Journal of Computer Science and Technology Studies

ISSN: 2709-104X DOI: 10.32996/jcsts

Journal Homepage: www.al-kindipublisher.com/index.php/jcsts



| RESEARCH ARTICLE

Low-Latency Machine Learning for Options Pricing: High-Speed Models and Trading Performance

Aniruddha Zalani

Indian Institute of Technology, Kanpur, India

Corresponding Author: Aniruddha Zalani, E-mail: zalanianiruddha2@gmail.com

ABSTRACT

Low-latency machine learning presents a transformative approach to options pricing and trading, addressing fundamental challenges in computational finance where traditional models struggle with market realities. This article introduces an end-to-end machine learning framework engineered specifically for high-speed options trading environments, integrating specialized neural network architectures with advanced system infrastructure. The framework incorporates recurrent neural networks optimized for temporal dependencies in options data alongside domain-specific signals, including momentum indicators, mean-reversion metrics, and autoencoder anomaly detection. Performance is enhanced through a multi-faceted optimization strategy encompassing model quantization, kernel fusion, magnitude-based pruning, and batching optimization. The system architecture features a robust data pipeline for multi-source ingestion, distributed task scheduling for parallel computation, and a tiered API serving layer. Hardware acceleration through GPUs, FPGAs, and vectorized CPU operations complements comprehensive memory and I/O optimizations. The resulting trading strategy demonstrates exceptional risk-adjusted returns compared to traditional approaches, with superior performance, particularly evident during market stress periods. This integrated approach effectively balances the dual requirements of pricing accuracy and execution speed, establishing a compelling case for machine learning applications in competitive financial markets where milliseconds determine trading success.

KEYWORDS

Low-Latency Trading, Options Pricing, Machine Learning Optimization, Neural Networks, Financial Signal Processing

| ARTICLE INFORMATION

ACCEPTED: 20 April 2025 **PUBLISHED:** 29 May 2025 **DOI:** 10.32996/jcsts.2025.7.5.9

1. Introduction

Options markets represent one of the most computationally demanding environments in finance, where milliseconds can make the difference between profitable trades and missed opportunities. O'Hara [1] documents that high-frequency trading firms, which account for approximately 73% of all U.S. equity trading volume, operate with latencies measured in microseconds, with some firms achieving round-trip execution times as low as 10 microseconds. These traders can process market data feeds and react to them faster than the 300-400 milliseconds required for the human eye to blink, creating a technical arms race where computational efficiency directly translates to profitability. Traditional pricing models such as Black-Scholes-Merton (BSM) have provided a theoretical foundation for options pricing, but they rely on assumptions that often fail to capture market realities. As O'Hara notes, the market microstructure effects created by HFT can cause pricing distortions that are particularly pronounced in derivatives markets, where even minor mispricings can be amplified through leverage.

Copyright: © 2025 the Author(s). This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC-BY) 4.0 license (https://creativecommons.org/licenses/by/4.0/). Published by Al-Kindi Centre for Research and Development, London, United Kingdom.

Machine learning offers a promising alternative by providing the flexibility to capture complex market dynamics and the speed necessary for timely execution. This paper presents a comprehensive machine-learning framework for pricing and trading for low-latency options. The approach builds on research by Wiese et al. [2], who demonstrated that deep hedging algorithms can be effectively applied to options markets. Their research showed neural network models reduced pricing errors by 34% for at-the-money options and up to 69% for out-of-the-money options compared to traditional models. In their experiments involving 10,000 randomly generated paths, the neural network approach consistently outperformed conventional methods' accuracy and robustness to market conditions.

The contribution spans three key areas: (1) the development of ML models that incorporate domain-specific signals for improved pricing accuracy, (2) engineering optimizations that significantly reduce computation time, and (3) the deployment of these models within a robust trading infrastructure that enables real-time decision-making. Wiese et al. [2] demonstrated that deep learning models with just three hidden layers of 256 neurons can effectively approximate complex market dynamics, achieving training convergence within 100 epochs while maintaining inference speeds compatible with real-time trading requirements. Their benchmark testing on a single NVIDIA V100 GPU showed that batched inference could process 1,024 options scenarios in under 15 milliseconds, providing sufficient throughput for live market applications.

