (pp. 205–217). Springer Netherlands. https://doi.org/10.1007/978-94-015-8769-3_10
- Box, G. E. P., Jenkins, G. M., & Reinsel, G. C. (1970). *Time series analysis: Forecasting and control*. Holden-Day.
- Campbell, J., Lo, A., MacKinlay, A., & Whitelaw, R. (1998). The econometrics of financial market. *Macroeconomic Dynamics*, 2, 559–562. <https://doi.org/10.1017/S1365100598009092>
- Christoffersen, P. F. (1998). Evaluating interval forecasts. *International Economic Review*, 39(4), 841–862. <https://doi.org/10.2307/2527341>
- Cont, R. (2001). Empirical properties of asset returns: Stylized facts and statistical issues. *Quantitative Finance*, 1(2), 223–236. <https://doi.org/10.1080/713665670>
- Corsi, F. (2009). A simple approximate long-memory model of realized volatility. *Journal of Financial Econometrics*, 7(2), 174–196. <https://doi.org/10.1093/jjfinec/nbp001>
- Dickey, D., & Fuller, W. (1979). Distribution of the estimators for autoregressive time series with a unit root. *JASA. Journal of the American Statistical Association*, 74. <https://doi.org/10.2307/2286348>
- Ding, Z., Granger, C. W. J., & Engle, R. F. (1993). A long memory property of stock market returns and a new model. *Journal of Empirical Finance*, 1(1), 83–106. [https://doi.org/10.1016/0927-5398\(93\)90006-D](https://doi.org/10.1016/0927-5398(93)90006-D)
- Engle, R. F. (1982). Autoregressive conditional heteroskedasticity with estimates of the variance of the united kingdom inflation. *Econometrica*, 50(4), 987–1007. <https://doi.org/10.2307/1912773>
- Engle, R. F., Lilien, D. M., & Robins, R. P. (1987). Estimating time varying risk premia in the term structure: The arch-m model. *Econometrica*, 55(2), 391–407. <https://doi.org/10.2307/1913242>
- Fischer, T., & Krauss, C. (2018). Deep learning with long short-term memory networks for financial market predictions. *European Journal of Operational Research*, 270, 654–669. <https://doi.org/10.1016/j.ejor.2017.11.054>
- Giot, P., & Laurent, S. (2005). Market risk in commodity markets: Value-at-risk measurement and extreme events modelling. *Applied Financial Economics*, 15(8), 597–610. <https://doi.org/10.1080/09603100500153030>
- Glosten, L. R., Jagannathan, R., & Runkle, D. E. (1993). On the relation between the expected value and the volatility of the nominal excess return on stocks. *Journal of Finance*, 48(5), 1779–1801.
- Hamilton, J. (1994). *Time series analysis*. Princeton University Press. https://books.google.co.jp/books?id=B8_1UBmqVUoC
- Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735–1780. <https://doi.org/10.1162/neco.1922.214.171.1245>
- Hoseinzade, E., & Haratizadeh, S. (2019). Cnnpred: Cnn-based stock market prediction using a diverse set of variables. *Expert Systems with Applications*, 129, 273–285. <https://doi.org/10.1016/j.eswa.2019.03.029>
- Jorion, P. (2007). *Value at risk: The new benchmark for managing financial risk* (3rd, Vol. 3). McGraw-Hill.
- Kakade, A., et al. (2022). Ensemble learning garch-lstm models for var forecasting. *Quantitative Finance*.

