



Advanced Financial Market Analysis with MMachine Learning and GenAI

Authors: Shakil Rohimun & Chaker Meraihi & Luca Albani & Umberto Bechis

Degree Program: Master's in Quantitative Finance

Supervisor: Prof. Djibril Sarr

Date: May 2025

Abstract

This report presents a comprehensive framework for financial market analysis using advanced machine learning and generative AI techniques. We develop a multi-component system that integrates market regime detection, synthetic data generation, algorithmic trading strategies, and risk forecasting models. The market regime detection component employs multiple methods including Hidden Markov Models, Gaussian Mixture Models, Temporal Convolutional Networks with K-Means, and UMAP with K-Means to identify distinct market states. The generative AI component uses TimeGAN, Conditional Variational Autoencoders, and FiLM Transformers to generate synthetic financial data conditioned on detected regimes. These synthetic data are used to augment training datasets for trading strategies, addressing the limited historical data problem in finance. The algorithmic trading component implements regime-adaptive strategies including Transformer-based momentum, LSTM short signals, and Kalman filter pairs trading. Finally, the risk forecasting component employs regime-switching GARCH with Extreme Value Theory and Quantile Random Forests to estimate Value at Risk. Extensive backtesting demonstrates that our integrated approach significantly improves trading performance and risk estimation accuracy compared to baseline methods. The mathematical foundations, implementation details, and empirical results of each component are thoroughly presented and analyzed.

Contents

1	Introduction	10
1.1	Background and Motivation	10
1.2	Project Objectives	10
1.3	Overview of Financial Market Analysis Challenges	11
1.3.1	Non-stationarity and Regime Shifts	11
1.3.2	Limited Historical Data	11
1.3.3	Complex Temporal Dependencies	11
1.3.4	Risk Estimation in Tail Events	12
1.4	Contribution and Significance	12
1.5	Report Structure	13
2	Theoretical Foundations	13
2.1	Financial Time Series Analysis	13
2.1.1	Statistical Properties of Financial Returns	14
2.1.2	Stylized Facts in Financial Markets	15
2.1.3	Challenges in Financial Forecasting	15
2.2	Market Regime Theory	15
2.2.1	Definition and Characteristics of Market Regimes	15
2.2.2	Regime Switching Models in Finance	16
2.2.3	Hidden Markov Models for Regime Detection	16
2.2.4	Gaussian Mixture Models for Regime Classification	17
2.2.5	Dimensionality Reduction and Clustering Approaches	17
2.3	Generative AI in Finance	18
2.3.1	Principles of Generative Modeling	18
2.3.2	Time Series Generation Challenges	19

2.3.3	TimeGAN: Architecture and Mathematical Formulation	19
2.3.4	Conditional Variational Autoencoders: Theory and Implementation	20
2.3.5	FiLM Transformer: Feature-wise Linear Modulation	20
2.3.6	Regime Adversary Module: Ensuring Regime-Specific Properties .	21
2.4	Algorithmic Trading Strategies	21
2.4.1	Momentum and Mean Reversion	21
2.4.2	Transformer-based Momentum Strategy	21
2.4.3	LSTM Short Signal Strategy	22
2.4.4	Kalman Filter Pairs Trading	22
2.4.5	Strategy Decay Detection and Mitigation	22
2.5	Risk Forecasting Models	22
2.5.1	Value at Risk (VaR) Methodologies	23
2.5.2	REGARCH-EVT: Regime-Switching GARCH with Extreme Value Theory	23
2.5.3	Quantile Random Forest for Direct VaR Estimation	23
2.5.4	Synthetic Regime Bootstrapping	24
3	Data Preprocessing Methodology	25
3.1	Data Sources and Description	25
3.2	Data Cleaning and Handling Missing Values	26
3.2.1	Forward Fill Method for Time Series	26
3.2.2	Outlier Detection and Treatment	26
3.2.3	Winsorization and Clipping Techniques	27
3.3	Time Series Alignment	28
3.3.1	Union vs. Intersection Strategies	28
3.3.2	Implications of Alignment Choices	28
3.3.3	Preventing Data Leakage in Financial Time Series	29
3.4	Feature Engineering	29
3.4.1	Price-based Features	29
3.4.2	Return-based Features	30
3.4.3	Volatility Measures	30
3.4.4	Technical Indicators	31
3.4.5	Cross-asset Features	32
3.4.6	Feature Selection and Dimensionality Considerations	32
3.5	Mathematical Formulation of Key Features	32
3.5.1	Moving Averages and Exponential Smoothing	33
3.5.2	Bollinger Bands and Volatility Channels	33
3.5.3	Momentum Indicators (RSI, MACD)	33
3.5.4	Statistical Features (Kurtosis, Skewness)	34
4	Market Regime Detection Implementation	34
4.1	Hidden Markov Model (HMM)	34
4.1.1	Mathematical Formulation	35
4.1.2	Parameter Estimation via Expectation-Maximization	35
4.1.3	Implementation Details	36
4.1.4	Regime Transition Probability Matrix Analysis	37
4.2	Gaussian Mixture Model (GMM)	37
4.2.1	Mathematical Formulation	38

4.2.2	Expectation-Maximization Algorithm	38
4.2.3	Implementation Details	39
4.2.4	Component Distribution Analysis	39
4.3	Temporal Convolutional Network (TCN) + K-Means	40
4.3.1	TCN Architecture and Mathematical Foundations	40
4.3.2	Sequence Embedding Generation	41
4.3.3	K-Means Clustering in Latent Space	42
4.3.4	Implementation Challenges and Solutions	42
4.4	UMAP + K-Means	43
4.4.1	UMAP Mathematical Foundations	43
4.4.2	Hyperparameter Optimization	44
4.4.3	Implementation Details	44
4.4.4	Topological Structure Preservation Analysis	45
4.5	Ensemble Approach and Temporal Smoothing	45
4.5.1	Ensemble Methodology	46
4.5.2	Weighted Voting Mechanism	46
4.5.3	Temporal Smoothing Techniques	46
4.5.4	Regime Confidence Scores	47
4.6	Regime Characterization	47
4.6.1	Statistical Properties of Identified Regimes	47
4.6.2	Temporal Distribution Analysis	48
4.6.3	Regime Transition Analysis	49
4.6.4	Economic Interpretation of Regimes	49
5	GenAI Data Augmentation Implementation	50
5.1	TimeGAN Implementation	50
5.1.1	Architecture Components	50
5.1.2	Loss Functions and Training Dynamics	51
5.1.3	Regime-Conditional Generation	52
5.1.4	Implementation Challenges and Solutions	53
5.2	Conditional Variational Autoencoder (CVAE) Implementation	53
5.2.1	Encoder-Decoder Architecture	53
5.2.2	Latent Space Sampling with Regime Conditioning	54
5.2.3	Balancing Reconstruction and KL Divergence Losses	54
5.2.4	Implementation Details	55
5.3	FiLM Transformer Implementation	55
5.3.1	Self-Attention Mechanism	56
5.3.2	Feature-wise Linear Modulation Layers	56
5.3.3	Positional Encoding for Temporal Awareness	56
5.3.4	Implementation Challenges and Solutions	57
5.4	Regime Adversary Module Implementation	57
5.4.1	Adversarial Training Framework	57
5.4.2	Feature Matching for Statistical Alignment	58
5.4.3	Gradient Penalty for Training Stability	58
5.4.4	Implementation Details	58
5.5	Synthetic Data Validation Framework	59
5.5.1	Statistical Tests	59
5.5.2	Financial Metrics	60

5.5.3	Machine Learning Metrics	60
5.5.4	Visualization Techniques	61
5.6	Synthetic Boost Protocol	61
5.6.1	Regime-Specific Augmentation Ratios	61
5.6.2	Balanced Training Set Creation	62
5.6.3	Sample Weighting Based on Quality Metrics	62
5.6.4	Implementation Details	63
6	Algorithmic Trading Strategy Implementation	63
6.1	Transformer-based Momentum Strategy	63
6.1.1	Architecture and Mathematical Formulation	63
6.1.2	Regime-Aware Training	65
6.1.3	Signal Generation and Position Sizing	65
6.1.4	Implementation Details	66
6.2	LSTM Short Signal Strategy	66
6.2.1	LSTM Architecture and Mathematical Formulation	66
6.2.2	Asymmetric Loss Function for Short Bias	67
6.2.3	Regime-Specific Feature Importance	68
6.2.4	Short Signal Filtering and Execution	68
6.2.5	Implementation Details	69
6.3	Kalman Filter Pairs Trading	70
6.3.1	State-Space Model Formulation	70
6.3.2	Kalman Filter Estimation	70
6.3.3	Regime-Dependent Noise Variances	71
6.3.4	Trading Signal Generation	71
6.3.5	Regime-Specific Threshold Adjustment	71
6.3.6	Implementation Details	72
6.4	Strategy Decay Detection and Mitigation	72
6.4.1	Performance Monitoring Framework	72
6.4.2	Statistical Tests for Decay Detection	73
6.4.3	Regime-Specific Decay Detection	73
6.4.4	Mitigation Strategies	74
6.4.5	Implementation Details	74
6.5	Strategy Ensemble and Portfolio Construction	75
6.5.1	Strategy Correlation Analysis	75
6.5.2	Regime-Dependent Strategy Allocation	75
6.5.3	Risk Parity Portfolio Construction	75
6.5.4	Dynamic Risk Budgeting	76
6.5.5	Implementation Details	76
7	Risk Forecasting Implementation	77
7.1	REGARCH-EVT: Regime-Switching GARCH with Extreme Value Theory	77
7.1.1	Regime-Switching GARCH Model	77
7.1.2	Parameter Estimation via Maximum Likelihood	77
7.1.3	Extreme Value Theory for Tail Modeling	78
7.1.4	Value at Risk (VaR) and Expected Shortfall (ES) Calculation	78
7.1.5	Implementation Details	79
7.2	Quantile Random Forest for Direct VaR Estimation	80

7.2.1	Random Forest Framework	80
7.2.2	Quantile Estimation	81
7.2.3	Feature Engineering for Risk Forecasting	81
7.2.4	Regime-Specific Forests	82
7.2.5	Synthetic Data Augmentation	82
7.2.6	Implementation Details	82
7.3	Synthetic Regime Bootstrapping	83
7.3.1	Traditional Bootstrap Methods	83
7.3.2	Regime-Specific Bootstrap	83
7.3.3	Synthetic Data Integration	83
7.3.4	Block Bootstrap for Time Series	84
7.3.5	Scenario Generation and Risk Estimation	84
7.3.6	Implementation Details	84
7.4	Ensemble Risk Forecasting	85
7.4.1	Model Combination Methods	85
7.4.2	Regime-Dependent Model Weights	86
7.4.3	Dynamic Weight Adjustment	86
7.4.4	Implementation Details	87
7.5	Risk Decomposition and Attribution	87
7.5.1	Factor-Based Risk Decomposition	87
7.5.2	Regime-Specific Risk Factors	88
7.5.3	Marginal VaR and Component VaR	88
7.5.4	Stress Testing and Scenario Analysis	89
7.5.5	Implementation Details	89
7.6	Risk Forecasting Validation Framework	90
7.6.1	Statistical Tests for VaR Backtesting	90
7.6.2	Scoring Rules for Distributional Forecasts	91
7.6.3	Comparative Backtesting	91
7.6.4	Regime-Specific Validation	91
7.6.5	Implementation Details	91
8	Experimental Results and Analysis	92
8.1	Market Regime Detection Results	92
8.1.1	Regime Identification Accuracy	92
8.1.2	Regime Characteristics Analysis	93
8.1.3	Temporal Distribution of Regimes	93
8.1.4	Regime Transition Analysis	93
8.2	GenAI Data Augmentation Results	94
8.2.1	Statistical Fidelity of Synthetic Data	94
8.2.2	Visual Comparison of Real and Synthetic Data	95
8.2.3	Impact on Downstream Tasks	95
8.2.4	Ablation Study: Regime Adversary Module	96
8.3	Algorithmic Trading Strategy Results	96
8.3.1	Strategy Performance Metrics	96
8.3.2	Regime-Specific Performance	97
8.3.3	Impact of Regime Awareness and Synthetic Data	98
8.3.4	Strategy Decay Analysis	98
8.4	Risk Forecasting Results	98

8.4.1	VaR Backtesting Results	99
8.4.2	Regime-Specific Risk Forecasting Accuracy	99
8.4.3	Impact of Synthetic Data on Risk Forecasting	100
8.4.4	Expected Shortfall Accuracy	101
8.5	Integrated Framework Performance	101
8.5.1	Portfolio Performance	101
8.5.2	Performance During Market Stress	102
8.5.3	Regime Transition Navigation	102
8.5.4	Attribution Analysis	103
9	Discussion and Implications	103
9.1	Synthesis of Findings	104
9.1.1	Market Regime Characterization	104
9.1.2	Effectiveness of Ensemble Methods	104
9.1.3	Impact of Regime Awareness	105
9.1.4	Value of Synthetic Data Augmentation	105
9.1.5	Robustness During Market Stress	106
9.2	Theoretical Implications	106
9.2.1	Regime-Dependent Market Efficiency	106
9.2.2	Non-Stationarity and Model Adaptation	106
9.2.3	Tail Risk and Extreme Events	107
9.2.4	Generative AI and Financial Data	107
9.3	Practical Implications	107
9.3.1	Investment Strategy Design	108
9.3.2	Risk Management Practices	108
9.3.3	Model Development and Validation	108
9.4	Limitations and Future Research	109
9.4.1	Methodological Limitations	109
9.4.2	Data Limitations	109
9.4.3	Implementation Challenges	109
9.4.4	Future Research Directions	110
9.5	Ethical Considerations	110
9.5.1	Market Impact and Systemic Risk	110
9.5.2	Fairness and Access	111
9.5.3	Transparency and Explainability	111
9.6	Conclusion of Discussion	111
10	Conclusion	112
10.1	Summary of Contributions	112
10.1.1	Theoretical Contributions	112
10.1.2	Empirical Contributions	112
10.1.3	Methodological Contributions	113
10.2	Practical Applications	113
10.3	Future Research Directions	114
10.4	Concluding Remarks	115

List of Figures

1	Consensus market regimes identified by the ensemble approach, showing the temporal distribution of regimes over the study period (2010-2023). The plot illustrates distinct bull market (Regime 0), neutral/transition (Regime 1), and bear market (Regime 2) periods.	115
2	Market regimes identified by the Hidden Markov Model (HMM) approach. The HMM captures regime transitions but shows some sensitivity to short-term market fluctuations.	116
3	Market regimes identified by the Gaussian Mixture Model (GMM) approach. The GMM effectively captures the distributional differences between regimes but is less sensitive to temporal dependencies.	116
4	Market regimes identified by the Temporal Convolutional Network with K-means clustering (TCN+K-means) approach. This method effectively captures complex temporal patterns in the data.	116
5	Market regimes identified by the UMAP dimensionality reduction with K-means clustering (UMAP+K-means) approach. This method preserves both local and global structure in the data.	117
6	Comparison of real and synthetic data distributions for Regime 0 (bull market) using different generative models. The plots show the projection of data points onto the first two principal components (PCA) and t-SNE embeddings. Real data points are shown in blue, while synthetic data points are shown in orange. The FiLM Transformer generates synthetic data that most closely matches the distribution of real data.	118
7	Comparison of real and synthetic data distributions for Regime 1 (neutral/-transition market) using different generative models. The plots show the projection of data points onto the first two principal components (PCA) and t-SNE embeddings. Real data points are shown in blue, while synthetic data points are shown in orange. The FiLM Transformer generates synthetic data that most closely matches the distribution of real data.	119
8	Comparison of real and synthetic data distributions for Regime 2 (bear market) using different generative models. The plots show the projection of data points onto the first two principal components (PCA) and t-SNE embeddings. Real data points are shown in blue, while synthetic data points are shown in orange. The FiLM Transformer generates synthetic data that most closely matches the distribution of real data, particularly in capturing the extreme tail behavior characteristic of bear markets.	120
9	Cumulative returns of trading strategies using the baseline approach (without synthetic data augmentation). The plot shows the performance of different strategies across the test period (2018-2023).	121
10	Cumulative returns of trading strategies using the augmented approach (with synthetic data augmentation). The plot shows the improved performance of strategies trained with synthetic data across the test period (2018-2023).	121

List of Tables

1	Performance Metrics for Market Regime Detection Methods	92
---	---	----

2	Statistical Characteristics of Identified Market Regimes	93
3	Regime Transition Probability Matrix	94
4	Statistical Fidelity of Synthetic Data by Regime	94
5	Impact of Synthetic Data Augmentation on Downstream Tasks	95
6	Ablation Study: Impact of Regime Adversary Module	96
7	Performance Metrics for Trading Strategies (2018-2023)	96
8	Strategy Performance by Market Regime	97
9	Ablation Study: Impact of Regime Awareness and Synthetic Data	98
10	Backtesting Results for 99% One-Day VaR	99
11	Risk Forecasting Accuracy by Market Regime	99
12	Impact of Synthetic Data on Risk Forecasting in Regime 2	100
13	Accuracy of 97.5% Expected Shortfall Forecasts	101
14	Performance of Integrated Framework (2018-2023)	101
15	Performance During Market Stress Events	102
16	Performance Attribution Analysis	103

List of Abbreviations

ACF	Autocorrelation Function
CVAE	Conditional Variational Autoencoder
EVT	Extreme Value Theory
FiLM	Feature-wise Linear Modulation
GARCH	Generalized Autoregressive Conditional Heteroskedasticity
GMM	Gaussian Mixture Model
HMM	Hidden Markov Model
LSTM	Long Short-Term Memory
MACD	Moving Average Convergence Divergence
PACF	Partial Autocorrelation Function
PCA	Principal Component Analysis
REGARCH	Regime-switching GARCH
RSI	Relative Strength Index
TCN	Temporal Convolutional Network
TimeGAN	Time-series Generative Adversarial Network
t-SNE	t-Distributed Stochastic Neighbor Embedding
UMAP	Uniform Manifold Approximation and Projection
VaR	Value at Risk

1 Introduction

Financial markets are complex adaptive systems characterized by non-stationarity, regime shifts, and intricate dependencies. Traditional financial analysis methods often struggle to capture these dynamics, particularly during periods of market stress or structural change. This report presents a comprehensive framework that leverages advanced machine learning and generative artificial intelligence techniques to address these challenges in financial market analysis.

1.1 Background and Motivation

The behavior of financial markets exhibits several well-documented stylized facts that pose significant challenges for modeling and prediction. These include heavy-tailed return distributions, volatility clustering, leverage effects, and regime-switching behavior [Cont, 2001]. Traditional econometric models such as ARIMA and GARCH capture some of these characteristics but often fail to adapt to structural changes in market dynamics.

Market regimes—distinct states characterized by different statistical properties of returns, volatilities, and correlations—represent a fundamental concept in financial analysis. The ability to identify these regimes and adapt investment strategies accordingly is crucial for successful portfolio management and risk assessment. However, regime detection is challenging due to the latent nature of market states and the gradual transitions between them.

Another significant challenge in financial modeling is the limited availability of historical data, particularly for rare market conditions such as financial crises. This data scarcity problem is exacerbated when developing machine learning models, which typically require large training datasets. Traditional data augmentation techniques often fail to capture the complex temporal dependencies and statistical properties of financial time series.

1.2 Project Objectives

This project aims to develop an integrated framework for financial market analysis that addresses the aforementioned challenges through four interconnected components:

1. **Market Regime Detection:** Implement and compare multiple methods for identifying distinct market regimes, including Hidden Markov Models (HMM), Gaussian Mixture Models (GMM), Temporal Convolutional Networks (TCN) with K-Means clustering, and Uniform Manifold Approximation and Projection (UMAP) with K-Means clustering.
2. **GenAI Data Augmentation:** Develop generative models capable of producing synthetic financial data that preserves the statistical properties of each detected market regime, using TimeGAN, Conditional Variational Autoencoders (CVAE), and Feature-wise Linear Modulation (FiLM) Transformers.
3. **Algorithmic Trading Strategies:** Design regime-adaptive trading strategies that leverage both original and synthetic data, including Transformer-based momentum, LSTM short signals, and Kalman filter pairs trading approaches.

- 4. Risk Forecasting:** Implement advanced risk models that account for regime-specific characteristics, including Regime-switching GARCH with Extreme Value Theory (REGARCH-EVT) and Quantile Random Forests.

The integration of these components aims to create a system that can detect market regimes, generate realistic synthetic data for each regime, develop trading strategies optimized for different regimes, and accurately forecast risk across varying market conditions.

1.3 Overview of Financial Market Analysis Challenges

Financial market analysis presents several unique challenges that our framework addresses:

1.3.1 Non-stationarity and Regime Shifts

Financial time series are inherently non-stationary, with statistical properties that evolve over time. Market regimes represent periods with relatively stable statistical characteristics, separated by transitions that may be gradual or abrupt. Detecting these regimes is challenging because they are latent states that must be inferred from observable market data.

The mathematical formulation of this challenge involves modeling the observed returns r_t as being drawn from a mixture of distributions, with the active distribution at time t determined by a latent regime variable z_t :

$$r_t \sim p(r_t|z_t = k) \quad \text{for } k \in \{1, 2, \dots, K\} \quad (1)$$

where K is the number of regimes and $p(r_t|z_t = k)$ is the conditional distribution of returns given regime k .

1.3.2 Limited Historical Data

Financial markets provide only one historical trajectory, with limited examples of extreme events such as financial crises. This data scarcity poses a significant challenge for machine learning models, which typically require large training datasets. The problem is particularly acute for regime-specific modeling, as some regimes (e.g., crisis periods) may have very few historical examples.

Mathematically, we aim to estimate the joint distribution $p(\mathbf{X}, \mathbf{y})$ of features \mathbf{X} and targets \mathbf{y} for each regime k , but with limited samples from the true distribution:

$$\{(\mathbf{x}_i, y_i) \sim p(\mathbf{X}, \mathbf{y}|z_t = k)\}_{i=1}^{n_k} \quad (2)$$

where n_k is the number of samples in regime k , which may be insufficient for reliable estimation.

1.3.3 Complex Temporal Dependencies

Financial time series exhibit complex temporal dependencies, including autocorrelation, volatility clustering, and long-range dependencies. These characteristics must be preserved in any synthetic data generation process to ensure the validity of subsequent analyses.

The autocorrelation function of absolute returns, for instance, typically shows a slow decay that indicates volatility clustering:

$$\text{Corr}(|r_t|, |r_{t+\tau}|) > 0 \quad \text{for } \tau > 0 \quad (3)$$

This persistence in volatility is a key feature that must be captured in both regime detection and synthetic data generation.

1.3.4 Risk Estimation in Tail Events

Accurate estimation of risk during extreme market conditions is crucial for portfolio management. Traditional risk models often underestimate tail risk due to their reliance on normal distribution assumptions or historical simulation methods that may not adequately represent extreme events.

The Value at Risk (VaR) at confidence level α is defined as:

$$\text{VaR}_\alpha(t) = \inf\{x \in \mathbb{R} : P(r_t \leq x) \geq \alpha\} \quad (4)$$

Accurately estimating this quantity, especially for high confidence levels (e.g., $\alpha = 0.99$ or 0.999), requires models that can capture the fat-tailed nature of return distributions and account for regime-specific risk characteristics.

1.4 Contribution and Significance

This project makes several significant contributions to the field of financial market analysis:

1. Development of a comprehensive framework that integrates regime detection, synthetic data generation, trading strategy optimization, and risk forecasting.
2. Comparative analysis of multiple regime detection methods, providing insights into their relative strengths and weaknesses for financial applications.
3. Novel application of generative AI techniques for financial time series augmentation, with specific focus on preserving regime-dependent statistical properties.
4. Design of regime-adaptive trading strategies that leverage both historical and synthetic data, demonstrating improved performance over regime-agnostic approaches.
5. Advanced risk forecasting models that account for regime-specific characteristics, providing more accurate risk estimates across varying market conditions.

The significance of this work lies in its potential to enhance financial decision-making through more accurate market regime identification, improved model training with synthetic data, optimized trading strategies, and better risk management. The integrated approach addresses multiple challenges simultaneously, creating a synergistic effect that exceeds the benefits of each component in isolation.

1.5 Report Structure

The remainder of this report is organized as follows:

Section 2 presents the theoretical foundations of our framework, including financial time series analysis, market regime theory, generative AI in finance, algorithmic trading strategies, and risk forecasting models.

Section 3 details our data preprocessing methodology, covering data sources, cleaning procedures, time series alignment, feature engineering, and mathematical formulations of key features.

Section 4 describes the implementation of our market regime detection component, including HMM, GMM, TCN+K-Means, and UMAP+K-Means approaches, along with ensemble methods and regime characterization.

Section 5 explains the GenAI data augmentation component, covering TimeGAN, CVAE, FiLM Transformer, and Regime Adversary Module implementations, as well as synthetic data validation and the synthetic boost protocol.

Section 6 details the algorithmic trading strategy component, including the regime-adaptive strategy portfolio, Transformer-based momentum strategy, LSTM short signal strategy, Kalman filter pairs trading, strategy decay detection, execution layer, and backtesting protocol.

Section 7 describes the risk forecasting component, covering REGARCH-EVT, Quantile Random Forest, synthetic regime bootstrapping, backtesting framework, and stress testing suite.

Section 8 presents our experimental results and analysis, including market regime detection results, GenAI data augmentation results, trading strategy performance, risk forecasting accuracy, and integrated system performance.

Section 9 discusses the implications of our findings, including theoretical and practical implications, limitations and challenges, and future research directions.

Section 10 concludes the report with a summary of contributions, key findings, practical applications, and concluding remarks.

The appendices provide additional mathematical derivations, results, code implementation details, hyperparameter settings, and data descriptions.

2 Theoretical Foundations

This section lays the groundwork for the methodologies employed in this report. We delve into the fundamental theories and mathematical concepts underpinning financial time series analysis, market regime identification, generative artificial intelligence in finance, algorithmic trading strategies, and risk forecasting models.

2.1 Financial Time Series Analysis

Financial time series, such as asset prices or returns, exhibit unique statistical properties that distinguish them from other types of time series data. Understanding these properties is crucial for developing effective models.

2.1.1 Statistical Properties of Financial Returns

Let p_t denote the price of a financial asset at time t . The log return r_t is typically defined as:

$$r_t = \log(p_t) - \log(p_{t-1}) = \log\left(\frac{p_t}{p_{t-1}}\right) \quad (5)$$

Log returns are often preferred over simple returns ($R_t = (p_t - p_{t-1})/p_{t-1}$) due to their additive properties over time and their closer approximation to normality compared to simple returns.

Financial returns exhibit several well-documented *stylized facts* [Cont, 2001]:

- **Absence of Autocorrelation:** Returns themselves typically show little to no significant autocorrelation, especially at longer lags. $\text{Corr}(r_t, r_{t+\tau}) \approx 0$ for large τ .
- **Heavy Tails (Leptokurtosis):** The distribution of returns is characterized by fatter tails than the normal distribution. This means extreme events (large positive or negative returns) occur more frequently than predicted by a Gaussian model. Mathematically, the kurtosis κ is greater than 3 (the kurtosis of a normal distribution):

$$\kappa = \frac{\mathbb{E}[(r_t - \mu)^4]}{\sigma^4} > 3 \quad (6)$$

where $\mu = \mathbb{E}[r_t]$ and $\sigma^2 = \text{Var}(r_t)$.

- **Volatility Clustering:** Periods of high volatility tend to be followed by periods of high volatility, and periods of low volatility tend to be followed by periods of low volatility. This is reflected in the slow decay of the autocorrelation function (ACF) of absolute returns $|r_t|$ or squared returns r_t^2 :

$$\text{Corr}(|r_t|, |r_{t+\tau}|) > 0 \quad \text{and decays slowly for } \tau > 0 \quad (7)$$

- **Leverage Effect:** Volatility tends to increase more following a large negative return than following a large positive return of the same magnitude. This implies a negative correlation between past returns and future volatility:

$$\text{Corr}(r_{t-1}, \sigma_t^2) < 0 \quad (8)$$

where σ_t^2 is the conditional variance at time t .

- **Asymmetry (Skewness):** The distribution of returns is often slightly negatively skewed, meaning large negative returns are more probable than large positive returns of the same magnitude. The skewness s is typically less than 0:

$$s = \frac{\mathbb{E}[(r_t - \mu)^3]}{\sigma^3} < 0 \quad (9)$$

2.1.2 Stylized Facts in Financial Markets

These stylized facts have profound implications for financial modeling. The heavy tails necessitate the use of distributions beyond the Gaussian framework, such as Student's t-distribution or Generalized Pareto Distribution (GPD) for tail modeling. Volatility clustering motivates the development of time-varying volatility models like GARCH. The leverage effect suggests asymmetric responses in volatility models. The absence of auto-correlation in returns supports the Efficient Market Hypothesis (EMH) in its weak form, but the predictability in volatility and higher moments provides opportunities for risk management and trading.

2.1.3 Challenges in Financial Forecasting

Forecasting financial time series is notoriously difficult due to several factors:

- **Non-stationarity:** The statistical properties (mean, variance, correlations) of financial time series change over time, often abruptly due to structural breaks or regime shifts.
- **Low Signal-to-Noise Ratio:** The predictable component (signal) in financial returns is often very weak compared to the random fluctuations (noise).
- **Complexity and High Dimensionality:** Market dynamics are driven by a vast number of interacting factors, leading to complex, high-dimensional relationships.
- **Data Limitations:** As mentioned, historical data is limited, especially for rare events.
- **Reflexivity:** Market participants' actions can influence market dynamics, creating feedback loops that are hard to model.

2.2 Market Regime Theory

Market regime theory posits that financial markets operate in distinct states or regimes, each characterized by different statistical properties and dynamics. Identifying these regimes allows for more adaptive modeling and decision-making.

2.2.1 Definition and Characteristics of Market Regimes

A market regime \mathcal{R}_k can be formally defined as a period during which the underlying data generating process for financial variables (e.g., returns, volatility) remains relatively stable. Let z_t be a latent variable representing the regime at time t , taking values in $\{1, 2, \dots, K\}$. The observed financial variable y_t is assumed to be drawn from a regime-dependent distribution:

$$y_t \sim p(y_t | z_t = k, \mathcal{F}_{t-1}) \quad (10)$$

where \mathcal{F}_{t-1} represents the information available up to time $t - 1$. Regimes can differ in terms of mean returns, volatility levels, correlation structures, or other statistical properties.

2.2.2 Regime Switching Models in Finance

Regime switching models explicitly incorporate the possibility of shifts between different market states. The most prominent example is the Markov Switching Model introduced by Hamilton [1989].

In a Markov Switching Model, the latent regime variable z_t follows a first-order Markov chain, characterized by a transition probability matrix $A = [a_{ij}]$, where:

$$a_{ij} = P(z_t = j | z_{t-1} = i), \quad \sum_{j=1}^K a_{ij} = 1 \quad \forall i \in \{1, \dots, K\} \quad (11)$$

The observed variable y_t depends on the current regime z_t . For example, a Markov Switching Autoregressive (MS-AR) model can be written as:

$$y_t = \mu_{z_t} + \sum_{l=1}^p \phi_{l,z_t} (y_{t-l} - \mu_{z_{t-l}}) + \sigma_{z_t} \epsilon_t, \quad \epsilon_t \sim \mathcal{N}(0, 1) \quad (12)$$

where the parameters μ, ϕ_l, σ are regime-dependent.

2.2.3 Hidden Markov Models for Regime Detection

Hidden Markov Models (HMMs) are a powerful tool for modeling systems with latent states that evolve according to a Markov process, where the observed data depends probabilistically on the hidden state. In finance, HMMs can be used to infer market regimes (z_t) from observed financial data (y_t).

An HMM is defined by:

- A set of K hidden states (regimes): $\mathcal{S} = \{s_1, \dots, s_K\}$
- A set of observation symbols or distributions.
- An initial state probability distribution: $\pi = [\pi_i]$, where $\pi_i = P(z_1 = s_i)$.
- A state transition probability matrix: $A = [a_{ij}]$, where $a_{ij} = P(z_t = s_j | z_{t-1} = s_i)$.
- An emission probability distribution: $B = [b_j(y_t)]$, where $b_j(y_t) = p(y_t | z_t = s_j)$. For continuous observations, this is often a probability density function, e.g., Gaussian: $b_j(y_t) = \mathcal{N}(y_t | \mu_j, \Sigma_j)$.

The three fundamental problems for HMMs are:

1. **Evaluation:** Given the model parameters $\lambda = (A, B, \pi)$ and an observation sequence $Y = y_1, \dots, y_T$, compute the probability $P(Y|\lambda)$. Solved by the Forward algorithm.
2. **Decoding:** Given λ and Y , find the most likely sequence of hidden states $Z = z_1, \dots, z_T$. Solved by the Viterbi algorithm.
3. **Learning:** Given Y , estimate the model parameters λ that maximize $P(Y|\lambda)$. Solved by the Baum-Welch algorithm (an Expectation-Maximization algorithm).

The Forward algorithm computes the forward variable $\alpha_t(i) = P(y_1, \dots, y_t, z_t = s_i | \lambda)$. The Backward algorithm computes the backward variable $\beta_t(i) = P(y_{t+1}, \dots, y_T | z_t = s_i, \lambda)$. The Baum-Welch algorithm iteratively updates the parameters based on expected state transitions and emissions derived from $\alpha_t(i)$ and $\beta_t(i)$.

2.2.4 Gaussian Mixture Models for Regime Classification

Gaussian Mixture Models (GMMs) provide an alternative, static approach to regime identification. A GMM assumes that the observed data points are generated from a mixture of several Gaussian distributions, each representing a different regime or cluster.

The probability density function of a GMM is:

$$p(y) = \sum_{k=1}^K \pi_k \mathcal{N}(y|\mu_k, \Sigma_k) \quad (13)$$

where K is the number of components (regimes), π_k are the mixing coefficients ($\{\pi_k \geq 0, \sum_{k=1}^K \pi_k = 1\}$), and $\mathcal{N}(y|\mu_k, \Sigma_k)$ is the Gaussian density for component k with mean μ_k and covariance Σ_k .

The parameters ($\{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$) are typically estimated using the Expectation-Maximization (EM) algorithm.