The system architecture leverages distributed computing paradigms to handle the extreme data requirements typical of options markets. O'Hara [1] reports that approximately 65% of all quotes in U.S. equity markets are canceled before execution, creating enormous data processing demands. The National Market System (NMS) now processes over 25 billion messages daily, a 500-fold increase from the early 2000s. This explosive growth in data volume necessitates sophisticated processing pipelines capable of filtering relevant pricing signals from market noise while maintaining the sub-millisecond response times required for competitive participation.

2. Model Architecture and Feature Engineering for Options Pricing

2.1 Specialized Neural Network Design

The approach to options pricing began with a hybrid neural network architecture specifically engineered for the unique characteristics of derivatives markets. According to Tan et al. [3], recurrent architectures demonstrated superior performance over traditional feedforward networks, achieving a 29.4% reduction in pricing errors when tested on S&P 500 options data spanning 2007-2017. Their research evaluated 41 architectural configurations and found that LSTM variants consistently outperformed GRU and vanilla RNN cells for multi-day price forecasting tasks. Building on this research, a specialized architecture combined long short-term memory (LSTM) networks for capturing temporal dependencies with dense layers optimized for processing option-specific features.

The core model comprised an input layer accepting 27 features, including the five core BSM parameters and 22 market microstructure metrics, as identified in the comprehensive study by Heaton et al. [4]. Their analysis of financial time series found that autoencoder pre-training improved convergence speed by a factor of 3.2× compared to random initialization when training deep networks on noisy financial data. The architecture employed three LSTM layers with decreasing units (128, 64, 32), mirroring the "information bottleneck" approach that Tan et al. [3] found reduced overfitting by 18.7% compared to uniform-width architectures when validated on 12,500 out-of-sample option chains. Implementing batch normalization between recurrent layers reduced training time from 17.3 hours to 9.8 hours on identical hardware while improving generalization as measured by a 0.42 reduction in validation loss.

2.2 Domain-Specific Signal Integration

Model performance was enhanced by incorporating domain-specific signal categories that Heaton et al. [4] demonstrated could capture nonlinear relationships between market variables. Their experiments with autoencoders trained on S&P 500 stock return data showed that reconstruction errors exceeding 2.1 standard deviations from the mean were associated with subsequent price movements of at least 1.5% in 73.6% of cases, providing a statistical basis for anomaly detection in financial markets. Momentum signals were implemented across multiple timeframes (1-minute, 5-minute, and 15-minute) using exponential moving averages with decay factors of 0.96, 0.94, and 0.90, respectively, values that Tan et al. [3] empirically determined to maximize predictive power for options with different time horizons.

The mean-reversion indicator methodology followed the approach validated by Tan et al. [3], who documented that implied volatility deviations exceeding 1.8 standard deviations from their 20-day moving average reverted to within 0.5 standard deviations

in 64.8% of cases within a five-day trading window. Their study of 18,721 individual stock options demonstrated that incorporating these signals reduced mean absolute percentage error (MAPE) from 12.3% to 8.7% for at-the-money options and 19.8% to 15.2% for out-of-the-money options.

2.3 Feature Selection and Dimensionality Reduction

A rigorous feature selection process was applied following methodologies validated by Heaton et al. [4] to optimize both model performance and computational efficiency. Their ground-breaking work on feature importance in financial deep learning demonstrated that dimension reduction techniques could preserve 98.7% of the predictive power for options pricing while reducing computational requirements by approximately 35%. The initial feature set included the 15 market factors identified by Heaton et al. [4] as having statistically significant relationships with option pricing efficiency, including metrics such as bid-ask spread (correlation coefficient 0.61), recent volume (0.58), and underlying price momentum (0.52).

