

Intraday trading CUDA optimized

A Project Component for
Parallel and Distributing Computing (UCS645)

By

Sr	Name	Roll No
1	Arnav Gupta	102203590

Under the guidance of
Dr. Saif Nalband
(Assistant Professor, DCSE)



DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING
THAPAR INSTITUTE OF ENGINEERING AND TECHNOLOGY
(DEEMED TO BE UNIVERSITY)

PATIALA - 147004

MAY, 2025

Table of Contents

1	Introduction	1
1.1	Background	1
1.2	Introduction to Problem Statement	1
2	Literature Review	3
2.1	Technical Indicators in Trading Systems	3
2.2	Limitations of Traditional CPU-Based Systems	3
2.3	GPU and CUDA in Financial Computing	3
2.4	Visualization in Trading Systems	4
3	Research Gaps	6
3.1	Lack of Real-Time GPU-Accelerated Technical Indicator Systems	6
3.2	Integration Challenges in End-to-End Trading Pipelines	6
3.3	Lack of Evaluation on Lightweight, Resource-Constrained Systems	6
3.4	Visualization Often Treated as an Afterthought	7
3.5	Multi-Asset Scalability and Optimization Gaps	7
4	Problem Formulation	8
5	Objectives	9
6	Methodology	10
6.1	System Architecture Overview	10
6.2	Workflow Pipeline	10
6.3	Tools & Libraries	11
6.4	Equation	12
7	Datasets	16
8	Results and Discussion	17
8.1	Indicator Computation Performance	17
8.2	Trading Logic Execution	17

8.3 Discussion	18
References	19

List of Figures

List of Tables

1	Summary of Literature Review	4
2	Stock Price Dataset Used for Real-Time Technical Indicator Computation	16
3	Stock Price Dataset Used for Real-Time Technical Indicator Computation	17

1 Introduction

1.1 Background

In recent years, algorithmic trading has gained immense popularity, revolutionizing the financial markets by automating decision-making processes. The integration of advanced computational techniques, including parallel computing, artificial intelligence, and machine learning, has significantly enhanced trading efficiency and accuracy. One such area where computational power plays a crucial role is the use of technical indicators to guide stock trading decisions.

Technical indicators are mathematical calculations based on historical price, volume, or open interest data that traders use to predict future price movements. Examples of these include Moving Averages (SMA, EMA), Bollinger Bands, RSI (Relative Strength Index), and MACD (Moving Average Convergence Divergence). By analyzing these indicators, traders attempt to identify potential buying or selling points for a given asset.

However, the traditional approach to calculating these indicators is computationally expensive, especially when large datasets or high-frequency data are involved. This has driven the need for more efficient algorithms that can handle massive amounts of data in real time. CUDA (Compute Unified Device Architecture), a parallel computing platform and programming model developed by NVIDIA, allows for significant acceleration in such computations by leveraging the processing power of GPUs. This project focuses on harnessing the power of CUDA to build a stock trading model that calculates technical indicators and makes buy/sell decisions at high speeds, facilitating real-time trading.

This is how you cite a reference [?].

if there are subsection use subsection

1.2 Introduction to Problem Statement

In financial markets, speed and precision are critical. Traders and institutions rely on rapid data analysis and decision-making to gain a competitive edge. As more traders adopt algorithmic approaches, the computational load associated with analyzing large volumes of stock data in real-time has grown exponentially. Traditional CPU-bound

trading systems are becoming a bottleneck, especially when dealing with high-frequency data or when multiple technical indicators are involved in the decision-making process.

Furthermore, technical analysis—which forms the backbone of many trading strategies—requires the computation of complex indicators like MACD, RSI, and Bollinger Bands, often over sliding time windows. These calculations are inherently parallelizable but are typically executed sequentially in CPU-based systems. This inefficiency limits the scalability and responsiveness of such trading models.