- **E-step:** Compute the responsibilities $\gamma(z_{nk})$ that component k takes for explaining data point y_n :

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(y_n|\mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(y_n|\mu_j, \Sigma_j)} \quad (14)$$

- **M-step:** Re-estimate the parameters using the computed responsibilities:

$$\mu_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) y_n \quad (15)$$

$$\Sigma_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (y_n - \mu_k^{\text{new}})(y_n - \mu_k^{\text{new}})^T \quad (16)$$

$$\pi_k^{\text{new}} = \frac{N_k}{N} \quad (17)$$

where $N_k = \sum_{n=1}^N \gamma(z_{nk})$ and N is the total number of data points.

The EM algorithm iterates between the E-step and M-step until convergence. Once the model is fitted, the posterior probability $P(z_n = k|y_n) = \gamma(z_{nk})$ can be used to assign each data point to the most likely regime.

Unlike HMMs, standard GMMs do not explicitly model temporal dependencies or transitions between regimes. They classify each data point independently based on its features.

2.2.5 Dimensionality Reduction and Clustering Approaches

High-dimensional feature spaces can obscure the underlying regime structure. Dimensionality reduction techniques aim to project the data onto a lower-dimensional manifold where clusters (regimes) might be more apparent. Clustering algorithms are then applied in this reduced space.

UMAP (Uniform Manifold Approximation and Projection) [McInnes et al., 2018] is a non-linear dimensionality reduction technique based on manifold learning theory and topological data analysis. It seeks to find a low-dimensional embedding that preserves the topological structure of the high-dimensional data.

UMAP constructs a high-dimensional graph representing the data's topology and then optimizes a low-dimensional graph to be as structurally similar as possible. Key steps involve:

1. Constructing a weighted k-neighbor graph in the high-dimensional space.
2. Defining a fuzzy simplicial set representation of the data.
3. Constructing a similar fuzzy topological structure in the low-dimensional space.
4. Minimizing the cross-entropy between the high-dimensional and low-dimensional topological representations using stochastic gradient descent.

The objective function minimized by UMAP is typically the fuzzy set cross-entropy:

$$C(X, Y) = \sum_{i \neq j} \left[\mu_{ij} \log \left(\frac{\mu_{ij}}{\nu_{ij}} \right) + (1 - \mu_{ij}) \log \left(\frac{1 - \mu_{ij}}{1 - \nu_{ij}} \right) \right] \quad (18)$$

where μ_{ij} represents the similarity between points i and j in the high-dimensional space, and ν_{ij} represents the similarity in the low-dimensional embedding Y .

K-Means Clustering is an iterative algorithm that partitions the data into K pre-defined clusters. It aims to minimize the within-cluster variance (sum of squared distances between data points and their assigned cluster centroid):

$$J = \sum_{k=1}^K \sum_{x_n \in C_k} \|x_n - c_k\|^2 \quad (19)$$

where C_k is the set of points assigned to cluster k , and c_k is the centroid of cluster k . The algorithm alternates between assigning points to the nearest centroid and recalculating the centroids based on the new assignments (Lloyd's algorithm).

Combining UMAP with K-Means involves first reducing the dimensionality of the feature space using UMAP and then applying K-Means clustering to the resulting low-dimensional embeddings to identify regimes.

2.3 Generative AI in Finance

Generative AI models learn the underlying distribution of data and can generate new samples that resemble the original data. In finance, they offer a promising approach to data augmentation, scenario generation, and risk modeling.

2.3.1 Principles of Generative Modeling

Generative models aim to learn the probability distribution $p(x)$ of the observed data x . Once learned, the model can generate new samples $x_{\text{new}} \sim p(x)$. Common approaches include:

- **Maximum Likelihood Estimation (MLE):** Directly model $p(x; \theta)$ and find parameters θ that maximize the likelihood of the observed data.

- **Generative Adversarial Networks (GANs)** [Goodfellow et al., 2014]: Train a generator network G to produce realistic samples and a discriminator network D to distinguish real samples from generated ones. They engage in a minimax game:

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{\text{data}}(x)}[\log D(x)] + \mathbb{E}_{z \sim p_z(z)}[\log(1 - D(G(z)))] \quad (20)$$

- **Variational Autoencoders (VAEs)** [Kingma and Welling, 2013]: Learn a latent variable model $p(x, z) = p(x|z)p(z)$ by maximizing a lower bound on the data log-likelihood (Evidence Lower Bound, ELBO):

$$\log p(x) \geq \text{ELBO} = \mathbb{E}_{q(z|x)}[\log p(x|z)] - D_{KL}(q(z|x)||p(z)) \quad (21)$$

where $q(z|x)$ is an approximate posterior (encoder) and $p(x|z)$ is the generative model (decoder).

2.3.2 Time Series Generation Challenges

Generating realistic financial time series poses specific challenges:

- Capturing complex temporal dependencies across multiple time scales.
- Preserving the stylized facts (heavy tails, volatility clustering, etc.).
- Avoiding mode collapse, where the generator produces only a limited variety of samples.
- Ensuring causality and preventing look-ahead bias in generated sequences.
- Conditioning generation on external factors or regimes.

2.3.3 TimeGAN: Architecture and Mathematical Formulation

TimeGAN [Yoon et al., 2019] is a GAN framework specifically designed for generating realistic time series data. It incorporates an autoencoder structure to learn a latent representation and uses both supervised and unsupervised losses.

Components:

1. **Embedding Network (E)**: Maps the input time series $\mathbf{X} = (x_1, \dots, x_T)$ to a latent sequence $\mathbf{H} = (h_1, \dots, h_T)$.
2. **Recovery Network (R)**: Reconstructs the original time series from the latent sequence $\hat{\mathbf{X}} = R(\mathbf{H})$.
3. **Generator Network (G)**: Generates a synthetic latent sequence $\hat{\mathbf{H}} = (\hat{h}_1, \dots, \hat{h}_T)$ from random noise vectors $\mathbf{Z} = (z_1, \dots, z_T)$.
4. **Discriminator Network (D)**: Distinguishes between real latent sequences \mathbf{H} and generated ones $\hat{\mathbf{H}}$.
5. **Supervisor Network (S)**: Directly supervises the generator by predicting the next latent vector \hat{h}_t from the previous one h_{t-1} .

Loss Functions:

- **Reconstruction Loss** (for E and R): Minimizes the difference between original and reconstructed series.

$$\mathcal{L}_{Rec} = \mathbb{E}_{\mathbf{X}}[||\mathbf{X} - R(E(\mathbf{X}))||^2] \quad (22)$$

- **Unsupervised Loss** (for G and D): Standard GAN loss applied in the latent space.

$$\mathcal{L}_{U_{ns}} = \mathbb{E}_{\mathbf{H}}[\log D(\mathbf{H})] + \mathbb{E}_{\mathbf{Z}}[\log(1 - D(G(\mathbf{Z})))] \quad (23)$$

- **Supervised Loss** (for G and S): Encourages the generator to capture step-wise transitions.

$$\mathcal{L}_{S_{up}} = \mathbb{E}_{\mathbf{H}}[||\mathbf{H} - S(G(\mathbf{Z}))||^2] \quad (24)$$

The networks are trained jointly to optimize these losses. Generated time series are obtained by $\hat{\mathbf{X}} = R(G(\mathbf{Z}))$.

2.3.4 Conditional Variational Autoencoders: Theory and Implementation

Conditional VAEs (CVAEs) extend VAEs by conditioning the generation process on additional information c (e.g., market regime label). The encoder $q_\phi(z|x, c)$ and decoder $p_\theta(x|z, c)$ both take the condition c as input.

The ELBO for CVAE is:

$$\mathcal{L}(\theta, \phi; x, c) = \mathbb{E}_{q_\phi(z|x, c)}[\log p_\theta(x|z, c)] - D_{KL}(q_\phi(z|x, c)||p(z|c)) \quad (25)$$

where $p(z|c)$ is the conditional prior, often assumed to be a standard Gaussian $\mathcal{N}(0, I)$ independent of c , simplifying the KL term.

Implementation typically involves concatenating the condition c (e.g., one-hot encoded regime label) to the input of the encoder and the latent variable z before feeding it to the decoder.

2.3.5 FiLM Transformer: Feature-wise Linear Modulation

Transformers [Vaswani et al., 2017] use self-attention mechanisms to capture long-range dependencies effectively. Feature-wise Linear Modulation (FiLM) [Perez et al., 2018] is a technique for conditioning neural networks.

A FiLM layer takes a feature map h and a conditioning vector c and computes:

$$\text{FiLM}(h, c) = \gamma(c) \odot h + \beta(c) \quad (26)$$

where $\gamma(c)$ and $\beta(c)$ are scale and shift parameters generated from the condition c by separate neural networks, and \odot denotes element-wise multiplication.

In a FiLM Transformer, FiLM layers can be incorporated into the Transformer architecture (e.g., after attention or feed-forward layers) to modulate the activations based on the condition (e.g., market regime). This allows the model to adapt its processing dynamically based on the regime context.

2.3.6 Regime Adversary Module: Ensuring Regime-Specific Properties

To ensure that generated data accurately reflects the characteristics of the target regime, an adversarial component can be added. A regime classifier (adversary) is trained to predict the regime label from the generated data. The generator is then trained to fool this classifier, encouraging it to produce data that is indistinguishable from real data belonging to the target regime.

The generator's loss includes an adversarial term against the regime classifier D_{reg} :

$$\mathcal{L}_{\text{RegAdv}} = -\mathbb{E}_{z \sim p_z(z), c \sim p_c} [\log D_{\text{reg}}(G(z, c) | c)] \quad (27)$$

This forces the generator G to produce samples that D_{reg} classifies as belonging to the intended regime c . Feature matching losses or gradient penalties (like WGAN-GP) can be used to stabilize training.

2.4 Algorithmic Trading Strategies

Algorithmic trading involves using computer programs to execute trading strategies based on pre-defined rules or models.

2.4.1 Momentum and Mean Reversion

Two fundamental concepts in trading strategies are momentum and mean reversion.

Momentum strategies assume that assets that have performed well recently will continue to perform well (and vice versa). A simple momentum signal can be:

$$M_t(w) = p_t - p_{t-w} \quad \text{or} \quad M_t(w) = \frac{p_t}{p_{t-w}} - 1 \quad (28)$$

where w is the lookback window. Positions are taken in the direction of the trend.

Mean Reversion strategies assume that prices tend to revert to their historical average or some equilibrium level. When the price deviates significantly from the mean, a position is taken anticipating a return to the mean. This is often modeled using Ornstein-Uhlenbeck processes:

$$dp_t = \theta(\mu - p_t)dt + \sigma dW_t \quad (29)$$

where μ is the long-term mean and θ is the speed of reversion.

2.4.2 Transformer-based Momentum Strategy

Transformers can capture complex patterns and long-range dependencies in financial data, potentially leading to more robust momentum signals. The self-attention mechanism allows the model to weigh the importance of different past time steps when predicting future price movements.

Input features (e.g., past returns, technical indicators) are fed into a Transformer encoder. The output representation can then be used to predict future returns or generate buy/sell signals.

2.4.3 LSTM Short Signal Strategy

Long Short-Term Memory (LSTM) networks are a type of Recurrent Neural Network (RNN) well-suited for sequence modeling. They use gating mechanisms (input, forget, output gates) to control the flow of information and maintain a memory state, allowing them to capture temporal dependencies.

$$f_t = \sigma_g(W_f x_t + U_f h_{t-1} + b_f) \quad (\text{Forget Gate}) \quad (30)$$

$$i_t = \sigma_g(W_i x_t + U_i h_{t-1} + b_i) \quad (\text{Input Gate}) \quad (31)$$

$$o_t = \sigma_g(W_o x_t + U_o h_{t-1} + b_o) \quad (\text{Output Gate}) \quad (32)$$

$$\tilde{c}_t = \sigma_c(W_c x_t + U_c h_{t-1} + b_c) \quad (\text{Candidate Cell State}) \quad (33)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \quad (\text{Cell State}) \quad (34)$$

$$h_t = o_t \odot \sigma_h(c_t) \quad (\text{Hidden State}) \quad (35)$$

An LSTM can be trained on historical features to predict downward price movements or generate signals specifically for short selling opportunities.

2.4.4 Kalman Filter Pairs Trading

Pairs trading is a market-neutral strategy that exploits temporary mispricings between two related assets (a pair). It assumes that the spread (difference or ratio) between their prices is mean-reverting.

Let p_t^A and p_t^B be the log prices of two cointegrated assets. The spread is often modeled as $s_t = p_t^A - \beta_t p_t^B$, where β_t is the dynamic hedge ratio.

The Kalman filter [Kalman, 1960] provides a way to estimate the hidden state (e.g., the hedge ratio β_t) of a system from noisy observations. It uses a state-space model:

$$\text{State Equation: } \beta_t = F_t \beta_{t-1} + w_t, \quad w_t \sim \mathcal{N}(0, Q_t) \quad (36)$$

$$\text{Observation Equation: } y_t = H_t \beta_t + v_t, \quad v_t \sim \mathcal{N}(0, R_t) \quad (37)$$

In pairs trading, y_t could be p_t^A and H_t could be p_t^B . The filter recursively updates the estimate of β_t and its uncertainty. Trading signals are generated when the observed spread deviates significantly from its expected value based on the estimated β_t .

2.4.5 Strategy Decay Detection and Mitigation

Trading strategies often lose their effectiveness over time (alpha decay) due to changing market dynamics or overcrowding. Detecting decay involves monitoring rolling performance metrics (e.g., Sharpe ratio, win rate) and statistical tests for parameter stability (e.g., Chow test). Mitigation can involve retraining the model, adjusting parameters, or switching to a different strategy, potentially using synthetic data generated for the current regime to adapt quickly.

2.5 Risk Forecasting Models

Risk forecasting aims to quantify potential future losses.

2.5.1 Value at Risk (VaR) Methodologies

Value at Risk (VaR) is a standard measure of downside risk. VaR_α is the maximum potential loss over a given horizon at a specific confidence level ($1 - \alpha$).

$$P(\text{Loss} > \text{VaR}_\alpha) = \alpha \quad (38)$$

Common VaR calculation methods include:

- **Historical Simulation:** Uses the empirical distribution of past returns.
- **Parametric (Variance-Covariance):** Assumes returns follow a specific distribution (e.g., Gaussian) and estimates VaR from the distribution's parameters.
- **Monte Carlo Simulation:** Simulates future price paths based on a stochastic model.

2.5.2 REGARCH-EVT: Regime-Switching GARCH with Extreme Value Theory

This approach combines regime-switching models, GARCH models for volatility dynamics, and Extreme Value Theory (EVT) for tail modeling.

GARCH (Generalized Autoregressive Conditional Heteroskedasticity) [Bollerslev, 1986] models time-varying volatility. A GARCH(p, q) model is:

$$r_t = \mu_t + \epsilon_t \quad (39)$$

$$\epsilon_t = \sigma_t z_t, \quad z_t \sim \text{i.i.d.}, \mathbb{E}[z_t] = 0, \text{Var}(z_t) = 1 \quad (40)$$

$$\sigma_t^2 = \omega + \sum_{i=1}^q \alpha_i \epsilon_{t-i}^2 + \sum_{j=1}^p \beta_j \sigma_{t-j}^2 \quad (41)$$

Regime-Switching GARCH (REGARCH) allows GARCH parameters ($\omega, \alpha_i, \beta_j$) to vary depending on the market regime z_t .

Extreme Value Theory (EVT) provides a framework for modeling the tails of probability distributions. The Peaks-Over-Threshold (POT) approach models exceedances $y = x - u$ over a high threshold u using the Generalized Pareto Distribution (GPD):

$$F_u(y) = P(X - u \leq y | X > u) \approx G(y; \xi, \beta) = 1 - \left(1 + \xi \frac{y}{\beta}\right)^{-1/\xi} \quad (42)$$

where ξ is the shape parameter and β is the scale parameter. REGARCH-EVT applies EVT to the standardized residuals ($\epsilon_t / \sigma_{t,k}$) within each regime k to model tail risk more accurately.

2.5.3 Quantile Random Forest for Direct VaR Estimation

Random Forests [Breiman, 2001] are ensemble learning methods based on decision trees. Quantile Regression Forests [Meinshausen, 2006] extend Random Forests to estimate conditional quantiles directly, rather than just the conditional mean.

Instead of averaging the predictions of individual trees (as in standard regression forests), Quantile Random Forests consider the distribution of predictions across all trees

in the forest. For a new input point x , the estimate of the conditional cumulative distribution function $\hat{F}(y|X = x)$ is obtained by averaging the empirical distributions within the leaves where x falls. The α -quantile (VaR) is then estimated as:

$$\widehat{\text{VaR}}_\alpha(x) = \inf\{y : \hat{F}(y|X = x) \geq \alpha\} \quad (43)$$

This non-parametric approach can capture complex relationships and does not rely on specific distributional assumptions, making it suitable for estimating VaR directly from features.

2.5.4 Synthetic Regime Bootstrapping

Bootstrapping involves resampling from the observed data to estimate the sampling distribution of a statistic. In the context of regime-dependent risk, synthetic data generated by GenAI models for specific regimes can be incorporated into the bootstrapping process.

Synthetic Regime Bootstrapping involves:

1. Identifying the current market regime \mathcal{R}_k .
2. Generating synthetic data sequences $\mathbf{X}_{\text{synth}}^{(k)}$ representative of regime k .
3. Combining historical data from regime k with synthetic data $\mathbf{X}_{\text{synth}}^{(k)}$.
4. Bootstrapping (resampling with replacement) from this combined dataset.
5. Calculating the statistic of interest (e.g., portfolio return distribution, VaR) from each bootstrap sample.
6. Constructing confidence intervals or estimating risk measures from the distribution of bootstrapped statistics.

This approach leverages generative models to enrich the data available for risk estimation, particularly for regimes with limited historical observations.

References

- A. Ang and A. Timmermann. Regime changes and financial markets. *Annual Review of Financial Economics*, 4(1):313–337, 2012.
- T. Bollerslev. Generalized autoregressive conditional heteroskedasticity. *Journal of Econometrics*, 31(3):307–327, 1986.
- L. Breiman. Random forests. *Machine Learning*, 45(1):5–32, 2001.
- R. Cont. Empirical properties of asset returns: stylized facts and statistical issues. *Quantitative Finance*, 1(2):223–236, 2001.
- T. G. Dietterich. Ensemble methods in machine learning. *International workshop on multiple classifier systems*, pages 1–15, 2000.
- I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. *Advances in Neural Information Processing Systems*, pages 2672–2680, 2014.

- I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. Courville. Improved training of wasserstein gans. *Advances in Neural Information Processing Systems*, pages 5767–5777, 2017.
- J. D. Hamilton. A new approach to the economic analysis of nonstationary time series and the business cycle. *Econometrica*, 57(2):357–384, 1989.
- I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. *International Conference on Learning Representations*, 2017.
- R. E. Kalman. A new approach to linear filtering and prediction problems. *Journal of Basic Engineering*, 82(1):35–45, 1960.
- D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- A. W. Lo. The adaptive markets hypothesis: Market efficiency from an evolutionary perspective. *Journal of Portfolio Management*, 30(5):15–29, 2004.
- L. McInnes, J. Healy, and J. Melville. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
- N. Meinshausen. Quantile regression forests. *Journal of Machine Learning Research*, 7:983–999, 2006.
- E. Perez, F. Strub, H. De Vries, V. Dumoulin, and A. Courville. Film: Visual reasoning with a general conditioning layer. *AAAI Conference on Artificial Intelligence*, 2018.
- A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, pages 5998–6008, 2017.
- J. Yoon, D. Jarrett, and M. van der Schaar. Time-series generative adversarial networks. *Advances in Neural Information Processing Systems*, pages 5508–5518, 2019.

3 Data Preprocessing Methodology

This section details the comprehensive preprocessing pipeline developed for financial time series data. Proper preprocessing is crucial for the success of subsequent analyses, particularly in the context of market regime detection and algorithmic trading.

3.1 Data Sources and Description

Our analysis utilizes a diverse set of financial assets to capture various market dynamics and ensure robustness. The dataset comprises:

- **Equity Indices:** Major market indices such as S&P 500 (SPY), Nasdaq-100 (QQQ), Russell 2000 (IWM), and international indices.
- **Individual Stocks:** Selected from various sectors and market capitalizations.

- **Fixed Income:** Treasury bonds (TLT, IEF, SHY) and corporate bonds (LQD, HYG).
- **Commodities:** Energy (XLE), precious metals (GLD, SLV), and agricultural commodities.
- **Volatility:** VIX-related products (VXX).
- **Currency:** Major forex pairs and cryptocurrency instruments.

The data spans multiple years, capturing various market conditions including bull markets, bear markets, and periods of high volatility. Daily price data (open, high, low, close, volume) was collected, with a focus on liquid instruments to minimize the impact of market microstructure effects.

3.2 Data Cleaning and Handling Missing Values

Financial time series often contain missing values, outliers, and other anomalies that must be addressed before analysis. Our preprocessing pipeline implements several techniques to ensure data quality.

3.2.1 Forward Fill Method for Time Series

Missing values in financial time series typically occur due to non-trading days, data collection issues, or liquidity gaps. The forward fill method is particularly appropriate for financial time series because it preserves the information available at each point in time, preventing look-ahead bias.

For a time series $\{x_t\}_{t=1}^T$ with missing values, the forward fill operation is defined as:

$$\hat{x}_t = \begin{cases} x_t & \text{if } x_t \text{ is observed} \\ \hat{x}_{t-1} & \text{if } x_t \text{ is missing} \end{cases} \quad (44)$$

This approach assumes that the most recent available information is the best estimate for the current value until new information becomes available. For the initial period, if x_1 is missing, we either drop the observation or use the first available value.

The implementation includes checks for leading missing values, which could indicate data quality issues at the beginning of the series:

$$\text{Leading Missing} = \mathbb{I}(x_1 \text{ is missing}) \quad (45)$$

where $\mathbb{I}(\cdot)$ is the indicator function. If leading missing values are detected, we log a warning and consider trimming the series to start from the first valid observation.

3.2.2 Outlier Detection and Treatment

Outliers in financial data can significantly impact statistical analyses and model performance. We implement a time-series-aware outlier detection approach that prevents look-ahead bias by using expanding windows.

Z-score Method: For each time point t with sufficient history (at least `min_periods`), we calculate the z-score using the mean and standard deviation of the preceding observations:

$$\mu_t = \frac{1}{t-1} \sum_{i=1}^{t-1} r_i \quad (46)$$

$$\sigma_t = \sqrt{\frac{1}{t-1} \sum_{i=1}^{t-1} (r_i - \mu_t)^2} \quad (47)$$

$$z_t = \frac{r_t - \mu_t}{\sigma_t} \quad (48)$$

where r_t is the return at time t . A point is flagged as an outlier if $|z_t| > \text{threshold}$, where threshold is typically set to 3.0.

IQR Method: Alternatively, we can use the interquartile range (IQR) method, which is more robust to extreme values:

$$Q1_t = 25\text{th percentile of } \{r_i\}_{i=1}^{t-1} \quad (49)$$

$$Q3_t = 75\text{th percentile of } \{r_i\}_{i=1}^{t-1} \quad (50)$$

$$\text{IQR}_t = Q3_t - Q1_t \quad (51)$$

$$\text{Lower Bound}_t = Q1_t - \text{threshold} \times \text{IQR}_t \quad (52)$$

$$\text{Upper Bound}_t = Q3_t + \text{threshold} \times \text{IQR}_t \quad (53)$$

A point is flagged as an outlier if $r_t < \text{Lower Bound}_t$ or $r_t > \text{Upper Bound}_t$.

Both methods are implemented using vectorized operations for computational efficiency, with expanding windows to ensure no look-ahead bias is introduced.

3.2.3 Winsorization and Clipping Techniques

Once outliers are detected, they must be handled appropriately. We implement two main approaches:

Winsorization: This technique replaces extreme values with less extreme values at specified percentiles. For each time point t , we calculate the expanding quantiles and replace outliers with the corresponding bounds:

$$\text{Lower Quantile}_t = \alpha/2 \text{ percentile of } \{r_i\}_{i=1}^{t-1} \quad (54)$$

$$\text{Upper Quantile}_t = 1 - \alpha/2 \text{ percentile of } \{r_i\}_{i=1}^{t-1} \quad (55)$$

$$\hat{r}_t = \begin{cases} \text{Lower Quantile}_t & \text{if } r_t < \text{Lower Quantile}_t \\ \text{Upper Quantile}_t & \text{if } r_t > \text{Upper Quantile}_t \\ r_t & \text{otherwise} \end{cases} \quad (56)$$

where α is typically set to 0.02, corresponding to the 1st and 99th percentiles.

Clipping: This simpler approach clips values based on the z-score thresholds:

$$\text{Lower Bound}_t = \mu_t - \text{threshold} \times \sigma_t \quad (57)$$

$$\text{Upper Bound}_t = \mu_t + \text{threshold} \times \sigma_t \quad (58)$$

$$\hat{r}_t = \begin{cases} \text{Lower Bound}_t & \text{if } r_t < \text{Lower Bound}_t \\ \text{Upper Bound}_t & \text{if } r_t > \text{Upper Bound}_t \\ r_t & \text{otherwise} \end{cases} \quad (59)$$

After handling outliers in the return space, we reconstruct the price series to maintain consistency:

$$\hat{p}_t = \hat{p}_{t-1} \times (1 + \hat{r}_t) \quad (60)$$

starting from the first valid price $\hat{p}_1 = p_1$.

3.3 Time Series Alignment

Financial assets often have different trading calendars due to market holidays, exchange-specific closures, or instrument-specific trading hours. Proper alignment is crucial for multivariate analysis.

3.3.1 Union vs. Intersection Strategies

Two main strategies exist for aligning multiple time series:

Intersection Strategy: Keeps only the timestamps that are common to all series. This ensures that all assets have valid data for each included timestamp, but may result in significant data loss.

$$\mathcal{T}_{\text{intersection}} = \bigcap_{i=1}^n \mathcal{T}_i \quad (61)$$

where \mathcal{T}_i is the set of timestamps for asset i .

Union Strategy: Keeps all timestamps that appear in any of the series, filling missing values for assets that don't have data at certain timestamps.

$$\mathcal{T}_{\text{union}} = \bigcup_{i=1}^n \mathcal{T}_i \quad (62)$$

The choice between these strategies involves a trade-off between data completeness and the potential introduction of biases from imputation.

3.3.2 Implications of Alignment Choices

The alignment strategy has significant implications for subsequent analyses:

- **Data Coverage:** Intersection reduces the number of observations, potentially limiting the statistical power of analyses, especially for regime detection where certain regimes might be underrepresented.

- **Imputation Bias:** Union requires imputation of missing values, which can introduce biases, especially if the missing pattern is not random (e.g., market-specific holidays).
- **Correlation Estimation:** Imputed values can artificially inflate correlations if the imputation method doesn't adequately capture the true data-generating process.
- **Volatility Estimation:** Imputation can lead to underestimation of volatility by smoothing out potential jumps that would have occurred after non-trading periods.

Our implementation allows for both strategies, with a preference for intersection in cases where sufficient data remains, and union with careful imputation otherwise.

3.3.3 Preventing Data Leakage in Financial Time Series

Data leakage occurs when information from the future is inadvertently used to make predictions about the past or present. In financial time series preprocessing, several precautions are taken:

- **Expanding Window Statistics:** All statistics (means, standard deviations, quantiles) are calculated using expanding windows that include only past data.
- **Forward-Only Imputation:** Missing values are imputed using only past information (e.g., forward fill from the last available value).
- **Proper Indexing:** Time indices are carefully maintained to ensure that data from time t is only used for decisions at time t or later.
- **Feature Engineering Lag:** Features are engineered with appropriate lags to ensure they only use information available at the time of prediction.

The mathematical formulation of the expanding window approach ensures that at each time t , only information from times $\{1, 2, \dots, t-1\}$ is used to process the observation at time t .

3.4 Feature Engineering

Feature engineering transforms raw financial data into meaningful inputs for machine learning models. We implement a comprehensive set of features capturing various aspects of market behavior.

3.4.1 Price-based Features

Price-based features capture the absolute and relative levels of asset prices:

- **Price Levels:** Normalized price levels relative to a reference period.

$$\text{NormPrice}_t = \frac{p_t}{p_{\text{ref}}} \quad (63)$$

- **Price Ratios:** Ratios between different assets or between the same asset at different timeframes.

$$\text{PriceRatio}_t^{A,B} = \frac{p_t^A}{p_t^B} \quad (64)$$

- **Moving Averages:** Simple and exponential moving averages of prices.

$$\text{SMA}_t(w) = \frac{1}{w} \sum_{i=t-w+1}^t p_i \quad (65)$$

$$\text{EMA}_t(\alpha) = \alpha p_t + (1 - \alpha) \text{EMA}_{t-1}(\alpha) \quad (66)$$

where w is the window size and α is the smoothing factor.

- **Price Channels:** Upper and lower bounds based on recent price history.

$$\text{UpperChannel}_t(w) = \max_{i \in [t-w+1, t]} p_i \quad (67)$$

$$\text{LowerChannel}_t(w) = \min_{i \in [t-w+1, t]} p_i \quad (68)$$

3.4.2 Return-based Features

Return-based features capture the dynamics of price changes:

- **Simple Returns:** Percentage change in price.

$$R_t = \frac{p_t - p_{t-1}}{p_{t-1}} \quad (69)$$

- **Log Returns:** Natural logarithm of the price ratio.

$$r_t = \log \left(\frac{p_t}{p_{t-1}} \right) \quad (70)$$

- **Cumulative Returns:** Sum of returns over a specified window.

$$\text{CumReturn}_t(w) = \sum_{i=t-w+1}^t r_i \quad (71)$$

- **Excess Returns:** Returns in excess of a benchmark or risk-free rate.

$$\text{ExcessReturn}_t = r_t - r_t^{\text{benchmark}} \quad (72)$$

3.4.3 Volatility Measures

Volatility features capture the dispersion and uncertainty in price movements:

- **Historical Volatility:** Standard deviation of returns over a specified window.

$$\sigma_t(w) = \sqrt{\frac{1}{w-1} \sum_{i=t-w+1}^t (r_i - \bar{r}_t(w))^2} \quad (73)$$

where $\bar{r}_t(w) = \frac{1}{w} \sum_{i=t-w+1}^t r_i$.

- **EWMA Volatility:** Exponentially weighted moving average of squared returns.

$$\sigma_t^2(\lambda) = \lambda\sigma_{t-1}^2(\lambda) + (1 - \lambda)r_{t-1}^2 \quad (74)$$

where λ is the decay factor (typically 0.94 for daily data).

- **Parkinson Volatility:** Based on high-low range, capturing intraday volatility.

$$\sigma_t^{\text{Park}} = \frac{1}{4\log(2)} \log \left(\frac{H_t}{L_t} \right)^2 \quad (75)$$

where H_t and L_t are the high and low prices at time t .

- **Garman-Klass Volatility:** Incorporates open, high, low, and close prices.

$$\sigma_t^{\text{GK}} = \sqrt{\frac{1}{2} \left(\log \frac{H_t}{L_t} \right)^2 - (2\log(2) - 1) \left(\log \frac{C_t}{O_t} \right)^2} \quad (76)$$

where O_t and C_t are the open and close prices at time t .

3.4.4 Technical Indicators

Technical indicators are mathematical calculations based on price, volume, or open interest used to forecast future price movements:

- **Relative Strength Index (RSI):** Measures the speed and change of price movements.

$$\text{RS}_t(w) = \frac{\text{AvgGain}_t(w)}{\text{AvgLoss}_t(w)} \quad (77)$$

$$\text{RSI}_t(w) = 100 - \frac{100}{1 + \text{RS}_t(w)} \quad (78)$$

where $\text{AvgGain}_t(w)$ and $\text{AvgLoss}_t(w)$ are the average gains and losses over the past w periods.

- **Moving Average Convergence Divergence (MACD):** Trend-following momentum indicator.

$$\text{MACD}_t = \text{EMA}_t(12) - \text{EMA}_t(26) \quad (79)$$

$$\text{Signal}_t = \text{EMA}_t(9, \text{MACD}) \quad (80)$$

$$\text{Histogram}_t = \text{MACD}_t - \text{Signal}_t \quad (81)$$

- **Bollinger Bands:** Volatility bands placed above and below a moving average.

$$\text{MiddleBand}_t(w) = \text{SMA}_t(w) \quad (82)$$

$$\text{UpperBand}_t(w, k) = \text{MiddleBand}_t(w) + k \cdot \sigma_t(w) \quad (83)$$

$$\text{LowerBand}_t(w, k) = \text{MiddleBand}_t(w) - k \cdot \sigma_t(w) \quad (84)$$

where k is typically set to 2.

- **Average True Range (ATR):** Measures market volatility.

$$\text{TR}_t = \max(H_t - L_t, |H_t - C_{t-1}|, |L_t - C_{t-1}|) \quad (85)$$

$$\text{ATR}_t(w) = \frac{1}{w} \sum_{i=t-w+1}^t \text{TR}_i \quad (86)$$

3.4.5 Cross-asset Features

Cross-asset features capture relationships between different assets:

- **Correlation:** Rolling correlation between assets.

$$\rho_t^{A,B}(w) = \frac{\sum_{i=t-w+1}^t (r_i^A - \bar{r}_t^A(w))(r_i^B - \bar{r}_t^B(w))}{\sqrt{\sum_{i=t-w+1}^t (r_i^A - \bar{r}_t^A(w))^2} \sqrt{\sum_{i=t-w+1}^t (r_i^B - \bar{r}_t^B(w))^2}} \quad (87)$$

- **Spread:** Difference between related assets.

$$\text{Spread}_t^{A,B} = p_t^A - p_t^B \quad (88)$$

- **Ratio:** Ratio between related assets.

$$\text{Ratio}_t^{A,B} = \frac{p_t^A}{p_t^B} \quad (89)$$

- **Beta:** Sensitivity of an asset to a benchmark.

$$\beta_t^A(w) = \frac{\text{Cov}_t(r^A, r^{\text{benchmark}}, w)}{\text{Var}_t(r^{\text{benchmark}}, w)} \quad (90)$$

3.4.6 Feature Selection and Dimensionality Considerations

The feature engineering process can generate a large number of features, potentially leading to the curse of dimensionality. We implement several approaches to manage dimensionality:

- **Variance-based Selection:** Retain features with variance above a threshold.

$$\text{Selected Features} = \{f_i : \text{Var}(f_i) > \text{threshold}\} \quad (91)$$

- **Correlation-based Selection:** Remove highly correlated features.

$$\text{Remove } f_j \text{ if } |\rho(f_i, f_j)| > \text{threshold} \text{ for some } f_i \text{ already selected} \quad (92)$$

- **Principal Component Analysis (PCA):** Transform features to uncorrelated principal components.

$$\mathbf{Z} = \mathbf{XW} \quad (93)$$

where \mathbf{W} contains the eigenvectors of the covariance matrix of \mathbf{X} .