- LJUNG, G. M., & BOX, G. E. P. (1978). On a measure of lack of fit in time series models. *Biometrika*, 65(2), 297–303. <https://doi.org/10.1093/biomet/65.2.297>
- McAleer, M., & Medeiros, M. C. (2008). Realized volatility: A review. *Econometric Reviews*, 27(1-3), 10–45. <https://doi.org/10.1080/07474930701853509>
- Nelson, D. B. (1991). Conditional heteroskedasticity in asset returns: A new approach. *Econometrica*, 59(2), 347–370.
- Paszke, A., et al. (2019). Pytorch: An imperative style, high-performance deep learning library. 32, 8024–8035.
- Roszyk, N., & Ślepaczuk, R. (2024). The hybrid forecast of s&p 500 volatility ensembled from vix, garch and lstm models.
- Taylor, J. L. (2018). *arch: Arch models in Python*. arXiv preprint arXiv:1811.01825.
- Taylor, S. J. (1986). *Modelling financial time series*. John Wiley & Sons.
- Tetlock, P. C. (2007). Giving content to investor sentiment: The role of media in the stock market. *The Journal of Finance*, 62(3), 1139–1168. <https://doi.org/10.1111/j.1540-6261.2007.01232.x>
- Tsay, R. S. (2010). *Analysis of financial time series*. Wiley.
- Whaley, R. E. (2000). The investor fear gauge. *The Journal of Portfolio Management*, 26(3), 12–17. <https://doi.org/10.3905/jpm.2000.319754>
- Xu, Z., Liechty, J., Benthall, S., Skar-Gislinge, N., & McComb, C. (2024). Garch-informed neural networks for volatility prediction in financial markets.
- Zhou, H., Zhang, S., Peng, J., Zhang, S., Li, J., Xiong, H., & Zhang, W. (2021). Informer: Beyond efficient transformer for long sequence time-series forecasting. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(12), 11106–11115. <https://ojs.aaai.org/index.php/AAAI/article/view/17325>
- Zolfaghari, A., & Gholami, Z. (2021). Value-at-risk forecasting: A hybrid ensemble learning garch-lstm based approach. *Journal of Risk*.

Appendices

For more details, please visit the GitHub repository at <https://github.com/allenchan308/Volatility-Prediction-LSTM-GARCH>

Appendix A: Python Code for ADF Test

```
1 import yfinance as yf
2 import numpy as np
3 import pandas as pd
4 from statsmodels.tsa.stattools import adfuller
5 import matplotlib.pyplot as plt
6 from statsmodels.tsa.arima.model import ARIMA
7 from statsmodels.stats.diagnostic import acorr_ljungbox, het_arch
8 from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
9 from arch import arch_model
10
11 # Download S&P 500 historical data (ticker: ^GSPC)
12 data = yf.download('^GSPC', start='2000-01-01', end='2020-12-31')
13
14 # Drop rows with missing values
15 data = data.dropna()
16
17 # Calculate log returns from 'Adj Close' prices
18 data['LogReturn'] = np.log(data['Close'] / data['Close'].shift(1))
19
20 # Drop NaNs resulting from the shift
21 log_returns = data['LogReturn'].dropna()
22
23 import matplotlib.pyplot as plt
24
25 plt.plot(log_returns)
26 plt.title('Log Returns of S&P 500 (2000-2020)') # Add title
27 plt.xlabel('Date') # Add x-axis label
28 plt.ylabel('Log Return') # Add y-axis label
29 plt.savefig('log_returns_plot.png')
30
31 adf_result = adfuller(log_returns)
32
33 print('ADF Statistic: %f' % adf_result[0])
34 print('p-value: %f' % adf_result[1])
35 print('Critical Values:')
36 for key, value in adf_result[4].items():
37     print('\t%s: %.3f' % (key, value))
38
39 if adf_result[1] <= 0.05:
40     print("Reject the null hypothesis: Time series is stationary")
41 else:
42     print("Fail to reject the null hypothesis: Time series is non-stationary")
```