The core problem this project aims to tackle is the inefficiency of traditional technical analysis-based trading systems in processing data fast enough to respond to market movements in real time. By shifting computationally intensive parts of the trading system—especially indicator calculations and trade signal generation—to the GPU using CUDA, we can exploit parallelism and achieve significant performance gains.

This project seeks to answer the following research and implementation questions:

- How can CUDA be effectively utilized to speed up technical indicator calculations?
- What performance improvements can be observed over a CPU-only implementation?
- How does this acceleration impact the responsiveness and effectiveness of trading decisions?
- Can GPU parallelism enable near real-time processing of live stock market data?

Through this, the project aims to develop a high-performance, single-stock trading model that integrates:

- Real-time data ingestion
- GPU-accelerated technical analysis
- Actionable buy/sell signals
- Portfolio performance logging

Though this current iteration does not support multi-stock analysis or graphical visualization, it lays a solid foundation for extending into a full-fledged, scalable, real-time trading platform in future versions.

2 Literature Review

The convergence of finance and technology has led to significant advancements in algorithmic trading systems. These systems use predefined rules based on statistical models, machine learning, or technical indicators to execute trades automatically. A critical component in improving their efficiency and responsiveness is the use of high-performance computing (HPC) methods, particularly parallel computing using GPUs.

2.1 Technical Indicators in Trading Systems

Technical analysis has long been used as a tool for understanding market trends and identifying trading opportunities. Indicators like Simple Moving Average (SMA), Exponential Moving Average (EMA), Relative Strength Index (RSI), MACD, and Bollinger Bands have proven effective in various market conditions [1]. These indicators are computationally intensive, especially when applied to high-frequency or large-scale historical data.

Past works (e.g., Achelis, 2001 [2]) have thoroughly described the mathematical basis and empirical relevance of technical indicators. More recent studies explore how combining multiple indicators can reduce false signals and improve trading accuracy (Rundo, 2019 [3]).

2.2 Limitations of Traditional CPU-Based Systems

Traditional trading models implemented on CPUs often suffer from limitations in speed and scalability. As datasets grow in size and complexity, especially in real-time trading, sequential processing becomes a bottleneck. Several researchers have pointed out the inefficiency of using CPUs for real-time financial data analysis, particularly in high-frequency trading environments (Aldridge, 2013 [4]).

2.3 GPU and CUDA in Financial Computing

The use of GPU acceleration in financial computing has garnered increasing attention. CUDA, developed by NVIDIA, enables parallel execution of thousands of lightweight

threads, ideal for tasks that involve repetitive numerical computationsâsuch as calculating technical indicators over sliding windows.

In a comparative study by M. Dixon et al. (2013) [5], CUDA implementations of financial algorithms demonstrated a significant speedup (up to 100x) over traditional CPU implementations. Similarly, Rakhsha et al. (2017) [6] implemented GPU-accelerated Monte Carlo simulations for option pricing, showing a dramatic improvement in execution times.

These studies collectively highlight the potential for GPU-accelerated trading systems, especially when real-time responsiveness and high throughput are required.

2.4 Visualization in Trading Systems

While visualization tools (e.g., matplotlib, Plotly) are widely used for analyzing trading performance and market behavior, their integration into real-time systems is still evolving. Research by Shah et al. (2018) [7] emphasized the importance of clear visual indicators (such as candlestick charts, overlays of technical indicators, and signal markers) for human decision-making in hybrid trading systems.

Though not currently implemented in this project, visualization remains a critical area for future development to enhance interpretability and user feedback.