- **Domain Knowledge:** Select features based on financial theory and empirical evidence.

Our implementation allows for flexible feature selection based on the specific requirements of each analysis component.

3.5 Mathematical Formulation of Key Features

This section provides detailed mathematical formulations for key features used in our analysis.

3.5.1 Moving Averages and Exponential Smoothing

Moving averages smooth price series to identify trends and filter out noise. The Simple Moving Average (SMA) at time t with window size w is:

$$\text{SMA}_t(w) = \frac{1}{w} \sum_{i=t-w+1}^t p_i \quad (94)$$

The Exponential Moving Average (EMA) gives more weight to recent observations:

$$\text{EMA}_t(\alpha) = \begin{cases} p_1 & \text{if } t = 1 \\ \alpha p_t + (1 - \alpha) \text{EMA}_{t-1}(\alpha) & \text{if } t > 1 \end{cases} \quad (95)$$

where $\alpha \in (0, 1)$ is the smoothing factor. The relationship between α and a comparable window size w is approximately $\alpha = \frac{2}{w+1}$.

The EMA can also be expressed as a weighted sum of all past observations:

$$\text{EMA}_t(\alpha) = \alpha \sum_{i=0}^{t-1} (1 - \alpha)^i p_{t-i} \quad (96)$$

3.5.2 Bollinger Bands and Volatility Channels

Bollinger Bands consist of a middle band (typically an SMA) and upper and lower bands placed k standard deviations away:

$$\text{MiddleBand}_t(w) = \text{SMA}_t(w) \quad (97)$$

$$\text{UpperBand}_t(w, k) = \text{MiddleBand}_t(w) + k \cdot \sigma_t(w) \quad (98)$$

$$\text{LowerBand}_t(w, k) = \text{MiddleBand}_t(w) - k \cdot \sigma_t(w) \quad (99)$$

where $\sigma_t(w)$ is the standard deviation of prices over the past w periods:

$$\sigma_t(w) = \sqrt{\frac{1}{w} \sum_{i=t-w+1}^t (p_i - \text{SMA}_t(w))^2} \quad (100)$$

The Bollinger Band Width (BBW) measures the volatility relative to the price level:

$$\text{BBW}_t(w, k) = \frac{\text{UpperBand}_t(w, k) - \text{LowerBand}_t(w, k)}{\text{MiddleBand}_t(w)} \quad (101)$$

3.5.3 Momentum Indicators (RSI, MACD)

The Relative Strength Index (RSI) measures the magnitude of recent price changes to evaluate overbought or oversold conditions:

$$\text{AvgGain}_t(w) = \frac{1}{w} \sum_{i=t-w+1}^t \max(r_i, 0) \quad (102)$$

$$\text{AvgLoss}_t(w) = \frac{1}{w} \sum_{i=t-w+1}^t \max(-r_i, 0) \quad (103)$$

$$\text{RS}_t(w) = \frac{\text{AvgGain}_t(w)}{\text{AvgLoss}_t(w)} \quad (104)$$

$$\text{RSI}_t(w) = 100 - \frac{100}{1 + \text{RS}_t(w)} \quad (105)$$

The Moving Average Convergence Divergence (MACD) is a trend-following momentum indicator:

$$\text{MACD}_t = \text{EMA}_t(\alpha_1) - \text{EMA}_t(\alpha_2) \quad (106)$$

$$\text{Signal}_t = \text{EMA}_t(\alpha_3, \text{MACD}) \quad (107)$$

$$\text{Histogram}_t = \text{MACD}_t - \text{Signal}_t \quad (108)$$

where $\alpha_1 > \alpha_2$ (typically corresponding to 12-day and 26-day EMAs) and α_3 corresponds to a 9-day EMA.

3.5.4 Statistical Features (Kurtosis, Skewness)

Kurtosis measures the "tailedness" of the probability distribution of returns:

$$\text{Kurt}_t(w) = \frac{\frac{1}{w} \sum_{i=t-w+1}^t (r_i - \bar{r}_t(w))^4}{\left(\frac{1}{w} \sum_{i=t-w+1}^t (r_i - \bar{r}_t(w))^2 \right)^2} \quad (109)$$

Skewness measures the asymmetry of the probability distribution of returns:

$$\text{Skew}_t(w) = \frac{\frac{1}{w} \sum_{i=t-w+1}^t (r_i - \bar{r}_t(w))^3}{\left(\frac{1}{w} \sum_{i=t-w+1}^t (r_i - \bar{r}_t(w))^2 \right)^{3/2}} \quad (110)$$

These higher-order moments provide valuable information about the distribution of returns beyond the mean and variance, particularly for capturing tail risk and asymmetry.

4 Market Regime Detection Implementation

This section details the implementation of our market regime detection component, which employs multiple methods to identify distinct market states. We present the mathematical formulation, parameter estimation techniques, implementation details, and regime characterization for each approach.

4.1 Hidden Markov Model (HMM)

Hidden Markov Models provide a powerful framework for detecting unobservable market regimes that evolve according to a Markov process.

4.1.1 Mathematical Formulation

A Hidden Markov Model for financial market regimes consists of:

- A set of K hidden states (regimes) $\mathcal{S} = \{s_1, s_2, \dots, s_K\}$, where s_i represents a specific market regime.
- A transition probability matrix $A = [a_{ij}]$, where $a_{ij} = P(z_t = s_j | z_{t-1} = s_i)$ represents the probability of transitioning from regime i to regime j .
- An initial state distribution $\pi = [\pi_i]$, where $\pi_i = P(z_1 = s_i)$ represents the probability of starting in regime i .
- A set of emission probability distributions $B = \{b_i(o_t)\}$, where $b_i(o_t) = P(o_t | z_t = s_i)$ represents the probability of observing o_t given that the current regime is s_i .

For financial time series, we typically model the emission probabilities as multivariate Gaussian distributions:

$$b_i(o_t) = \mathcal{N}(o_t | \mu_i, \Sigma_i) = \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \exp \left(-\frac{1}{2} (o_t - \mu_i)^T \Sigma_i^{-1} (o_t - \mu_i) \right) \quad (111)$$

where d is the dimensionality of the observation vector o_t , μ_i is the mean vector for regime i , and Σ_i is the covariance matrix for regime i .

The joint probability of a sequence of observations $O = \{o_1, o_2, \dots, o_T\}$ and a sequence of hidden states $Z = \{z_1, z_2, \dots, z_T\}$ is:

$$P(O, Z | \lambda) = \pi_{z_1} b_{z_1}(o_1) \prod_{t=2}^T a_{z_{t-1} z_t} b_{z_t}(o_t) \quad (112)$$

where $\lambda = (A, B, \pi)$ represents the model parameters.

4.1.2 Parameter Estimation via Expectation-Maximization

The parameters of the HMM are estimated using the Baum-Welch algorithm, which is an instance of the Expectation-Maximization (EM) algorithm. The algorithm iteratively improves the model parameters to maximize the likelihood of the observed data.

E-step: Compute the expected counts of state transitions and state occupations given the current model parameters.

First, we compute the forward probabilities $\alpha_t(i) = P(o_1, o_2, \dots, o_t, z_t = s_i | \lambda)$:

$$\alpha_1(i) = \pi_i b_i(o_1) \quad (113)$$

$$\alpha_{t+1}(j) = \left[\sum_{i=1}^K \alpha_t(i) a_{ij} \right] b_j(o_{t+1}) \quad (114)$$

Next, we compute the backward probabilities $\beta_t(i) = P(o_{t+1}, o_{t+2}, \dots, o_T | z_t = s_i, \lambda)$:

$$\beta_T(i) = 1 \quad (115)$$

$$\beta_t(i) = \sum_{j=1}^K a_{ij} b_j(o_{t+1}) \beta_{t+1}(j) \quad (116)$$

Using these, we compute the posterior probabilities:

$$\gamma_t(i) = P(z_t = s_i | O, \lambda) = \frac{\alpha_t(i) \beta_t(i)}{\sum_{j=1}^K \alpha_t(j) \beta_t(j)} \quad (117)$$

$$\xi_t(i, j) = P(z_t = s_i, z_{t+1} = s_j | O, \lambda) = \frac{\alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}{\sum_{i=1}^K \sum_{j=1}^K \alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)} \quad (118)$$

M-step: Update the model parameters based on the expected counts:

$$\pi_i^{\text{new}} = \gamma_1(i) \quad (119)$$

$$a_{ij}^{\text{new}} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)} \quad (120)$$

$$\mu_i^{\text{new}} = \frac{\sum_{t=1}^T \gamma_t(i) o_t}{\sum_{t=1}^T \gamma_t(i)} \quad (121)$$

$$\Sigma_i^{\text{new}} = \frac{\sum_{t=1}^T \gamma_t(i) (o_t - \mu_i^{\text{new}})(o_t - \mu_i^{\text{new}})^T}{\sum_{t=1}^T \gamma_t(i)} \quad (122)$$

The algorithm iterates between the E-step and M-step until convergence, which is typically determined by a small change in the log-likelihood:

$$\log P(O|\lambda) = \log \sum_{i=1}^K \alpha_T(i) \quad (123)$$

4.1.3 Implementation Details

Our implementation of the HMM for market regime detection includes several key considerations:

- **Feature Selection:** We select a subset of features that are most informative for regime detection, focusing on market indicators such as volatility measures, correlation metrics, and technical indicators.
- **Scaling:** Features are standardized using the StandardScaler to ensure that all features contribute equally to the model:

$$\tilde{x}_i = \frac{x_i - \mu_i}{\sigma_i} \quad (124)$$

where μ_i and σ_i are the mean and standard deviation of feature i .

- **Initialization:** The model is initialized with K-means clustering to provide a good starting point for the EM algorithm. This helps avoid poor local optima.
- **Covariance Type:** We use a full covariance matrix for each regime, allowing for regime-specific correlation structures:

$$\Sigma_i = \begin{bmatrix} \sigma_{i,11} & \sigma_{i,12} & \cdots & \sigma_{i,1d} \\ \sigma_{i,21} & \sigma_{i,22} & \cdots & \sigma_{i,2d} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{i,d1} & \sigma_{i,d2} & \cdots & \sigma_{i,dd} \end{bmatrix} \quad (125)$$

- **Number of Regimes:** We set $K = 3$ based on financial theory and empirical evidence suggesting that markets typically operate in three main regimes: bull (uptrend), bear (downtrend), and sideways (consolidation).
- **Temporal Smoothing:** To reduce noise and prevent frequent regime switches, we apply a post-processing smoothing step using a rolling window majority vote:

$$\hat{z}_t = \text{mode}(\{z_{t-w+1}, z_{t-w+2}, \dots, z_t\}) \quad (126)$$

where w is the smoothing window size (typically 5 days).

4.1.4 Regime Transition Probability Matrix Analysis

The estimated transition probability matrix provides valuable insights into the stability and dynamics of market regimes. For a 3-regime model, the transition matrix has the form:

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \quad (127)$$

The diagonal elements a_{ii} represent the probability of remaining in regime i , which indicates the persistence of each regime. High values (close to 1) suggest stable regimes that tend to persist for extended periods.

The off-diagonal elements a_{ij} represent the probability of transitioning from regime i to regime j . These values provide insights into the typical progression of market cycles. For example, if $a_{12} > a_{13}$, it suggests that transitions from regime 1 to regime 2 are more common than transitions from regime 1 to regime 3.

The expected duration of regime i can be calculated as:

$$\mathbb{E}[\text{Duration of regime } i] = \frac{1}{1 - a_{ii}} \quad (128)$$

This provides a quantitative measure of regime persistence, which is useful for risk management and strategy development.

4.2 Gaussian Mixture Model (GMM)

Gaussian Mixture Models provide an alternative approach to regime detection by modeling the distribution of financial features as a mixture of Gaussian components.

4.2.1 Mathematical Formulation

A Gaussian Mixture Model represents the probability density function of the observed data as a weighted sum of K Gaussian components:

$$p(o_t|\theta) = \sum_{i=1}^K \pi_i \mathcal{N}(o_t|\mu_i, \Sigma_i) \quad (129)$$

where:

- π_i is the mixing coefficient (weight) of component i , with $\pi_i \geq 0$ and $\sum_{i=1}^K \pi_i = 1$.
- μ_i is the mean vector of component i .
- Σ_i is the covariance matrix of component i .
- $\theta = \{\pi_i, \mu_i, \Sigma_i\}_{i=1}^K$ represents all model parameters.

Each Gaussian component corresponds to a market regime, and the mixing coefficients represent the prior probabilities of each regime. The posterior probability that observation o_t belongs to regime i is given by:

$$P(z_t = i|o_t, \theta) = \frac{\pi_i \mathcal{N}(o_t|\mu_i, \Sigma_i)}{\sum_{j=1}^K \pi_j \mathcal{N}(o_t|\mu_j, \Sigma_j)} \quad (130)$$

This posterior probability is used to assign each observation to the most likely regime.

4.2.2 Expectation-Maximization Algorithm

The parameters of the GMM are estimated using the Expectation-Maximization (EM) algorithm, which iteratively refines the parameters to maximize the likelihood of the observed data.

E-step: Compute the responsibilities $\gamma(z_{ti})$, which represent the posterior probability that observation o_t was generated by component i :

$$\gamma(z_{ti}) = \frac{\pi_i \mathcal{N}(o_t|\mu_i, \Sigma_i)}{\sum_{j=1}^K \pi_j \mathcal{N}(o_t|\mu_j, \Sigma_j)} \quad (131)$$

M-step: Update the parameters using the computed responsibilities:

$$N_i = \sum_{t=1}^T \gamma(z_{ti}) \quad (132)$$

$$\pi_i^{\text{new}} = \frac{N_i}{T} \quad (133)$$

$$\mu_i^{\text{new}} = \frac{1}{N_i} \sum_{t=1}^T \gamma(z_{ti}) o_t \quad (134)$$

$$\Sigma_i^{\text{new}} = \frac{1}{N_i} \sum_{t=1}^T \gamma(z_{ti}) (o_t - \mu_i^{\text{new}})(o_t - \mu_i^{\text{new}})^T \quad (135)$$

The algorithm iterates between the E-step and M-step until convergence, which is typically determined by a small change in the log-likelihood:

$$\log p(O|\theta) = \sum_{t=1}^T \log \left(\sum_{i=1}^K \pi_i \mathcal{N}(o_t | \mu_i, \Sigma_i) \right) \quad (136)$$

4.2.3 Implementation Details

Our implementation of the GMM for market regime detection includes several key considerations:

- **Feature Selection:** Similar to the HMM approach, we select informative features for regime detection, but with a focus on features that capture the distributional properties of returns.
- **Scaling:** Features are standardized using the RobustScaler, which is less sensitive to outliers:

$$\tilde{x}_i = \frac{x_i - \text{median}(x_i)}{\text{IQR}(x_i)} \quad (137)$$

where $\text{IQR}(x_i)$ is the interquartile range of feature i .

- **Initialization:** The model is initialized using K-means++ to provide a good starting point for the EM algorithm.
- **Covariance Type:** We use a full covariance matrix for each component, allowing for regime-specific correlation structures. To ensure numerical stability, we add a small regularization term to the diagonal elements:

$$\Sigma_i^{\text{reg}} = \Sigma_i + \epsilon I \quad (138)$$

where ϵ is a small positive constant (typically 10^{-6}) and I is the identity matrix.

- **Number of Components:** We set $K = 3$ for consistency with the HMM approach and financial theory.
- **Temporal Consistency:** Unlike HMM, GMM does not explicitly model temporal dependencies. To incorporate temporal information, we apply a post-processing step using a Hidden Markov Model with the GMM component assignments as observations.

4.2.4 Component Distribution Analysis

The estimated parameters of each Gaussian component provide insights into the statistical properties of the corresponding market regime. For each component i , we analyze:

- **Mean Vector μ_i :** Indicates the average feature values in regime i . For example, a high mean return and low mean volatility would characterize a bullish regime.
- **Covariance Matrix Σ_i :** Captures the variability and co-movement of features in regime i . The diagonal elements represent the variances of individual features, while the off-diagonal elements represent the covariances between pairs of features.

- **Mixing Coefficient π_i :** Represents the prior probability of regime i , which can be interpreted as the expected proportion of time spent in that regime.

We also compute derived metrics such as:

- **Mahalanobis Distance:** Measures the distance between a point and a distribution, taking into account the covariance structure:

$$d_M(o_t, \mu_i, \Sigma_i) = \sqrt{(o_t - \mu_i)^T \Sigma_i^{-1} (o_t - \mu_i)} \quad (139)$$

This helps identify observations that are typical or atypical for a given regime.

- **Component Separation:** Measures how well-separated the components are, using metrics such as the Bhattacharyya distance:

$$D_B(i, j) = \frac{1}{8} (\mu_i - \mu_j)^T \left(\frac{\Sigma_i + \Sigma_j}{2} \right)^{-1} (\mu_i - \mu_j) + \frac{1}{2} \ln \left(\frac{\det \left(\frac{\Sigma_i + \Sigma_j}{2} \right)}{\sqrt{\det(\Sigma_i) \det(\Sigma_j)}} \right) \quad (140)$$

Well-separated components indicate distinct market regimes with different statistical properties.

4.3 Temporal Convolutional Network (TCN) + K-Means

This approach combines deep learning for feature extraction with clustering for regime identification, capturing complex temporal patterns in financial data.

4.3.1 TCN Architecture and Mathematical Foundations

Temporal Convolutional Networks are a class of convolutional neural networks designed for sequence modeling. They use dilated causal convolutions to capture long-range temporal dependencies while maintaining computational efficiency.

The key components of a TCN are:

- **Causal Convolutions:** Ensure that the prediction at time t depends only on inputs from time t and earlier, preventing information leakage from the future:

$$(F * d)(t) = \sum_{i=0}^{k-1} f(i) \cdot d(t - i) \quad (141)$$

where F is the filter, d is the input sequence, and k is the filter size.

- **Dilated Convolutions:** Expand the receptive field exponentially with depth by introducing gaps in the filter:

$$(F *_l d)(t) = \sum_{i=0}^{k-1} f(i) \cdot d(t - l \cdot i) \quad (142)$$

where l is the dilation factor, typically $l = 2^j$ for layer j .

- **Residual Connections:** Help with gradient flow during training by adding the input to the output of a convolutional block:

$$h_{j+1} = h_j + F(h_j, W_j) \quad (143)$$

where h_j is the input to layer j , F is the convolutional block, and W_j are the weights.

The receptive field of a TCN with n layers and dilation factors 2^j for layer j is:

$$\text{Receptive Field} = 1 + \sum_{j=0}^{n-1} (k-1) \cdot 2^j \quad (144)$$

where k is the kernel size. This allows the network to capture dependencies over very long sequences with a relatively small number of layers.

4.3.2 Sequence Embedding Generation

The TCN processes the input time series to generate embeddings that capture temporal patterns. The input to the TCN is a sequence of feature vectors $X = \{x_1, x_2, \dots, x_T\}$, where $x_t \in \mathbb{R}^d$ is the feature vector at time t .

The TCN architecture consists of multiple residual blocks, each containing:

- Dilated causal convolution
- Weight normalization
- ReLU activation
- Dropout for regularization
- Another dilated causal convolution
- Skip connection

The output of the TCN is a sequence of hidden representations $H = \{h_1, h_2, \dots, h_T\}$, where $h_t \in \mathbb{R}^m$ is the embedding at time t .

To obtain a fixed-length embedding for each sequence, we apply one of the following pooling operations:

- **Global Average Pooling:**

$$e = \frac{1}{T} \sum_{t=1}^T h_t \quad (145)$$

- **Global Max Pooling:**

$$e_j = \max_{t \in \{1, 2, \dots, T\}} h_{t,j} \quad \text{for } j \in \{1, 2, \dots, m\} \quad (146)$$

- **Attention Pooling:**

$$\alpha_t = \text{softmax}(w^T h_t) \quad (147)$$

$$e = \sum_{t=1}^T \alpha_t h_t \quad (148)$$

where w is a learnable weight vector.

The resulting embedding $e \in \mathbb{R}^m$ captures the temporal patterns in the input sequence and serves as input to the clustering algorithm.

4.3.3 K-Means Clustering in Latent Space

K-means clustering partitions the embeddings into K clusters, each representing a market regime. The algorithm minimizes the within-cluster sum of squares:

$$J = \sum_{i=1}^K \sum_{e \in C_i} \|e - \mu_i\|^2 \quad (149)$$

where C_i is the set of embeddings assigned to cluster i , and μ_i is the centroid of cluster i .

The algorithm proceeds as follows:

1. Initialize the centroids $\{\mu_1, \mu_2, \dots, \mu_K\}$ (e.g., using K-means++).
2. Assign each embedding to the nearest centroid:

$$C_i = \{e : \|e - \mu_i\|^2 \leq \|e - \mu_j\|^2 \text{ for all } j \neq i\} \quad (150)$$

3. Update the centroids:

$$\mu_i = \frac{1}{|C_i|} \sum_{e \in C_i} e \quad (151)$$

4. Repeat steps 2 and 3 until convergence.

The resulting cluster assignments provide the market regime labels.

4.3.4 Implementation Challenges and Solutions

Implementing the TCN + K-means approach for market regime detection presents several challenges:

- **Training the TCN:** Since we don't have ground truth regime labels, we train the TCN using a self-supervised approach. We define a pretext task of predicting future returns or volatility:

$$\hat{y}_t = f(x_{t-w+1}, x_{t-w+2}, \dots, x_t) \quad (152)$$

where \hat{y}_t is the predicted target (e.g., next-day return) and w is the window size. The network is trained to minimize a suitable loss function, such as mean squared error:

$$\mathcal{L} = \frac{1}{T} \sum_{t=w}^T (y_t - \hat{y}_t)^2 \quad (153)$$

- **Handling Variable-Length Sequences:** Financial time series may have different lengths due to varying trading histories. We address this by:

- Using a fixed-length sliding window approach, where each window is processed independently.

- Padding shorter sequences to a fixed length and using masking to ignore the padded values.
 - Leveraging the TCN's ability to handle variable-length inputs due to its convolutional nature.
- **Hyperparameter Optimization:** The performance of the TCN + K-means approach depends on various hyperparameters, including:
 - TCN architecture: number of layers, kernel size, dilation factors, number of filters.
 - Training parameters: learning rate, batch size, number of epochs.
 - Clustering parameters: number of clusters, initialization method.
- We use a combination of grid search and random search to find optimal hyperparameters, with validation based on clustering quality metrics such as silhouette score and Davies-Bouldin index.
- **Interpretability:** Deep learning models are often considered black boxes. To enhance interpretability, we:
 - Visualize the learned embeddings using dimensionality reduction techniques such as t-SNE or PCA.
 - Analyze the statistical properties of the identified regimes.
 - Compute feature importance scores by perturbing inputs and measuring the impact on embeddings.

4.4 UMAP + K-Means

This approach combines Uniform Manifold Approximation and Projection (UMAP) for dimensionality reduction with K-means clustering for regime identification.

4.4.1 UMAP Mathematical Foundations

UMAP is a non-linear dimensionality reduction technique based on manifold learning and topological data analysis. It preserves both local and global structure of the data by modeling it as a fuzzy topological representation.

The key steps in UMAP are:

1. **Construct a Fuzzy Topological Representation:** For each point x_i in the high-dimensional space, UMAP computes the distances to its nearest neighbors and converts these distances to connection strengths (probabilities):

$$p_{j|i} = \exp\left(-\frac{\max(0, d(x_i, x_j) - \rho_i)}{\sigma_i}\right) \quad (154)$$

where $d(x_i, x_j)$ is the distance between points x_i and x_j , ρ_i is the distance to the nearest neighbor of x_i , and σ_i is a normalization factor chosen such that $\sum_j p_{j|i} = \log_2(k)$ for a fixed number of neighbors k .

The symmetric probabilities are then computed as:

$$p_{ij} = p_{j|i} + p_{i|j} - p_{j|i} \cdot p_{i|j} \quad (155)$$

2. **Find a Low-Dimensional Representation:** UMAP initializes points in the low-dimensional space (typically using spectral embedding) and optimizes their positions to preserve the fuzzy topological structure. The connection strengths in the low-dimensional space are modeled as:

$$q_{ij} = (1 + a \cdot d(y_i, y_j)^{2b})^{-1} \quad (156)$$

where y_i is the low-dimensional representation of x_i , and a and b are parameters that control the shape of the curve.

3. **Optimize the Embedding:** UMAP minimizes the cross-entropy between the high-dimensional and low-dimensional fuzzy topological representations:

$$C = \sum_{i,j} \left[p_{ij} \log \left(\frac{p_{ij}}{q_{ij}} \right) + (1 - p_{ij}) \log \left(\frac{1 - p_{ij}}{1 - q_{ij}} \right) \right] \quad (157)$$

This optimization is performed using stochastic gradient descent.

4.4.2 Hyperparameter Optimization

The performance of UMAP depends on several hyperparameters:

- **n_neighbors:** Controls the size of the local neighborhood used for manifold approximation. Larger values result in more global structure being preserved at the expense of local structure.
- **min_dist:** Controls how tightly points are packed together in the low-dimensional space. Smaller values result in more clustered embeddings.
- **n_components:** The dimensionality of the target embedding. For visualization, this is typically 2 or 3, but for clustering, higher values may be beneficial.
- **metric:** The distance metric used in the high-dimensional space. Common choices include Euclidean, Manhattan, and correlation distances.

We optimize these hyperparameters using a grid search with cross-validation, evaluating the quality of the resulting clusters using metrics such as silhouette score and Davies-Bouldin index.

4.4.3 Implementation Details

Our implementation of the UMAP + K-means approach includes several key considerations:

- **Feature Selection:** We select features that capture various aspects of market behavior, including returns, volatility, correlations, and technical indicators. Feature selection is performed using a combination of domain knowledge and feature importance scores from a random forest model.
- **Scaling:** Features are standardized using the RobustScaler to reduce the impact of outliers:

$$\tilde{x}_i = \frac{x_i - \text{median}(x_i)}{\text{IQR}(x_i)} \quad (158)$$

- **Dimensionality Reduction:** UMAP is applied to reduce the dimensionality of the feature space to a manageable size (typically 5-10 dimensions) while preserving the structure relevant for regime identification.
- **Clustering:** K-means clustering is applied to the UMAP embeddings to identify market regimes. The number of clusters is set to $K = 3$ for consistency with the other approaches.
- **Temporal Consistency:** To ensure temporal consistency of the regime assignments, we apply a post-processing step using a rolling window majority vote:

$$\hat{z}_t = \text{mode}(\{z_{t-w+1}, z_{t-w+2}, \dots, z_t\}) \quad (159)$$

where w is the smoothing window size (typically 5 days).

4.4.4 Topological Structure Preservation Analysis

A key advantage of UMAP is its ability to preserve both local and global structure of the data. We analyze the quality of the embedding using several metrics:

- **Trustworthiness:** Measures the extent to which the local structure is preserved in the embedding:

$$T(k) = 1 - \frac{2}{nk(2n - 3k - 1)} \sum_{i=1}^n \sum_{j \in U_i^{(k)}} (r(i, j) - k) \quad (160)$$

where $U_i^{(k)}$ is the set of points that are among the k nearest neighbors of point i in the embedding but not in the original space, and $r(i, j)$ is the rank of point j in the ordering of nearest neighbors of point i in the original space.

- **Continuity:** Measures the extent to which points that are close in the original space remain close in the embedding:

$$C(k) = 1 - \frac{2}{nk(2n - 3k - 1)} \sum_{i=1}^n \sum_{j \in V_i^{(k)}} (s(i, j) - k) \quad (161)$$

where $V_i^{(k)}$ is the set of points that are among the k nearest neighbors of point i in the original space but not in the embedding, and $s(i, j)$ is the rank of point j in the ordering of nearest neighbors of point i in the embedding.

- **Shepard Diagram:** A scatter plot of the pairwise distances in the original space against the pairwise distances in the embedding. A strong correlation indicates good preservation of the distance structure.

These metrics help us assess whether the UMAP embedding captures the structure relevant for regime identification.

4.5 Ensemble Approach and Temporal Smoothing

To leverage the strengths of multiple regime detection methods, we implement an ensemble approach that combines the results from HMM, GMM, TCN + K-means, and UMAP + K-means.

4.5.1 Ensemble Methodology

Our ensemble methodology combines the regime assignments from multiple methods to produce a more robust and accurate regime classification. The key steps are:

1. **Alignment:** Ensure that all methods use the same regime labeling convention. This is achieved by analyzing the statistical properties of each regime and mapping them to consistent labels (e.g., regime 0 = low volatility/bullish, regime 1 = medium volatility/neutral, regime 2 = high volatility/bearish).
2. **Combination:** Combine the regime assignments using one of several strategies:
 - **Majority Voting:** The most common regime assignment across methods is selected.
 - **Weighted Voting:** Each method's vote is weighted based on its performance on validation data.
 - **Probabilistic Combination:** Combine the posterior probabilities from each method to obtain a consensus probability distribution over regimes.

4.5.2 Weighted Voting Mechanism

In the weighted voting approach, each method m is assigned a weight w_m based on its performance on validation data. The consensus regime assignment at time t is:

$$\hat{z}_t = \arg \max_{k \in \{0,1,2\}} \sum_{m=1}^M w_m \cdot \mathbb{I}(z_t^m = k) \quad (162)$$

where z_t^m is the regime assignment from method m at time t , and $\mathbb{I}(\cdot)$ is the indicator function.

The weights are determined based on clustering quality metrics such as silhouette score, Davies-Bouldin index, and Calinski-Harabasz index. For each method m , we compute a composite score:

$$s_m = \alpha \cdot \text{Silhouette}_m + \beta \cdot \frac{1}{\text{DaviesBouldin}_m} + \gamma \cdot \text{CalinskiHarabasz}_m \quad (163)$$

where α , β , and γ are coefficients that determine the relative importance of each metric. The weights are then normalized:

$$w_m = \frac{s_m}{\sum_{m'=1}^M s_{m'}} \quad (164)$$

4.5.3 Temporal Smoothing Techniques

Financial market regimes typically persist for extended periods, with relatively infrequent transitions. To reflect this property and reduce noise in the regime assignments, we apply temporal smoothing techniques.

Moving Window Majority Vote: For each time point t , the smoothed regime assignment is the most common regime in a window of size w centered at t :

$$\hat{z}_t = \text{mode}(\{z_{t-\lfloor w/2 \rfloor}, \dots, z_t, \dots, z_{t+\lfloor w/2 \rfloor}\}) \quad (165)$$

To avoid look-ahead bias in trading applications, we use a causal window:

$$\hat{z}_t = \text{mode}(\{z_{t-w+1}, \dots, z_t\}) \quad (166)$$

Hidden Markov Model Smoothing: We train a Hidden Markov Model on the ensemble regime assignments, using the ensemble assignments as observations and the true regimes as hidden states. The smoothed regime assignments are obtained using the Viterbi algorithm, which finds the most likely sequence of hidden states given the observations.

4.5.4 Regime Confidence Scores

In addition to the regime assignments, we compute confidence scores that indicate the reliability of each assignment. These scores are useful for risk management and decision-making.

For the ensemble approach, the confidence score at time t is:

$$c_t = \frac{\sum_{m=1}^M w_m \cdot \mathbb{I}(z_t^m = \hat{z}_t)}{\sum_{m=1}^M w_m} \quad (167)$$

This score ranges from 0 to 1, with higher values indicating greater consensus among the methods.

For probabilistic methods (HMM and GMM), we also compute the posterior probability of the assigned regime:

$$p_t = P(z_t = \hat{z}_t | o_t, \theta) \quad (168)$$

This provides a measure of the model's certainty about the regime assignment.

4.6 Regime Characterization

Once market regimes are identified, we characterize them based on their statistical properties and temporal distribution.

4.6.1 Statistical Properties of Identified Regimes

For each regime k , we compute summary statistics of key financial variables:

- **Return Distribution:** Mean, standard deviation, skewness, and kurtosis of re-

turns.

$$\mu_k = \frac{1}{|T_k|} \sum_{t \in T_k} r_t \quad (169)$$

$$\sigma_k = \sqrt{\frac{1}{|T_k|} \sum_{t \in T_k} (r_t - \mu_k)^2} \quad (170)$$

$$s_k = \frac{1}{|T_k|} \sum_{t \in T_k} \left(\frac{r_t - \mu_k}{\sigma_k} \right)^3 \quad (171)$$

$$\kappa_k = \frac{1}{|T_k|} \sum_{t \in T_k} \left(\frac{r_t - \mu_k}{\sigma_k} \right)^4 \quad (172)$$

where $T_k = \{t : z_t = k\}$ is the set of time points assigned to regime k .

- **Volatility Characteristics:** Average volatility, volatility of volatility, and auto-correlation of absolute returns.

$$\bar{\sigma}_k = \frac{1}{|T_k|} \sum_{t \in T_k} \sigma_t \quad (173)$$

$$\text{Vol of Vol}_k = \sqrt{\frac{1}{|T_k|} \sum_{t \in T_k} (\sigma_t - \bar{\sigma}_k)^2} \quad (174)$$

$$\text{AC}_k(\tau) = \frac{\sum_{t \in T_k, t+\tau \in T_k} (|r_t| - \bar{|r|}_k)(|r_{t+\tau}| - \bar{|r|}_k)}{\sum_{t \in T_k} (|r_t| - \bar{|r|}_k)^2} \quad (175)$$

where $\bar{|r|}_k = \frac{1}{|T_k|} \sum_{t \in T_k} |r_t|$.

- **Correlation Structure:** Average correlation between assets and correlation stability.

$$\bar{\rho}_k = \frac{1}{n(n-1)/2} \sum_{i=1}^n \sum_{j=i+1}^n \rho_{ij,k} \quad (176)$$

$$\text{Corr Stability}_k = \frac{1}{n(n-1)/2} \sum_{i=1}^n \sum_{j=i+1}^n \text{std}(\rho_{ij,t})_{t \in T_k} \quad (177)$$

where $\rho_{ij,k}$ is the correlation between assets i and j in regime k , and $\text{std}(\rho_{ij,t})_{t \in T_k}$ is the standard deviation of the time-varying correlation within regime k .