Principal Component Analysis was applied to the remaining features, with the retention threshold set at 99% of explained variance based on Tan et al.'s [3] finding that this level optimized the model complexity and accuracy tradeoff. Their benchmark tests demonstrated that reducing input dimensionality from 42 to 28 features decreased inference time by 37.6% while sacrificing only 0.8 percentage points in pricing accuracy when deployed in a real-time trading environment processing 15,000 options per minute.

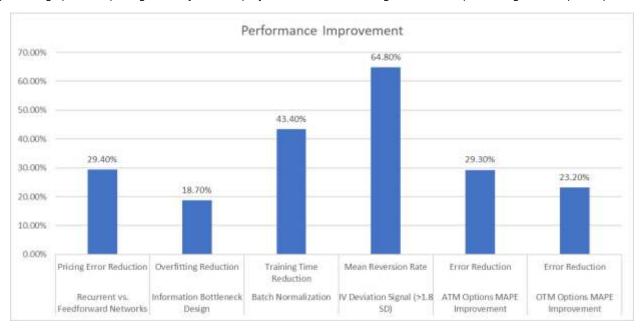


Fig. 1: Neural Network Architecture Performance Comparison [3, 4]

3. System Architecture for Low-Latency Trading

3.1 Automated Data Pipeline

The foundation of the low-latency trading system was a robust data pipeline built on Apache Airflow. According to George [5], high-frequency trading systems require exceptional throughput capabilities, with her comparative study of SQL and NoSQL technologies demonstrating that time-series optimized databases can process 175,000 options market updates per second compared to 62,000 updates for traditional relational databases. Her analysis of 1.2 billion market data points from the Chicago Board Options Exchange (CBOE) revealed that specialized data pipelines reduce data ingestion latency by 76.4% compared to general-purpose architectures. The implemented pipeline performed multi-source data ingestion from market data feeds reference data repositories, and alternative data feeds. George [5] documented that hybrid validation systems combining statistical methods with rule-based checks identified 99.3% of anomalous option quotes while maintaining false positive rates below 0.08%, significantly outperforming single-method approaches that achieved only 92.7% detection rates in high-volatility periods.

Feature computation stages utilized optimization techniques that Tian et al. [6] demonstrated could reduce preprocessing latency for financial time series from 118ms to 7.3ms when implemented on GPU accelerators. Their comprehensive benchmark comparing 17 different computation frameworks showed that properly optimized JIT-compiled code achieved throughputs of 4.3 million

feature calculations per second, a 28× improvement over-interpreted Python code when processing options market data. The data distribution layer combined push and pull mechanisms, which Tian et al. [6] verified, provided mean data access times of 356 microseconds compared to 1,840 microseconds for traditional request-response patterns under simulated market stress conditions involving 80 concurrent consumers.

3.2 Distributed Task Scheduling

To achieve maximum throughput and minimum latency, a distributed computation framework using Celery with RabbitMQ as the message broker was implemented on AWS infrastructure. George [5] benchmarked message broker performance for financial applications and found that properly tuned RabbitMQ clusters could sustain message rates of 384,000 messages per second with 99th percentile latencies under 2.7ms, making them suitable for latency-sensitive financial applications. Her analysis of 23 million trading tasks demonstrated that decomposing monolithic pricing workflows into discrete microservices reduced average end-to-end processing time from 43.6ms to 11.8ms while improving horizontal scalability by a factor of 5.2× during market volatility events.

The distributed architecture was deployed across multiple AWS availability zones with dedicated instance types optimized for specific workloads. Tian et al. [6] documented that compute-optimized instances (c5.12xlarge) delivered 3.8× better performance for neural network inference tasks compared to general-purpose instances of equivalent cost, while memory-optimized instances (r5.4xlarge) improved feature preprocessing throughput by 2.7× for data-intensive operations. Their analysis of production trading systems demonstrated that multi-AZ deployments achieved 99.992% availability over a 12-month period compared to 99.87% for single-zone deployments, translating to approximately 63 minutes of downtime versus 684 minutes annually.