Table 1: Summary of Literature Review

References	Concepts Highlighted	Technique(s) used	Significance/ Proposal	Limitations
[1]	Technical analysis basics and its role in trading	Use of SMA, EMA, RSI, MACD, Bollinger Bands	Foundation for technical indicator-based trading strategies	Lacks implementation-level detail; not performance-oriented

[2]	Mathematical definitions and variations of indicators	Algorithmic formulas for technical indicators	In-depth explanation of how indicators work	No optimization or real-time focus
[3]	Combining indicators to improve accuracy	Hybrid indicator strategies	Reduces false signals in trading	May overfit in volatile markets; lacks execution optimization
[4]	Drawbacks of CPU-based trading systems	Sequential processing in real-time systems	Highlights the inefficiencies in conventional trading setups	No mention of parallel or GPU-based solutions
[5]	GPU use in financial computing	CUDA implementation of financial models	Shows significant speedup (up to 100x) over CPU	Focuses more on option pricing, not technical indicators
[6]	GPU-accelerated simulations	CUDA for Monte Carlo simulations	Demonstrates practical use of GPUs in finance	Narrow scope; not directly linked to trading signal systems
[7]	Role of visualization in trading	Use of charts and overlays	Enhances interpretability for decision-making	Visualization is not real-time focused; lacks GPU integration

3 Research Gaps

Despite the increasing integration of computational techniques in finance, several critical research gaps persist at the intersection of technical analysis, high-performance computing, and real-time trading systems. The existing literature, while robust in foundational theory and simulation studies, reveals limitations in practical implementation, especially for GPU-accelerated trading logic. These gaps form the basis for the motivation and direction of the current project.

3.1 Lack of Real-Time GPU-Accelerated Technical Indicator Systems

While CUDA and GPU computing have been effectively applied to financial simulations and option pricing models, there is a noticeable scarcity of research that directly applies GPU acceleration to technical indicator-based trading systems. Most existing works focus on batch simulations or theoretical performance benchmarks, rather than real-time trade decision-making in volatile markets.

3.2 Integration Challenges in End-to-End Trading Pipelines

Studies often isolate performance improvements in either data processing or model computation, but rarely consider a full pipeline—from data fetching, preprocessing, indicator calculation, to trade signal execution. There is limited research on tightly integrating GPU-accelerated components into a real-time, modular trading framework.

3.3 Lack of Evaluation on Lightweight, Resource-Constrained Systems

Many implementations assume access to high-end servers or multi-GPU clusters. Few works explore how much performance gain can be achieved using CUDA on consumer-grade GPUs, making them less applicable to individual traders or small-scale setups.

3.4 Visualization Often Treated as an Afterthought

While some literature acknowledges the importance of data visualization for human decision-making, it is often treated as a separate task rather than being integrated into the system pipeline. Real-time systems that support visualization of both indicators and trading decisions are still underexplored—especially those that interact with GPU-processed data efficiently.

3.5 Multi-Asset Scalability and Optimization Gaps

Most academic studies and small-scale systems focus on a single asset or ticker, with limited provisions for extending to a multi-asset portfolio analysis. Optimizing technical analysis and trading decisions concurrently across multiple stocks—especially using GPU resources efficiently—is still an open research problem.

4 Problem Formulation

In the realm of algorithmic trading, technical indicators are extensively used to guide buy and sell decisions. These indicators often rely on historical price data and are computed over sliding windows, resulting in high computational overhead. When multiple indicators are used concurrently—as is common in robust trading strategies—the processing time increases significantly, especially as the volume of stock data grows. This latency poses a major problem in environments where real-time decision-making is crucial, such as intraday or high-frequency trading.

Most conventional trading systems are implemented using CPU-based architectures, which process indicator computations sequentially. While these systems can be accurate, they struggle with speed and scalability. As a result, traders relying on these systems may miss time-sensitive opportunities or suffer from lag in execution, particularly in volatile market conditions.

At the same time, advancements in GPU programming, particularly through NVIDIA’s CUDA framework, offer a promising solution. GPU architectures are well-suited for parallel computation and can execute thousands of threads simultaneously, making them ideal for the repetitive, data-parallel operations inherent in technical indicator calculations.