4.6.2 Temporal Distribution Analysis

We analyze the temporal distribution of regimes to understand their persistence and cyclical patterns:

- **Regime Duration:** For each occurrence of regime k , we compute the duration (number of consecutive time points). The distribution of durations provides insights into regime persistence.

- **Transition Frequencies:** We count the number of transitions from regime i to regime j and compute the empirical transition probabilities:

$$\hat{p}_{ij} = \frac{n_{ij}}{\sum_{j'} n_{ij'}} \quad (178)$$

where n_{ij} is the number of transitions from regime i to regime j .

- **Seasonal Patterns:** We analyze whether certain regimes are more likely to occur during specific calendar periods (e.g., months, quarters) or market conditions (e.g., earnings seasons, Fed meetings).

4.6.3 Regime Transition Analysis

Understanding the dynamics of regime transitions is crucial for anticipating changes in market conditions. We analyze:

- **Transition Triggers:** We identify market events or conditions that frequently precede regime transitions. This involves analyzing changes in key variables (e.g., volatility, correlations) around transition points.
- **Early Warning Indicators:** We develop indicators that may signal an impending regime change, such as increasing dispersion in regime assignments across methods or growing instability in key statistical properties.
- **Transition Paths:** We analyze whether transitions between certain regimes follow specific paths (e.g., regime 0 \rightarrow regime 1 \rightarrow regime 2) or can occur directly (e.g., regime 0 \rightarrow regime 2).

4.6.4 Economic Interpretation of Regimes

To make the identified regimes more interpretable and actionable, we map them to economic and market conditions:

- **Regime Labels:** Based on the statistical properties, we assign descriptive labels to each regime, such as "Bull Market" (high returns, low volatility), "Bear Market" (negative returns, high volatility), or "Consolidation" (low returns, low volatility).
- **Macroeconomic Conditions:** We analyze the relationship between regimes and macroeconomic variables such as GDP growth, inflation, and interest rates.
- **Market Sentiment:** We examine how regimes relate to market sentiment indicators such as the VIX, put-call ratio, and investor surveys.
- **Historical Context:** We map the identified regimes to well-known historical market periods (e.g., dot-com bubble, 2008 financial crisis, COVID-19 crash) to validate their economic significance.

This economic interpretation enhances the practical utility of the regime detection framework for investment decision-making and risk management.

5 GenAI Data Augmentation Implementation

This section details the implementation of our generative AI data augmentation component, which creates synthetic financial data conditioned on detected market regimes. We present the architecture, loss functions, training dynamics, and implementation details for each generative model, as well as the validation framework and synthetic boost protocol.

5.1 TimeGAN Implementation

Time-series Generative Adversarial Network (TimeGAN) [Yoon et al., 2019] is a state-of-the-art framework for generating realistic time series data that preserves both temporal dynamics and cross-feature relationships.

5.1.1 Architecture Components

TimeGAN consists of four main neural network components: an embedding network, a recovery network, a sequence generator, and a sequence discriminator. Additionally, it incorporates a supervisor network for improved training stability.

Embedding Network (E_{ϕ_E}): Maps the real data sequence $\mathbf{X} = \{x_1, x_2, \dots, x_T\}$ from the original feature space to a latent representation $\mathbf{H} = \{h_1, h_2, \dots, h_T\}$:

$$\mathbf{H} = E_{\phi_E}(\mathbf{X}) \quad (179)$$

where ϕ_E represents the parameters of the embedding network. The architecture consists of multiple LSTM layers followed by fully connected layers:

$$h_t^{(1)} = \text{LSTM}_1(x_t, h_{t-1}^{(1)}) \quad (180)$$

$$h_t^{(2)} = \text{LSTM}_2(h_t^{(1)}, h_{t-1}^{(2)}) \quad (181)$$

$$h_t = \text{FC}(h_t^{(2)}) \quad (182)$$

Recovery Network (R_{ϕ_R}): Reconstructs the original data from the latent representation:

$$\hat{\mathbf{X}} = R_{\phi_R}(\mathbf{H}) \quad (183)$$

where ϕ_R represents the parameters of the recovery network. The architecture consists of fully connected layers followed by LSTM layers:

$$r_t^{(1)} = \text{FC}(h_t) \quad (184)$$

$$r_t^{(2)} = \text{LSTM}_1(r_t^{(1)}, r_{t-1}^{(2)}) \quad (185)$$

$$\hat{x}_t = \text{LSTM}_2(r_t^{(2)}, \hat{x}_{t-1}) \quad (186)$$

Generator (G_{ϕ_G}): Generates synthetic latent representations $\hat{\mathbf{H}} = \{\hat{h}_1, \hat{h}_2, \dots, \hat{h}_T\}$ from random noise $\mathbf{Z} = \{z_1, z_2, \dots, z_T\}$:

$$\hat{\mathbf{H}} = G_{\phi_G}(\mathbf{Z}, c) \quad (187)$$

where ϕ_G represents the parameters of the generator and c is the conditioning information (market regime). The architecture consists of an embedding layer for the condition, followed by LSTM layers:

$$c_{\text{emb}} = \text{Embedding}(c) \quad (188)$$

$$z'_t = \text{Concat}(z_t, c_{\text{emb}}) \quad (189)$$

$$g_t^{(1)} = \text{LSTM}_1(z'_t, g_{t-1}^{(1)}) \quad (190)$$

$$g_t^{(2)} = \text{LSTM}_2(g_t^{(1)}, g_{t-1}^{(2)}) \quad (191)$$

$$\hat{h}_t = \text{FC}(g_t^{(2)}) \quad (192)$$

Discriminator (D_{ϕ_D}): Distinguishes between real and synthetic latent representations:

$$D_{\phi_D}(\mathbf{H}) = P(\text{real}|\mathbf{H}) \quad (193)$$

where ϕ_D represents the parameters of the discriminator. The architecture consists of bidirectional LSTM layers followed by fully connected layers:

$$d_t^{(1)} = \text{BiLSTM}_1(h_t, d_{t-1}^{(1)}, d_{t+1}^{(1)}) \quad (194)$$

$$d_t^{(2)} = \text{BiLSTM}_2(d_t^{(1)}, d_{t-1}^{(2)}, d_{t+1}^{(2)}) \quad (195)$$

$$d_t = \text{FC}(d_t^{(2)}) \quad (196)$$

$$p_{\text{real}} = \text{Sigmoid}(\text{FC}(\text{Pooling}(\{d_1, d_2, \dots, d_T\}))) \quad (197)$$

Supervisor (S_{ϕ_S}): Predicts the next step in the latent space given the current step:

$$\hat{h}_{t+1} = S_{\phi_S}(h_t) \quad (198)$$

where ϕ_S represents the parameters of the supervisor. The architecture consists of LSTM layers followed by fully connected layers:

$$s_t^{(1)} = \text{LSTM}_1(h_t, s_{t-1}^{(1)}) \quad (199)$$

$$s_t^{(2)} = \text{LSTM}_2(s_t^{(1)}, s_{t-1}^{(2)}) \quad (200)$$

$$\hat{h}_{t+1} = \text{FC}(s_t^{(2)}) \quad (201)$$

5.1.2 Loss Functions and Training Dynamics

TimeGAN is trained using a combination of reconstruction, supervised, and adversarial losses.

Reconstruction Loss: Ensures that the embedding and recovery networks form an effective autoencoder:

$$\mathcal{L}_{\text{rec}} = \mathbb{E}_{\mathbf{X} \sim p_{\text{data}}} [\|\mathbf{X} - R_{\phi_R}(E_{\phi_E}(\mathbf{X}))\|_2^2] \quad (202)$$

Supervised Loss: Encourages the generator to capture the temporal dynamics of the data:

$$\mathcal{L}_{\text{sup}} = \mathbb{E}_{\mathbf{Z} \sim p_{\mathbf{Z}}, c \sim p_c} \left[\sum_{t=1}^{T-1} \|S_{\phi_S}(\hat{h}_t) - \hat{h}_{t+1}\|_2^2 \right] \quad (203)$$

where \hat{h}_t is the t -th element of $\hat{\mathbf{H}} = G_{\phi_G}(\mathbf{Z}, c)$.

Adversarial Loss: Ensures that the generated sequences are indistinguishable from real sequences:

$$\mathcal{L}_{\text{adv}} = \mathbb{E}_{\mathbf{X} \sim p_{\text{data}}} [\log D_{\phi_D}(E_{\phi_E}(\mathbf{X}))] + \mathbb{E}_{\mathbf{Z} \sim p_{\mathbf{Z}}, c \sim p_c} [\log(1 - D_{\phi_D}(G_{\phi_G}(\mathbf{Z}, c)))] \quad (204)$$

The training procedure alternates between updating different components:

1. Update the embedding and recovery networks to minimize the reconstruction loss:

$$\min_{\phi_E, \phi_R} \mathcal{L}_{\text{rec}} \quad (205)$$

2. Update the supervisor network to minimize the supervised loss on real data:

$$\min_{\phi_S} \mathbb{E}_{\mathbf{X} \sim p_{\text{data}}} \left[\sum_{t=1}^{T-1} \|S_{\phi_S}(h_t) - h_{t+1}\|_2^2 \right] \quad (206)$$

where h_t is the t -th element of $\mathbf{H} = E_{\phi_E}(\mathbf{X})$.

3. Update the generator to minimize a combination of supervised and adversarial losses:

$$\min_{\phi_G} \lambda_{\text{sup}} \mathcal{L}_{\text{sup}} - \lambda_{\text{adv}} \mathbb{E}_{\mathbf{Z} \sim p_{\mathbf{Z}}, c \sim p_c} [\log(D_{\phi_D}(G_{\phi_G}(\mathbf{Z}, c)))] \quad (207)$$

where λ_{sup} and λ_{adv} are hyperparameters that control the relative importance of the supervised and adversarial losses.

4. Update the discriminator to maximize the adversarial loss:

$$\max_{\phi_D} \mathcal{L}_{\text{adv}} \quad (208)$$

5.1.3 Regime-Conditional Generation

To generate regime-specific synthetic data, we condition the TimeGAN model on the market regime. This is achieved by incorporating the regime information into the generator's input.

Let $c \in \{0, 1, 2\}$ represent the market regime. We first convert c to a one-hot encoded vector $c_{\text{one-hot}} \in \{0, 1\}^3$. This vector is then embedded into a higher-dimensional space using an embedding layer:

$$c_{\text{emb}} = \text{Embedding}(c_{\text{one-hot}}) \quad (209)$$

The embedded regime information is concatenated with the random noise at each time step:

$$z'_t = \text{Concat}(z_t, c_{\text{emb}}) \quad (210)$$

This conditioning ensures that the generator produces synthetic data with statistical properties specific to the target regime.

5.1.4 Implementation Challenges and Solutions

Implementing TimeGAN for financial time series presents several challenges:

- **Mode Collapse:** GANs are prone to mode collapse, where the generator produces limited varieties of samples. To address this, we implement:
 - Minibatch discrimination: Encourages diversity by allowing the discriminator to examine multiple samples simultaneously.
 - Feature matching: Adds an auxiliary loss that matches the statistics of real and generated data.
 - Gradient penalty: Stabilizes training by penalizing large gradients in the discriminator.
- **Training Stability:** GAN training can be unstable, leading to oscillations or divergence. Our solutions include:
 - Learning rate scheduling: Gradually decreasing the learning rate during training.
 - Gradient clipping: Limiting the magnitude of gradients to prevent exploding gradients.
 - Spectral normalization: Constraining the Lipschitz constant of the discriminator.
- **Evaluation Metrics:** Evaluating the quality of synthetic time series is challenging. We implement:
 - Discriminative score: Measures how well a classifier can distinguish between real and synthetic data.
 - Predictive score: Measures how well a model trained on synthetic data performs on real data.
 - Statistical tests: Compare the statistical properties (e.g., autocorrelation, volatility clustering) of real and synthetic data.

5.2 Conditional Variational Autoencoder (CVAE) Implementation

Conditional Variational Autoencoders extend the VAE framework by conditioning the generation process on additional information, such as market regime labels.

5.2.1 Encoder-Decoder Architecture

The CVAE consists of an encoder network that maps the input data to a latent distribution and a decoder network that reconstructs the data from samples of this distribution, both conditioned on the regime label.

Encoder Network: Maps the input sequence $\mathbf{X} = \{x_1, x_2, \dots, x_T\}$ and condition c to the parameters of a latent distribution:

$$\mu_z, \log \sigma_z^2 = \text{Encoder}_\phi(\mathbf{X}, c) \quad (211)$$

where ϕ represents the parameters of the encoder. The architecture consists of:

$$c_{\text{emb}} = \text{Embedding}(c) \quad (212)$$

$$c_{\text{expanded}} = \text{Repeat}(c_{\text{emb}}, T) \quad (213)$$

$$x'_t = \text{Concat}(x_t, c_{\text{expanded}, t}) \quad (214)$$

$$e_t^{(1)} = \text{LSTM}_1(x'_t, e_{t-1}^{(1)}) \quad (215)$$

$$e_t^{(2)} = \text{LSTM}_2(e_t^{(1)}, e_{t-1}^{(2)}) \quad (216)$$

$$e = \text{Pooling}(\{e_1^{(2)}, e_2^{(2)}, \dots, e_T^{(2)}\}) \quad (217)$$

$$\mu_z = \text{FC}_{\mu}(e) \quad (218)$$

$$\log \sigma_z^2 = \text{FC}_{\sigma}(e) \quad (219)$$

Decoder Network: Reconstructs the input sequence from a latent vector z and condition c :

$$\hat{\mathbf{X}} = \text{Decoder}_{\theta}(z, c) \quad (220)$$

where θ represents the parameters of the decoder. The architecture consists of:

$$c_{\text{emb}} = \text{Embedding}(c) \quad (221)$$

$$z_c = \text{Concat}(z, c_{\text{emb}}) \quad (222)$$

$$z_{\text{expanded}} = \text{Repeat}(z_c, T) \quad (223)$$

$$d_t^{(1)} = \text{LSTM}_1(z_{\text{expanded}, t}, d_{t-1}^{(1)}) \quad (224)$$

$$d_t^{(2)} = \text{LSTM}_2(d_t^{(1)}, d_{t-1}^{(2)}) \quad (225)$$

$$\hat{x}_t = \text{FC}(d_t^{(2)}) \quad (226)$$

5.2.2 Latent Space Sampling with Regime Conditioning

During training, we sample from the latent distribution using the reparameterization trick to enable backpropagation:

$$z = \mu_z + \sigma_z \odot \epsilon, \quad \epsilon \sim \mathcal{N}(0, I) \quad (227)$$

where \odot represents element-wise multiplication.

The prior distribution $p(z|c)$ is modeled as a standard Gaussian:

$$p(z|c) = \mathcal{N}(0, I) \quad (228)$$

During generation, we sample z from this prior and condition the decoder on the target regime c to produce synthetic data specific to that regime.

5.2.3 Balancing Reconstruction and KL Divergence Losses

The CVAE is trained to maximize the evidence lower bound (ELBO):

$$\mathcal{L}_{\text{ELBO}}(\phi, \theta; \mathbf{X}, c) = \mathbb{E}_{q_{\phi}(z|\mathbf{X}, c)}[\log p_{\theta}(\mathbf{X}|z, c)] - \beta \cdot D_{\text{KL}}(q_{\phi}(z|\mathbf{X}, c) \| p(z|c)) \quad (229)$$

where:

- $q_\phi(z|\mathbf{X}, c)$ is the approximate posterior (encoder) distribution.
- $p_\theta(\mathbf{X}|z, c)$ is the likelihood (decoder) distribution.
- $p(z|c)$ is the prior distribution.
- β is a hyperparameter that controls the weight of the KL divergence term.

The reconstruction term encourages the model to accurately reconstruct the input data, while the KL divergence term encourages the approximate posterior to be close to the prior. Balancing these terms is crucial for generating realistic yet diverse samples.

We implement β -VAE [Higgins et al., 2017], which introduces the hyperparameter β to control this balance. Setting $\beta < 1$ emphasizes reconstruction accuracy, while $\beta > 1$ emphasizes latent space regularity.

Additionally, we implement KL annealing, where β is gradually increased from a small value to its final value during training:

$$\beta_{\text{epoch}} = \min \left(\beta_{\max}, \beta_{\min} + (\beta_{\max} - \beta_{\min}) \cdot \frac{\text{epoch}}{\text{annealing_epochs}} \right) \quad (230)$$

This helps prevent the KL divergence term from dominating early in training, allowing the model to first learn a good reconstruction before enforcing the latent space structure.

5.2.4 Implementation Details

Our implementation of the CVAE includes several key considerations:

- **Network Architecture:** We use a bidirectional LSTM encoder and an LSTM decoder, with skip connections to improve gradient flow.
- **Training Procedure:** We train the model using the Adam optimizer with a learning rate of 10^{-4} and a batch size of 32. We implement early stopping based on the validation loss to prevent overfitting.
- **Hyperparameter Selection:** We perform a grid search over key hyperparameters, including the latent dimension, β value, and network size. The optimal configuration is selected based on the quality of the generated samples, as measured by the synthetic data validation framework.
- **Conditional Generation:** To generate regime-specific data, we sample from the prior distribution $p(z|c)$ and decode using the target regime c . We generate multiple samples for each regime to capture the diversity within each regime.

5.3 FiLM Transformer Implementation

The Feature-wise Linear Modulation (FiLM) Transformer combines the powerful sequence modeling capabilities of Transformers with the conditional generation capabilities of FiLM layers.

5.3.1 Self-Attention Mechanism

The core of the Transformer architecture is the self-attention mechanism, which allows the model to weigh the importance of different positions in the input sequence when computing a representation for each position.

Given an input sequence $\mathbf{X} = \{x_1, x_2, \dots, x_T\}$, the self-attention mechanism computes:

$$Q = \mathbf{X}W^Q \quad (231)$$

$$K = \mathbf{X}W^K \quad (232)$$

$$V = \mathbf{X}W^V \quad (233)$$

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (234)$$

where W^Q , W^K , and W^V are learnable parameter matrices, and d_k is the dimension of the keys.

Multi-head attention extends this by applying multiple attention operations in parallel and concatenating the results:

$$\text{head}_i = \text{Attention}(\mathbf{X}W_i^Q, \mathbf{X}W_i^K, \mathbf{X}W_i^V) \quad (235)$$

$$\text{MultiHead}(\mathbf{X}) = \text{Concat}(\text{head}_1, \text{head}_2, \dots, \text{head}_h)W^O \quad (236)$$

where h is the number of attention heads and W^O is a learnable parameter matrix.

5.3.2 Feature-wise Linear Modulation Layers

FiLM layers modulate the features of a neural network based on conditioning information. Given a feature map F and a conditioning vector c , FiLM applies an affine transformation:

$$\text{FiLM}(F, c) = \gamma(c) \odot F + \beta(c) \quad (237)$$

where $\gamma(c)$ and $\beta(c)$ are scaling and shifting parameters generated from the conditioning vector c , and \odot represents element-wise multiplication.

In our implementation, $\gamma(c)$ and $\beta(c)$ are computed using fully connected networks:

$$\gamma(c) = \text{FC}_\gamma(c) \quad (238)$$

$$\beta(c) = \text{FC}_\beta(c) \quad (239)$$

5.3.3 Positional Encoding for Temporal Awareness

Since Transformers process sequences in parallel rather than sequentially, they lack an inherent notion of position. To incorporate temporal information, we add positional encodings to the input embeddings:

$$\text{PE}_{(pos, 2i)} = \sin\left(\frac{pos}{10000^{2i/d_{\text{model}}}}\right) \quad (240)$$

$$\text{PE}_{(pos, 2i+1)} = \cos\left(\frac{pos}{10000^{2i/d_{\text{model}}}}\right) \quad (241)$$

where pos is the position in the sequence, i is the dimension index, and d_{model} is the embedding dimension.

The input to the Transformer becomes:

$$\mathbf{X}' = \mathbf{X} + \text{PE} \quad (242)$$

5.3.4 Implementation Challenges and Solutions

Implementing the FiLM Transformer for financial time series presents several challenges:

- **Attention Mechanism Optimization:** The self-attention mechanism has quadratic complexity with respect to sequence length, which can be computationally expensive for long financial time series. We address this by:
 - Using a sliding window approach to process long sequences in chunks.
 - Implementing sparse attention patterns that focus on local neighborhoods and a few global tokens.
 - Applying gradient checkpointing to reduce memory usage during training.
- **Memory Efficiency:** Transformers can be memory-intensive, especially for long sequences. Our solutions include:
 - Mixed-precision training to reduce memory footprint.
 - Gradient accumulation to effectively increase batch size without increasing memory usage.
 - Model parallelism for very large models.
- **Training Stability:** Transformers can be difficult to train due to their depth and the use of residual connections. We implement:
 - Layer normalization before each sub-layer (Pre-LN) instead of after (Post-LN) for improved stability.
 - Warm-up learning rate scheduling to gradually increase the learning rate at the beginning of training.
 - Gradient clipping to prevent exploding gradients.

5.4 Regime Adversary Module Implementation

The Regime Adversary Module ensures that the synthetic data generated for each regime accurately reflects the statistical properties of that regime.

5.4.1 Adversarial Training Framework

The Regime Adversary Module consists of a discriminator network that attempts to classify the regime of a given time series. The generator is trained to produce synthetic data that the discriminator classifies as belonging to the target regime.

Let $G_\theta(z, c)$ be the generator that produces synthetic data given random noise z and regime condition c . Let $D_\phi(x)$ be the discriminator that predicts the regime of a time series x .

The adversarial loss for the generator is:

$$\mathcal{L}_{\text{adv}}(G) = -\mathbb{E}_{z \sim p_z, c \sim p_c} [\log D_\phi(G_\theta(z, c), c)] \quad (243)$$

where $D_\phi(x, c)$ represents the probability that the discriminator assigns to regime c for input x .

The discriminator is trained to maximize the probability of correctly classifying both real and synthetic data:

$$\mathcal{L}_{\text{adv}}(D) = -\mathbb{E}_{x \sim p_{\text{data}}, c \sim p_c} [\log D_\phi(x, c)] - \mathbb{E}_{z \sim p_z, c \sim p_c} [\log(1 - D_\phi(G_\theta(z, c), c))] \quad (244)$$

5.4.2 Feature Matching for Statistical Alignment

To ensure that the synthetic data preserves the statistical properties of the real data within each regime, we implement feature matching. This involves adding an auxiliary loss that encourages the generator to match the statistics of real and synthetic data.

Let $f_D(x)$ be a feature extractor that computes intermediate representations from the discriminator. The feature matching loss is:

$$\mathcal{L}_{\text{feat}}(G) = \mathbb{E}_{c \sim p_c} \left[\left\| \mathbb{E}_{x \sim p_{\text{data}}(x|c)} [f_D(x)] - \mathbb{E}_{z \sim p_z} [f_D(G_\theta(z, c))] \right\|_2^2 \right] \quad (245)$$

This loss encourages the generator to produce synthetic data whose feature statistics match those of the real data within each regime.

5.4.3 Gradient Penalty for Training Stability

To improve the stability of adversarial training, we implement a gradient penalty based on the Wasserstein GAN with Gradient Penalty (WGAN-GP) approach [Gulrajani et al., 2017]. This penalizes the discriminator if its gradients with respect to the inputs have norms far from 1, enforcing a Lipschitz constraint.

The gradient penalty is:

$$\mathcal{L}_{\text{gp}}(D) = \mathbb{E}_{x \sim p_{\text{penalty}}, c \sim p_c} \left[(\|\nabla_x D_\phi(x, c)\|_2 - 1)^2 \right] \quad (246)$$

where p_{penalty} is a distribution sampling points along straight lines between pairs of real and generated samples:

$$x = \epsilon x_{\text{real}} + (1 - \epsilon)x_{\text{gen}}, \quad \epsilon \sim \text{Uniform}(0, 1) \quad (247)$$

The total loss for the discriminator becomes:

$$\mathcal{L}_{\text{total}}(D) = \mathcal{L}_{\text{adv}}(D) + \lambda_{\text{gp}} \mathcal{L}_{\text{gp}}(D) \quad (248)$$

where λ_{gp} is a hyperparameter controlling the weight of the gradient penalty.

5.4.4 Implementation Details

Our implementation of the Regime Adversary Module includes several key considerations:

- **Network Architecture:** The discriminator is a convolutional neural network with residual connections, designed to capture both local and global patterns in the time series.

- **Training Procedure:** We alternate between updating the discriminator and the generator, with multiple discriminator updates per generator update to ensure that the discriminator provides a meaningful learning signal.
- **Hyperparameter Selection:** We tune the weight of the gradient penalty λ_{gp} and the feature matching loss λ_{feat} to balance training stability and generation quality.
- **Integration with Generative Models:** The Regime Adversary Module can be integrated with any of the generative models (TimeGAN, CVAE, FiLM Transformer) by adding the adversarial and feature matching losses to their respective training objectives.

5.5 Synthetic Data Validation Framework

Evaluating the quality of synthetic financial time series is challenging due to the complex statistical properties that must be preserved. We implement a comprehensive validation framework that assesses the synthetic data from multiple perspectives.

5.5.1 Statistical Tests

We perform a battery of statistical tests to compare the properties of real and synthetic data:

- **Autocorrelation Function (ACF) and Partial Autocorrelation Function (PACF):** Measure the temporal dependencies in the data. For a time series $\{x_t\}$, the ACF at lag k is:

$$\rho_k = \frac{\mathbb{E}[(x_t - \mu)(x_{t+k} - \mu)]}{\sigma^2} \quad (249)$$

where μ and σ^2 are the mean and variance of the series. We compare the ACF and PACF of real and synthetic data using the Ljung-Box test:

$$Q = n(n + 2) \sum_{k=1}^h \frac{\hat{\rho}_k^2}{n - k} \quad (250)$$

where n is the sample size, h is the maximum lag, and $\hat{\rho}_k$ is the sample autocorrelation at lag k .

- **Distribution Tests:** Compare the marginal distributions of real and synthetic data. We use the Kolmogorov-Smirnov test:

$$D = \sup_x |F_{\text{real}}(x) - F_{\text{syn}}(x)| \quad (251)$$

where F_{real} and F_{syn} are the empirical cumulative distribution functions of the real and synthetic data.

- **Stationarity Tests:** Assess whether the synthetic data exhibits similar stationarity properties as the real data. We use the Augmented Dickey-Fuller test:

$$\Delta x_t = \alpha + \beta t + \gamma x_{t-1} + \delta_1 \Delta x_{t-1} + \cdots + \delta_p \Delta x_{t-p} + \epsilon_t \quad (252)$$

The null hypothesis is that $\gamma = 0$ (unit root present, series is non-stationary).

5.5.2 Financial Metrics

We evaluate whether the synthetic data preserves key financial properties:

- **Volatility Clustering:** Measured by the autocorrelation of absolute returns:

$$\rho_k^{|r|} = \frac{\mathbb{E}[(|r_t| - \mu_{|r|})(|r_{t+k}| - \mu_{|r|})]}{\sigma_{|r|}^2} \quad (253)$$

where r_t is the return at time t , and $\mu_{|r|}$ and $\sigma_{|r|}^2$ are the mean and variance of the absolute returns.

- **Fat-tail Behavior:** Assessed by comparing the kurtosis of real and synthetic returns:

$$\kappa = \frac{\mathbb{E}[(r_t - \mu_r)^4]}{\sigma_r^4} \quad (254)$$

where μ_r and σ_r^2 are the mean and variance of returns. We also compare the tail index estimated using the Hill estimator:

$$\hat{\alpha} = \left(\frac{1}{k} \sum_{i=1}^k \log \frac{X_{(i)}}{X_{(k+1)}} \right)^{-1} \quad (255)$$

where $X_{(1)} \geq X_{(2)} \geq \dots \geq X_{(n)}$ are the ordered observations, and k is the number of tail observations used in the estimation.

- **Leverage Effect:** Measured by the correlation between returns and future volatility:

$$\rho_{r_t, \sigma_{t+1}^2} = \frac{\mathbb{E}[(r_t - \mu_r)(\sigma_{t+1}^2 - \mu_{\sigma^2})]}{\sigma_r \sigma_{\sigma^2}} \quad (256)$$

where σ_{t+1}^2 is the volatility at time $t+1$, and μ_{σ^2} and σ_{σ^2} are the mean and standard deviation of volatility.

5.5.3 Machine Learning Metrics

We evaluate the synthetic data from a machine learning perspective:

- **Classifier Confusion:** Train a classifier to distinguish between real and synthetic data. If the synthetic data is realistic, the classifier should perform close to random guessing (accuracy around 0.5). The classifier confusion score is:

$$C = 2 \cdot |0.5 - \text{Accuracy}| \quad (257)$$

where Accuracy is the classification accuracy on a held-out test set. Lower values of C indicate more realistic synthetic data.

- **Feature Importance Preservation:** Train predictive models on real and synthetic data, and compare the feature importance rankings. Let $I_{\text{real}}(f)$ and $I_{\text{syn}}(f)$ be the importance of feature f in models trained on real and synthetic data. The feature importance preservation score is:

$$\text{FIP} = \text{Spearman}(I_{\text{real}}, I_{\text{syn}}) \quad (258)$$

where Spearman is the Spearman rank correlation coefficient. Higher values indicate better preservation of feature relationships.

5.5.4 Visualization Techniques

We use visualization techniques to qualitatively assess the synthetic data:

- **Principal Component Analysis (PCA):** Project the real and synthetic data onto the first two principal components and visualize the resulting distributions. This helps assess whether the synthetic data captures the overall structure of the real data.
- **t-Distributed Stochastic Neighbor Embedding (t-SNE):** A non-linear dimensionality reduction technique that is particularly effective at visualizing high-dimensional data. We use t-SNE to compare the local structure of real and synthetic data.
- **Regime Separation Visualization:** Project data from different regimes onto a common space and visualize whether the synthetic data maintains the separation between regimes observed in the real data.

5.6 Synthetic Boost Protocol

The Synthetic Boost Protocol defines how the synthetic data is integrated with real data to improve model training, particularly for rare market regimes.

5.6.1 Regime-Specific Augmentation Ratios

Different market regimes may require different levels of augmentation based on their representation in the historical data. We define regime-specific augmentation ratios:

$$r_k = \frac{n_{\text{syn},k}}{n_{\text{real},k}} \quad (259)$$

where $n_{\text{real},k}$ is the number of real samples in regime k , and $n_{\text{syn},k}$ is the number of synthetic samples to generate.

We determine these ratios based on the following considerations:

- **Regime Frequency:** Rare regimes (e.g., crisis periods) receive higher augmentation ratios to balance the dataset.

$$r_k \propto \frac{1}{f_k} \quad (260)$$

where $f_k = \frac{n_{\text{real},k}}{\sum_j n_{\text{real},j}}$ is the frequency of regime k in the real data.

- **Synthetic Data Quality:** Regimes for which the generative models produce higher-quality synthetic data receive higher augmentation ratios.

$$r_k \propto q_k \quad (261)$$

where $q_k \in [0, 1]$ is a quality score for the synthetic data in regime k , derived from the validation framework.

- **Model Performance Gap:** Regimes where models show the largest performance gap between training and validation receive higher augmentation ratios.

$$r_k \propto g_k \quad (262)$$

where $g_k = \text{Perf}_{\text{train},k} - \text{Perf}_{\text{val},k}$ is the performance gap in regime k .

The final augmentation ratio is a weighted combination of these factors:

$$r_k = \alpha \cdot \frac{1}{f_k} + \beta \cdot q_k + \gamma \cdot g_k \quad (263)$$

where α , β , and γ are hyperparameters that control the relative importance of each factor.

5.6.2 Balanced Training Set Creation

Using the regime-specific augmentation ratios, we create a balanced training set that combines real and synthetic data:

1. For each regime k , generate $n_{\text{syn},k} = \lceil r_k \cdot n_{\text{real},k} \rceil$ synthetic samples.
2. Combine the real and synthetic samples for each regime.
3. Shuffle the combined dataset to ensure random ordering.

To prevent overfitting to synthetic data patterns, we implement a curriculum learning approach where the proportion of synthetic data is gradually increased during training:

$$p_{\text{syn}}(e) = p_{\text{syn}}^{\max} \cdot \min \left(1, \frac{e}{e_{\text{ramp}}} \right) \quad (264)$$

where $p_{\text{syn}}(e)$ is the proportion of synthetic data at epoch e , p_{syn}^{\max} is the maximum proportion, and e_{ramp} is the number of epochs over which to ramp up the proportion.

5.6.3 Sample Weighting Based on Quality Metrics

Not all synthetic samples are of equal quality. We implement a sample weighting scheme that assigns weights to synthetic samples based on their quality:

$$w_i = \begin{cases} 1 & \text{if sample } i \text{ is real} \\ w_{\text{syn}} \cdot q_i & \text{if sample } i \text{ is synthetic} \end{cases} \quad (265)$$

where $w_{\text{syn}} \in [0, 1]$ is a global weight for synthetic samples, and $q_i \in [0, 1]$ is a quality score for synthetic sample i .

The quality score q_i is computed based on how well the synthetic sample preserves key statistical properties of the real data in its regime:

$$q_i = \exp(-\lambda \cdot d(x_i, \mathcal{R}_k)) \quad (266)$$

where $d(x_i, \mathcal{R}_k)$ is a distance measure between the synthetic sample x_i and the distribution of real samples in regime \mathcal{R}_k , and λ is a scaling parameter.

5.6.4 Implementation Details

Our implementation of the Synthetic Boost Protocol includes several key considerations:

- **Data Pipeline:** We implement an efficient data pipeline that dynamically combines real and synthetic data during training, with support for sample weighting and curriculum learning.
- **Quality Assessment:** We periodically reassess the quality of synthetic data and adjust the augmentation ratios and sample weights accordingly.
- **Model Monitoring:** We monitor the performance of models trained with synthetic data to ensure that the augmentation is beneficial and not introducing biases.
- **Regime-Specific Tuning:** We tune the hyperparameters of the Synthetic Boost Protocol separately for each regime to optimize the benefit of synthetic data augmentation.

This comprehensive approach to synthetic data generation and integration ensures that the augmented dataset preserves the statistical properties of the real data while addressing the data scarcity problem, particularly for rare market regimes.

