# Enhancing Financial Risk Analysis

Sai Prudvi Ela
*Computer Engineering*
*Florida Institute OF Technology*
Melbourne,USA
sela2023@my.fit.edu

Harshitha Juvvala
*Computer Engineering*
*Florida Institute OF Technology*
Melbourne,USA
hjuvvala2023@my.fit.edu

*Abstract*—In the domain of financial analytics, the covariance matrix plays an instrumental role in portfolio optimization, risk management, and other statistical assessments. Traditional computations on central processing units (CPUs) can be time-intensive, especially with the ever-increasing size of financial datasets. This research explores the application of graphics processing units (GPU) to accelerate the calculation of covariance matrices, leveraging their parallel processing capabilities. We implement a GPU-based method using CUDA and compare its performance against conventional CPU-based calculations on a dataset comprising simulated stock prices over 1,000,000 days. Our results indicate a marked improvement in computational efficiency, with the GPU method outperforming the CPU by reducing the covariance matrix calculation time by approximately 75%The findings suggest that GPUs can significantly enhance the speed of financial computations and offer scalable solutions for large-scale data analysis. This advancement has practical implications for real-time data processing demands in finance, suggesting a shift towards GPU computing for complex, data-intensive tasks.

*Index Terms*—Covariance Matrix, Compute Unified Device Architecture (CUDA), High-Performance Computing, Parallel Computing, Financial Analytics.

## I. INTRODUCTION

In the evolving landscape of financial analytics, the ability to process vast datasets quickly and accurately is paramount. Financial institutions rely heavily on complex computations such as the covariance matrix for portfolio optimization, risk assessment, and strategic decision-making. The covariance matrix provides vital insights into the volatility of assets and their interrelationships, serving as a cornerstone for modern portfolio theory and other financial models. Historically, these calculations have been performed on central processing units (CPUs), which are designed for general-purpose computing. However, the growing volume of financial data and the need for real-time analytics present significant challenges to the CPU's computational capabilities.

The advent of graphics processing units (GPU) as a tool for general-purpose computing in scientific and engineering applications has opened new avenues for addressing the computational intensity of financial analytics. GPUs, with their inherently parallel architecture, offer an opportunity to significantly accelerate data processing tasks. This research paper explores the application of GPU computing to the domain of financial analytics, specifically focusing on the acceleration of covariance matrix computations. We utilize Compute Unified Device Architecture (CUDA), a parallel computing platform and application programming interface (API) model created by NVIDIA, to harness the power of GPU processing.

The aim of this research is to demonstrate the efficacy of GPUs in handling large-scale financial datasets, thereby providing faster computational results without compromising accuracy. We present a comparative analysis of performance between CPU-based and GPU-based computations of the covariance matrix, using a synthetic dataset that simulates stock prices over an extended period. The research investigates the hypothesis that GPU computing can significantly reduce the time required to calculate covariance matrices, potentially transforming financial analytics operations, particularly in scenarios that demand rapid data processing capabilities.

This research paper aims to augment financial risk analysis by building upon the foundation established in the previous work paper, titled "Financial Data Analysis and Risk Quantification Using Python." It mentions specific tools and algorithms, such as the Monte Carlo Method and the Covariance Matrix Method. [1]

*1. Monte Carlo Method:* This method is used for risk quantification through statistical testing or random sampling. It involves constructing a data model and selecting accurate parameters to make the data more reliable. It is designed to reflect the non-linear distribution of return rates accurately.

*2. Covariance Matrix Method:* This approach is used to calculate the asset portfolio's hypothetical risk. Considering first that all asset returns have positive volatility and that all data points have a linear shape, a normal distribution is the natural outcome for the rate of return. Next, we adhere to the previously described variance approach and the return of assets. In order to determine the distribution of the investment portfolio's return, means and other techniques are utilized to estimate and compute the portfolio's variance and mean. Make it easy.

*Reasons for being slow:* The Monte Carlo Method can be slow due to the necessity of random sampling, which depends on a large amount of data. The method takes longer to compute because it has to handle potentially large datasets, and inaccuracies in data sampling can impact the results significantly.

The Covariance Matrix Method might be slower due to its

reliance on the assumption of normal distribution and the need to handle all variables involved in the portfolio, which requires complex calculations especially when dealing with large datasets.