This project addresses the problem of slow, sequential indicator processing in traditional trading systems by formulating a CUDA-accelerated trading model. The model focuses on offloading computationally intensive indicator calculations—such as SMA, EMA, RSI, MACD, and Bollinger Bands—to the GPU for parallel execution. It uses real-time stock data (via Alpha Vantage API), parses it, computes indicators in parallel, and makes trading decisions based on predefined logic.

By accelerating these operations using CUDA, the model aims to significantly reduce processing latency, enabling faster, more responsive trading. Although the current implementation is focused on a single asset and does not yet include live visualization or multi-asset handling, it lays the groundwork for a scalable and performance-oriented trading engine suitable for real-world deployment.

5 Objectives

The main objective of this project is to design and implement a high-performance, GPU-accelerated trading system that leverages technical indicators for real-time stock trading decisions.

1. Primary Objectives

- To develop a CUDA-based computation engine for technical indicators such as SMA, EMA, RSI, MACD, and Bollinger Bands, optimizing them for parallel execution on the GPU.
- To build an end-to-end stock trading model that integrates real-time data fetching, indicator computation, trade decision logic, and result logging.
- To minimize latency in decision-making, enabling more responsive trading actions compared to traditional CPU-bound systems.
- To evaluate performance improvements in indicator computation between GPU (CUDA) and CPU implementations.

2. Secondary Objectives

- To implement a robust trade signal generation mechanism based on multiple technical indicators.
- To maintain a structured portfolio log capturing trade actions, execution prices, and value over time via CSV for later evaluation.
- To assess the resource efficiency of running the model on consumer-grade GPUs, enhancing accessibility for small-scale or individual traders.

6 Methodology

This project implements a CUDA-accelerated stock trading system in C++, integrating real-time data processing, technical indicator computation, and trading logic. The methodology involves building an efficient pipeline that shifts computationally expensive operations to the GPU using CUDA, while leveraging CPU resources for I/O and logic orchestration.

6.1 System Architecture Overview

The system is composed of the following modules:

- **Data Acquisition Module:** Uses Alpha Vantage API to fetch intraday stock price data in CSV format.
- **Preprocessing Module:** Cleans, parses, and structures the raw data into arrays for indicator calculation.
- **CUDA Indicator Kernel Module:** Implements GPU-parallelized computation of technical indicators (SMA, EMA, RSI, MACD, Bollinger Bands).
- **Trading Logic Module:** Applies decision rules based on indicator thresholds to generate buy/sell signals.
- **Portfolio Management Module:** Records trade actions and portfolio value over time into CSV logs.

6.2 Workflow Pipeline

1. **Fetch Stock Data:** Pull historical OHLC (Open, High, Low, Close) data using Alpha Vantage API. Store as CSV.
2. **Preprocess Data:** Read and parse the CSV. Convert it to suitable numerical arrays (e.g., float arrays) for GPU processing.
3. **GPU Computation of Indicators:** Launch CUDA kernels to calculate:

- SMA (Simple Moving Average)
- EMA (Exponential Moving Average)
- RSI (Relative Strength Index)
- MACD (Moving Average Convergence Divergence)
- Bollinger Bands
- Stochastic Oscillator
- ATR

Each kernel processes multiple data points in parallel using CUDA threads, exploiting the SIMD architecture of GPUs.

4. Trade Signal Generation: Based on user-defined logic (e.g., RSI ≥ 30 = Buy, MACD crossover = Buy/Sell), generate signals. This step is handled by the CPU after collecting computed indicators from the GPU.
5. Portfolio Update & Logging: Execute trades virtually and log:
 - Action (Buy/Sell)
 - Execution price
 - Updated portfolio value
 - Timestamp

6.3 Tools & Libraries

- C++: Core programming language for the trading logic and data handling.
- CUDA: NVIDIA's parallel computing platform for accelerating indicator computation.
- Alpha Vantage API: For real-time stock market data.
- CSV & Json Parsing Libraries: For efficient reading and writing of trade logs and price data. (Curl, nlohmann)

for equation you can use AI Tools to get the code or use this URL <https://latexeditor.lagrida.com/>