Appendix B: Python Code for Checking Volatility Clustering

```
1 # --- 1. ARMA Model Fitting ---
2 # Choose ARMA(1,0,1) as example; adjust (p,d,q) as needed
3 arma_order = (1, 0, 1)
4 arma_model = ARIMA(log_returns, order=arma_order).fit()
5 print("\nARMA Model Summary:")
6 print(arma_model.summary())
7
8 # Extract ARMA parameters for table
9 arma_params = arma_model.params
10 arma_bse = arma_model.bse
11 arma_tvalues = arma_model.tvalues
12 arma_pvalues = arma_model.pvalues
13
14 arma_table = pd.DataFrame({
15     "Parameter": arma_params.index,
16     "Estimate": arma_params.values,
17     "Std. Error": arma_bse.values,
18     "t-value": arma_tvalues.values,
19     "p-value": arma_pvalues.values
20 })
21
22 # --- 2. Residual Diagnostics ---
23 # Ljung-Box test for residual autocorrelation
24 ljung_box = acorr_ljungbox(arma_model.resid, lags=[10], return_df=True)
25 print("\nLjung-Box test for residual autocorrelation:")
26 print(ljung_box)
27
28 # ARCH-LM test for ARCH effects
29 arch_test_stat, arch_test_pvalue, _, _ = het_arch(arma_model.resid)
30 print(f"\nARCH-LM Test Statistic: {arch_test_stat:.4f}, p-value: {arch_test_pvalue:.4f}")
31
32 arch_table = pd.DataFrame({
33     "Test Statistic": [arch_test_stat],
34     "p-value": [arch_test_pvalue]
35 })
36
37 # --- 3. GARCH(1,1) Model Fitting ---
38 # Fit GARCH(1,1) with zero mean (assuming ARMA captured mean)
39 garch_model = arch_model(log_returns, vol='Garch', p=1, q=1, mean='Zero', dist='normal')
40 garch_fit = garch_model.fit(update_freq=5, disp='off')
41 print("\nGARCH(1,1) Model Summary:")
42 print(garch_fit.summary())
43
44 # Extract GARCH parameters for table
45 garch_params = garch_fit.params
46 garch_bse = garch_fit.std_err
47 garch_tvalues = garch_fit.tvalues
48 garch_pvalues = garch_fit.pvalues
```



```

49
50 garch_table = pd.DataFrame({
51     "Parameter": garch_params.index,
52     "Estimate": garch_params.values,
53     "Std. Error": garch_bse.values,
54     "t-value": garch_tvalues.values,
55     "p-value": garch_pvalues.values
56 })
57
58 # --- 4. Plotting ---
59
60 # a) Time Series Plot of Data
61 plt.figure(figsize=(12, 4))
62 plt.plot(log_returns)
63 plt.title('Time Series Plot of Stationary Data')
64 plt.xlabel('Time')
65 plt.ylabel('Value')
66 plt.tight_layout()
67 plt.show()
68
69 # b) ACF and PACF plots
70 fig, ax = plt.subplots(1, 2, figsize=(14, 4))
71 plot_acf(log_returns, lags=40, ax=ax[0])
72 ax[0].set_title('ACF of Data')
73 plot_pacf(log_returns, lags=40, ax=ax[1])
74 ax[1].set_title('PACF of Data')
75 plt.tight_layout()
76 plt.savefig('acf_pacf_plot.png')
77
78 # c) Residual and Squared Residual Plots from ARMA
79 fig, ax = plt.subplots(2, 1, figsize=(12, 6))
80 ax[0].plot(arma_model.resid)
81 ax[0].set_title('ARMA Model Residuals')
82 ax[0].set_xlabel('Time')
83 ax[0].set_ylabel('Residuals')
84
85 ax[1].plot(arma_model.resid**2)
86 ax[1].set_title('Squared Residuals from ARMA Model')
87 ax[1].set_xlabel('Time')
88 ax[1].set_ylabel('Squared Residuals')
89
90 plt.tight_layout()
91 plt.savefig('arma_resid_sqresid.png')
92
93 # d) Conditional Volatility Plot from GARCH Model
94 plt.figure(figsize=(12, 4))
95 plt.plot(garch_fit.conditional_volatility)
96 plt.title('Conditional Volatility from GARCH(1,1) Model')
97 plt.xlabel('Time')

```

```

98 plt.ylabel('Volatility')
99 plt.tight_layout()
100 plt.savefig('garch_cond_vol.png')
101
102 # e) Standardized Residuals and Squared Standardized Residuals from GARCH
103 std_resid = garch_fit.std_resid
104
105 fig, ax = plt.subplots(2, 1, figsize=(12, 6))
106 ax[0].plot(std_resid)
107 ax[0].set_title('Standardized Residuals from GARCH Model')
108 ax[0].set_xlabel('Time')
109 ax[0].set_ylabel('Std Residuals')
110
111 ax[1].plot(std_resid**2)
112 ax[1].set_title('Squared Standardized Residuals from GARCH Model')
113 ax[1].set_xlabel('Time')
114 ax[1].set_ylabel('Squared Std Residuals')
115
116 plt.tight_layout()
117 plt.savefig('garch_std_resid.png')