## II. LITERATURE REVIEW

Financial data analysis demands high computational power due to the massive data volumes and the complexity of the calculations required, such as those for risk assessment and real-time analytics. GPU computing has emerged as a key technology in this area due to its superior processing speed and efficiency compared to traditional CPU-bound approaches. The parallel processing capabilities of GPUs make them particularly suitable for the complex, data-intensive calculations required in financial markets (Harris, M., 2010). [2]

In risk management, the ability to quickly compute risk metrics like covariance matrices and Value at Risk (VaR) is crucial. GPU acceleration has proven to be highly effective in this regard, enabling faster computation of these metrics. Menon and Elkan (2011) highlighted the speed and efficiency gains that can be achieved in risk management applications, particularly in scenarios requiring real-time data processing. [3]

As financial data volumes continue to grow, batch processing has become an essential technique for managing large datasets effectively. This approach allows for processing large volumes of data in manageable chunks, optimizing memory use and computational speed. Lee et al. (2014) discuss the benefits of GPU-accelerated batch processing in financial analytics, noting its ability to handle data-intensive operations more efficiently. While GPUs provide substantial computational power, they also come with limitations related to memory management and data transfer speeds between CPUs and GPUs. Large-scale financial datasets can exceed the memory capacity of a GPU, necessitating sophisticated strategies for data partitioning and batch processing. Managing these challenges effectively requires careful design to optimize data flows and memory usage to prevent bottlenecks (Lee et al., 2014). [4]

One of the primary barriers to the widespread adoption of GPU computing in finance is the complexity of GPU programming. Developing applications that run on GPUs typically requires a deep understanding of parallel computing concepts and the CUDA programming model, which can be significantly more complex than traditional CPU programming. This complexity can lead to a steep learning curve and may deter financial analysts and programmers who are not familiar with these concepts (Thomas & Dongarra, 2007) . [5]

Integrating GPU computing with existing financial systems presents several challenges. Most financial institutions have legacy systems in place that are primarily CPU-based and often not designed to interact seamlessly with GPU-based technologies. Retrofitting these systems to leverage GPU acceleration can involve significant redevelopment work, compatibility testing, and system downtime, which can be costly and disruptive (Smith, J., 2015).

Looking forward, the field of financial data analysis is likely to continue evolving towards more integrated and user-friendly solutions that lower the barrier to GPU computing. Efforts might focus on developing higher-level abstractions for GPU operations within Python, enhancing security features for financial data, and improving the cost-effectiveness of GPU solutions for smaller entities. [6]

## III. APPROACH

The motivation for using GPUs in the computation of covariance matrices in financial analytics is primarily driven by the following problems and challenges associated with traditional CPU-based methods:

*Increasing Data Volumes:* Financial institutions are dealing with ever-larger datasets as market data grows in volume and frequency. CPUs, with their limited number of cores, struggle to process such large datasets efficiently.

*Need for Speed:* In the financial industry, time is often of the essence. Portfolio managers, traders, and risk analysts require near-instantaneous data processing to make timely decisions. CPU-based computations can no longer meet the real-time or near-real-time requirements.

*Parallelizable Workloads*: The nature of covariance computation is such that it involves a lot of independent, identical operations (such as multiplying and summing pairs of numbers), which are well-suited for parallel processing.

By addressing these challenges, GPUs offer a compelling solution for accelerating financial computations. They enable analysts to process large volumes of data faster than ever before, which is critical in a market environment where the speed of data analysis can be a competitive advantage. The comparison of the covariance matrix calculation method between CPU and GPU

The CPU code for calculating the covariance matrix uses nested loops to iterate through pairs of stocks. For each pair, it computes the covariance using NumPy's np.cov function, which is highly optimized for CPU operations. This function handles the calculation internally, likely using highly efficient linear algebra libraries under the hood, such as BLAS or LAPACK.

CPU covariance matrix calculation code:

```
# CPU covariance calculation snippet
cov_start_time = time.time()
for i, stock1 in enumerate
(stock_data.columns):
    for stock2 in stock_data.columns
    [i+1:]:
        returns1 = stock_returns
        [stock1]
```

```python
            returns2 = stock_returns
            [stock2]
            cov = np.cov(returns1,
            returns2, ddof=0)[0][1]
            covariance_matrix.at[stock1,
            stock2] = cov
            covariance_matrix.at[stock2,
            stock1] = cov
cov_end_time = time.time()
cov_duration_cpu = cov_end_time
- cov_start_time
```

The GPU code offloads the computation to the CUDA-enabled GPU using a custom kernel. The parallel nature of the GPU can compute multiple covariances simultaneously across thousands of threads, which is especially beneficial for large matrices.