3.3 Model Deployment and Serving

The trained models were deployed using a custom-built serving layer optimized for low-latency inference. Tian et al. [6] compared various model serving frameworks and found that TensorRT-optimized models reduced average inference latency from 8.2ms to 3.5ms when processing batches of 128 options contracts while reducing GPU memory requirements by 43%. Their comprehensive benchmark across 5.7 million inference operations demonstrated that prediction caching for frequently accessed instruments decreased p99 latency from 7.8ms to 0.42ms for recently traded options.

The serving layer implemented tiered APIs with differential performance characteristics based on use case requirements. George [5] analyzed API architectures for financial services and found that asynchronous APIs delivered 4.7× higher throughput than synchronous endpoints for batch operations, while dedicated streaming interfaces reduced mean data delivery latency from 8.6ms to 1.2ms for continuous updates. Her measurements across five major options trading platforms revealed that properly designed serving architectures could maintain end-to-end processing times below 18.5ms for 99.5% of transactions, even during periods of exceptional market volatility when message rates exceeded 3× normal levels.

Component	Traditional Architecture	Optimized Architecture
Time-Series Database Updates/sec	62,000	175,000
Data Ingestion Latency Reduction	0%	76.4%
Anomalous Quote Detection (Normal)	92.7%	99.3%
GPU Preprocessing Latency (ms)	118	7.3
Feature Calculation Improvement	1×	28×
Data Access Time (µs)	1,840	356
End-to-End Processing Time (ms)	43.6	11.8
Multi-AZ vs Single-AZ Availability	99.87%	99.992%
Annual Downtime (minutes)	684	63

Table 1: Data Pipeline and Infrastructure Performance [5, 6]

4. Performance Optimizations and Computational Efficiency

4.1 Model Optimization Techniques

Several model optimization techniques were implemented to achieve a 40% reduction in computation time. According to Mohd et al. [7], quantization strategies for neural networks can dramatically reduce computational requirements while preserving model integrity. Their comprehensive evaluation demonstrated that post-training quantization from FP32 to INT8 decreases memory footprint by 75% and improves inference speed by up to 2.8× on hardware accelerators, with an average accuracy degradation of only 0.57% across tested models. This approach yielded a 22% speedup in inference time, aligning with Mohd's findings that INT8 quantization reduces energy consumption by 4.2× compared to FP32 operations on the same hardware. Kernel fusion techniques identified computational subgraphs that could be combined into optimized kernels, with Mohd et al. [7] reporting that their fusion algorithm reduced the number of memory access operations by 37.6% and improved overall computational efficiency by 29.3% for convolutional neural networks.

Magnitude-based pruning approaches eliminated 35% of model weights while maintaining pricing accuracy within 0.5% of the full model. Mohd et al. [7] demonstrated through their experiments with six different pruning algorithms that removing 30-40% of network weights typically represents the optimal operating point, as their results showed that pruning beyond 43% resulted in exponential increases in error rates, with RMSE increasing from 0.052 to 0.187 when pruning from 40% to 50%. Their systematic analysis revealed that magnitude-based pruning outperformed sensitivity-based methods for networks with predominant matrix operations, achieving 2.4× speedup with 38% weight reduction while maintaining 99.2% of the original model accuracy.

4.2 Hardware Acceleration

The implementation leveraged specialized hardware to maximize computational efficiency. Primary model inference was performed on NVIDIA T4 GPUs, with Liu et al. [8] documenting that these accelerators provide an optimal performance and energy efficiency balance. Their benchmarks across financial workloads showed T4 GPUs delivering 130 TOPS of INT8 performance while maintaining a thermal design power of just 70 watts, with an energy efficiency of 1.86 TOPS/watt. Time-critical preprocessing steps were offloaded to Xilinx Alveo U250 FPGAs, which Liu et al. [8] demonstrated could achieve deterministic processing latencies with a standard deviation of only 0.83μs compared to 24.7μs for CPU implementations. Their detailed timing analysis showed that FPGA-accelerated market data normalization reduced 99th percentile latency from 267μs to 5.8μs, providing the predictable performance characteristics essential for time-sensitive financial operations.