```

Appendix C: Python Code for LSTM Model

```
1 !pip install arch
2 !pip install yfinance --upgrade --no-cache-dir
3 import yfinance as yf
4 import pandas as pd
5 import numpy as np
6 import torch
7 import torch.nn as nn
8 from torch.utils.data import DataLoader, TensorDataset
9 from arch import arch_model
10 from sklearn.preprocessing import MinMaxScaler
11 from sklearn.metrics import mean_absolute_error, mean_squared_error
12 import matplotlib.pyplot as plt
13 from tqdm import tqdm # For progress bars
14 import math # For sqrt
15
16 # --- Parameters ---
17 # Data Parameters
18 TICKER_SPX = '^GSPC'
19 TICKER_VIX = '^VIX'
20 START_DATE = '2000-01-01'
21 END_DATE = '2020-12-31'
22
23 # Preprocessing Parameters
24 REALIZED_VOL_WINDOW = 5
25 ANNUALIZATION_FACTOR = np.sqrt(252)
26
27 # Model Parameters
28 LSTM_INPUT_SEQ_LEN = 120
29 LSTM_HIDDEN_SIZE = 24
30 LSTM_NUM_LAYERS = 1
31 LSTM_DROPOUT = 0.3
32 GARCH_P = 1
33 GARCH_Q = 1
34
35 # Training & Rolling Window Parameters
36 TRAIN_TEST_SPLIT_RATIO = 0.8
37 LSTM_EPOCHS = 30
38 LSTM_BATCH_SIZE = 32
39 LSTM_LEARNING_RATE = 0.001
40 PERIODIC_RETRAIN = False # Keep False unless implementing re-training logic
41 RETRAIN_INTERVAL = 252
42 history = {'loss': []}
43
44 # --- Data Fetching and Preprocessing ---
45 def fetch_data(ticker, start, end):
46     """Fetches historical adjusted closing prices from Yahoo Finance."""
47     try:
48         data = yf.download(ticker, start=start, end=end, progress=False)
```

```

49     if data is None or data.empty:
50         raise ValueError(f"No data fetched for {ticker}")
51     print(f"Fetched {len(data)} rows for {ticker}")
52     return data['Close']
53 except Exception as e:
54     print(f"Error fetching data for {ticker}: {type(e).__name__}, {e}")
55     return None
56
57 spx_price = fetch_data(TICKER_SPX, START_DATE, END_DATE)
58 vix_price = fetch_data(TICKER_VIX, START_DATE, END_DATE)
59
60 if spx_price is None or vix_price is None:
61     raise SystemExit("Failed to fetch necessary data. Exiting.")
62
63 df = pd.concat([spx_price, vix_price], axis=1)
64 df.columns = ['SPX', 'VIX']
65 df.ffill(inplace=True)
66 df.dropna(inplace=True)
67 print(f"Combined data shape after initial NaN handling: {df.shape}")
68
69 # Calculate Features and Targets
70 df['SPX_LogRet'] = np.log(df['SPX'] / df['SPX'].shift(1))
71 df['VIX_Level'] = df['VIX'] # Feature for LSTM
72
73 # Calculate Realized Volatility (Target for LSTM and Evaluation)
74 df['Realized_Vol'] = df['SPX_LogRet'].rolling(window=REALIZED_VOL_WINDOW).std() *
    ANNUALIZATION_FACTOR
75 df['Realized_Vol_Target'] = df['Realized_Vol'].shift(-1) # Predict vol for t+1 using info
    up to t
76
77 initial_rows = len(df)
78 df.dropna(inplace=True)
79 print(f"Data shape after calculating returns/vol and dropping NaNs: {df.shape} (dropped {
    initial_rows - len(df)} rows)")
80
81 print("Data preprocessed. Sample:")
82 print(df.head())
83
84 # --- Train/Test Split ---
85 n_obs = len(df)
86 n_train = int(n_obs * TRAIN_TEST_SPLIT_RATIO)
87 train_df = df.iloc[:n_train].copy()
88 test_df = df.iloc[n_train:].copy()
89 test_indices = df.index[n_train:]
90
91 print(f"\nData Split:")
92 print(f"Training set size: {len(train_df)}")
93 print(f"Test set size: {len(test_df)}")
94