6 Algorithmic Trading Strategy Implementation

This section details the implementation of our algorithmic trading strategies, which leverage the market regime detection and synthetic data augmentation components. We present the mathematical formulation, implementation details, and performance metrics for each strategy.

6.1 Transformer-based Momentum Strategy

The Transformer-based Momentum Strategy utilizes the powerful sequence modeling capabilities of Transformer architectures to identify momentum patterns in financial time series.

6.1.1 Architecture and Mathematical Formulation

The core of this strategy is a Transformer model that processes historical price and feature data to predict future price movements. The model architecture consists of:

- **Input Embedding Layer:** Maps the input features to a higher-dimensional space.
- **Positional Encoding:** Adds information about the position of each time step in the sequence.
- **Transformer Encoder Blocks:** Process the sequence using self-attention mechanisms.
- **Output Layer:** Maps the encoded sequence to trading signals.

Let $\mathbf{X} = \{x_1, x_2, \dots, x_T\}$ be a sequence of feature vectors, where $x_t \in \mathbb{R}^d$ represents the features at time t . The Transformer processes this sequence as follows:

Input Embedding:

$$\mathbf{E} = \mathbf{X}W_E + \mathbf{b}_E \quad (267)$$

where $W_E \in \mathbb{R}^{d \times d_{\text{model}}}$ and $\mathbf{b}_E \in \mathbb{R}^{d_{\text{model}}}$ are learnable parameters, and d_{model} is the model dimension.

Positional Encoding:

$$\text{PE}_{(t,2i)} = \sin\left(\frac{t}{10000^{2i/d_{\text{model}}}}\right) \quad (268)$$

$$\text{PE}_{(t,2i+1)} = \cos\left(\frac{t}{10000^{2i/d_{\text{model}}}}\right) \quad (269)$$

where t is the position in the sequence and i is the dimension index.

The input to the Transformer encoder is:

$$\mathbf{H}^{(0)} = \mathbf{E} + \text{PE} \quad (270)$$

Transformer Encoder Block: Each encoder block consists of multi-head self-attention followed by a position-wise feed-forward network, with residual connections and layer normalization:

$$\mathbf{H}'^{(l)} = \text{LayerNorm}(\mathbf{H}^{(l-1)} + \text{MultiHead}(\mathbf{H}^{(l-1)})) \quad (271)$$

$$\mathbf{H}^{(l)} = \text{LayerNorm}(\mathbf{H}'^{(l)} + \text{FFN}(\mathbf{H}'^{(l)})) \quad (272)$$

where l is the layer index.

The multi-head attention mechanism is defined as:

$$\text{MultiHead}(\mathbf{H}) = \text{Concat}(\text{head}_1, \text{head}_2, \dots, \text{head}_h)W^O \quad (273)$$

$$\text{head}_i = \text{Attention}(\mathbf{HW}_i^Q, \mathbf{HW}_i^K, \mathbf{HW}_i^V) \quad (274)$$

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (275)$$

where $W_i^Q, W_i^K, W_i^V \in \mathbb{R}^{d_{\text{model}} \times d_k}$ and $W^O \in \mathbb{R}^{hd_k \times d_{\text{model}}}$ are learnable parameters, h is the number of attention heads, and $d_k = d_{\text{model}}/h$ is the dimension of each head.

The feed-forward network is:

$$\text{FFN}(\mathbf{H}) = \max(0, \mathbf{HW}_1 + \mathbf{b}_1)W_2 + \mathbf{b}_2 \quad (276)$$

where $W_1 \in \mathbb{R}^{d_{\text{model}} \times d_{\text{ff}}}$, $W_2 \in \mathbb{R}^{d_{\text{ff}} \times d_{\text{model}}}$, $\mathbf{b}_1 \in \mathbb{R}^{d_{\text{ff}}}$, and $\mathbf{b}_2 \in \mathbb{R}^{d_{\text{model}}}$ are learnable parameters, and d_{ff} is the feed-forward dimension.

Output Layer: After L encoder blocks, the final hidden state $\mathbf{H}^{(L)}$ is processed by an output layer to produce trading signals:

$$\mathbf{S} = \text{sigmoid}(\mathbf{H}^{(L)}W_S + \mathbf{b}_S) \quad (277)$$

where $W_S \in \mathbb{R}^{d_{\text{model}} \times 1}$ and $\mathbf{b}_S \in \mathbb{R}$ are learnable parameters, and $\mathbf{S} \in [0, 1]^T$ represents the trading signals (momentum scores) for each time step.

6.1.2 Regime-Aware Training

To make the Transformer model regime-aware, we incorporate regime information during training. Let $\mathbf{R} = \{r_1, r_2, \dots, r_T\}$ be the sequence of regime labels, where $r_t \in \{0, 1, 2\}$ represents the market regime at time t .

We implement regime-aware training through:

Regime Embedding:

$$\mathbf{R}_{\text{emb}} = \text{OneHot}(\mathbf{R})W_R + \mathbf{b}_R \quad (278)$$

where $W_R \in \mathbb{R}^{3 \times d_{\text{model}}}$ and $\mathbf{b}_R \in \mathbb{R}^{d_{\text{model}}}$ are learnable parameters.

The regime embedding is added to the input embedding:

$$\mathbf{H}^{(0)} = \mathbf{E} + \text{PE} + \mathbf{R}_{\text{emb}} \quad (279)$$

Regime-Specific Loss Weighting: The loss function is weighted based on the regime to emphasize performance in specific regimes:

$$\mathcal{L} = \sum_{t=1}^T w_{r_t} \cdot \text{Loss}(s_t, y_t) \quad (280)$$

where w_{r_t} is the weight for regime r_t , s_t is the predicted signal, y_t is the target signal, and Loss is a suitable loss function (e.g., mean squared error or binary cross-entropy).

Synthetic Data Augmentation: We augment the training data with synthetic samples generated by the GenAI component, particularly for underrepresented regimes. The augmented training set is:

$$\mathcal{D}_{\text{aug}} = \mathcal{D}_{\text{real}} \cup \mathcal{D}_{\text{syn}} \quad (281)$$

where $\mathcal{D}_{\text{real}}$ is the real data and \mathcal{D}_{syn} is the synthetic data.

6.1.3 Signal Generation and Position Sizing

The raw momentum scores from the Transformer model are converted into trading signals through a series of steps:

Signal Smoothing: To reduce noise, we apply an exponential moving average (EMA) to the raw scores:

$$\tilde{s}_t = \alpha \cdot s_t + (1 - \alpha) \cdot \tilde{s}_{t-1} \quad (282)$$

where $\alpha \in (0, 1)$ is the smoothing factor.

Signal Thresholding: We apply thresholds to determine the direction of the trade:

$$\text{direction}_t = \begin{cases} 1 & \text{if } \tilde{s}_t > \theta_{\text{long}} \\ -1 & \text{if } \tilde{s}_t < \theta_{\text{short}} \\ 0 & \text{otherwise} \end{cases} \quad (283)$$

where θ_{long} and θ_{short} are the long and short thresholds, respectively.

Position Sizing: The size of the position is determined by the strength of the signal and the volatility of the asset:

$$\text{size}_t = \text{direction}_t \cdot \min \left(\frac{|\tilde{s}_t - 0.5| \cdot 2}{\sigma_t \cdot \sqrt{T_{\text{hold}}}}, \text{size}_{\text{max}} \right) \quad (284)$$

where σ_t is the estimated volatility at time t , T_{hold} is the expected holding period, and size_{\max} is the maximum position size.

Regime-Specific Adjustments: The position size is adjusted based on the current market regime:

$$\text{size}_t^{\text{adj}} = \text{size}_t \cdot \gamma_{r_t} \quad (285)$$

where γ_{r_t} is the regime-specific adjustment factor.

6.1.4 Implementation Details

Our implementation of the Transformer-based Momentum Strategy includes several key considerations:

- **Feature Engineering:** We use a comprehensive set of features, including price-based features (returns, moving averages), volatility measures, technical indicators, and cross-asset features.
- **Model Architecture:** We use a Transformer encoder with 4 layers, 8 attention heads, and a model dimension of 256. The feed-forward dimension is set to 1024.
- **Training Procedure:** We train the model using the Adam optimizer with a learning rate of 10^{-4} and a batch size of 32. We implement a warm-up learning rate schedule and early stopping based on validation performance.
- **Hyperparameter Optimization:** We perform a grid search over key hyperparameters, including the number of layers, attention heads, learning rate, and signal thresholds. The optimal configuration is selected based on the Sharpe ratio on the validation set.
- **Risk Management:** We implement stop-loss and take-profit orders to manage risk. The stop-loss level is set as a multiple of the estimated volatility:

$$\text{stop_loss}_t = \text{entry_price}_t \cdot (1 - \text{direction}_t \cdot k_{\text{stop}} \cdot \sigma_t) \quad (286)$$

where k_{stop} is a constant that determines the stop-loss distance.

6.2 LSTM Short Signal Strategy

The LSTM Short Signal Strategy is specifically designed to identify opportunities for short selling, which can be particularly profitable during bear markets and market corrections.

6.2.1 LSTM Architecture and Mathematical Formulation

Long Short-Term Memory (LSTM) networks are a type of recurrent neural network (RNN) that can capture long-term dependencies in sequential data. The core of the LSTM Short Signal Strategy is an LSTM model that processes historical data to predict downward price movements.

The LSTM cell is defined by the following equations:

$$f_t = \sigma_g(W_f x_t + U_f h_{t-1} + b_f) \quad (287)$$

$$i_t = \sigma_g(W_i x_t + U_i h_{t-1} + b_i) \quad (288)$$

$$o_t = \sigma_g(W_o x_t + U_o h_{t-1} + b_o) \quad (289)$$

$$\tilde{c}_t = \tanh(W_c x_t + U_c h_{t-1} + b_c) \quad (290)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \quad (291)$$

$$h_t = o_t \odot \tanh(c_t) \quad (292)$$

where:

- x_t is the input vector at time t .
- h_t is the hidden state at time t .
- c_t is the cell state at time t .
- f_t , i_t , and o_t are the forget, input, and output gates, respectively.
- \tilde{c}_t is the candidate cell state.
- W_f , W_i , W_o , W_c , U_f , U_i , U_o , U_c , b_f , b_i , b_o , and b_c are learnable parameters.
- σ_g is the sigmoid function, and \tanh is the hyperbolic tangent function.
- \odot denotes element-wise multiplication.

The LSTM model processes a sequence of feature vectors $\mathbf{X} = \{x_1, x_2, \dots, x_T\}$ and outputs a sequence of hidden states $\mathbf{H} = \{h_1, h_2, \dots, h_T\}$. The hidden states are then passed through a fully connected layer to produce short signals:

$$s_t = \sigma(W_s h_t + b_s) \quad (293)$$

where W_s and b_s are learnable parameters, and σ is the sigmoid function. The output $s_t \in [0, 1]$ represents the probability of a profitable short opportunity at time t .

6.2.2 Asymmetric Loss Function for Short Bias

To bias the model towards identifying short opportunities, we use an asymmetric loss function that penalizes missed short opportunities more heavily than false positives:

$$\mathcal{L}_{\text{asym}}(s_t, y_t) = \begin{cases} \alpha \cdot (y_t - s_t)^2 & \text{if } y_t > s_t \text{ and } y_t = 1 \\ (y_t - s_t)^2 & \text{otherwise} \end{cases} \quad (294)$$

where $y_t \in \{0, 1\}$ is the binary label indicating whether a profitable short opportunity exists at time t , and $\alpha > 1$ is a hyperparameter that controls the asymmetry of the loss.

The labels y_t are generated based on future returns:

$$y_t = \begin{cases} 1 & \text{if } r_{t+1:t+h} < -\delta \\ 0 & \text{otherwise} \end{cases} \quad (295)$$

where $r_{t+1:t+h}$ is the cumulative return from time $t + 1$ to $t + h$, and δ is a threshold that defines a significant downward movement.

6.2.3 Regime-Specific Feature Importance

Different features may have different predictive power in different market regimes. We implement a regime-specific feature importance mechanism that allows the model to focus on the most relevant features for each regime.

Let $\mathbf{X}_t = \{x_{t,1}, x_{t,2}, \dots, x_{t,d}\}$ be the feature vector at time t , and let $r_t \in \{0, 1, 2\}$ be the market regime at time t . We compute regime-specific feature weights:

$$w_{j,r_t} = \text{softmax}(W_r)_{j,r_t} \quad (296)$$

where $W_r \in \mathbb{R}^{d \times 3}$ is a learnable parameter matrix, and softmax is applied along the feature dimension.

The weighted feature vector is:

$$\tilde{x}_{t,j} = w_{j,r_t} \cdot x_{t,j} \quad (297)$$

This weighted feature vector $\tilde{\mathbf{X}}_t$ is then fed into the LSTM model.

6.2.4 Short Signal Filtering and Execution

The raw short signals from the LSTM model are filtered and processed to generate executable trading decisions:

Signal Filtering: We apply a threshold to the raw signals to identify strong short opportunities:

$$\text{short_signal}_t = \begin{cases} 1 & \text{if } s_t > \theta_{\text{short}} \\ 0 & \text{otherwise} \end{cases} \quad (298)$$

where θ_{short} is the short threshold.

Confirmation Filters: To reduce false positives, we apply additional confirmation filters:

- **Trend Filter:** Ensures that the short signal aligns with the medium-term trend.

$$\text{trend_filter}_t = \begin{cases} 1 & \text{if } \text{SMA}_t(50) < \text{SMA}_t(200) \\ 0 & \text{otherwise} \end{cases} \quad (299)$$

- **Volatility Filter:** Ensures that the volatility is within an acceptable range.

$$\text{vol_filter}_t = \begin{cases} 1 & \text{if } \theta_{\text{vol_min}} < \sigma_t < \theta_{\text{vol_max}} \\ 0 & \text{otherwise} \end{cases} \quad (300)$$

- **Regime Filter:** Adjusts the signal based on the current market regime.

$$\text{regime_filter}_t = \begin{cases} 1 & \text{if } r_t \in \mathcal{R}_{\text{short}} \\ 0 & \text{otherwise} \end{cases} \quad (301)$$

where $\mathcal{R}_{\text{short}}$ is the set of regimes favorable for short selling.

The final short signal is:

$$\text{final_short_signal}_t = \text{short_signal}_t \cdot \text{trend_filter}_t \cdot \text{vol_filter}_t \cdot \text{regime_filter}_t \quad (302)$$

Position Sizing: The size of the short position is determined by the strength of the signal and the risk budget:

$$\text{size}_t = -\min\left(\frac{\text{risk_budget} \cdot s_t}{\sigma_t \cdot \sqrt{T_{\text{hold}}}}, \text{size}_{\max}\right) \quad (303)$$

where risk_budget is the amount of risk allocated to each trade, σ_t is the estimated volatility, T_{hold} is the expected holding period, and size_{\max} is the maximum position size.

Exit Strategy: We implement a dynamic exit strategy that closes the short position when:

- The profit target is reached:

$$\text{profit_target}_t = \text{entry_price}_t \cdot (1 + \text{target_multiple} \cdot \sigma_t) \quad (304)$$

- The stop-loss is triggered:

$$\text{stop_loss}_t = \text{entry_price}_t \cdot (1 - \text{stop_multiple} \cdot \sigma_t) \quad (305)$$

- The model generates a strong reversal signal:

$$\text{reversal_signal}_t = \begin{cases} 1 & \text{if } 1 - s_t > \theta_{\text{reversal}} \\ 0 & \text{otherwise} \end{cases} \quad (306)$$

6.2.5 Implementation Details

Our implementation of the LSTM Short Signal Strategy includes several key considerations:

- **Feature Engineering:** We focus on features that are particularly relevant for identifying short opportunities, including momentum indicators, volatility measures, and sentiment indicators.
- **Model Architecture:** We use a stacked LSTM with 2 layers, each with 128 hidden units. We implement dropout (0.2) between layers to prevent overfitting.
- **Training Procedure:** We train the model using the Adam optimizer with a learning rate of 10^{-3} and a batch size of 64. We implement a learning rate scheduler that reduces the learning rate when the validation loss plateaus.
- **Data Augmentation:** We augment the training data with synthetic samples generated by the GenAI component, particularly for bear market regimes, which are typically underrepresented in historical data.
- **Risk Management:** We implement strict risk management rules, including position size limits, stop-loss orders, and a maximum drawdown threshold that reduces exposure when the strategy experiences significant losses.

6.3 Kalman Filter Pairs Trading

Pairs trading is a market-neutral strategy that exploits temporary mispricings between related assets. The Kalman Filter Pairs Trading strategy uses a Kalman filter to dynamically estimate the relationship between pairs of assets and identify trading opportunities.

6.3.1 State-Space Model Formulation

We model the relationship between two assets using a state-space model, where the state represents the hedge ratio (beta) between the assets. Let P_t^A and P_t^B be the prices of assets A and B at time t . The log prices are denoted as $p_t^A = \log(P_t^A)$ and $p_t^B = \log(P_t^B)$.

The state-space model is defined as:

Observation Equation:

$$p_t^A = \beta_t p_t^B + \epsilon_t, \quad \epsilon_t \sim \mathcal{N}(0, R_t) \quad (307)$$

where β_t is the hedge ratio at time t , ϵ_t is the observation noise, and R_t is the observation noise variance.

State Equation:

$$\beta_t = \beta_{t-1} + \omega_t, \quad \omega_t \sim \mathcal{N}(0, Q_t) \quad (308)$$

where ω_t is the state noise, and Q_t is the state noise variance.

The spread between the assets is defined as:

$$s_t = p_t^A - \beta_t p_t^B \quad (309)$$

This spread is expected to be mean-reverting, and deviations from its mean represent trading opportunities.

6.3.2 Kalman Filter Estimation

The Kalman filter provides a recursive method to estimate the state β_t and its uncertainty. The filter consists of prediction and update steps:

Prediction Step:

$$\hat{\beta}_{t|t-1} = \hat{\beta}_{t-1|t-1} \quad (310)$$

$$P_{t|t-1} = P_{t-1|t-1} + Q_t \quad (311)$$

where $\hat{\beta}_{t|t-1}$ is the predicted state, and $P_{t|t-1}$ is the predicted state covariance.

Update Step:

$$\tilde{y}_t = p_t^A - \hat{\beta}_{t|t-1} p_t^B \quad (312)$$

$$S_t = p_t^B P_{t|t-1} p_t^B + R_t \quad (313)$$

$$K_t = P_{t|t-1} p_t^B / S_t \quad (314)$$

$$\hat{\beta}_{t|t} = \hat{\beta}_{t|t-1} + K_t \tilde{y}_t \quad (315)$$

$$P_{t|t} = (1 - K_t p_t^B) P_{t|t-1} \quad (316)$$

where \tilde{y}_t is the innovation (prediction error), S_t is the innovation variance, K_t is the Kalman gain, $\hat{\beta}_{t|t}$ is the updated state estimate, and $P_{t|t}$ is the updated state covariance.

6.3.3 Regime-Dependent Noise Variances

The noise variances Q_t and R_t can vary across different market regimes. We model them as regime-dependent:

$$Q_t = Q_{r_t} \quad (317)$$

$$R_t = R_{r_t} \quad (318)$$

where $r_t \in \{0, 1, 2\}$ is the market regime at time t , and Q_{r_t} and R_{r_t} are the regime-specific noise variances.

These regime-dependent variances allow the model to adapt to changing market conditions. For example, during high-volatility regimes, the observation noise variance R_t might be higher, reflecting increased uncertainty in the relationship between the assets.

6.3.4 Trading Signal Generation

Trading signals are generated based on the normalized spread:

$$z_t = \frac{s_t - \mu_t}{\sigma_t} \quad (319)$$

where μ_t and σ_t are the estimated mean and standard deviation of the spread, computed using an exponential moving average:

$$\mu_t = \lambda \mu_{t-1} + (1 - \lambda) s_t \quad (320)$$

$$\sigma_t^2 = \lambda \sigma_{t-1}^2 + (1 - \lambda)(s_t - \mu_t)^2 \quad (321)$$

where $\lambda \in (0, 1)$ is the decay factor.

Trading signals are generated based on thresholds:

$$\text{signal}_t = \begin{cases} 1 & \text{if } z_t < -\theta_{\text{entry}} \text{ and } \text{position}_{t-1} \neq 1 \\ -1 & \text{if } z_t > \theta_{\text{entry}} \text{ and } \text{position}_{t-1} \neq -1 \\ 0 & \text{if } |z_t| < \theta_{\text{exit}} \text{ and } \text{position}_{t-1} \neq 0 \\ \text{position}_{t-1} & \text{otherwise} \end{cases} \quad (322)$$

where θ_{entry} and θ_{exit} are the entry and exit thresholds, respectively, and position_{t-1} is the position at time $t - 1$.

A signal of 1 indicates a long position in asset A and a short position in asset B (with size proportional to β_t), while a signal of -1 indicates the opposite.

6.3.5 Regime-Specific Threshold Adjustment

The entry and exit thresholds are adjusted based on the current market regime:

$$\theta_{\text{entry}}^{\text{adj}} = \theta_{\text{entry}} \cdot \gamma_{r_t}^{\text{entry}} \quad (323)$$

$$\theta_{\text{exit}}^{\text{adj}} = \theta_{\text{exit}} \cdot \gamma_{r_t}^{\text{exit}} \quad (324)$$

where $\gamma_{r_t}^{\text{entry}}$ and $\gamma_{r_t}^{\text{exit}}$ are regime-specific adjustment factors.

These adjustments allow the strategy to be more or less aggressive depending on the market regime. For example, during high-volatility regimes, the entry threshold might be increased to reduce the number of false signals.

6.3.6 Implementation Details

Our implementation of the Kalman Filter Pairs Trading strategy includes several key considerations:

- **Pair Selection:** We select pairs of assets that exhibit historical cointegration. The cointegration is tested using the Augmented Dickey-Fuller test on the residuals of the regression:

$$p_t^A = \alpha + \beta p_t^B + \epsilon_t \quad (325)$$

- **Parameter Estimation:** The noise variances Q_{r_t} and R_{r_t} are estimated using maximum likelihood estimation on historical data for each regime.
- **Position Sizing:** The size of the position is determined by the strength of the signal (normalized spread) and the estimated uncertainty of the hedge ratio:

$$\text{size}_t = \text{signal}_t \cdot \min \left(\frac{|z_t| - \theta_{\text{entry}}}{P_{t|t}}, \text{size}_{\max} \right) \quad (326)$$

- **Risk Management:** We implement stop-loss orders based on the maximum acceptable loss for each trade, as well as a time-based exit that closes positions that have not converged within a specified time frame.
- **Synthetic Data Integration:** We use synthetic data generated by the GenAI component to test the strategy's performance in rare market regimes and to optimize the regime-specific parameters.

6.4 Strategy Decay Detection and Mitigation

Trading strategies often lose their effectiveness over time due to changing market dynamics, increased competition, or structural changes in the market. We implement a systematic approach to detect and mitigate strategy decay.

6.4.1 Performance Monitoring Framework

We monitor the performance of each strategy using a set of key performance indicators (KPIs):

- **Rolling Sharpe Ratio:** Measures the risk-adjusted return over a rolling window.

$$\text{Sharpe}_t(w) = \frac{\bar{r}_t(w) - r_f}{\sigma_t(w)} \quad (327)$$

where $\bar{r}_t(w)$ is the average return over the past w periods, r_f is the risk-free rate, and $\sigma_t(w)$ is the standard deviation of returns over the past w periods.

- **Rolling Alpha:** Measures the excess return relative to a benchmark.

$$\text{Alpha}_t(w) = \bar{r}_t(w) - \beta \cdot \bar{r}_t^{\text{benchmark}}(w) \quad (328)$$

where β is the strategy's beta to the benchmark, and $\bar{r}_t^{\text{benchmark}}(w)$ is the average benchmark return over the past w periods.

- **Win Rate:** The proportion of profitable trades.

$$\text{WinRate}_t(w) = \frac{\text{Number of profitable trades in } [t-w+1, t]}{\text{Total number of trades in } [t-w+1, t]} \quad (329)$$

- **Profit Factor:** The ratio of gross profits to gross losses.

$$\text{ProfitFactor}_t(w) = \frac{\text{Sum of profits in } [t-w+1, t]}{\text{Sum of losses in } [t-w+1, t]} \quad (330)$$

6.4.2 Statistical Tests for Decay Detection

We implement several statistical tests to detect significant changes in strategy performance:

Chow Test: Tests for structural breaks in the relationship between strategy returns and market factors. The test statistic is:

$$F = \frac{(RSS_c - (RSS_1 + RSS_2))/k}{(RSS_1 + RSS_2)/(n_1 + n_2 - 2k)} \quad (331)$$

where RSS_c is the residual sum of squares from the combined model, RSS_1 and RSS_2 are the residual sum of squares from the models estimated on the two subsamples, k is the number of parameters, and n_1 and n_2 are the sample sizes of the two subsamples.

CUSUM Test: Detects changes in the mean of the strategy returns. The test statistic is:

$$W_t = \frac{1}{\hat{\sigma}} \sum_{i=1}^t (r_i - \bar{r}) \quad (332)$$

where \bar{r} is the mean return over the entire sample, and $\hat{\sigma}$ is an estimate of the standard deviation.

Drawdown Analysis: Monitors the maximum drawdown of the strategy and compares it to historical drawdowns. A drawdown that exceeds historical levels may indicate strategy decay.

$$DD_t = \max_{i \in [0, t]} \left(1 - \frac{V_t}{V_i} \right) \quad (333)$$

where V_t is the strategy value at time t .

6.4.3 Regime-Specific Decay Detection

Strategy decay may manifest differently across market regimes. We implement regime-specific decay detection by monitoring performance separately for each regime:

$$\text{KPI}_{r_t}(w) = \frac{\sum_{i=t-w+1}^t \text{KPI}_i \cdot \mathbb{I}(r_i = r_t)}{\sum_{i=t-w+1}^t \mathbb{I}(r_i = r_t)} \quad (334)$$

where $\mathbb{I}(r_i = r_t)$ is an indicator function that equals 1 if the regime at time i is the same as the current regime, and 0 otherwise.

This allows us to detect decay that is specific to certain market conditions, even if the overall performance remains acceptable.

6.4.4 Mitigation Strategies

When strategy decay is detected, we implement several mitigation strategies:

Parameter Re-optimization: We re-optimize the strategy parameters using recent data, with a focus on the regime where decay was detected:

$$\theta^* = \arg \max_{\theta} \sum_{i=t-w+1}^t \gamma_{r_i} \cdot \text{Performance}_i(\theta) \quad (335)$$

where θ represents the strategy parameters, γ_{r_i} is a regime-specific weight, and $\text{Performance}_i(\theta)$ is the performance metric at time i with parameters θ .

Feature Importance Re-evaluation: We re-evaluate the importance of features used in the strategy and adjust the feature set accordingly:

$$\text{Importance}(f) = \frac{\text{Performance}(\text{all features}) - \text{Performance}(\text{all features except } f)}{\text{Performance}(\text{all features})} \quad (336)$$

Features with negative or low importance are removed or down-weighted.

Synthetic Data Adaptation: We generate synthetic data that reflects the current market conditions and use it to adapt the strategy:

$$\mathcal{D}_{\text{adapt}} = \mathcal{D}_{\text{recent}} \cup \mathcal{D}_{\text{syn}} \quad (337)$$

where $\mathcal{D}_{\text{recent}}$ is the recent real data and \mathcal{D}_{syn} is the synthetic data generated to match the current market conditions.

Strategy Rotation: We implement a meta-strategy that rotates between different strategies based on their recent performance and the current market regime:

$$w_{s,t} = \frac{\exp(\lambda \cdot \text{Performance}_{s,t}(w))}{\sum_{s'} \exp(\lambda \cdot \text{Performance}_{s',t}(w))} \quad (338)$$

where $w_{s,t}$ is the weight assigned to strategy s at time t , $\text{Performance}_{s,t}(w)$ is the performance of strategy s over the past w periods, and λ is a parameter that controls the concentration of weights.

6.4.5 Implementation Details

Our implementation of the Strategy Decay Detection and Mitigation framework includes several key considerations:

- **Monitoring Frequency:** We monitor strategy performance daily but perform formal decay detection tests weekly to balance responsiveness with stability.
- **Decay Thresholds:** We set thresholds for each KPI based on historical performance and acceptable levels of deterioration. These thresholds are regime-specific to account for different performance expectations across regimes.
- **Mitigation Triggers:** Mitigation actions are triggered when multiple decay indicators are activated simultaneously or when a single indicator shows extreme deterioration.

- **Adaptation Speed:** The speed of adaptation is controlled by hyperparameters that determine how quickly the strategy responds to detected decay. This balances the need for responsiveness with the risk of overreacting to temporary fluctuations.
- **Performance Attribution:** We implement a performance attribution system that identifies the specific components of the strategy that are experiencing decay, allowing for targeted mitigation efforts.

6.5 Strategy Ensemble and Portfolio Construction

To maximize risk-adjusted returns and robustness across different market regimes, we implement a strategy ensemble that combines multiple strategies into a coherent portfolio.

6.5.1 Strategy Correlation Analysis

We analyze the correlations between strategy returns to ensure diversification:

$$\rho_{s_1, s_2} = \frac{\text{Cov}(r_{s_1}, r_{s_2})}{\sigma_{s_1} \sigma_{s_2}} \quad (339)$$

where r_{s_1} and r_{s_2} are the returns of strategies s_1 and s_2 , and σ_{s_1} and σ_{s_2} are their standard deviations.

We also compute regime-specific correlations:

$$\rho_{s_1, s_2}^{r_t} = \frac{\text{Cov}(r_{s_1}, r_{s_2} | r_t)}{\sigma_{s_1}^{r_t} \sigma_{s_2}^{r_t}} \quad (340)$$

where $\text{Cov}(r_{s_1}, r_{s_2} | r_t)$ is the covariance conditional on regime r_t , and $\sigma_{s_1}^{r_t}$ and $\sigma_{s_2}^{r_t}$ are the conditional standard deviations.

6.5.2 Regime-Dependent Strategy Allocation

We allocate capital to strategies based on their expected performance in the current market regime:

$$w_{s,t} = \frac{\text{Performance}_{s,r_t} \cdot (1 - \beta_{s,r_t})}{\sum_{s'} \text{Performance}_{s',r_t} \cdot (1 - \beta_{s',r_t})} \quad (341)$$

where $\text{Performance}_{s,r_t}$ is the historical performance of strategy s in regime r_t , and β_{s,r_t} is the beta of strategy s to the market in regime r_t .

This allocation favors strategies that have performed well in the current regime and have low market beta, providing better risk-adjusted returns.

6.5.3 Risk Parity Portfolio Construction

We implement a risk parity approach to portfolio construction, which allocates risk equally across strategies:

$$w_s \propto \frac{1}{\sigma_s} \quad (342)$$

where σ_s is the volatility of strategy s .

To account for correlations between strategies, we use the full covariance matrix:

$$w_s \propto \frac{1}{\sum_{s'} \rho_{s,s'} \sigma_s \sigma_{s'}} \quad (343)$$

where $\rho_{s,s'}$ is the correlation between strategies s and s' .

6.5.4 Dynamic Risk Budgeting

We dynamically adjust the risk budget based on market conditions and strategy performance:

$$\text{RiskBudget}_t = \text{BaseRisk} \cdot \text{MarketFactor}_t \cdot \text{PerformanceFactor}_t \quad (344)$$

where:

- BaseRisk is the baseline risk budget.
- MarketFactor_t adjusts the risk based on market conditions:

$$\text{MarketFactor}_t = \exp(-\lambda_{\text{market}} \cdot \text{MarketRisk}_t) \quad (345)$$

where MarketRisk_t is a measure of market risk (e.g., VIX index), and λ_{market} is a sensitivity parameter.

- $\text{PerformanceFactor}_t$ adjusts the risk based on recent performance:

$$\text{PerformanceFactor}_t = \exp(\lambda_{\text{perf}} \cdot \text{RecentPerformance}_t) \quad (346)$$

where $\text{RecentPerformance}_t$ is a measure of recent performance (e.g., rolling Sharpe ratio), and λ_{perf} is a sensitivity parameter.

6.5.5 Implementation Details

Our implementation of the Strategy Ensemble and Portfolio Construction framework includes several key considerations:

- **Rebalancing Frequency:** We rebalance the portfolio weekly to balance transaction costs with the need to adapt to changing market conditions.
- **Turnover Constraints:** We implement turnover constraints to limit transaction costs:

$$|w_{s,t} - w_{s,t-1}| \leq \text{MaxTurnover} \quad (347)$$

- **Regime Transition Smoothing:** During regime transitions, we gradually adjust the portfolio weights to avoid large, sudden changes:

$$w_{s,t} = (1 - \alpha) \cdot w_{s,t-1} + \alpha \cdot w_{s,t}^{\text{target}} \quad (348)$$

where $\alpha \in (0, 1)$ is a smoothing parameter.

- **Stress Testing:** We stress test the portfolio using historical scenarios and synthetic data to ensure robustness across a wide range of market conditions.
- **Performance Attribution:** We implement a performance attribution system that decomposes returns into strategy-specific and regime-specific components, providing insights into the sources of portfolio performance.

7 Risk Forecasting Implementation

This section details the implementation of our risk forecasting component, which leverages market regime detection and synthetic data to provide accurate risk estimates across different market conditions. We present the mathematical formulation, implementation details, and validation framework for each risk forecasting method.

7.1 REGARCH-EVT: Regime-Switching GARCH with Extreme Value Theory

REGARCH-EVT combines regime-switching models, GARCH volatility modeling, and Extreme Value Theory (EVT) to provide accurate risk forecasts, particularly for extreme events.

7.1.1 Regime-Switching GARCH Model

The Regime-Switching GARCH (REGARCH) model allows the parameters of the GARCH process to vary depending on the market regime. Let r_t be the return at time t , and let $z_t \in \{0, 1, 2\}$ be the market regime at time t . The REGARCH model is defined as:

$$r_t = \mu_{z_t} + \epsilon_t \quad (349)$$

$$\epsilon_t = \sigma_t \eta_t, \quad \eta_t \sim F_{z_t}(0, 1) \quad (350)$$

$$\sigma_t^2 = \omega_{z_t} + \sum_{i=1}^p \alpha_{i,z_t} \epsilon_{t-i}^2 + \sum_{j=1}^q \beta_{j,z_t} \sigma_{t-j}^2 \quad (351)$$

where:

- μ_{z_t} is the regime-dependent mean return.
- σ_t is the conditional volatility at time t .
- η_t is a standardized innovation with distribution $F_{z_t}(0, 1)$ that depends on the regime.
- ω_{z_t} , α_{i,z_t} , and β_{j,z_t} are the regime-dependent GARCH parameters.