Algorithm 1 PPO

```

1: for  $iteration = 1, 2, \dots$  do
2:   for  $actor = 1, 2, \dots, N$  do
3:     Run policy  $\pi_{\theta_{old}}$  in environment for  $T$  time steps
4:     Compute advantage estimates  $\hat{A}_1, \dots, \hat{A}_T$ 
5:   end for
6:   Optimize surrogate  $L$  wrt.  $\theta$ , with  $K$  epochs and minibatch size  $M \leq NT$ 
7:    $\theta_{old} \leftarrow \theta$ 
8: end for

```

6.4 Equation

1. Simple Moving Average (SMA)

Formula:

$$SMA_t = \frac{1}{n} \sum_{i=0}^{n-1} P_{t-i} \quad (1)$$

Where:

- P_t is the closing price at time t
- n is the period of the moving average

CUDA Adaptation:

Each thread computes the SMA for a single time step by summing the last n values in parallel. Thread divergence is minimal, making this ideal for GPU.

2. Exponential Moving Average (EMA)

Formula:

$$EMA_t = \alpha \cdot P_t + (1 - \alpha) \cdot EMA_{t-1}, \quad \alpha = \frac{2}{n+1} \quad (2)$$

CUDA Adaptation:

Since EMA is recursive, a hybrid approach is used. The first few values are calculated on

the CPU, and then parallelized approximations (e.g., initializing with SMA) are executed on the GPU for subsequent steps.

3. Relative Strength Index (RSI)

Formula:

$$RSI = 100 - \left(\frac{100}{1 + RS} \right), \quad RS = \frac{\text{Avg Gain}}{\text{Avg Loss}} \quad (3)$$

CUDA Adaptation:

Threads independently compute gains and losses across intervals. Shared memory is used for intermediate reduction operations to calculate average gain/loss efficiently.

4. Moving Average Convergence Divergence (MACD)

Formula:

$$MACD = EMA_{12} - EMA_{26} \quad (4)$$

$$\text{Signal Line} = EMA_9(MACD) \quad (5)$$

CUDA Adaptation:

Multiple EMA streams are calculated in parallel using separate thread blocks, allowing fast calculation of the MACD and its signal line.

5. Bollinger Bands

Formula:

$$\text{Upper Band} = SMA + (k \cdot \sigma) \quad (6)$$

$$\text{Lower Band} = SMA - (k \cdot \sigma) \quad (7)$$

Where:

- σ is the standard deviation over the window
- k is a multiplier (typically 2)

CUDA Adaptation:

Threads compute both the SMA and variance over a moving window in parallel, followed by $\sqrt{}$ operations for standard deviation.

6. Stochastic Oscillator

The Stochastic Oscillator is a momentum indicator comparing a specific closing price to a range of prices over a given period.

$$\%K = \frac{(C - L_n)}{(H_n - L_n)} \times 100 \quad (8)$$

$$\%D = \text{SMA of } \%K \text{ over 3 periods} \quad (9)$$

Where:

- C = Current closing price
- L_n = Lowest low over the last n periods
- H_n = Highest high over the last n periods
- $\%K$ = Current stochastic value
- $\%D$ = Smoothed signal line (typically a 3-period SMA of $\%K$)

CUDA Adaptation: Each CUDA thread processes a sliding window to compute L_n and H_n in parallel using shared memory. Once the high-low range is determined, the $\%K$ value is calculated. The $\%D$ line is then generated using a parallel SMA kernel applied to $\%K$ values.

subsection*7. Average True Range (ATR)

The ATR measures market volatility based on recent price movements.