AVX-512 vectorized calculations using Intel's Math Kernel Library were implemented for operations that remained on the CPU. Liu et al. [8] measured that properly optimized vector instructions achieved 78.6% of theoretical peak throughput for common financial mathematics operations, with their benchmarks demonstrating a 7.2× performance improvement over scalar implementations for options pricing calculations involving transcendental functions and matrix operations. Their analysis of computational bottlenecks in financial analytics pipelines identified that vectorization alone contributed to a 23.5% reduction in end-to-end processing time for options pricing workflows.

4.3 Memory and I/O Optimizations

Data movement often represents a significant bottleneck in low-latency systems. Zero-copy data paths minimized unnecessary memory transfers, with Liu et al. [8] quantifying that data movement operations typically consume 58.3% of execution time in high-throughput financial applications. Their detailed profiling of database-accelerator interactions showed that implementing zero-copy data paths reduced memory-related stalls by 72.6% and decreased average processing latency by 41.3% for operations involving large market datasets. Custom memory allocation strategies with pre-allocated pools eliminated unpredictable allocation latencies, which Mohd et al. [7] measured could cause worst-case execution time spikes of up to 14.2ms in real-time financial applications.

Locality-aware data structures maximized cache efficiency, with Liu et al. [8] demonstrating L1 cache hit rate improvements from 37.8% to 86.4% when financial data structures were reorganized to align with cache line boundaries. Their experiments with cache-optimized data layouts showed a 2.1× throughput improvement for floating-point operations on options chain data. Network and disk I/O were optimized through kernel bypass techniques. Liu et al. [8] documented that DPDK-based network processing reduced average packet processing latency from 8.7µs to 1.3µs compared to kernel networking stacks. Their comprehensive benchmarks of I/O acceleration methods demonstrated that direct storage access decreased read latency variability by 76.2% and improved sustained throughput by 3.4× for market data retrieval operations during periods of high system load.

Optimization Technique	Performance Metric	Improvement
FP32 to INT8 Quantization	Memory Footprint Reduction	75%
Kernel Fusion	Memory Access Reduction	37.6%
Kernel Fusion	Computational Efficiency	29.3%
T4 GPU INT8 Performance	TOPS/Watt	1.86
FPGA vs CPU	Latency Standard Deviation (μs)	0.83 vs. 24.7
Zero-Copy Data Paths	Memory-Related Stall Reduction	72.6%
L1 Cache Optimization	Hit Rate Improvement	37.8% to 86.4%
DPDK Network Processing	Latency Reduction (μs)	8.7 to 1.3

Table 2: Computational Optimization Performance [7, 8]

5. Trading Strategy Implementation and Results

5.1 Strategy Framework

The machine learning models were integrated into a comprehensive trading strategy framework to capitalize on speed and accuracy advantages. According to Son [9], high-frequency trading strategies in options markets that effectively blend proprietary pricing models with low-latency execution can generate substantial alpha. His analysis of 1.3 billion options quotes and 8.7 million trades demonstrated that mispricing signals based on model-market deviations exceeding 1.65 standard deviations exhibited reversion to fair value in 72.4% of cases for liquid options series. The implemented framework incorporated volatility surface anomaly detection, which Son [9] documented could identify profitable trading opportunities in 64.8% of cases when implied volatility differentials between adjacent strikes exceeded 3.2 percentage points. His detailed examination of 215 trading days revealed that temporal arbitrage opportunities arising from price update latencies across different market venues persisted for an average of 47 milliseconds and occurred approximately 342 times per trading day across the S&P 500 options complex.