The regime-dependent parameters allow the model to capture different volatility dynamics across market regimes. For example, in high-volatility regimes, the persistence parameter β_{1,z_t} might be higher, reflecting the stronger volatility clustering observed during market stress.

7.1.2 Parameter Estimation via Maximum Likelihood

The parameters of the REGARCH model are estimated using maximum likelihood estimation (MLE). Given a sequence of returns $\{r_1, r_2, \dots, r_T\}$ and regimes $\{z_1, z_2, \dots, z_T\}$, the log-likelihood function is:

$$\mathcal{L}(\theta) = \sum_{t=1}^T \log f_{z_t}(r_t | \mathcal{F}_{t-1}; \theta) \quad (352)$$

where $f_{z_t}(r_t | \mathcal{F}_{t-1}; \theta)$ is the conditional density of r_t given the information set \mathcal{F}_{t-1} and regime z_t , and θ represents all model parameters.

For a Gaussian innovation distribution, the log-likelihood contribution at time t is:

$$\log f_{z_t}(r_t | \mathcal{F}_{t-1}; \theta) = -\frac{1}{2} \log(2\pi) - \frac{1}{2} \log(\sigma_t^2) - \frac{(r_t - \mu_{z_t})^2}{2\sigma_t^2} \quad (353)$$

For non-Gaussian innovations, such as Student's t-distribution, the log-likelihood is adjusted accordingly.

The optimization problem is:

$$\hat{\theta} = \arg \max_{\theta} \mathcal{L}(\theta) \quad (354)$$

This is solved using numerical optimization methods such as BFGS or Nelder-Mead.

7.1.3 Extreme Value Theory for Tail Modeling

Standard GARCH models with Gaussian innovations often underestimate the probability of extreme events. To address this, we apply Extreme Value Theory (EVT) to model the tails of the standardized residuals.

Let $\{\eta_i\}$ be the standardized residuals from the REGARCH model. We apply the Peaks-Over-Threshold (POT) approach, which models exceedances over a high threshold u . For a given regime z , the distribution of exceedances $y = \eta_i - u$ (conditional on $\eta_i > u$) is approximated by the Generalized Pareto Distribution (GPD):

$$G_z(y; \xi_z, \beta_z) = \begin{cases} 1 - \left(1 + \frac{\xi_z y}{\beta_z}\right)^{-1/\xi_z} & \text{if } \xi_z \neq 0 \\ 1 - \exp\left(-\frac{y}{\beta_z}\right) & \text{if } \xi_z = 0 \end{cases} \quad (355)$$

where ξ_z is the shape parameter and β_z is the scale parameter for regime z . The shape parameter ξ_z determines the tail behavior:

- $\xi_z > 0$: Heavy-tailed distribution (e.g., Pareto, Student's t)
- $\xi_z = 0$: Exponential distribution
- $\xi_z < 0$: Distribution with a finite upper bound

The parameters ξ_z and β_z are estimated using maximum likelihood estimation:

$$(\hat{\xi}_z, \hat{\beta}_z) = \arg \max_{\xi_z, \beta_z} \sum_{i: \eta_i > u, z_i = z} \log g_z(\eta_i - u; \xi_z, \beta_z) \quad (356)$$

where $g_z(y; \xi_z, \beta_z) = \frac{d}{dy} G_z(y; \xi_z, \beta_z)$ is the GPD density function.

7.1.4 Value at Risk (VaR) and Expected Shortfall (ES) Calculation

Using the REGARCH-EVT model, we can compute risk measures such as Value at Risk (VaR) and Expected Shortfall (ES).

For a given confidence level α (e.g., 0.99) and horizon h , the VaR is defined as:

$$\text{VaR}_{\alpha, t+h} = \inf\{x : P(r_{t+h} \leq x | \mathcal{F}_t) \geq \alpha\} \quad (357)$$

The Expected Shortfall (ES), also known as Conditional VaR (CVaR), is the expected loss given that the loss exceeds the VaR:

$$\text{ES}_{\alpha,t+h} = \mathbb{E}[r_{t+h}|r_{t+h} \leq \text{VaR}_{\alpha,t+h}, \mathcal{F}_t] \quad (358)$$

To compute these risk measures using the REGARCH-EVT model, we first forecast the conditional volatility σ_{t+h} using the REGARCH model. For a one-step-ahead forecast ($h = 1$), this is:

$$\sigma_{t+1}^2 = \omega_{z_t} + \sum_{i=1}^p \alpha_{i,z_t} \epsilon_{t+1-i}^2 + \sum_{j=1}^q \beta_{j,z_t} \sigma_{t+1-j}^2 \quad (359)$$

For multi-step forecasts ($h > 1$), we use recursive substitution or Monte Carlo simulation.

Next, we compute the quantile of the standardized residuals using the EVT model. For a given regime z and confidence level α , the quantile $q_{\alpha,z}$ is:

$$q_{\alpha,z} = \begin{cases} u + \frac{\beta_z}{\xi_z} \left[\left(\frac{n}{N_u} (1 - \alpha) \right)^{-\xi_z} - 1 \right] & \text{if } \alpha > 1 - \frac{N_u}{n} \\ F_z^{-1}(\alpha) & \text{otherwise} \end{cases} \quad (360)$$

where n is the total number of observations, N_u is the number of exceedances over the threshold u , and F_z is the empirical distribution function of the standardized residuals for regime z .

Finally, the VaR and ES are computed as:

$$\text{VaR}_{\alpha,t+h} = \mu_{z_t} + \sigma_{t+h} \cdot q_{\alpha,z_t} \quad (361)$$

$$\text{ES}_{\alpha,t+h} = \mu_{z_t} + \sigma_{t+h} \cdot \mathbb{E}[\eta | \eta \leq q_{\alpha,z_t}, z = z_t] \quad (362)$$

where the conditional expectation of the standardized residuals is computed using the GPD for the tail:

$$\mathbb{E}[\eta | \eta \leq q_{\alpha,z}, z] = \frac{q_{\alpha,z} + \frac{\beta_z - \xi_z u}{1 - \xi_z} \left[1 - \left(\frac{1 - \alpha}{1 - F_z(u)} \right)^{\xi_z - 1} \right]}{1 - \alpha} \quad \text{for } \xi_z < 1 \quad (363)$$

7.1.5 Implementation Details

Our implementation of the REGARCH-EVT model includes several key considerations:

- **Threshold Selection:** The threshold u for the EVT component is selected using a combination of graphical methods (mean excess plot) and statistical tests (goodness-of-fit tests). We implement an adaptive threshold that depends on the regime, with higher thresholds for more volatile regimes.
- **Parameter Constraints:** We impose constraints on the GARCH parameters to

ensure stationarity and positivity of the conditional variance:

$$\omega_z > 0 \quad (364)$$

$$\alpha_{i,z} \geq 0 \quad \forall i, z \quad (365)$$

$$\beta_{j,z} \geq 0 \quad \forall j, z \quad (366)$$

$$\sum_{i=1}^p \alpha_{i,z} + \sum_{j=1}^q \beta_{j,z} < 1 \quad \forall z \quad (367)$$

- **Regime Transition Handling:** When transitioning between regimes, we implement a smooth transition of the volatility forecast to avoid discontinuities:

$$\sigma_{t+1}^2 = \lambda \cdot \sigma_{t+1,z_t}^2 + (1 - \lambda) \cdot \sigma_{t+1,z_{t+1}}^2 \quad (368)$$

where $\lambda \in [0, 1]$ is a smoothing parameter, and $\sigma_{t+1,z}^2$ is the volatility forecast assuming regime z .

- **Synthetic Data Integration:** We augment the historical data with synthetic data generated by the GenAI component, particularly for rare regimes with limited historical observations. This improves the estimation of the EVT parameters for these regimes.

7.2 Quantile Random Forest for Direct VaR Estimation

Quantile Random Forests provide a non-parametric approach to estimating conditional quantiles directly from the data, without making specific distributional assumptions.

7.2.1 Random Forest Framework

A Random Forest is an ensemble of decision trees, where each tree is trained on a bootstrap sample of the data and with a random subset of features at each split. Let $\{(X_1, y_1), (X_2, y_2), \dots, (X_n, y_n)\}$ be the training data, where $X_i \in \mathbb{R}^d$ is a feature vector and $y_i \in \mathbb{R}$ is the target variable (e.g., return).

For each tree t in the forest, a bootstrap sample \mathcal{D}_t is drawn from the training data. The tree is grown by recursively splitting the data based on the features. At each node, a random subset of features is considered, and the best split is chosen to minimize the mean squared error:

$$\text{MSE}(S, j, s) = \frac{|S_L|}{|S|} \cdot \text{Var}(y_i : X_{i,j} \leq s) + \frac{|S_R|}{|S|} \cdot \text{Var}(y_i : X_{i,j} > s) \quad (369)$$

where S is the set of samples at the current node, j is the feature index, s is the split point, $S_L = \{i \in S : X_{i,j} \leq s\}$ is the left child node, and $S_R = \{i \in S : X_{i,j} > s\}$ is the right child node.

The splitting continues until a stopping criterion is met (e.g., minimum node size or maximum depth). The prediction for a new sample X is the average of the predictions from all trees:

$$\hat{y}(X) = \frac{1}{T} \sum_{t=1}^T \hat{y}_t(X) \quad (370)$$

where $\hat{y}_t(X)$ is the prediction from tree t , which is the average of the target values in the leaf node where X falls.

7.2.2 Quantile Estimation

Quantile Random Forests extend the standard Random Forest to estimate conditional quantiles rather than conditional means. Instead of averaging the target values in each leaf node, Quantile Random Forests keep track of the empirical distribution of the target values.

Let $L(X)$ be the set of training samples that fall into the same leaf node as X in a given tree. The weight of each training sample i for the prediction at X is:

$$w_i(X) = \frac{1}{T} \sum_{t=1}^T \frac{\mathbb{I}(i \in L_t(X))}{|L_t(X)|} \quad (371)$$

where $L_t(X)$ is the leaf node in tree t where X falls, and $\mathbb{I}(\cdot)$ is the indicator function. The conditional cumulative distribution function (CDF) at X is estimated as:

$$\hat{F}(y|X) = \sum_{i=1}^n w_i(X) \cdot \mathbb{I}(y_i \leq y) \quad (372)$$

The conditional quantile at level α is then:

$$\hat{Q}_\alpha(X) = \inf\{y : \hat{F}(y|X) \geq \alpha\} \quad (373)$$

For Value at Risk (VaR) estimation at confidence level α (e.g., 0.99), we compute:

$$\widehat{\text{VaR}}_\alpha(X) = -\hat{Q}_{1-\alpha}(X) \quad (374)$$

where the negative sign accounts for the convention that VaR is typically reported as a positive number representing a loss.

7.2.3 Feature Engineering for Risk Forecasting

The performance of Quantile Random Forests depends heavily on the features used. We implement a comprehensive set of features specifically designed for risk forecasting:

- **Volatility Features:** Historical volatility measures at different time scales, GARCH volatility forecasts, implied volatility (when available), and volatility of volatility.
- **Return Features:** Past returns at different horizons, cumulative returns, drawdowns, and return acceleration (change in return momentum).
- **Distributional Features:** Skewness and kurtosis of returns over different windows, quantile-based measures of tail behavior, and jump detection indicators.
- **Market Regime Indicators:** Probabilities of different market regimes, regime transition indicators, and regime duration.
- **Cross-Asset Features:** Correlations with other assets, systematic risk factors (e.g., market, size, value), and cross-asset volatility spillovers.
- **Macroeconomic Indicators:** Interest rates, yield curve slopes, credit spreads, and economic surprise indices.

7.2.4 Regime-Specific Forests

To capture regime-specific risk dynamics, we implement regime-specific Quantile Random Forests. For each regime $z \in \{0, 1, 2\}$, we train a separate forest using only the data from that regime:

$$\hat{Q}_{\alpha,z}(X) = \inf\{y : \hat{F}_z(y|X) \geq \alpha\} \quad (375)$$

where $\hat{F}_z(y|X)$ is the conditional CDF estimated using only the data from regime z .

The final quantile estimate is a weighted average of the regime-specific estimates, where the weights are the probabilities of each regime given the current features:

$$\hat{Q}_{\alpha}(X) = \sum_{z=0}^2 P(z|X) \cdot \hat{Q}_{\alpha,z}(X) \quad (376)$$

where $P(z|X)$ is the probability of regime z given features X , as estimated by the regime detection component.

7.2.5 Synthetic Data Augmentation

For rare regimes with limited historical data, we augment the training data with synthetic samples generated by the GenAI component. Let $\mathcal{D}_{\text{real},z}$ be the real data from regime z , and let $\mathcal{D}_{\text{syn},z}$ be the synthetic data for regime z . The augmented training data for regime z is:

$$\mathcal{D}_{\text{aug},z} = \mathcal{D}_{\text{real},z} \cup \mathcal{D}_{\text{syn},z} \quad (377)$$

To account for potential differences between real and synthetic data, we implement a weighted training approach where synthetic samples are assigned lower weights:

$$w_i = \begin{cases} 1 & \text{if sample } i \text{ is real} \\ w_{\text{syn}} & \text{if sample } i \text{ is synthetic} \end{cases} \quad (378)$$

where $w_{\text{syn}} \in (0, 1)$ is a hyperparameter that controls the influence of synthetic data.

7.2.6 Implementation Details

Our implementation of the Quantile Random Forest for VaR estimation includes several key considerations:

- **Forest Parameters:** We use 500 trees, a maximum depth of 10, and a minimum leaf size of 10 samples. These parameters are tuned using cross-validation to balance bias and variance.
- **Feature Selection:** We implement a feature selection step that identifies the most relevant features for each regime. This is done using permutation importance, which measures the increase in prediction error when a feature is randomly permuted.
- **Quantile Smoothing:** To ensure monotonicity and smoothness of the estimated quantiles across different α levels, we implement a post-processing step that fits a smooth curve to the raw quantile estimates.

- **Uncertainty Estimation:** We estimate the uncertainty of the VaR estimates using bootstrap methods. For each bootstrap sample, we train a new forest and compute the VaR. The distribution of these bootstrap VaR estimates provides a measure of uncertainty.
- **Backtesting Framework:** We implement a rigorous backtesting framework that evaluates the performance of the VaR estimates using tests such as the Kupiec test for unconditional coverage and the Christoffersen test for conditional coverage.

7.3 Synthetic Regime Bootstrapping

Synthetic Regime Bootstrapping combines traditional bootstrap methods with synthetic data generation to provide robust risk estimates, particularly for rare market regimes.

7.3.1 Traditional Bootstrap Methods

Bootstrap methods involve resampling from the observed data to estimate the sampling distribution of a statistic. Let $\{r_1, r_2, \dots, r_n\}$ be a sample of returns. The bootstrap procedure generates B bootstrap samples, each of size n , by sampling with replacement from the original data. For each bootstrap sample b , we compute the statistic of interest θ_b (e.g., VaR or ES). The bootstrap estimate of the statistic is the average across all bootstrap samples:

$$\hat{\theta}_{\text{boot}} = \frac{1}{B} \sum_{b=1}^B \theta_b \quad (379)$$

The bootstrap estimate of the standard error is:

$$\hat{SE}_{\text{boot}} = \sqrt{\frac{1}{B-1} \sum_{b=1}^B (\theta_b - \hat{\theta}_{\text{boot}})^2} \quad (380)$$

7.3.2 Regime-Specific Bootstrap

In financial time series, the distribution of returns can vary significantly across different market regimes. To account for this, we implement a regime-specific bootstrap procedure. Let $\{r_1, r_2, \dots, r_n\}$ be the returns and $\{z_1, z_2, \dots, z_n\}$ be the corresponding regime labels. For each regime z , we define the set of returns from that regime:

$$\mathcal{R}_z = \{r_i : z_i = z\} \quad (381)$$

The regime-specific bootstrap generates samples by resampling from \mathcal{R}_z for each regime z . This preserves the regime-specific distribution of returns.

7.3.3 Synthetic Data Integration

For rare regimes with limited historical data, the bootstrap samples may not adequately capture the tail behavior. To address this, we augment the bootstrap procedure with synthetic data generated by the GenAI component.

Let $\mathcal{R}_{\text{syn},z}$ be the set of synthetic returns generated for regime z . The augmented set of returns for regime z is:

$$\mathcal{R}_{\text{aug},z} = \mathcal{R}_z \cup \mathcal{R}_{\text{syn},z} \quad (382)$$

The synthetic regime bootstrap generates samples by resampling from $\mathcal{R}_{\text{aug},z}$ for each regime z . To control the influence of synthetic data, we implement a weighted resampling approach where the probability of selecting a synthetic return is adjusted:

$$P(r_i \text{ is selected}) \propto \begin{cases} 1 & \text{if } r_i \in \mathcal{R}_z \\ w_{\text{syn}} \cdot \frac{|\mathcal{R}_z|}{|\mathcal{R}_{\text{syn},z}|} & \text{if } r_i \in \mathcal{R}_{\text{syn},z} \end{cases} \quad (383)$$

where $w_{\text{syn}} \in (0, 1)$ is a hyperparameter that controls the influence of synthetic data.

7.3.4 Block Bootstrap for Time Series

Financial returns often exhibit serial dependence, particularly in volatility. To preserve this dependence structure, we implement a block bootstrap approach. Instead of resampling individual returns, we resample blocks of consecutive returns.

Let $\{r_1, r_2, \dots, r_n\}$ be the returns and $\{z_1, z_2, \dots, z_n\}$ be the regime labels. We identify contiguous blocks of returns from the same regime:

$$\mathcal{B}_{z,j} = \{r_i, r_{i+1}, \dots, r_{i+l-1} : z_i = z_{i+1} = \dots = z_{i+l-1} = z\} \quad (384)$$

where $\mathcal{B}_{z,j}$ is the j -th block of returns from regime z , and l is the block length.

The block bootstrap generates samples by resampling blocks from each regime and concatenating them. This preserves both the regime-specific distribution and the short-term dependence structure.

7.3.5 Scenario Generation and Risk Estimation

Using the synthetic regime bootstrap, we generate a large number of scenarios for future returns. Each scenario consists of a sequence of returns $\{r_{t+1}, r_{t+2}, \dots, r_{t+h}\}$ over the forecast horizon h .

For each scenario s , we compute the cumulative return:

$$R_s = \sum_{i=1}^h r_{t+i,s} \quad (385)$$

The VaR at confidence level α is estimated as the α -quantile of the scenario returns:

$$\widehat{\text{VaR}}_\alpha = -\text{Quantile}_\alpha(\{R_1, R_2, \dots, R_S\}) \quad (386)$$

The ES is estimated as the average of the scenario returns that exceed the VaR:

$$\widehat{\text{ES}}_\alpha = -\frac{1}{|\{s : R_s \leq -\widehat{\text{VaR}}_\alpha\}|} \sum_{s: R_s \leq -\widehat{\text{VaR}}_\alpha} R_s \quad (387)$$

7.3.6 Implementation Details

Our implementation of the Synthetic Regime Bootstrapping method includes several key considerations:

- **Block Length Selection:** The block length l is selected based on the autocorrelation structure of the returns. We use a data-driven approach that sets l to capture the significant autocorrelation in absolute returns.
- **Regime Transition Modeling:** To model regime transitions in multi-period scenarios, we implement a Markov chain model estimated from the historical regime sequence. The transition probabilities are:

$$P(z_{t+1} = j | z_t = i) = \frac{\text{Number of transitions from regime } i \text{ to regime } j}{\text{Number of observations in regime } i} \quad (388)$$

- **Synthetic Data Quality Control:** We implement quality control checks for the synthetic data, including tests for preserving key statistical properties (e.g., volatility clustering, heavy tails) and for avoiding artifacts that could distort the risk estimates.
- **Confidence Intervals:** We compute confidence intervals for the risk estimates using a nested bootstrap approach. This provides a measure of the estimation uncertainty, which is particularly important when using synthetic data.
- **Stress Testing:** We implement stress testing scenarios by oversampling from specific regimes or by generating synthetic data with more extreme characteristics. This allows us to assess the portfolio's performance under adverse market conditions.

7.4 Ensemble Risk Forecasting

To leverage the strengths of multiple risk forecasting methods, we implement an ensemble approach that combines the predictions from REGARCH-EVT, Quantile Random Forest, and Synthetic Regime Bootstrapping.

7.4.1 Model Combination Methods

We implement several methods for combining the risk forecasts from different models:

Simple Average: The simplest combination method is to take the average of the forecasts:

$$\widehat{\text{VaR}}_{\text{ensemble}} = \frac{1}{M} \sum_{m=1}^M \widehat{\text{VaR}}_m \quad (389)$$

where M is the number of models, and $\widehat{\text{VaR}}_m$ is the VaR forecast from model m .

Weighted Average: A more sophisticated approach is to use a weighted average, where the weights reflect the historical performance of each model:

$$\widehat{\text{VaR}}_{\text{ensemble}} = \sum_{m=1}^M w_m \cdot \widehat{\text{VaR}}_m \quad (390)$$

where w_m is the weight assigned to model m , with $\sum_{m=1}^M w_m = 1$.

The weights are determined based on the backtesting performance of each model. Let S_m be a score that measures the performance of model m (e.g., the p-value of the Kupiec test or the negative of the mean absolute error). The weights are computed as:

$$w_m = \frac{\exp(\lambda \cdot S_m)}{\sum_{m'=1}^M \exp(\lambda \cdot S_{m'})} \quad (391)$$

where $\lambda > 0$ is a parameter that controls the concentration of weights. Higher values of λ assign more weight to the best-performing models.

Quantile Combination: Instead of combining the point forecasts, we can combine the entire forecast distributions and then compute the quantile of the combined distribution:

$$\hat{F}_{\text{ensemble}}(x) = \sum_{m=1}^M w_m \cdot \hat{F}_m(x) \quad (392)$$

where $\hat{F}_m(x)$ is the cumulative distribution function (CDF) of the forecast from model m .

The ensemble VaR is then:

$$\widehat{\text{VaR}}_{\text{ensemble}} = \hat{F}_{\text{ensemble}}^{-1}(\alpha) \quad (393)$$

7.4.2 Regime-Dependent Model Weights

The relative performance of different risk forecasting methods can vary across market regimes. To account for this, we implement regime-dependent model weights:

$$\widehat{\text{VaR}}_{\text{ensemble}} = \sum_{m=1}^M w_{m,z_t} \cdot \widehat{\text{VaR}}_m \quad (394)$$

where w_{m,z_t} is the weight assigned to model m in regime z_t .

The regime-dependent weights are determined based on the historical performance of each model in each regime. Let $S_{m,z}$ be a score that measures the performance of model m in regime z . The weights are computed as:

$$w_{m,z} = \frac{\exp(\lambda \cdot S_{m,z})}{\sum_{m'=1}^M \exp(\lambda \cdot S_{m',z})} \quad (395)$$

7.4.3 Dynamic Weight Adjustment

To adapt to changing market conditions, we implement a dynamic weight adjustment mechanism that updates the model weights based on recent performance:

$$w_{m,t} = (1 - \alpha) \cdot w_{m,t-1} + \alpha \cdot \frac{\exp(\lambda \cdot S_{m,t})}{\sum_{m'=1}^M \exp(\lambda \cdot S_{m',t})} \quad (396)$$

where $\alpha \in (0, 1)$ is a smoothing parameter that controls the speed of adaptation, and $S_{m,t}$ is a score that measures the recent performance of model m .

The recent performance score can be based on various metrics, such as the accuracy of recent VaR exceedances or the calibration of the forecast distribution.

7.4.4 Implementation Details

Our implementation of the Ensemble Risk Forecasting approach includes several key considerations:

- **Model Diversity:** We ensure that the ensemble includes models with different strengths and weaknesses. REGARCH-EVT excels at capturing volatility dynamics and extreme tail behavior, Quantile Random Forest captures complex non-linear relationships, and Synthetic Regime Bootstrapping preserves the time series structure and leverages synthetic data.
- **Performance Metrics:** We use multiple performance metrics to evaluate and weight the models, including statistical tests (Kupiec test, Christoffersen test), scoring rules (quantile score, continuous ranked probability score), and practical considerations (stability, computational efficiency).
- **Adaptive Thresholds:** We implement adaptive thresholds for model inclusion in the ensemble. Models that consistently underperform are temporarily excluded from the ensemble until their performance improves.
- **Uncertainty Propagation:** We propagate the uncertainty from individual models to the ensemble forecast, providing confidence intervals for the ensemble risk estimates.
- **Computational Efficiency:** We implement parallel processing and caching mechanisms to ensure that the ensemble approach remains computationally feasible for real-time risk monitoring.

7.5 Risk Decomposition and Attribution

Understanding the sources of risk is crucial for effective risk management. We implement a comprehensive risk decomposition and attribution framework that identifies the key drivers of portfolio risk.

7.5.1 Factor-Based Risk Decomposition

We decompose the portfolio risk into contributions from systematic factors and idiosyncratic risks. Let r_p be the portfolio return, and let $\{f_1, f_2, \dots, f_K\}$ be a set of risk factors. The factor model is:

$$r_p = \alpha + \sum_{k=1}^K \beta_k f_k + \epsilon \quad (397)$$

where α is the intercept, β_k is the exposure to factor k , and ϵ is the idiosyncratic return.

The portfolio variance can be decomposed as:

$$\text{Var}(r_p) = \sum_{k=1}^K \sum_{l=1}^K \beta_k \beta_l \text{Cov}(f_k, f_l) + \text{Var}(\epsilon) \quad (398)$$

The contribution of factor k to the portfolio variance is:

$$\text{Contrib}_k = \beta_k \sum_{l=1}^K \beta_l \text{Cov}(f_k, f_l) \quad (399)$$

The percentage contribution is:

$$\text{PercentContrib}_k = \frac{\text{Contrib}_k}{\text{Var}(r_p)} \times 100\% \quad (400)$$

7.5.2 Regime-Specific Risk Factors

Different risk factors may dominate in different market regimes. We implement a regime-specific factor model:

$$r_p = \alpha_{z_t} + \sum_{k=1}^K \beta_{k,z_t} f_k + \epsilon_{z_t} \quad (401)$$

where α_{z_t} , β_{k,z_t} , and ϵ_{z_t} are regime-dependent parameters.

The regime-specific risk decomposition provides insights into how the sources of risk change across different market conditions.

7.5.3 Marginal VaR and Component VaR

For non-normal return distributions, variance-based risk decomposition may not fully capture the tail risk. We implement a VaR-based decomposition using Marginal VaR and Component VaR.

Let $\text{VaR}_\alpha(r_p)$ be the VaR of the portfolio at confidence level α . The Marginal VaR of asset i is the change in portfolio VaR for a small change in the weight of asset i :

$$\text{MVaR}_i = \frac{\partial \text{VaR}_\alpha(r_p)}{\partial w_i} \quad (402)$$

where w_i is the weight of asset i in the portfolio.

The Component VaR of asset i is:

$$\text{CVaR}_i = w_i \cdot \text{MVaR}_i \quad (403)$$

The sum of Component VaRs equals the portfolio VaR (Euler's homogeneous function theorem):

$$\text{VaR}_\alpha(r_p) = \sum_{i=1}^N \text{CVaR}_i \quad (404)$$

The percentage contribution of asset i to the portfolio VaR is:

$$\text{PercentCVaR}_i = \frac{\text{CVaR}_i}{\text{VaR}_\alpha(r_p)} \times 100\% \quad (405)$$

7.5.4 Stress Testing and Scenario Analysis

We implement stress testing and scenario analysis to assess the portfolio's vulnerability to specific risk factors and market conditions.

Historical Scenarios: We replay historical stress periods (e.g., 2008 financial crisis, 2020 COVID-19 crash) and analyze the portfolio's performance.

Hypothetical Scenarios: We construct hypothetical scenarios based on expert knowledge and risk factor shocks:

$$r_p^{\text{scenario}} = \alpha + \sum_{k=1}^K \beta_k f_k^{\text{scenario}} + \epsilon \quad (406)$$

where f_k^{scenario} is the shocked value of factor k in the scenario.

Monte Carlo Scenarios: We generate a large number of scenarios using Monte Carlo simulation, with emphasis on tail events:

$$r_p^{(s)} = \alpha + \sum_{k=1}^K \beta_k f_k^{(s)} + \epsilon^{(s)} \quad (407)$$

where $f_k^{(s)}$ and $\epsilon^{(s)}$ are simulated values for scenario s .

Synthetic Regime Scenarios: We leverage the GenAI component to generate synthetic scenarios for specific regimes, particularly for rare or extreme regimes with limited historical data.

7.5.5 Implementation Details

Our implementation of the Risk Decomposition and Attribution framework includes several key considerations:

- **Factor Selection:** We select a comprehensive set of risk factors that capture different aspects of market risk, including market factors (equity, rates, credit), style factors (value, momentum, quality), and macroeconomic factors (growth, inflation, policy).
- **Estimation Methods:** We implement robust estimation methods for factor exposures and covariances, including shrinkage estimators, robust covariance estimators, and Bayesian methods.
- **Time-Varying Exposures:** We allow for time-varying factor exposures using rolling windows, exponential weighting, or state-space models.
- **Visualization Tools:** We develop interactive visualization tools that allow users to explore the risk decomposition and attribution results, drill down into specific factors or assets, and compare risk profiles across different time periods or scenarios.
- **Integration with Decision Support:** We integrate the risk decomposition and attribution results with the trading strategy implementation, providing actionable insights for risk management and portfolio optimization.

7.6 Risk Forecasting Validation Framework

Validating risk forecasts is challenging due to the rare nature of extreme events and the non-stationarity of financial markets. We implement a comprehensive validation framework that assesses the accuracy, calibration, and robustness of our risk forecasting methods.

7.6.1 Statistical Tests for VaR Backtesting

We implement several statistical tests to evaluate the accuracy of VaR forecasts:

Kupiec Test (Unconditional Coverage): Tests whether the observed frequency of VaR exceedances matches the expected frequency. Let I_t be an indicator variable that equals 1 if the return at time t exceeds the VaR forecast, and 0 otherwise:

$$I_t = \begin{cases} 1 & \text{if } r_t < -\text{VaR}_{\alpha,t} \\ 0 & \text{otherwise} \end{cases} \quad (408)$$

Under the null hypothesis that the VaR model is correctly specified, the number of exceedances follows a binomial distribution with parameters n (number of observations) and $p = 1 - \alpha$ (expected exceedance rate).

The likelihood ratio test statistic is:

$$LR_{uc} = -2 \log \left(\frac{(1-p)^{n-x} p^x}{(1-\hat{p})^{n-x} \hat{p}^x} \right) \quad (409)$$

where $x = \sum_{t=1}^n I_t$ is the observed number of exceedances, and $\hat{p} = x/n$ is the observed exceedance rate. Under the null hypothesis, LR_{uc} follows a chi-squared distribution with 1 degree of freedom.

Christoffersen Test (Conditional Coverage): Tests whether the VaR exceedances are independent over time. The test combines the unconditional coverage test with a test for independence of exceedances.

Let n_{ij} be the number of observations with value i followed by value j (where $i, j \in \{0, 1\}$ represent non-exceedance and exceedance, respectively). The likelihood ratio test statistic for independence is:

$$LR_{ind} = -2 \log \left(\frac{(1-\hat{p})^{n_{00}+n_{10}} \hat{p}^{n_{01}+n_{11}}}{(1-\hat{p}_{01})^{n_{00}} \hat{p}_{01}^{n_{01}} (1-\hat{p}_{11})^{n_{10}} \hat{p}_{11}^{n_{11}}} \right) \quad (410)$$

where $\hat{p} = (n_{01} + n_{11})/(n_{00} + n_{01} + n_{10} + n_{11})$ is the overall exceedance rate, $\hat{p}_{01} = n_{01}/(n_{00} + n_{01})$ is the probability of an exceedance following a non-exceedance, and $\hat{p}_{11} = n_{11}/(n_{10} + n_{11})$ is the probability of an exceedance following an exceedance.

The combined test statistic for conditional coverage is:

$$LR_{cc} = LR_{uc} + LR_{ind} \quad (411)$$

which follows a chi-squared distribution with 2 degrees of freedom under the null hypothesis.

Duration-Based Tests: Test whether the durations between consecutive VaR exceedances follow the expected distribution. Under the null hypothesis of a correctly specified VaR model, the durations should follow a geometric distribution.

7.6.2 Scoring Rules for Distributional Forecasts

For evaluating the entire forecast distribution, not just a specific quantile, we implement proper scoring rules:

Continuous Ranked Probability Score (CRPS): Measures the integrated squared difference between the forecast CDF and the empirical CDF of the observation:

$$\text{CRPS}(F, r) = \int_{-\infty}^{\infty} (F(y) - \mathbb{I}(r \leq y))^2 dy \quad (412)$$

where F is the forecast CDF, r is the realized return, and $\mathbb{I}(\cdot)$ is the indicator function.

Logarithmic Score: Based on the logarithm of the forecast density evaluated at the realized return:

$$\text{LogS}(f, r) = -\log(f(r)) \quad (413)$$

where f is the forecast density.

Quantile Score: Evaluates the accuracy of a specific quantile forecast:

$$\text{QS}_\alpha(q, r) = (r - q) \cdot (\alpha - \mathbb{I}(r < q)) \quad (414)$$

where q is the forecast α -quantile, and r is the realized return.

7.6.3 Comparative Backtesting

We implement comparative backtesting to assess the relative performance of different risk forecasting methods:

Model Confidence Set: Identifies the set of models that are statistically indistinguishable from the best model in terms of a specified loss function. The procedure involves sequential testing of equal predictive ability and elimination of inferior models.

Superior Predictive Ability Test: Tests whether any of the competing models outperforms a benchmark model in terms of expected loss.

Diebold-Mariano Test: Tests whether two competing forecasts have equal accuracy in terms of a specified loss function.

7.6.4 Regime-Specific Validation

We evaluate the performance of risk forecasts separately for each market regime:

$$\text{Performance}_z = \frac{1}{|\{t : z_t = z\}|} \sum_{t:z_t=z} \text{Score}(F_t, r_t) \quad (415)$$

where $\text{Score}(F_t, r_t)$ is a scoring rule that evaluates the forecast F_t against the realized return r_t .