```

```

95 # --- Feature Scaling ---
96 # Scale features used by LSTM: Log Returns and VIX Level
97 scaler_logret = MinMaxScaler(feature_range=(-1, 1))
98 scaler_vix = MinMaxScaler(feature_range=(-1, 1))
99
100 # Fit scalers ONLY on training data
101 train_df['SPX_LogRet_Scaled'] = scaler_logret.fit_transform(train_df[['SPX_LogRet']])
102 train_df['VIX_Level_Scaled'] = scaler_vix.fit_transform(train_df[['VIX_Level']])
103
104 # Transform test data using the FITTED scalers
105 test_df['SPX_LogRet_Scaled'] = scaler_logret.transform(test_df[['SPX_LogRet']])
106 test_df['VIX_Level_Scaled'] = scaler_vix.transform(test_df[['VIX_Level']])
107
108 scaled_df = pd.concat([train_df, test_df])
109
110 # --- Model Definitions ---
111 # LSTM Model Definition (Predicts Volatility)
112 class LSTMVolModel(nn.Module):
113     def __init__(self, input_size, hidden_size, num_layers, output_size, dropout):
114         super(LSTMVolModel, self).__init__()
115         self.hidden_size = hidden_size
116         self.num_layers = num_layers
117         self.lstm = nn.LSTM(input_size, hidden_size, num_layers,
118                             batch_first=True, dropout=dropout if num_layers > 1 else 0)
119         self.fc = nn.Linear(hidden_size, output_size)
120
121     def forward(self, x):
122         h0 = torch.zeros(self.num_layers, x.size(0), self.hidden_size).to(x.device)
123         c0 = torch.zeros(self.num_layers, x.size(0), self.hidden_size).to(x.device)
124         out, _ = self.lstm(x, (h0, c0))
125         out = self.fc(out[:, -1, :]) # Use output of last time step
126         return out
127
128 # --- Data Preparation Functions ---
129 def create_lstm_sequences(data_logret, data_vix, target_data, seq_len):
130     """Prepares sequences for LSTM training/prediction."""
131     xs, ys = [], []
132     data_logret_np = np.array(data_logret)
133     data_vix_np = np.array(data_vix)
134     target_data_np = np.array(target_data)
135
136     for i in range(len(data_logret_np) - seq_len):
137         x_logret = data_logret_np[i:(i + seq_len)]
138         x_vix = data_vix_np[i:(i + seq_len)]
139         x = np.stack((x_logret, x_vix), axis=1) # Features: LogRet, VIX
140         y = target_data_np[i + seq_len] # Target: Value at t+1
141
142         if not np.isnan(y): # Check target validity
143             xs.append(x)

```

```

144         ys.append(y)
145
146     return np.array(xs), np.array(ys)
147
148 # --- Training Function ---
149 def train_lstm(model, dataloader, epochs, learning_rate, device, model_name="LSTM"):
150     """ Trains the PyTorch LSTM model. """
151     criterion = nn.MSELoss() # Use MSE for volatility prediction
152     optimizer = torch.optim.Adam(model.parameters(), lr=learning_rate)
153     model.to(device)
154     model.train()
155
156     print(f"\nTraining {model_name} on {device}...")
157     for epoch in range(epochs):
158         epoch_loss = 0.0
159         pbar_batch = tqdm(dataloader, desc=f'Epoch {epoch+1}/{epochs}', leave=False)
160         for inputs, targets in pbar_batch:
161             inputs, targets = inputs.to(device), targets.to(device)
162             outputs = model(inputs)
163             loss = criterion(outputs.squeeze(), targets.squeeze()) # Ensure shapes match
164
165             optimizer.zero_grad()
166             loss.backward()
167             optimizer.step()
168
169             epoch_loss += loss.item()
170             pbar_batch.set_postfix({'loss': loss.item()})
171
172         avg_epoch_loss = epoch_loss / len(dataloader)
173         print(f'Epoch [{epoch+1}/{epochs}], Average Loss: {avg_epoch_loss:.6f}')
174         history['loss'].append(avg_epoch_loss)
175
176     print(f"{model_name} Training finished.")
177     model.eval()
178     return model
179
180 # --- Initial Model Training ---
181 print("\n--- Initial Model Training ---")
182 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
183 print(f"Using device: {device}")
184
185 # 1. Train LSTM (Predicts Volatility Directly)
186 print("Preparing data for LSTM (Volatility Prediction)...")
187 # Target is the Realized Volatility (not scaled by default, predict actual value)
188 X_lstm_train_np, y_lstm_train_np = create_lstm_sequences(
189     train_df['SPX_LogRet_Scaled'], # Input feature 1 (scaled)
190     train_df['VIX_Level_Scaled'], # Input feature 2 (scaled)
191     train_df['Realized_Vol_Target'], # Target is the actual realized volatility
192     LSTM_INPUT_SEQ_LEN