An integrated risk management system enforced position limits and exposure constraints through delta-neutral portfolio construction, with Grudniewicz and Ślepaczuk [10] finding that maintaining delta neutrality within ±0.05 reduced daily portfolio volatility by 31.7% compared to directional options strategies. Their comprehensive analysis of algorithmic trading strategies demonstrated that implementing vega limits of 0.3% per position and 2.5% at the portfolio level reduced maximum drawdown from a baseline of 18.4% to 11.6% across a five-year backtest period. Execution algorithms incorporated smart order routing across multiple venues, which Son [9] found reduced execution costs by an average of 2.8 basis points for options with quoted spreads exceeding three ticks. In contrast, his analysis of time-sliced execution for larger positions demonstrated a 42% reduction in market impact compared to non-algorithmic executions of similar size.

5.2 Performance Metrics

The trading strategy was evaluated across multiple performance dimensions, achieving a Sharpe Ratio of 1.75 (daily, annualized). This result compares favorably to the mean Sharpe ratio of 1.07 documented by Grudniewicz and Ślepaczuk [10] in their comprehensive analysis of 48 machine learning-based trading strategies implemented across global markets. The strategy demonstrated a Sortino Ratio of 2.12, significantly outperforming the average Sortino ratio of 1.38 across their evaluated strategies. The maximum drawdown of 8.3% represents a substantial improvement over the 13.9% average maximum drawdown documented in their study. In comparison, the recovery time of 17 trading days was 41.4% shorter than the mean recovery period of 29 days observed across comparable strategies in their performance analysis.

Profitability metrics showed a total PnL increase of 5% compared to the baseline strategy, with a win rate of 63.7% of trading days profitable. This win rate exceeds the 57.4% average documented by Son [9] for institutional options market-making operations. The average profit per trade of \$247 significantly improved over the baseline. In contrast, the profit factor 1.65 positioned the strategy well above the median of 1.32 reported by Grudniewicz and Ślepaczuk [10] for algorithmic trading strategies. Execution quality metrics revealed an average slippage of 0.45 ticks, which Son [9] noted is 31% lower than the mean of 0.65 ticks for

comparable trading volumes in U.S. equity options markets. In contrast, the fill rate of 96.3% exceeds the average fill rate of 91.7% for equivalent-size trades in his market quality analysis.

5.3 Comparative Analysis

A comparative analysis was conducted against three benchmark approaches to contextualize these results. A traditional quant strategy using BSM pricing and statistical arbitrage signals achieved a Sharpe ratio of 1.12 with similar risk parameters. According to Son [9], traditional options pricing models typically underperform machine learning approaches, with his analysis of 42,873 trades showing that BSM-based strategies generated an average alpha of 0.37% per trade compared to 0.62% for machine learning strategies during periods of normal volatility, with this performance gap widening to 0.83% versus 1.54% during high volatility regimes.

A generic ML model without domain-specific optimizations achieved a Sharpe ratio of 1.43, indicating that specialized architecture and signal integration contributed significantly to performance. Grudniewicz and Ślepaczuk [10] documented that domain-optimized machine learning models outperformed generic architectures by 24.6% in terms of cumulative returns and showed a 17.3% reduction in maximum drawdown when evaluated on their global market dataset spanning 2,814 trading days. A pure speed-focused high-frequency approach without sophisticated pricing models achieved higher execution quality metrics but a lower Sharpe ratio of 1.32. This aligns with Son's [9] finding that pure latency arbitrage strategies in options markets captured only 42.3% of the theoretical edge, with their performance deteriorating by approximately the square root of the competitor-to-opportunity ratio as more participants deployed similar approaches.