Step 1: True Range (TR)

$$TR = \max(H_t - L_t, |H_t - C_{t-1}|, |L_t - C_{t-1}|) \quad (10)$$

Step 2: Average True Range (ATR)

$$ATR_t = \frac{1}{n} \sum_{i=0}^{n-1} TR_{t-i} \quad (11)$$

Where:

- H_t = High price at time t
- L_t = Low price at time t
- C_{t-1} = Previous closing price
- n = Period for averaging (commonly 14)

CUDA Adaptation: Threads compute all three components of TR in parallel for each time step. The maximum operation is performed using built-in CUDA device functions. Once TR values are obtained, an SMA kernel is used to compute the ATR across the window using shared memory to store intermediate sums.

These formulas serve as the core of the trading model’s decision-making logic, and their parallel implementation drastically reduces computation time, enabling faster and more frequent trading evaluations.

7 Datasets

Table 2: Stock Price Dataset Used for Real-Time Technical Indicator Computation

Name of Datasets	Different Features	Size	Year	Organization
Alpha Vantage Intraday Stock Price Data	Open, High, Low, Close, Volume, Timestamp (5-min interval); supports thousands of global stock tickers	Approx. 4,000 rows per stock per 30 days (5-min granularity)	2025 (Real-time)	Alpha Vantage Inc.

8 Results and Discussion

The CUDA-accelerated trading model was tested using intraday stock data retrieved via the Alpha Vantage API, processed on a NVIDIA GPU (NVIDIA RTX 3050). The results demonstrate the feasibility and performance gains of GPU acceleration in real-time technical analysis.

8.1 Indicator Computation Performance

Benchmark comparisons between CPU and GPU computations of indicators showed significant performance improvements:

Table 3: Stock Price Dataset Used for Real-Time Technical Indicator Computation

Simple Moving Average (SMA)	3.3656 ms	0.5535ms	6x
Exponential Moving Average (EMA)	2.8091 ms	0.4576ms	6x
Relative Strength Index (RSI)	6.4405 ms	1.1934ms	5.3x
Bollinger Bands	10.4532 ms	3.2598ms	3x
MACD	11.0155 ms	4.6283	2.5x
Stochastic Oscillator	12.7368 ms	3.2894	4x
Average True Range (ATR)	5.7527 ms	2.3421	2.5x

These results indicate that technical indicators—particularly those based on rolling windows—are highly parallelizable and benefit substantially from GPU execution.

8.2 Trading Logic Execution

Based on predefined conditions (e.g., RSI \geq 30 triggers Buy, MACD crossover triggers Sell), the model successfully generated buy/sell signals. Each action was logged into a CSV file along with execution price and updated portfolio value.

Although currently limited to a single asset, the logic framework is modular and designed to support multiple stock tickers with minor extensions.

8.3 Discussion

This project validates that GPU acceleration via CUDA significantly enhances the speed of technical indicator computations. While real-time performance and accuracy are promising, extending the system to handle real broker integration, error handling, and visualization would be key next steps. With such enhancements, this model could evolve into a lightweight, scalable trading assistant for retail or institutional use.

References

- [1] J. J. Murphy, *Technical Analysis of the Financial Markets*. New York Institute of Finance, 1999.
- [2] S. B. Achelis, *Technical Analysis from A to Z*. McGraw Hill, 2001.
- [3] F. Rundo, “A survey on technical indicators in algorithmic trading,” *International Journal of Computer Applications*, 2019.
- [4] I. Aldridge, *High-Frequency Trading: A Practical Guide to Algorithmic Strategies and Trading Systems*. Wiley, 2013.
- [5] M. Dixon, I. Halperin, and P. Bilokon, *Machine Learning in Finance: From Theory to Practice*. Wiley, 2013.
- [6] R. Rakhsha, T. T. Nguyen, and J. Wang, “Cuda-based high-performance option pricing,” *Journal of Computational Finance*, 2017.
- [7] D. Shah, A. Braverman, and J. Tsitsiklis, “Visual analytics in financial modeling,” *IEEE Transactions on Visualization and Computer Graphics*, 2018.