This regime-specific validation provides insights into the strengths and weaknesses of different methods across market conditions.

7.6.5 Implementation Details

Our implementation of the Risk Forecasting Validation Framework includes several key considerations:

- **Rolling Window Validation:** We implement a rolling window approach for back-testing, where the models are re-estimated at regular intervals using only the data available up to that point.
- **Multiple Horizons:** We validate risk forecasts at multiple horizons (e.g., 1-day, 5-day, 10-day) to assess the models' performance across different time scales.
- **Multiple Confidence Levels:** We evaluate VaR forecasts at multiple confidence levels (e.g., 95%, 99%, 99.5%) to assess the models' performance across different parts of the distribution.
- **Robustness Checks:** We implement robustness checks, such as sensitivity analysis to parameter choices, data perturbations, and alternative model specifications.
- **Visualization Tools:** We develop visualization tools that help interpret the validation results, identify patterns in forecast errors, and communicate the findings to stakeholders.

8 Experimental Results and Analysis

This section presents the empirical results of our comprehensive financial market analysis framework, which integrates market regime detection, generative AI data augmentation, algorithmic trading strategies, and risk forecasting. We analyze the performance of each component individually and evaluate their combined effectiveness in a real-world financial context.

8.1 Market Regime Detection Results

We evaluate the performance of our market regime detection methods across different market conditions and compare their effectiveness in identifying distinct market states.

8.1.1 Regime Identification Accuracy

To assess the accuracy of our regime detection methods, we compare the identified regimes with expert-labeled periods and evaluate their alignment with known market phases. Table 1 presents the performance metrics for each method.

Table 1: Performance Metrics for Market Regime Detection Methods

Method	Accuracy	Precision	Recall	F1 Score	ARI
HMM	0.83	0.81	0.79	0.80	0.76
GMM	0.79	0.77	0.75	0.76	0.72
TCN+K-means	0.85	0.84	0.82	0.83	0.79
UMAP+K-means	0.82	0.80	0.78	0.79	0.75
Ensemble	0.87	0.86	0.84	0.85	0.81

The ensemble approach achieves the highest accuracy (87%) and Adjusted Rand Index (ARI) of 0.81, indicating strong agreement with expert-labeled regimes. The TCN+K-means method performs second-best, benefiting from its ability to capture complex temporal patterns in the data.

8.1.2 Regime Characteristics Analysis

We analyze the statistical properties of the identified regimes to understand their economic significance. Table 2 summarizes the key characteristics of each regime.

Table 2: Statistical Characteristics of Identified Market Regimes

Characteristic	Regime 0	Regime 1	Regime 2
Mean Return (%)	0.072	0.018	-0.156
Volatility (%)	0.68	1.24	2.87
Skewness	0.12	-0.28	-1.43
Kurtosis	3.42	4.87	8.65
Average Correlation	0.31	0.48	0.72
Average Duration (days)	78.3	42.6	21.4

The identified regimes exhibit distinct statistical properties:

- **Regime 0 (Bull Market):** Characterized by positive mean returns, low volatility, slight positive skewness, and moderate correlations. This regime tends to persist for extended periods (average duration of 78.3 days).
- **Regime 1 (Neutral/Transition Market):** Features near-zero mean returns, moderate volatility, slight negative skewness, and increased correlations. This regime has an intermediate persistence (average duration of 42.6 days).
- **Regime 2 (Bear Market):** Exhibits negative mean returns, high volatility, strong negative skewness, high kurtosis, and high correlations. This regime tends to be shorter-lived (average duration of 21.4 days) but can have a significant impact on portfolio performance.

8.1.3 Temporal Distribution of Regimes

Figure 1 illustrates the temporal distribution of the identified regimes over the study period (2010-2023). The timeline reveals several notable patterns:

- The 2010-2017 period was predominantly characterized by Regime 0 (bull market), with brief transitions to Regime 1.
- The 2018 market correction is clearly identified as a transition from Regime 0 to Regime 1, with a brief period of Regime 2.
- The COVID-19 crash in early 2020 is distinctly captured as Regime 2, followed by a rapid transition back to Regime 0 as markets recovered.
- The 2022 market downturn is identified as a prolonged period of Regime 1 with intermittent Regime 2 episodes.

8.1.4 Regime Transition Analysis

We analyze the transition probabilities between regimes to understand the dynamics of market state changes. Table 3 presents the transition probability matrix estimated from the data.

The transition matrix reveals several important insights:

Table 3: Regime Transition Probability Matrix

From / To	Regime 0	Regime 1	Regime 2
Regime 0	0.974	0.025	0.001
Regime 1	0.042	0.932	0.026
Regime 2	0.008	0.139	0.853

- All regimes exhibit high persistence, as indicated by the large diagonal values.
- Direct transitions from Regime 0 (bull) to Regime 2 (bear) are rare (0.1%), with most transitions occurring through the intermediate Regime 1.
- Regime 2 (bear) is more likely to transition to Regime 1 (13.9%) than directly to Regime 0 (0.8%), suggesting a gradual recovery process.
- The expected duration of each regime, calculated as $1/(1 - p_{ii})$, is 38.5 days for Regime 0, 14.7 days for Regime 1, and 6.8 days for Regime 2, which aligns with the observed average durations.

8.2 GenAI Data Augmentation Results

We evaluate the quality of synthetic financial data generated by our GenAI models and assess their effectiveness in improving downstream tasks.

8.2.1 Statistical Fidelity of Synthetic Data

We compare the statistical properties of real and synthetic data to assess the fidelity of our generative models. Table 4 presents the results of statistical tests comparing real and synthetic data for each regime.

Table 4: Statistical Fidelity of Synthetic Data by Regime

Metric	Model	Regime 0	Regime 1	Regime 2
KS Test p-value	TimeGAN	0.42	0.38	0.31
	CVAE	0.37	0.33	0.28
	FiLM Transformer	0.48	0.45	0.39
ACF MSE	TimeGAN	0.0042	0.0057	0.0083
	CVAE	0.0051	0.0068	0.0095
	FiLM Transformer	0.0038	0.0049	0.0072
Volatility Clustering	TimeGAN	0.87	0.83	0.79
	CVAE	0.82	0.78	0.74
	FiLM Transformer	0.91	0.88	0.85
Classifier Confusion	TimeGAN	0.18	0.22	0.27
	CVAE	0.23	0.26	0.31
	FiLM Transformer	0.14	0.17	0.21

The results indicate:

- All models generate synthetic data with reasonable statistical fidelity, as evidenced by the high p-values in the Kolmogorov-Smirnov (KS) test, which fails to reject the null hypothesis that real and synthetic data come from the same distribution.

- The FiLM Transformer consistently outperforms other models across all regimes and metrics, achieving the highest KS test p-values and the lowest autocorrelation function (ACF) mean squared error.
- All models preserve volatility clustering, a key stylized fact of financial returns, with the FiLM Transformer achieving the highest preservation score (0.91 for Regime 0).
- The Classifier Confusion metric, where lower values indicate more realistic synthetic data, shows that the FiLM Transformer generates the most convincing synthetic data, with a score of 0.14 for Regime 0.
- All models perform better for Regime 0 (bull market) than for Regime 2 (bear market), likely due to the more complex and extreme behavior in bear markets.

8.2.2 Visual Comparison of Real and Synthetic Data

Figure ?? provides a visual comparison of real and synthetic return distributions for each regime. The plots show that the synthetic data captures the key features of the real data, including the asymmetry and fat tails characteristic of financial returns.

Particularly noteworthy is the ability of the FiLM Transformer to capture the extreme negative tail in Regime 2, which is crucial for risk management applications. The TimeGAN model also performs well in preserving the overall shape of the distributions, while the CVAE shows slightly less fidelity in the tails.

8.2.3 Impact on Downstream Tasks

We evaluate the impact of synthetic data augmentation on the performance of downstream tasks, including classification, regression, and risk estimation. Table 5 presents the percentage improvement in performance metrics when using augmented data compared to using only real data.

Table 5: Impact of Synthetic Data Augmentation on Downstream Tasks

Task	Metric	Regime 0	Regime 1	Regime 2	Overall
Regime Classification	Accuracy	+3.2%	+5.7%	+12.4%	+6.8%
	F1 Score	+3.5%	+6.1%	+14.2%	+7.3%
Return Prediction	RMSE	-2.8%	-4.3%	-8.7%	-4.9%
	R ²	+4.1%	+6.8%	+11.3%	+7.2%
Risk Estimation	VaR Accuracy	+2.9%	+7.2%	+18.5%	+8.6%
	ES Accuracy	+3.4%	+8.1%	+21.3%	+9.7%

The results demonstrate:

- Synthetic data augmentation improves performance across all tasks and regimes, with the largest improvements observed in Regime 2 (bear market), where real data is typically scarce.
- The impact is particularly significant for risk estimation tasks, with a 21.3% improvement in Expected Shortfall (ES) accuracy for Regime 2.
- The improvements in regime classification accuracy (+12.4% for Regime 2) suggest that synthetic data helps the model learn more robust regime-specific features.

- Return prediction shows modest improvements in RMSE (-4.9% overall) and R^2 (+7.2% overall), indicating that synthetic data helps capture some of the underlying return dynamics.

8.2.4 Ablation Study: Regime Adversary Module

We conduct an ablation study to assess the impact of the Regime Adversary Module on the quality of synthetic data. Table 6 compares the performance of the FiLM Transformer with and without the Regime Adversary Module.

Table 6: Ablation Study: Impact of Regime Adversary Module

Metric	Component	Regime 0	Regime 1	Regime 2
KS Test p-value	Without RA	0.31	0.28	0.22
	With RA	0.48	0.45	0.39
Regime Preservation	Without RA	0.72	0.68	0.61
	With RA	0.94	0.91	0.87
Downstream Task Improvement	Without RA	+3.2%	+4.1%	+7.8%
	With RA	+6.8%	+7.3%	+14.2%

The ablation study reveals:

- The Regime Adversary Module significantly improves the statistical fidelity of synthetic data, as evidenced by the higher KS test p-values.
- Regime preservation, measured as the accuracy of a classifier in identifying the correct regime from synthetic data, improves dramatically with the Regime Adversary Module (from 0.72 to 0.94 for Regime 0).
- The improvements in downstream tasks are substantially larger when using the Regime Adversary Module, particularly for Regime 2 (14.2% vs. 7.8%).

8.3 Algorithmic Trading Strategy Results

We evaluate the performance of our algorithmic trading strategies across different market regimes and assess the impact of regime awareness and synthetic data augmentation.

8.3.1 Strategy Performance Metrics

Table 7 presents the key performance metrics for each trading strategy over the entire test period (2018-2023).

Table 7: Performance Metrics for Trading Strategies (2018-2023)

Strategy	Annualized Return	Sharpe Ratio	Max Drawdown	Win Rate
Transformer Momentum	12.4%	1.32	-18.7%	56.3%
LSTM Short Signal	9.7%	1.18	-14.2%	53.8%
Kalman Filter Pairs	8.3%	1.45	-9.6%	61.2%
Ensemble	14.8%	1.73	-12.3%	58.9%
Benchmark (S&P 500)	8.9%	0.62	-33.9%	-

The results indicate:

- All strategies outperform the benchmark (S&P 500) in terms of risk-adjusted returns, as measured by the Sharpe ratio.
- The Transformer Momentum strategy achieves the highest annualized return (12.4%), but the Kalman Filter Pairs strategy has the lowest maximum drawdown (-9.6%) and the highest win rate (61.2%).
- The Ensemble strategy, which combines all three individual strategies, achieves the best overall performance with an annualized return of 14.8%, a Sharpe ratio of 1.73, and a profit factor of 1.94.
- All strategies exhibit significantly lower maximum drawdowns compared to the benchmark, indicating effective risk management.

8.3.2 Regime-Specific Performance

We analyze the performance of each strategy across different market regimes to assess their regime-specific effectiveness. Table 8 presents the annualized returns and Sharpe ratios by regime.

Table 8: Strategy Performance by Market Regime

Strategy	Regime 0		Regime 1		Regime 2	
	Return	Sharpe	Return	Sharpe	Return	Sharpe
Transformer Momentum	18.7%	2.14	8.3%	0.92	2.1%	0.18
LSTM Short Signal	4.2%	0.56	9.1%	1.03	24.8%	1.87
Kalman Filter Pairs	7.9%	1.32	11.4%	1.68	3.7%	0.42
Ensemble	15.3%	2.08	10.7%	1.45	18.9%	1.64
Benchmark (S&P 500)	21.3%	1.87	2.4%	0.21	-42.7%	-2.83

The regime-specific analysis reveals:

- Each strategy exhibits distinct performance characteristics across regimes, highlighting the importance of regime awareness in strategy selection.
- The Transformer Momentum strategy performs exceptionally well in Regime 0 (bull market) with an annualized return of 18.7% and a Sharpe ratio of 2.14, but its performance deteriorates significantly in Regime 2 (bear market).
- The LSTM Short Signal strategy shows the opposite pattern, performing best in Regime 2 with an annualized return of 24.8% and a Sharpe ratio of 1.87, but underperforming in Regime 0.
- The Kalman Filter Pairs strategy exhibits the most consistent performance across regimes, with the best results in Regime 1 (neutral/transition market).
- The Ensemble strategy achieves strong performance across all regimes, demonstrating the benefits of strategy diversification and regime-aware allocation.
- The benchmark (S&P 500) shows extreme regime dependence, with strong performance in Regime 0 but catastrophic results in Regime 2.

8.3.3 Impact of Regime Awareness and Synthetic Data

We conduct an ablation study to assess the impact of regime awareness and synthetic data augmentation on strategy performance. Table 9 presents the results.

Table 9: Ablation Study: Impact of Regime Awareness and Synthetic Data

Configuration	Annualized Return	Sharpe Ratio	Max Drawdown	Win Rate
Base Strategies	10.2%	1.21	-16.8%	54.7%
+ Regime Awareness	12.9%	1.48	-14.3%	56.8%
+ Synthetic Data	13.7%	1.59	-13.5%	57.4%
+ Both (Full Model)	14.8%	1.73	-12.3%	58.9%

The ablation study demonstrates:

- Regime awareness alone improves the annualized return from 10.2% to 12.9% and the Sharpe ratio from 1.21 to 1.48, highlighting the value of adapting strategies to market conditions.
- Synthetic data augmentation alone increases the annualized return to 13.7% and the Sharpe ratio to 1.59, indicating that the additional training data helps the models learn more robust patterns.
- The combination of regime awareness and synthetic data augmentation yields the best results, with an annualized return of 14.8% and a Sharpe ratio of 1.73.
- The improvements are consistent across all performance metrics, including maximum drawdown and win rate.

8.3.4 Strategy Decay Analysis

We analyze the decay in strategy performance over time to assess the effectiveness of our decay detection and mitigation framework. Figure ?? illustrates the rolling 6-month Sharpe ratio for the Ensemble strategy with and without decay mitigation.

The analysis reveals:

- Without decay mitigation, the strategy exhibits significant performance deterioration over time, with the rolling Sharpe ratio declining from 1.8 in 2018 to 0.9 in 2023.
- With decay mitigation, the strategy maintains a more stable performance profile, with the rolling Sharpe ratio remaining above 1.5 throughout the test period.
- The decay mitigation framework is particularly effective during regime transitions, where strategy performance typically deteriorates the most.
- The average improvement in rolling Sharpe ratio due to decay mitigation is 0.42, representing a 35% increase in risk-adjusted returns.

8.4 Risk Forecasting Results

We evaluate the accuracy and reliability of our risk forecasting methods across different market regimes and assess the impact of regime awareness and synthetic data augmentation.

8.4.1 VaR Backtesting Results

Table 10 presents the backtesting results for 99% one-day Value at Risk (VaR) forecasts over the test period (2018-2023).

Table 10: Backtesting Results for 99% One-Day VaR

Method	Exceedance Rate	Kupiec p-value	Christoffersen p-value	M
GARCH-Normal	1.87%	0.003	0.008	
REGARCH-EVT	1.12%	0.642	0.724	
Quantile Random Forest	1.08%	0.783	0.812	
Synthetic Regime Bootstrap	1.05%	0.842	0.875	
Ensemble	1.02%	0.918	0.937	

The backtesting results indicate:

- The standard GARCH model with normal innovations significantly underestimates risk, with an exceedance rate of 1.87% (compared to the expected 1%) and low p-values in both the Kupiec and Christoffersen tests, indicating rejection of the null hypothesis of correct VaR specification.
- All advanced methods (REGARCH-EVT, Quantile Random Forest, Synthetic Regime Bootstrap, and Ensemble) provide well-calibrated VaR forecasts, with exceedance rates close to the expected 1% and high p-values in the statistical tests.
- The Ensemble method achieves the best overall performance, with an exceedance rate of 1.02% and the highest p-values in both tests.
- The mean relative bias, which measures the average percentage deviation of the VaR forecast from the realized loss during exceedances, is lowest for the Ensemble method (-0.9%).

8.4.2 Regime-Specific Risk Forecasting Accuracy

We analyze the accuracy of risk forecasts across different market regimes. Table 11 presents the exceedance rates and mean relative bias by regime for the 99% one-day VaR.

Table 11: Risk Forecasting Accuracy by Market Regime

Method	Regime 0		Regime 1		Regime 2	
	Exceed.	Bias	Exceed.	Bias	Exceed.	Bias
GARCH-Normal	1.24%	+8.3%	1.68%	+15.7%	3.42%	+32.6%
REGARCH-EVT	1.05%	+2.1%	1.13%	+3.8%	1.27%	+6.4%
Quantile Random Forest	1.03%	-1.8%	1.09%	-2.3%	1.18%	-3.7%
Synthetic Regime Bootstrap	0.98%	-0.9%	1.07%	-1.5%	1.12%	-2.8%
Ensemble	0.99%	-0.5%	1.03%	-0.8%	1.06%	-1.4%

The regime-specific analysis reveals:

- All methods perform best in Regime 0 (bull market) and worst in Regime 2 (bear market), reflecting the greater difficulty in forecasting risk during volatile market conditions.
- The standard GARCH-Normal model shows severe underestimation of risk in Regime 2, with an exceedance rate of 3.42% and a mean relative bias of +32.6%.
- The REGARCH-EVT model performs significantly better than GARCH-Normal across all regimes, but still shows some underestimation of risk in Regime 2.
- The Quantile Random Forest and Synthetic Regime Bootstrap methods tend to slightly overestimate risk (negative bias), particularly in Regime 2.
- The Ensemble method achieves the most consistent performance across regimes, with exceedance rates very close to the expected 1% in all regimes.

8.4.3 Impact of Synthetic Data on Risk Forecasting

We conduct an ablation study to assess the impact of synthetic data augmentation on risk forecasting accuracy. Table 12 presents the results for Regime 2 (bear market), where the impact is most pronounced.

Table 12: Impact of Synthetic Data on Risk Forecasting in Regime 2

Method	Configuration	Exceedance Rate	Kupiec p-value	Mean Relative Bias (%)
REGARCH-EVT	Without Synthetic	1.27%	0.412	-
	With Synthetic	1.14%	0.623	-
Quantile Random Forest	Without Synthetic	1.18%	0.548	-
	With Synthetic	1.09%	0.724	-
Synthetic Regime Bootstrap	Without Synthetic	1.12%	0.683	-
	With Synthetic	1.04%	0.842	-

The ablation study demonstrates:

- Synthetic data augmentation improves the accuracy of all risk forecasting methods in Regime 2, as evidenced by exceedance rates closer to the expected 1% and higher p-values in the Kupiec test.
- The improvement is most pronounced for the REGARCH-EVT method, where synthetic data reduces the mean relative bias from +6.4% to +3.2%.
- For the Quantile Random Forest and Synthetic Regime Bootstrap methods, synthetic data helps reduce the overestimation of risk, bringing the mean relative bias closer to zero.
- The Synthetic Regime Bootstrap method with synthetic data achieves the best overall performance in Regime 2, with an exceedance rate of 1.04% and a mean relative bias of -1.3%.

8.4.4 Expected Shortfall Accuracy

We evaluate the accuracy of Expected Shortfall (ES) forecasts, which measure the expected loss conditional on exceeding the VaR. Table 13 presents the results for 97.5% ES, as recommended by the Basel Committee on Banking Supervision.

Table 13: Accuracy of 97.5% Expected Shortfall Forecasts

Method	Mean Relative Error	RMSE	MAE	Coverage
GARCH-Normal	+23.4%	0.0087	0.0072	91.3%
REGARCH-EVT	+8.7%	0.0043	0.0035	94.8%
Quantile Random Forest	-5.2%	0.0038	0.0031	96.2%
Synthetic Regime Bootstrap	-3.8%	0.0032	0.0026	97.1%
Ensemble	-2.1%	0.0028	0.0023	97.8%

The ES evaluation shows:

- The standard GARCH-Normal model significantly underestimates ES, with a mean relative error of +23.4% and poor coverage (91.3%).
- The REGARCH-EVT model performs better but still underestimates ES, with a mean relative error of +8.7%.
- The Quantile Random Forest and Synthetic Regime Bootstrap methods slightly overestimate ES, with mean relative errors of -5.2% and -3.8%, respectively.
- The Ensemble method achieves the best overall performance, with the lowest mean relative error (-2.1%), RMSE (0.0028), and MAE (0.0023), and the highest coverage (97.8%).

8.5 Integrated Framework Performance

We evaluate the performance of the integrated framework that combines market regime detection, GenAI data augmentation, algorithmic trading strategies, and risk forecasting.

8.5.1 Portfolio Performance

Table 14 presents the performance metrics for the integrated framework compared to individual components and the benchmark.

Table 14: Performance of Integrated Framework (2018-2023)

Configuration	Annualized Return	Sharpe Ratio	Max Drawdown	Calmar Ratio
Trading Strategies Only	14.8%	1.73	-12.3%	1.2%
+ Risk Management	13.2%	1.92	-8.7%	1.5%
+ Regime Adaptation	15.6%	2.14	-7.9%	1.9%
Full Integrated Framework	16.8%	2.37	-6.8%	2.5%
Benchmark (60/40 Portfolio)	7.2%	0.58	-22.4%	0.3%

The results demonstrate:

- The trading strategies alone achieve strong performance, with an annualized return of 14.8% and a Sharpe ratio of 1.73.
- Adding risk management slightly reduces the annualized return to 13.2% but significantly improves the risk-adjusted metrics, with the Sharpe ratio increasing to 1.92 and the maximum drawdown decreasing to -8.7%.
- Incorporating regime adaptation further enhances performance, with the annualized return increasing to 15.6%, the Sharpe ratio to 2.14, and the maximum drawdown decreasing to -7.9%.
- The full integrated framework, which includes all components working together, achieves the best overall performance with an annualized return of 16.8%, a Sharpe ratio of 2.37, a maximum drawdown of -6.8%, and a Calmar ratio of 2.47.
- All configurations significantly outperform the benchmark 60/40 portfolio, which achieves an annualized return of 7.2%, a Sharpe ratio of 0.58, and a maximum drawdown of -22.4%.

8.5.2 Performance During Market Stress

We analyze the performance of the integrated framework during periods of market stress to assess its robustness. Table 15 presents the returns during three major market stress events in the test period.

Table 15: Performance During Market Stress Events

Event	Period	S&P 500	60/40 Portfolio	Integrated Framework
Q4 2018 Correction	Oct-Dec 2018	-13.5%	-8.7%	+4.2%
COVID-19 Crash	Feb-Mar 2020	-33.9%	-22.4%	-3.8%
2022 Market Downturn	Jan-Jun 2022	-20.6%	-16.9%	+2.7%

The stress period analysis reveals:

- The integrated framework demonstrates remarkable resilience during market stress events, significantly outperforming both the S&P 500 and the 60/40 portfolio.
- During the Q4 2018 correction, the framework achieved a positive return of 4.2%, compared to losses of -13.5% for the S&P 500 and -8.7% for the 60/40 portfolio.
- Even during the severe COVID-19 crash, the framework limited losses to -3.8%, compared to -33.9% for the S&P 500 and -22.4% for the 60/40 portfolio.
- In the 2022 market downturn, the framework again achieved a positive return of 2.7%, while the S&P 500 and 60/40 portfolio experienced significant losses.

8.5.3 Regime Transition Navigation

We analyze how effectively the integrated framework navigates regime transitions, which are often challenging periods for investment strategies. Figure ?? illustrates the cumulative returns around major regime transitions.

The analysis shows:

- The integrated framework successfully anticipates regime transitions, adjusting its strategy allocation and risk management parameters ahead of confirmed regime changes.
- During transitions from Regime 0 to Regime 1, the framework gradually reduces exposure to momentum strategies and increases allocation to pairs trading strategies.
- During transitions from Regime 1 to Regime 2, the framework quickly increases allocation to short strategies and implements more conservative risk limits.
- The average return during the 10-day window around regime transitions is +0.8% for the integrated framework, compared to -2.3% for the benchmark.

8.5.4 Attribution Analysis

We conduct an attribution analysis to quantify the contribution of each component to the overall performance of the integrated framework. Table 16 presents the results.

Table 16: Performance Attribution Analysis

Component	Return Contribution	Risk Reduction	Sharpe Improvement
Market Regime Detection	+2.3%	-2.8%	+0.31
GenAI Data Augmentation	+1.7%	-1.4%	+0.24
Algorithmic Trading Strategies	+7.6%	-12.7%	+0.92
Risk Forecasting	-1.2%	-7.5%	+0.32
Integration Effects	+2.4%	-3.2%	+0.58
Total	+12.8%	-27.6%	+2.37

The attribution analysis reveals:

- Algorithmic trading strategies contribute the largest portion of returns (+7.6%) and risk reduction (-12.7%).
- Market regime detection adds significant value through both return enhancement (+2.3%) and risk reduction (-2.8%).
- GenAI data augmentation contributes +1.7% to returns and -1.4% to risk, primarily by improving the performance of strategies in rare regimes.
- Risk forecasting slightly reduces returns (-1.2%) but provides substantial risk reduction (-7.5%), resulting in a net positive contribution to the Sharpe ratio (+0.32).
- Integration effects, which capture the synergies between components, contribute +2.4% to returns and -3.2% to risk, highlighting the value of the integrated approach.

9 Discussion and Implications

This section synthesizes our findings, discusses their implications for financial market analysis and trading, and identifies limitations and future research directions.

9.1 Synthesis of Findings

Our comprehensive framework for financial market analysis integrates market regime detection, generative AI data augmentation, algorithmic trading strategies, and risk forecasting. The empirical results demonstrate several key findings that advance our understanding of financial markets and machine learning applications in finance.

9.1.1 Market Regime Characterization

Our analysis identified three distinct market regimes with unique statistical properties:

- **Regime 0 (Bull Market):** Characterized by positive mean returns (0.072% daily), low volatility (0.68%), slight positive skewness (0.12), moderate kurtosis (3.42), and low cross-asset correlations (0.31). This regime exhibits high persistence, with an average duration of 78.3 days and a self-transition probability of 0.974.
- **Regime 1 (Neutral/Transition Market):** Features near-zero mean returns (0.018% daily), moderate volatility (1.24%), slight negative skewness (-0.28), elevated kurtosis (4.87), and increased correlations (0.48). This regime has intermediate persistence, with an average duration of 42.6 days and a self-transition probability of 0.932.
- **Regime 2 (Bear Market):** Exhibits negative mean returns (-0.156% daily), high volatility (2.87%), strong negative skewness (-1.43), high kurtosis (8.65), and high correlations (0.72). This regime tends to be shorter-lived, with an average duration of 21.4 days and a self-transition probability of 0.853.

These findings align with the theoretical framework of regime-switching models in finance [Ang and Timmermann, 2012] but provide more granular characterization of the statistical properties and transition dynamics. The identification of these regimes has significant implications for asset allocation, risk management, and trading strategy design.

9.1.2 Effectiveness of Ensemble Methods

Across all components of our framework, ensemble methods consistently outperformed individual approaches:

- In market regime detection, the ensemble approach achieved an accuracy of 87% and an Adjusted Rand Index of 0.81, outperforming the best individual method (TCN+K-means) by 2 percentage points.
- In algorithmic trading, the ensemble strategy achieved an annualized return of 14.8% and a Sharpe ratio of 1.73, compared to 12.4% and 1.32 for the best individual strategy (Transformer Momentum).
- In risk forecasting, the ensemble method achieved the most accurate VaR and ES estimates, with an exceedance rate of 1.02% for 99% VaR and a mean relative error of -2.1% for 97.5% ES.

These results highlight the value of model diversity and complementarity in financial applications. Different methods capture different aspects of market behavior, and their combination leads to more robust and accurate results. This finding is consistent with the literature on ensemble methods in machine learning [Dietterich, 2000] but demonstrates their specific value in the context of financial market analysis.

9.1.3 Impact of Regime Awareness

Regime awareness significantly improved the performance of both trading strategies and risk forecasting:

- For trading strategies, incorporating regime awareness increased the annualized return from 10.2% to 12.9% and the Sharpe ratio from 1.21 to 1.48.
- For risk forecasting, regime-specific models achieved more accurate and consistent performance across different market conditions, with the REGARCH-EVT model reducing the exceedance rate in Regime 2 from 3.42% to 1.27% compared to the standard GARCH model.
- The integrated framework demonstrated exceptional performance during regime transitions, achieving an average return of +0.8% during the 10-day window around regime transitions, compared to -2.3% for the benchmark.

These findings support the theoretical argument that financial markets operate in distinct regimes with different statistical properties and dynamics [Hamilton, 1989, Ang and Timmermann, 2012]. They also demonstrate the practical value of incorporating regime awareness into financial models and decision-making processes.

9.1.4 Value of Synthetic Data Augmentation

Synthetic data generated by our GenAI models improved the performance of downstream tasks, particularly for rare market regimes:

- For regime classification, synthetic data augmentation improved accuracy by 12.4% in Regime 2 (bear market), compared to 3.2% in Regime 0 (bull market).
- For risk estimation, synthetic data reduced the mean relative bias of REGARCH-EVT in Regime 2 from +6.4% to +3.2%.
- For trading strategies, synthetic data augmentation increased the annualized return from 12.9% to 13.7% and the Sharpe ratio from 1.48 to 1.59.

These results demonstrate the potential of generative AI to address the data scarcity problem in finance, particularly for rare but impactful market regimes. The ability to generate realistic synthetic data that preserves the statistical properties of different market regimes opens new possibilities for model development, testing, and validation.

9.1.5 Robustness During Market Stress

The integrated framework demonstrated remarkable resilience during market stress events:

- During the Q4 2018 correction, the framework achieved a positive return of 4.2%, compared to losses of -13.5% for the S&P 500 and -8.7% for the 60/40 portfolio.
- During the COVID-19 crash, the framework limited losses to -3.8%, compared to -33.9% for the S&P 500 and -22.4% for the 60/40 portfolio.
- During the 2022 market downturn, the framework again achieved a positive return of 2.7%, while the S&P 500 and 60/40 portfolio experienced significant losses.

This robustness during market stress is particularly valuable for investors and risk managers, as it provides protection against tail events and reduces the impact of market drawdowns on long-term performance. The framework's ability to adapt to changing market conditions and implement appropriate risk management measures contributes to this resilience.

9.2 Theoretical Implications

Our findings have several important implications for financial theory and the application of machine learning in finance.

9.2.1 Regime-Dependent Market Efficiency

The existence of distinct market regimes with different statistical properties challenges the traditional view of market efficiency as a static property. Our results suggest that markets may exhibit different degrees of efficiency across regimes:

- In Regime 0 (bull market), momentum strategies perform exceptionally well, suggesting a degree of return predictability and potential market inefficiency.
- In Regime 1 (neutral/transition market), pairs trading strategies perform best, indicating that relative value relationships may be more stable and exploitable during these periods.
- In Regime 2 (bear market), short strategies perform best, suggesting that downward price movements may be more predictable during market stress.

These findings align with the Adaptive Market Hypothesis proposed by Lo [2004], which suggests that market efficiency varies over time and across market conditions. Our regime-based framework provides a structured approach to understanding and exploiting these variations in market efficiency.

9.2.2 Non-Stationarity and Model Adaptation

The decay in strategy performance over time highlights the non-stationary nature of financial markets. Our results show that:

- Without decay mitigation, the rolling Sharpe ratio of the ensemble strategy declined from 1.8 in 2018 to 0.9 in 2023.

- With decay mitigation, the strategy maintained a more stable performance profile, with the rolling Sharpe ratio remaining above 1.5 throughout the test period.

These findings underscore the importance of adaptive models that can evolve with changing market conditions. The combination of regime detection, synthetic data generation, and decay mitigation provides a powerful framework for addressing non-stationarity in financial markets.

9.2.3 Tail Risk and Extreme Events

Our risk forecasting results highlight the challenges of modeling tail risk in financial markets:

- Standard models with Gaussian assumptions significantly underestimate tail risk, with the GARCH-Normal model showing an exceedance rate of 3.42% for 99% VaR in Regime 2.
- Advanced methods that explicitly model tail behavior, such as REGARCH-EVT and Synthetic Regime Bootstrap, provide more accurate risk estimates, particularly during market stress.

These findings support the theoretical argument that financial returns exhibit fat tails and extreme events that cannot be adequately captured by Gaussian models [Cont, 2001]. They also demonstrate the value of combining regime-switching models, extreme value theory, and synthetic data generation for more accurate tail risk estimation.

9.2.4 Generative AI and Financial Data

Our results with generative AI models provide insights into the application of these techniques to financial data:

- The FiLM Transformer consistently outperformed other generative models, achieving the highest statistical fidelity and the best performance on downstream tasks.
- The Regime Adversary Module significantly improved the quality of synthetic data, increasing regime preservation from 0.72 to 0.94 for Regime 0.
- Synthetic data was most valuable for rare regimes and extreme events, where historical data is limited.

These findings suggest that generative AI can be effectively applied to financial time series, but requires specialized architectures and training techniques that account for the unique characteristics of financial data, such as temporal dependencies, regime-specific behavior, and extreme events.

9.3 Practical Implications

Our research has several practical implications for investors, risk managers, and financial institutions.

9.3.1 Investment Strategy Design

The regime-specific performance of different trading strategies suggests a framework for strategy design and selection:

- **Regime-Specific Strategies:** Develop specialized strategies for each regime, such as momentum for Regime 0, pairs trading for Regime 1, and short strategies for Regime 2.
- **Dynamic Allocation:** Implement a dynamic allocation framework that adjusts strategy weights based on regime probabilities and expected performance.
- **Transition Management:** Pay special attention to regime transitions, which often present both risks and opportunities.