```

```

193 )
194
195 # Convert to PyTorch tensors
196 X_lstm_train = torch.tensor(X_lstm_train_np, dtype=torch.float32)
197 y_lstm_train = torch.tensor(y_lstm_train_np, dtype=torch.float32).unsqueeze(1)
198
199 # Create DataLoader
200 lstm_dataset = TensorDataset(X_lstm_train, y_lstm_train)
201 lstm_dataloader = DataLoader(lstm_dataset, batch_size=LSTM_BATCH_SIZE, shuffle=True)
202
203 # Initialize and train the LSTM model
204 lstm_model = LSTMVolModel(input_size=2, # Two features: LogRet, VIX
205                             hidden_size=LSTM_HIDDEN_SIZE,
206                             num_layers=LSTM_NUM_LAYERS,
207                             output_size=1, # Output is 1 (predicted volatility)
208                             dropout=LSTM_DROPOUT).to(device)
209
210 lstm_model = train_lstm(lstm_model, lstm_dataloader, LSTM_EPOCHS, LSTM_LEARNING_RATE,
211                          device, "LSTM (Volatility)")
212
213 # 2. Estimate Initial GARCH Parameters
214 print("\nFitting initial GARCH model on log returns...")
215 try:
216     # Use log returns directly from the training period
217     train_log_returns = train_df['SPX_LogRet'].dropna()
218     # Use arch_model for GARCH(1,1)
219     garch_model_spec = arch_model(train_log_returns * 100, vol='Garch', p=GARCH_P, q=
220                                   GARCH_Q,
221                                   mean='Zero', dist='Normal') # Scale by 100 for stability
222     res_garch_initial = garch_model_spec.fit(disp='off', show_warning=False)
223     print(res_garch_initial.summary())
224
225     # Extract parameters needed for rolling forecast
226     omega_garch = res_garch_initial.params['omega']
227     alpha_garch = res_garch_initial.params['alpha[1]']
228     beta_garch = res_garch_initial.params['beta[1]']
229
230     # Get last variance and residual for rolling forecast start (unscale variance/residual
231     )
232     last_h_garch = res_garch_initial.conditional_volatility[-1]**2 / (100**2) # Variance =
233     vol^2, unscale
234     last_resid_garch = train_log_returns.iloc[-1] # Use the actual last log return as the
235     last 'residual' for standard GARCH
236     print("GARCH parameters estimated.")
237 except Exception as e:
238     print(f"ERROR: Initial GARCH fitting failed: {e}")
239     raise SystemExit("GARCH fitting failed.")

```

```

237 # --- Rolling Forecast ---
238 print(f"\n--- Starting Rolling Forecast from {test_indices[0].date()} ---")
239
240 predictions = {
241     'Date': [],
242     'Actual': [],
243     'GARCH': [],          # Renamed from GARCH_Baseline
244     'LSTM': [],           # Renamed from LSTM_Baseline
245     'Combined_LSTM_GARCH': [] # New combined forecast
246 }
247
248 # Set models to evaluation mode
249 lstm_model.eval()
250
251 # Iterate through the test set for rolling predictions
252 for i in tqdm(range(len(test_indices) - 1), desc="Rolling Forecast"):
253     t_index = n_train + i
254     current_date = df.index[t_index]
255     next_date = df.index[t_index + 1] # Date for which volatility is being predicted
256
257     # 1. Prepare Input Data for step t
258     start_idx = t_index - LSTM_INPUT_SEQ_LEN
259     end_idx = t_index # Exclusive
260     if start_idx < 0: continue # Skip if not enough history
261
262     # Get scaled features for LSTM input
263     input_data_logret_scaled = scaled_df['SPX_LogRet_Scaled'].iloc[start_idx:end_idx].
        values
264     input_data_vix_scaled = scaled_df['VIX_Level_Scaled'].iloc[start_idx:end_idx].values
265
266     # Get actual log return at time t needed for GARCH update
267     actual_log_ret_t = df['SPX_LogRet'].iloc[t_index]
268
269     # Create LSTM input tensor (shape: [1, seq_len, num_features])
270     x_input_np = np.stack((input_data_logret_scaled, input_data_vix_scaled), axis=1).
        reshape(1, LSTM_INPUT_SEQ_LEN, 2)
271     x_input = torch.tensor(x_input_np, dtype=torch.float32).to(device)
272
273     # 2. Get Prediction from LSTM Model (Predicts Volatility sigma_hat_{t+1})
274     with torch.no_grad():
275         vol_hat_lstm_tplus1 = lstm_model(x_input).cpu().numpy().flatten()[0]
276         # Ensure prediction is non-negative
277         vol_hat_lstm_tplus1 = max(vol_hat_lstm_tplus1, 0)
278
279     # 3. Calculate GARCH Volatility for t+1
280     # Use actual log return at t as the 'innovation' epsilon_t for standard GARCH
281     epsilon_t_garch = actual_log_ret_t
282
283     # Calculate variance for day t: h_t = omega + alpha * epsilon_{t-1}^2 + beta * h_{t-1}