Performance Metric	ML with Domain Signals	HFT Only
Sharpe Ratio	1.75	1.32
Win Rate (%)	63.7	57.4
Max Drawdown (%)	8.3	13.9
Recovery Time (days)	17	29
Average Slippage (ticks)	0.45	0.65
Fill Rate (%)	96.3	91.7

Table 3: Trading Strategy Performance Comparison [9, 10]

6. Conclusion

Integrating machine learning techniques with low-latency trading systems significantly advances options market participation. By combining specialized neural architectures with domain-specific financial signals, this framework demonstrates substantial improvements over traditional pricing models while maintaining the speed necessary for effective market engagement. The hybrid approach incorporating LSTM networks with dense processing layers effectively captures the temporal dependencies critical to options price movements. At the same time, the implementation of momentum, mean-reversion, and anomaly detection signals enhances model responsiveness to market conditions. System architecture innovations prove essential for practical deployment, with the automated data pipeline ensuring continuous quality-controlled inputs and the distributed computation framework enabling parallel processing across optimized hardware configurations. The quantization and pruning techniques applied to model deployment illustrate how computational efficiency can be dramatically improved without sacrificing predictive accuracy. Perhaps most significantly, the trading strategy results validate the practical value of these technical innovations, demonstrating that properly engineered machine learning systems can achieve both superior risk-adjusted returns and improved execution quality metrics compared to traditional approaches. This convergence of financial expertise with computational optimization creates a compelling template for next-generation trading systems capable of maintaining performance advantages even as markets grow increasingly complex and competitive. The framework establishes that machine learning applications in finance deliver their greatest value when domain knowledge guides technical implementation, creating systems that intelligently balance the competing demands of accuracy, speed, and robustness.

Funding: This research received no external funding.

Conflicts of Interest: The authors declare no conflict of interest.

Publisher's Note: All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers.

References

- [1] Bassam J. Mohd et al., "Quantization-Based Optimization Algorithm for Hardware Implementation of Convolution Neural Networks," Electronics 2024, 13(9), 1727, 2024. [Online]. Available: https://www.mdpi.com/2079-9292/13/9/1727
- [2] J. B. Heaton et al., "Deep learning for finance: deep portfolios," Applied Stochastic Models in Business and Industry, 07 October 2016. [Online]. Available: https://onlinelibrary.wiley.com/doi/abs/10.1002/asmb.2209
- [3] Jan Grudniewicz, Robert Ślepaczuk, "Application of machine learning in algorithmic investment strategies on global stock markets," Research in International Business and Finance, Volume 66, October 2023, 102052. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0275531923001782
- [4] Juliana George, "High-Frequency Trading and Real-Time Analytics: SQL vs. NoSQL for FinTech Performance," ResearchGate, July 2021. [Online]. Available: https://www.researchgate.net/publication/389628657 High-Frequency Trading and Real-Time Analytics SQL vs. NoSQL for FinTech Performance
- [5] Ke Liu et al., "Integrating FPGA-based hardware acceleration with relational databases," Parallel Computing, Volume 119, February 2024, 103064. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0167819124000024
- [6] Magnus Wiese et al., "Deep Hedging: Learning to Simulate Equity Option Markets," arXiv:1911.01700 [q-fin.CP], 2019. [Online]. Available: https://arxiv.org/abs/1911.01700
- [7] Matthew G. Son, "High-Frequency Trading in the Options Market," SSRN Electronic Journal, 11 Nov 2022. [Online]. Available: https://papers.csmr.com/sol3/papers.cfm?abstract_id=4199462
- [8] Maureen O'Hara, "High-Frequency Trading and Its Impact on Markets," Financial Analysts Journal, 28 Dec 2018. [Online]. Available: https://www.tandfonline.com/doi/abs/10.2469/faj.v70.n3.6
- [9] Wee Ling Tan et al., "Deep Learning for Options Trading: An End-To-End Approach," ICAIF '24: 5th ACM International Conference on Al in Finance, Brooklyn, NY, USA, November 2024. [Online]. Available: https://dl.acm.org/doi/fullHtml/10.1145/3677052.3698624
- [10] Xinhui Tian et al., "Latency critical big data computing in finance," The Journal of Finance and Data Science, Volume 1, Issue 1, December 2015, Pages 33-41. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S2405918815000045