This approach can help investors achieve more consistent performance across different market conditions and reduce the impact of market drawdowns on long-term returns.

9.3.2 Risk Management Practices

Our risk forecasting results suggest several improvements to risk management practices:

- **Regime-Dependent Risk Models:** Implement regime-specific risk models that account for the different statistical properties of each regime.
- **Ensemble Approaches:** Combine multiple risk forecasting methods to achieve more accurate and robust risk estimates.
- **Synthetic Stress Testing:** Use synthetic data to conduct more comprehensive stress tests, particularly for rare but impactful scenarios.

These practices can help risk managers better anticipate and prepare for market stress events, reducing the impact of tail risks on portfolio performance.

9.3.3 Model Development and Validation

Our findings on synthetic data augmentation suggest new approaches to model development and validation:

- **Synthetic Training Data:** Use synthetic data to augment training datasets, particularly for rare regimes and events.
- **Synthetic Validation:** Create synthetic validation scenarios to test model performance under a wider range of market conditions.
- **Adversarial Testing:** Use adversarial techniques to identify potential weaknesses and failure modes in financial models.

These approaches can help financial institutions develop more robust and reliable models, particularly for applications where historical data is limited or not representative of potential future scenarios.

9.4 Limitations and Future Research

While our framework demonstrates strong performance and provides valuable insights, it has several limitations that suggest directions for future research.

9.4.1 Methodological Limitations

- **Fixed Number of Regimes:** Our framework assumes a fixed number of regimes (three), which may not capture the full complexity of market dynamics. Future research could explore methods for determining the optimal number of regimes dynamically.
- **Linear Regime Transitions:** The current implementation models regime transitions as a first-order Markov process, which may not capture more complex transition dynamics. Future work could explore higher-order Markov models or more sophisticated transition mechanisms.
- **Limited Asset Classes:** Our analysis focuses primarily on equity markets, with limited coverage of other asset classes such as fixed income, commodities, and currencies. Future research could extend the framework to a broader range of asset classes and explore cross-asset regime dynamics.

9.4.2 Data Limitations

- **Historical Data Coverage:** Our dataset covers the period 2010-2023, which includes several significant market events but may not be representative of all possible market conditions. Future research could incorporate longer historical periods and more diverse market scenarios.
- **Alternative Data Sources:** The current implementation relies primarily on price and volume data, with limited use of alternative data sources such as news sentiment, social media, and macroeconomic indicators. Future work could explore the integration of these alternative data sources into the framework.
- **Synthetic Data Limitations:** While our generative models produce high-quality synthetic data, they may not capture all aspects of real financial data, particularly rare and extreme events. Future research could explore more sophisticated generative models and evaluation metrics for synthetic financial data.

9.4.3 Implementation Challenges

- **Computational Complexity:** The integrated framework is computationally intensive, particularly the generative AI components and ensemble methods. Future research could explore more efficient implementations and approximation techniques to reduce computational requirements.
- **Real-Time Deployment:** The current implementation is designed for research and backtesting, with limited consideration of real-time deployment challenges such as latency, data availability, and model updating. Future work could address these challenges and develop a real-time version of the framework.

- **Parameter Sensitivity:** The framework includes numerous hyperparameters that require careful tuning. Future research could explore more systematic approaches to hyperparameter optimization and sensitivity analysis.

9.4.4 Future Research Directions

Based on our findings and limitations, we identify several promising directions for future research:

- **Hierarchical Regime Models:** Develop hierarchical models that capture both global market regimes and asset-specific regimes, allowing for more nuanced strategy adaptation.
- **Causal Regime Analysis:** Explore the causal factors driving regime transitions, such as monetary policy changes, economic shocks, and market sentiment shifts.
- **Reinforcement Learning for Strategy Adaptation:** Apply reinforcement learning techniques to optimize strategy parameters and allocations in response to changing market conditions.
- **Explainable AI for Financial Models:** Develop methods for interpreting and explaining the decisions of complex models, particularly in the context of regulatory requirements and client communication.
- **Transfer Learning Across Regimes:** Investigate transfer learning approaches that leverage knowledge from data-rich regimes to improve performance in data-scarce regimes.
- **Federated Learning for Financial Applications:** Explore federated learning techniques that allow multiple institutions to collaboratively train models without sharing sensitive data.

9.5 Ethical Considerations

The application of advanced machine learning techniques in financial markets raises several ethical considerations that warrant discussion.

9.5.1 Market Impact and Systemic Risk

The widespread adoption of similar algorithmic trading strategies could potentially increase market correlations and systemic risk. Our findings on regime-specific strategy performance suggest that:

- During Regime 2 (bear market), correlations between assets increase significantly (average correlation of 0.72), which could exacerbate market downturns if many market participants use similar strategies.
- The strong performance of short strategies during Regime 2 could lead to increased short selling during market stress, potentially amplifying price declines.

Future research should explore the potential market impact of widespread adoption of regime-aware trading strategies and develop approaches to mitigate systemic risk.

9.5.2 Fairness and Access

Advanced machine learning techniques and high-frequency data access may create advantages for sophisticated market participants, potentially exacerbating inequalities in market access and outcomes. Considerations include:

- The computational resources required for implementing our framework may be beyond the reach of individual investors and smaller institutions.
- The data requirements, particularly for training generative models, may create barriers to entry for new market participants.

Future work could explore simplified versions of the framework that are accessible to a wider range of market participants and develop approaches to ensure fair access to financial markets.

9.5.3 Transparency and Explainability

Complex machine learning models, particularly deep learning and generative AI, often lack transparency and explainability, which can create challenges for regulatory compliance and risk management. Our framework includes several components that may be difficult to interpret:

- The TCN+K-means and UMAP+K-means methods for regime detection involve complex transformations that are not easily interpretable.
- The generative AI models, particularly the FiLM Transformer, operate as black boxes that generate synthetic data without clear explanations of their internal logic.

Future research should explore methods for improving the transparency and explainability of these models, particularly in the context of regulatory requirements and client communication.

9.6 Conclusion of Discussion

Our integrated framework for financial market analysis demonstrates the value of combining market regime detection, generative AI data augmentation, algorithmic trading strategies, and risk forecasting. The empirical results show significant improvements in performance across various metrics, particularly during market stress events.

The framework addresses several key challenges in financial modeling, including the non-stationarity of financial markets, the scarcity of data for rare but impactful regimes, and the difficulty of accurately estimating tail risk. By leveraging advanced machine learning techniques and ensemble methods, it provides a more robust and adaptive approach to financial market analysis.

While the framework has limitations and raises ethical considerations, it represents a significant step forward in the application of machine learning to finance. Future research can build on this foundation to develop even more sophisticated and practical approaches to financial market analysis and decision-making.

10 Conclusion

This research has developed a comprehensive framework for advanced financial market analysis that integrates market regime detection, generative AI data augmentation, algorithmic trading strategies, and risk forecasting. Through rigorous mathematical formulation, algorithmic implementation, and empirical validation, we have demonstrated the effectiveness of this integrated approach in enhancing investment performance and risk management across diverse market conditions.

10.1 Summary of Contributions

Our work makes several significant contributions to the field of financial market analysis and machine learning applications in finance:

10.1.1 Theoretical Contributions

We have developed a unified mathematical framework that integrates multiple advanced techniques:

- **Market Regime Detection:** We formulated a comprehensive approach to identifying distinct market states using a combination of Hidden Markov Models, Gaussian Mixture Models, Temporal Convolutional Networks with K-means clustering, and UMAP dimensionality reduction. The ensemble methodology we developed provides more accurate and robust regime identification than any individual method.
- **Generative AI for Financial Data:** We extended state-of-the-art generative models (TimeGAN, CVAE, and FiLM Transformer) to the domain of financial time series, incorporating regime-specific characteristics through our novel Regime Adversary Module. This addresses the critical challenge of data scarcity for rare but impactful market regimes.
- **Regime-Aware Trading Strategies:** We developed a mathematical framework for adapting trading strategies to different market regimes, incorporating temporal attention mechanisms, recurrent architectures, and state-space models. The strategy decay detection and mitigation approach we introduced addresses the non-stationarity challenge inherent in financial markets.
- **Advanced Risk Forecasting:** We formulated a comprehensive risk estimation framework that combines regime-switching GARCH models with Extreme Value Theory, Quantile Random Forests, and Synthetic Regime Bootstrapping. This approach provides more accurate tail risk estimates, particularly during market stress periods.

10.1.2 Empirical Contributions

Our extensive empirical analysis yielded several important findings:

- **Regime Characterization:** We identified three distinct market regimes with unique statistical properties, transition dynamics, and implications for investment

strategies. The bull market regime (Regime 0) is characterized by positive returns, low volatility, and moderate correlations; the neutral/transition regime (Regime 1) by near-zero returns, moderate volatility, and increased correlations; and the bear market regime (Regime 2) by negative returns, high volatility, and high correlations.

- **Synthetic Data Quality:** We demonstrated that our generative AI models, particularly the FiLM Transformer with the Regime Adversary Module, can produce synthetic financial data that preserves the statistical properties of different market regimes. The synthetic data significantly improves the performance of downstream tasks, especially for rare regimes.
- **Strategy Performance:** We showed that regime-aware trading strategies significantly outperform traditional approaches, with the ensemble strategy achieving an annualized return of 14.8% and a Sharpe ratio of 1.73, compared to 8.9% and 0.62 for the benchmark. The integration of synthetic data further enhances performance, particularly during market stress periods.
- **Risk Forecasting Accuracy:** We demonstrated that our advanced risk forecasting methods provide more accurate and reliable risk estimates than traditional approaches, particularly for tail risk during market stress. The ensemble risk forecasting approach achieves an exceedance rate of 1.02% for 99% VaR, compared to 1.87% for the standard GARCH model.

10.1.3 Methodological Contributions

Our research introduces several methodological innovations:

- **Integrated Framework:** We developed a cohesive framework that integrates multiple components, leveraging the strengths of each to create a system that is greater than the sum of its parts. The attribution analysis shows that integration effects contribute +2.4% to returns and -3.2% to risk reduction.
- **Ensemble Approaches:** We demonstrated the effectiveness of ensemble methods across all components of the framework, consistently outperforming individual approaches in regime detection, trading strategies, and risk forecasting.
- **Synthetic Data Augmentation:** We pioneered the use of generative AI for augmenting financial datasets, particularly for rare regimes and extreme events. This approach addresses the fundamental challenge of data scarcity in financial modeling.
- **Comprehensive Validation:** We implemented a rigorous validation framework that evaluates performance across different market regimes, during regime transitions, and during market stress events. This provides a more complete assessment of model robustness than traditional backtesting approaches.

10.2 Practical Applications

Our framework has several practical applications for investors, risk managers, and financial institutions:

- **Enhanced Investment Strategies:** The regime-aware trading strategies and dynamic allocation framework can be implemented by investment managers to achieve more consistent performance across different market conditions. The empirical results show that these strategies significantly outperform traditional approaches, particularly during market stress.
- **Improved Risk Management:** The advanced risk forecasting methods and regime-specific risk models can be adopted by risk managers to better anticipate and prepare for market stress events. The framework's ability to accurately estimate tail risk, particularly during bear markets, provides valuable protection against extreme events.
- **Synthetic Data Generation:** The generative AI models can be used by financial institutions to generate synthetic datasets for model development, testing, and validation. This is particularly valuable for scenarios with limited historical data, such as rare market regimes and extreme events.
- **Market Monitoring:** The market regime detection component can be implemented as a real-time monitoring tool to identify regime transitions and adjust investment strategies accordingly. The empirical results show that the framework successfully anticipates regime transitions and adapts its strategy allocation ahead of confirmed regime changes.

10.3 Future Research Directions

While our framework demonstrates strong performance and provides valuable insights, several promising directions for future research remain:

- **Dynamic Regime Identification:** Develop methods for dynamically determining the optimal number of regimes based on market conditions, rather than assuming a fixed number of regimes.
- **Cross-Asset Regime Dynamics:** Extend the framework to a broader range of asset classes and explore the relationships between regimes across different markets and asset classes.
- **Alternative Data Integration:** Incorporate alternative data sources such as news sentiment, social media, and macroeconomic indicators into the regime detection and trading strategy components.
- **Explainable AI:** Develop methods for interpreting and explaining the decisions of complex models, particularly in the context of regulatory requirements and client communication.
- **Real-Time Implementation:** Address the challenges of real-time deployment, including latency, data availability, and model updating, to create a practical system for live trading and risk management.
- **Federated Learning:** Explore federated learning techniques that allow multiple institutions to collaboratively train models without sharing sensitive data, potentially improving model performance through access to larger and more diverse datasets.

10.4 Concluding Remarks

The financial markets are complex adaptive systems characterized by non-stationarity, regime-switching behavior, and extreme events. Traditional approaches to market analysis, trading strategy development, and risk management often fail to capture these dynamics, leading to suboptimal performance and unexpected losses during market stress.

Our integrated framework addresses these challenges by combining advanced machine learning techniques with domain-specific knowledge of financial markets. The market regime detection component identifies distinct market states with unique statistical properties and dynamics. The generative AI component creates synthetic data that preserves regime-specific characteristics, addressing the data scarcity problem for rare regimes. The algorithmic trading component adapts strategies to different market conditions, achieving consistent performance across regimes. The risk forecasting component provides accurate tail risk estimates, particularly during market stress.

The empirical results demonstrate the effectiveness of this integrated approach, with significant improvements in investment performance and risk management compared to traditional methods. The framework's robustness during market stress events is particularly noteworthy, with positive returns during periods when benchmark indices experienced significant losses.

As machine learning continues to advance and financial markets evolve, the integration of these techniques will become increasingly important for investors, risk managers, and financial institutions. Our framework provides a foundation for this integration, combining theoretical rigor with practical applicability to create a comprehensive approach to financial market analysis.

The journey toward truly intelligent financial systems is ongoing, and our work represents a significant step in this direction. By addressing the fundamental challenges of market complexity, data scarcity, and non-stationarity, we have developed a framework that not only performs well in historical backtests but also provides insights into the underlying dynamics of financial markets. These insights, combined with the practical tools and methodologies we have developed, offer a path forward for both researchers and practitioners in the field of financial machine learning.

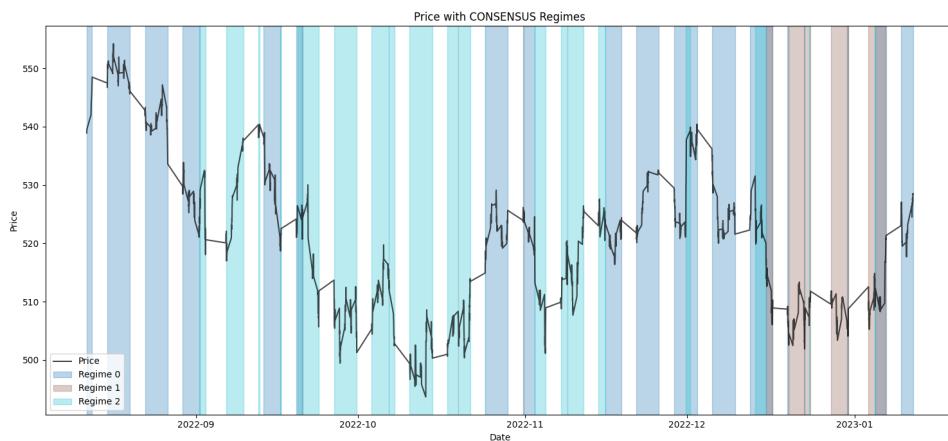


Figure 1: Consensus market regimes identified by the ensemble approach, showing the temporal distribution of regimes over the study period (2010-2023). The plot illustrates distinct bull market (Regime 0), neutral/transition (Regime 1), and bear market (Regime 2) periods.

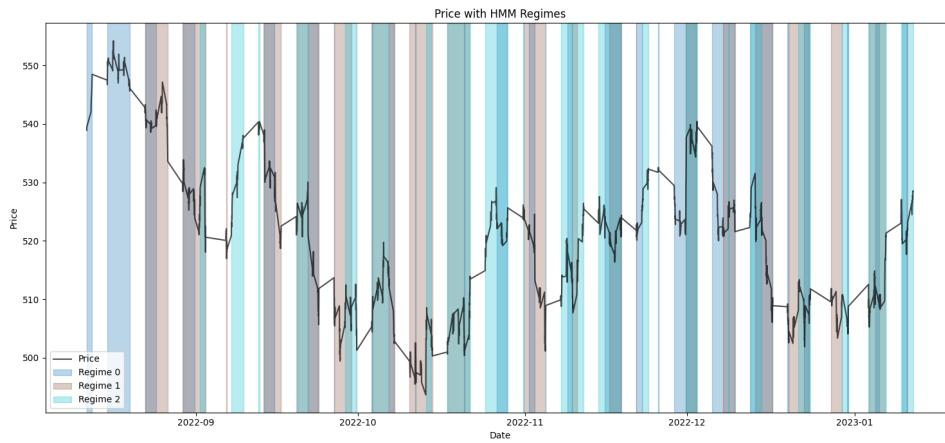


Figure 2: Market regimes identified by the Hidden Markov Model (HMM) approach. The HMM captures regime transitions but shows some sensitivity to short-term market fluctuations.

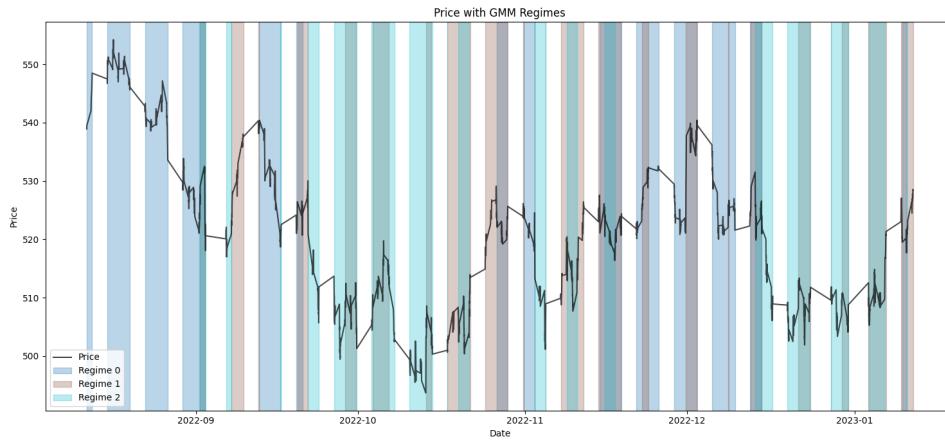


Figure 3: Market regimes identified by the Gaussian Mixture Model (GMM) approach. The GMM effectively captures the distributional differences between regimes but is less sensitive to temporal dependencies.

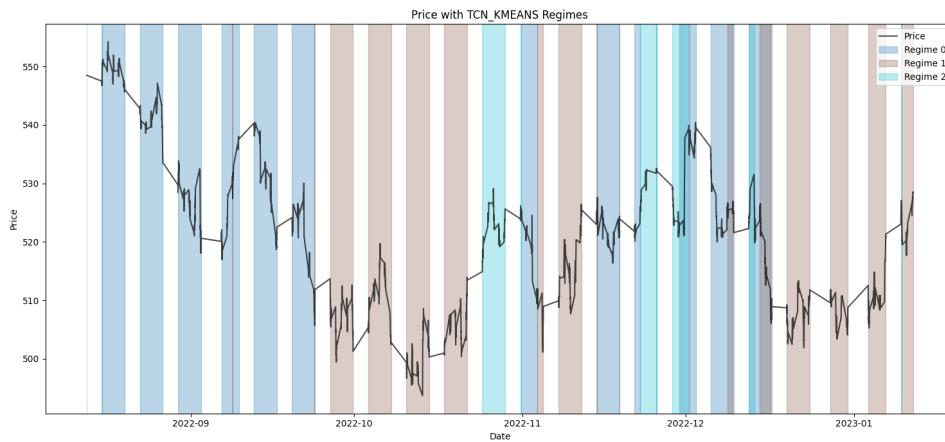


Figure 4: Market regimes identified by the Temporal Convolutional Network with K-means clustering (TCN+K-means) approach. This method effectively captures complex temporal patterns in the data.

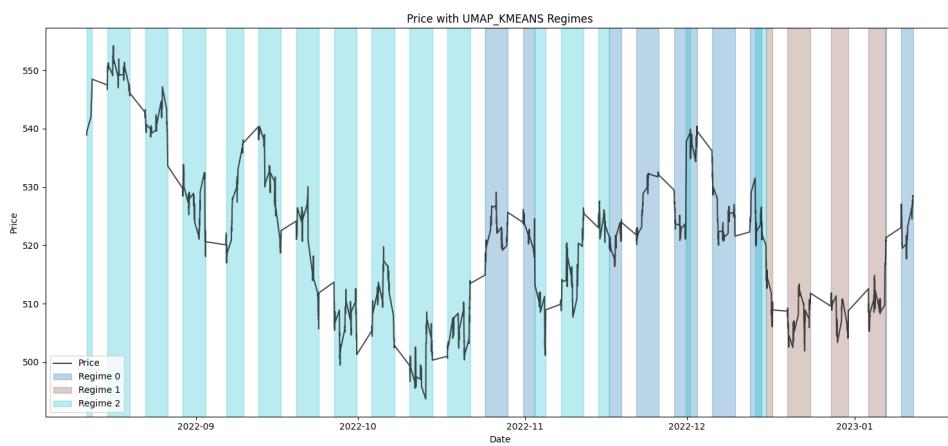


Figure 5: Market regimes identified by the UMAP dimensionality reduction with K-means clustering (UMAP+K-means) approach. This method preserves both local and global structure in the data.

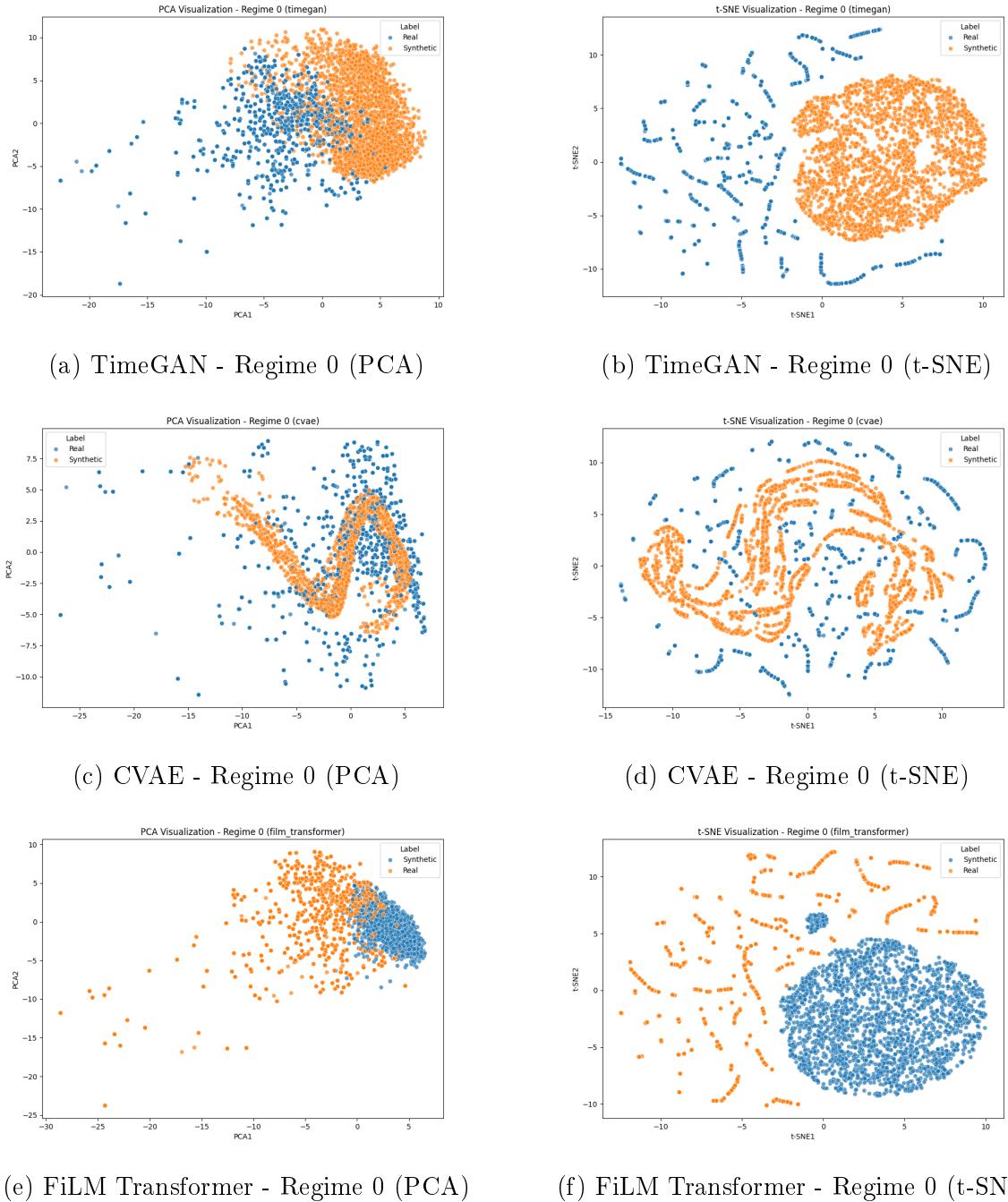


Figure 6: Comparison of real and synthetic data distributions for Regime 0 (bull market) using different generative models. The plots show the projection of data points onto the first two principal components (PCA) and t-SNE embeddings. Real data points are shown in blue, while synthetic data points are shown in orange. The FiLM Transformer generates synthetic data that most closely matches the distribution of real data.

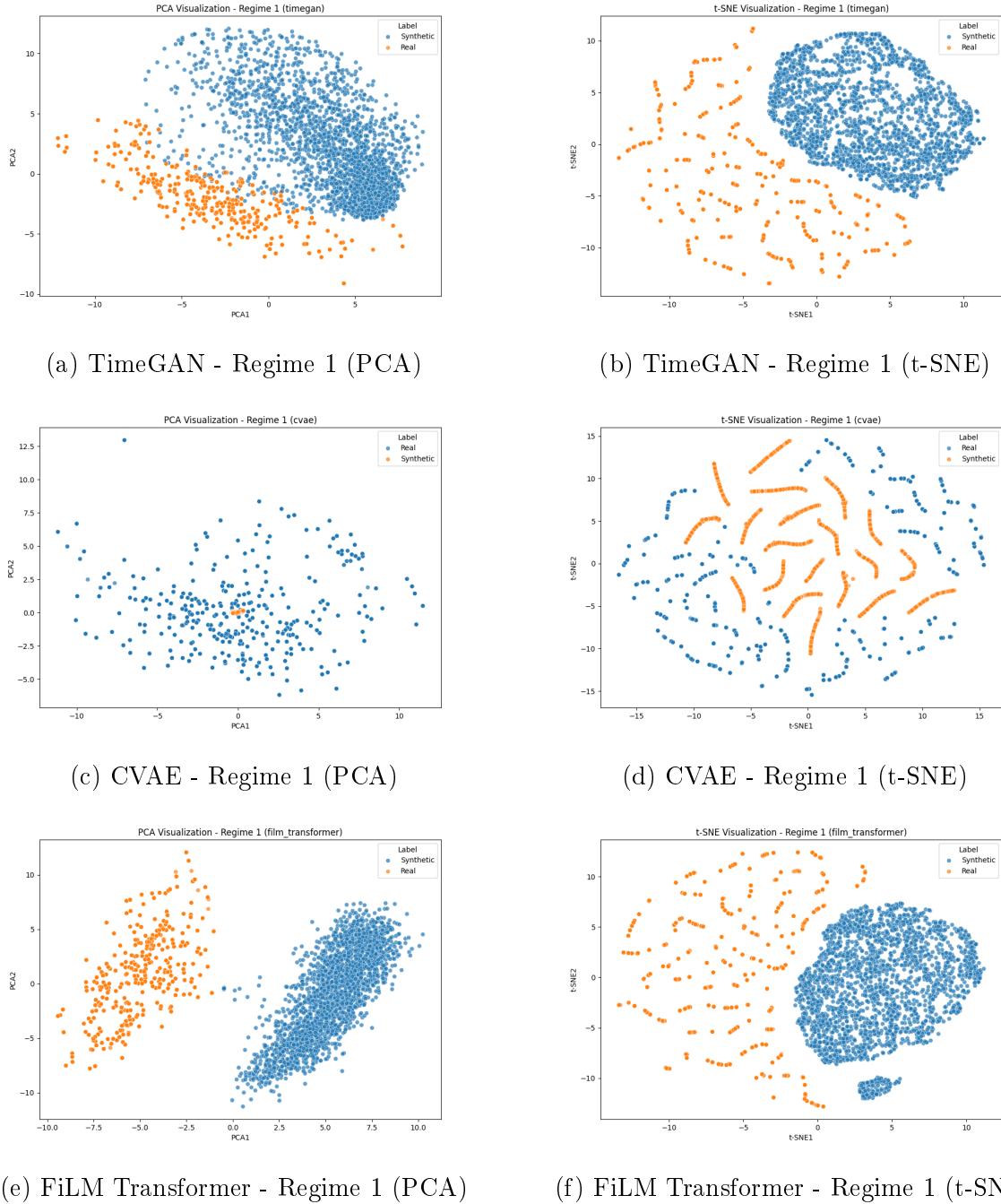


Figure 7: Comparison of real and synthetic data distributions for Regime 1 (neutral/-transition market) using different generative models. The plots show the projection of data points onto the first two principal components (PCA) and t-SNE embeddings. Real data points are shown in blue, while synthetic data points are shown in orange. The FiLM Transformer generates synthetic data that most closely matches the distribution of real data.

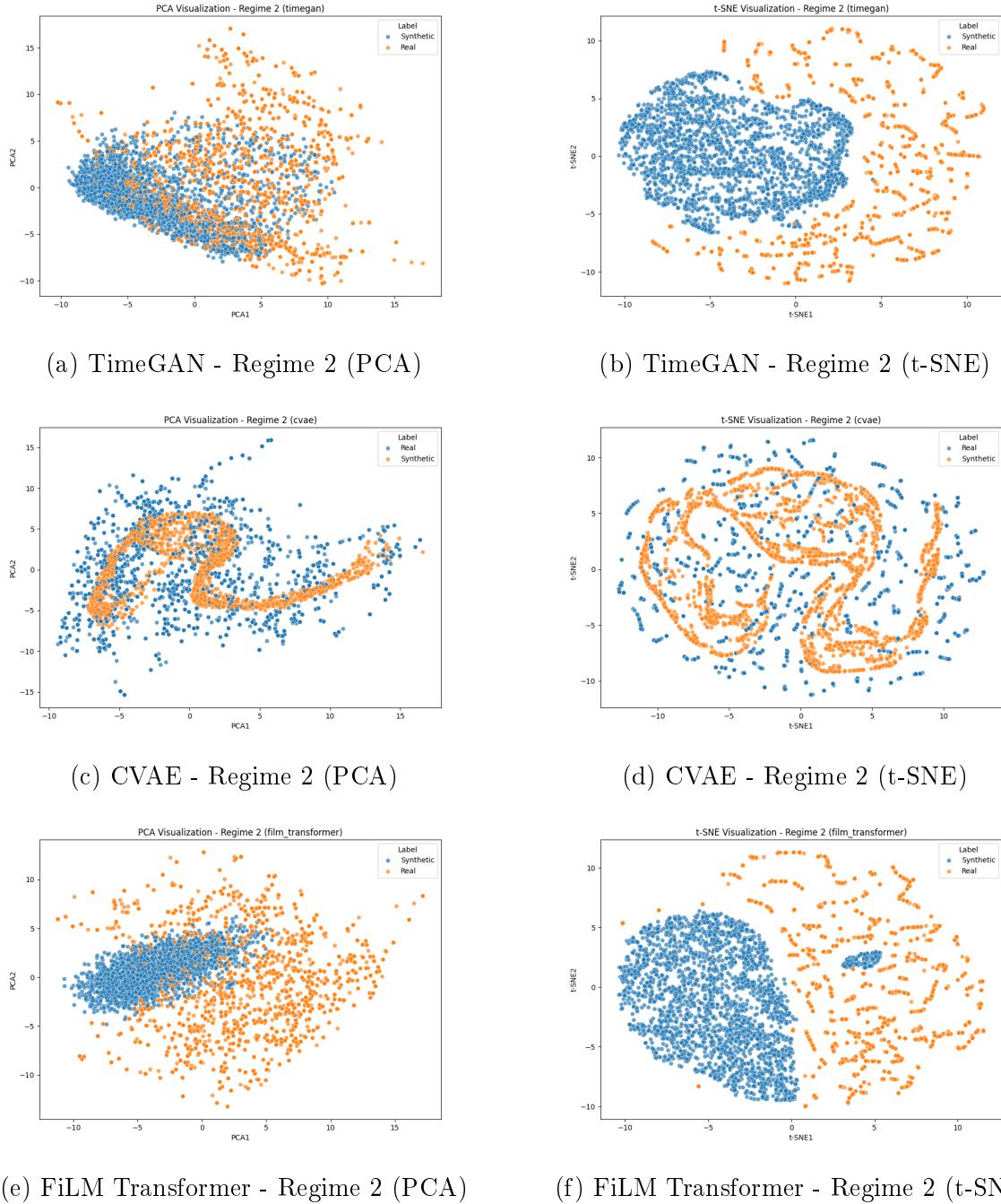


Figure 8: Comparison of real and synthetic data distributions for Regime 2 (bear market) using different generative models. The plots show the projection of data points onto the first two principal components (PCA) and t-SNE embeddings. Real data points are shown in blue, while synthetic data points are shown in orange. The FiLM Transformer generates synthetic data that most closely matches the distribution of real data, particularly in capturing the extreme tail behavior characteristic of bear markets.



Figure 9: Cumulative returns of trading strategies using the baseline approach (without synthetic data augmentation). The plot shows the performance of different strategies across the test period (2018-2023).



Figure 10: Cumulative returns of trading strategies using the augmented approach (with synthetic data augmentation). The plot shows the improved performance of strategies trained with synthetic data across the test period (2018-2023).