```



```

284 # Use last_resid_garch (logret_{t-1}) and last_h_garch (h_{t-1})
285 h_t_garch = omega_garch / (100**2) + alpha_garch * (last_resid_garch**2) + beta_garch
      * last_h_garch
286 # Note: omega was estimated on scaled returns, so unscale it here. alpha/beta are
      scale invariant.
287
288 # Forecast variance for day t+1: h_{t+1} = omega + alpha * epsilon_t^2 + beta * h_t
289 h_tplus1_garch = omega_garch / (100**2) + alpha_garch * (epsilon_t_garch**2) +
      beta_garch * h_t_garch
290
291 # Ensure variance is positive and calculate annualized volatility
292 h_tplus1_garch = max(h_tplus1_garch, 1e-12) # Prevent sqrt(0) or negative
293 vol_hat_garch_tplus1 = np.sqrt(h_tplus1_garch) * ANNUALIZATION_FACTOR
294
295 # Update GARCH state for next iteration (t+1)
296 last_h_garch = h_t_garch # h_t becomes the next h_{t-1}
297 last_resid_garch = epsilon_t_garch # epsilon_t becomes the next epsilon_{t-1}
298
299 # 4. Calculate Combined Forecast (Simple Average)
300 vol_hat_combined_tplus1 = (vol_hat_lstm_tplus1 + vol_hat_garch_tplus1) / 2.0
301
302 # 5. Store Predictions and Actuals for day t+1
303 # Actual realized volatility for day t+1 is stored at index t in 'Realized_Vol_Target'
304 actual_vol_tplus1 = df['Realized_Vol_Target'].iloc[t_index]
305
306 if not pd.isna(actual_vol_tplus1):
307     predictions['Date'].append(next_date) # Store the date the prediction is FOR
308     predictions['Actual'].append(actual_vol_tplus1)
309     predictions['GARCH'].append(vol_hat_garch_tplus1)
310     predictions['LSTM'].append(vol_hat_lstm_tplus1)
311     predictions['Combined_LSTM_GARCH'].append(vol_hat_combined_tplus1)
312 # else: # If actual is NaN, skip storing this step's predictions
313
314
315 print("Rolling forecast finished.")
316
317 # --- Evaluation ---
318 print("\n--- Evaluating Model Performance ---")
319
320 # Create DataFrame for results
321 results_df = pd.DataFrame(predictions)
322 results_df.set_index('Date', inplace=True)
323
324 # Drop any rows with NaN predictions or actuals
325 results_df.dropna(inplace=True)
326
327 if results_df.empty:
328     print("No valid results found for evaluation. Check data or prediction loop.")
329 else:

```

```

330 # Calculate Metrics: MAE, RMSE, MSE
331 metrics = {}
332 print("Evaluation Metrics (Test Set):")
333 # Evaluate GARCH, LSTM, and the new Combined model
334 for model_name in ['GARCH', 'LSTM', 'Combined_LSTM_GARCH']:
335     pred = results_df[model_name]
336     actual = results_df['Actual']
337     mae = mean_absolute_error(actual, pred)
338     mse = mean_squared_error(actual, pred)
339     rmse = np.sqrt(mse)
340     metrics[model_name] = {'MAE': mae, 'MSE': mse, 'RMSE': rmse}
341     print(f"\n--- {model_name} ---")
342     print(f"MAE: {mae:.6f}")
343     print(f"MSE: {mse:.6f}")
344     print(f"RMSE: {rmse:.6f}")
345
346 # Check Objective: Compare Combined model against baseline GARCH
347 if 'GARCH' in metrics and 'Combined_LSTM_GARCH' in metrics:
348     baseline_rmse = metrics['GARCH']['RMSE']
349     combined_rmse = metrics['Combined_LSTM_GARCH']['RMSE']
350     if baseline_rmse > 0: # Avoid division by zero
351         improvement = ((baseline_rmse - combined_rmse) / baseline_rmse) * 100
352         print(f"\n--- Objective Check ---")
353         print(f"Baseline GARCH RMSE: {baseline_rmse:.6f}")
354         print(f"Combined LSTM-GARCH RMSE: {combined_rmse:.6f}")
355         print(f"RMSE Reduction vs GARCH: {improvement:.2f}%")
356         # Set your target improvement percentage here
357         target_improvement = 15
358         if improvement >= target_improvement:
359             print(f"Objective MET: Combined model achieved >= {target_improvement}%
360                 RMSE reduction vs GARCH.")
361         else:
362             print(f"Objective NOT MET: Combined model achieved < {target_improvement
363                }% RMSE reduction vs GARCH.")
364     else:
365         print("\n--- Objective Check ---")
366         print("Baseline GARCH RMSE is zero, cannot calculate percentage improvement.")
367
368 else:
369     print("\nCould not perform objective check due to missing model results.")
370
371 # --- Plotting ---
372 if not results_df.empty:
373     # Plot actual vs predicted volatility
374     plt.figure(figsize=(15, 7))
375     plt.plot(results_df.index, results_df['Actual'], label='Actual Realized Volatility',
376             alpha=0.7, linewidth=1.5, color='blue')
377     plt.plot(results_df.index, results_df['GARCH'], label='GARCH(1,1) Forecast', alpha
378             =0.7, linewidth=1, linestyle='--', color='orange')

```

```

374 plt.plot(results_df.index, results_df['LSTM'], label='LSTM Forecast', alpha=0.7,
375           linewidth=1, linestyle=':', color='green')
376 plt.plot(results_df.index, results_df['Combined_LSTM_GARCH'], label='Combined (Avg)
377           Forecast', alpha=0.9, color='red', linewidth=1.2) # Highlight combined
378 plt.title('S&P 500 Volatility Forecast Comparison (Test Set)')
379 plt.xlabel('Date')
380 plt.ylabel('Annualized Volatility')
381 plt.legend()
382 plt.grid(True, linestyle='--', alpha=0.5)
383 plt.tight_layout()
384 plt.show()
385
386 # Plot forecast errors
387 plt.figure(figsize=(15, 7))
388 plt.plot(results_df.index, results_df['Actual'] - results_df['GARCH'], label='GARCH
389           Error', alpha=0.7)
390 plt.plot(results_df.index, results_df['Actual'] - results_df['LSTM'], label='LSTM
391           Error', alpha=0.7)
392 plt.plot(results_df.index, results_df['Actual'] - results_df['Combined_LSTM_GARCH'],
393           label='Combined (Avg) Error', alpha=0.7, color='red') # Highlight combined error
394 plt.title('Volatility Forecast Errors (Actual - Predicted) on Test Set')
395 plt.xlabel('Date')
396 plt.ylabel('Forecast Error')
397 plt.axhline(0, color='black', linestyle='--', linewidth=0.8)
398 plt.legend()
399 plt.grid(True, linestyle='--', alpha=0.5)
400 plt.tight_layout()
401 plt.show()
402
403 else:
404     print("Skipping plots as results are empty.")

```

Appendix D: Comparison with Traditional Models

```

1  from statsmodels.tsa.arima.model import ARIMA
2  from sklearn.model_selection import train_test_split
3  from sklearn.metrics import mean_squared_error
4
5  # 1. Split the data
6  train_data, test_data = train_test_split(df, test_size=0.2, shuffle=False)
7
8  # 2. Prepare train and test sets
9  # Exogenous variables for train and test
10 train_exog = train_data[['SPX_LogRet', 'VIX_Level']]
11 test_exog = test_data[['SPX_LogRet', 'VIX_Level']]
12
13 # Target variable for train and test
14 train_target = train_data['Realized_Vol_Target']
15 test_target = test_data['Realized_Vol_Target']
16

```

```
17 # 3. Fit the model on the training data
18 model = ARIMA(train_target, exog=train_exog, order=(1, 0, 1))
19 model_fit = model.fit()
20
21 # 4. Make predictions on the test data
22 predictions = model_fit.predict(start=0, end=len(test_data)-1, exog=test_exog)
23
24 # 5. Evaluate the model on the test data
25 rmse = np.sqrt(mean_squared_error(test_target, predictions))
26 print(f"RMSE on Test Set: {rmse}")
```