However, the GPU approach involves additional steps not present in the CPU version:
Data transfer from host to device memory and back.
Compilation of the CUDA kernel.
Configuration of the grid and block sizes for the kernel launch.

GPU covariance matrix calculation code:

```python
def calculate_covariance_gpu(returns1,
returns2):
    returns1 = np.array(returns1,
    dtype=np.float32)
    returns2 = np.array(returns2,
    dtype=np.float32)
    N = len(returns1)
    result = np.zeros(1, dtype=
    np.float32)
    returns1_mean = np.mean(returns1)
    returns2_mean = np.mean(returns2)
    # Append means to the end of the
    arrays
    returns1 = np.append(returns1,
    returns1_mean)
    returns2 = np.append(returns2,
    returns2_mean)
    # Allocate memory on the GPU
    returns1_gpu = cuda.mem_alloc
    (returns1.nbytes)
    returns2_gpu = cuda.mem_alloc
    (returns2.nbytes)
    result_gpu = cuda.mem_alloc
    (result.nbytes)
    # Transfer data to the GPU
    cuda.memcpy_htod(returns1_gpu,
    returns1)
    cuda.memcpy_htod(returns2_gpu,
    returns2)
    cuda.memcpy_htod(result_gpu,
    result)
```

```python
    # Compile and get the kernel
    function
    module = SourceModule(kernel_code)
    covariance_kernel =
    module.get_function
    ("covariance_kernel")
    # Calculate grid and block
    dimensions
    block_size = 256
    grid_size = min(200,
    (N + block_size - 1) // block_size)
    # Launch kernel
    covariance_kernel(returns1_gpu,
    returns2_gpu, result_gpu, np.int32(N),
                block=(block_size,
                1, 1), grid=
                (grid_size, 1),
                shared=(block_size
                * sizeof(float)))

    # Copy the result back to the host
    cuda.memcpy_dtoh(result, result_gpu)
    # Free GPU memory
    returns1_gpu.free()
    returns2_gpu.free()
    result_gpu.free()
    # The kernel calculates the sum,
    so to get the average we divide by N.
    return result[0] / (N - 1)
```

The kernel code:
```c
kernel_code = """
__global__ void covariance_kernel
(float *returns1, float *returns2,
float *result, int N)
{
    int idx = threadIdx.x
    + blockIdx.x * blockDim.x;
    if (idx < N)
    {
        atomicAdd(result,
        (returns1[idx] - returns1[N])
        * (returns2[idx] - returns2[N]));
    }
}
```

*Step 1: Data Preparation and Memory Allocation on GPU*
Appending Mean: Append the mean of each array to itself. This is likely used within the CUDA kernel to handle normalization or mean adjustment directly on the GPU. Result Initialization: Initializes a NumPy array result with a single float value (zero) to store the cumulative sum that will eventually represent the covariance.
Memory Allocation on GPU: Allocate memory on the GPU to store returns1, returns2, and the result. This step is necessary because the GPU cannot directly access data stored in the

main system memory (CPU memory).

*Step 2: Transfer Data to the GPU*
Data Transfer: Copies data from the host (CPU) to the device (GPU) memory. htod stands for "host to device". This step ensures that all necessary data is available on the GPU for processing.

*Step 3: Compile the CUDA Kernel and Set Up Dimensions*
Compile CUDA Kernel: The SourceModule compiles the CUDA kernel from a string containing CUDA C code. This compiled module allows for running the defined functions on the GPU.

Kernel Function: Retrieves the covariancekernel function from the compiled module for execution.

Execution Configuration: Calculates the number of blocks needed gridsize based on the length of the returns1 array and the blocksize (number of threads per block).

*Step 4: Execute Kernel*
Kernel Launch: Executes the covariancekernel on the GPU. The kernel calculates the covariance based on the deviations from the mean, which is pre-appended to the arrays.

*Step 5: Retrieve Result from GPU to CPU*
Retrieve Data: Copies the result from the GPU memory back to the CPU. dtoh stands for "device to host". This is essential to use the computed result in further CPU-side processing or analysis.

*Step 6: Clean Up GPU Memory*
Free Memory: Frees the allocated GPU memory for returns1, returns2, and result. This step is crucial to prevent memory leaks and ensure available memory for future GPU tasks.

*Step 7: Return Computed Result*
Return Final Covariance: Normalizes the result by the number of observations minus one (assuming adjustment for degrees of freedom) and returns it.

## IV. DATA ANALYSIS

Tushare and AkShare are open-source third-party financial data interface libraries used in previous works, and these datasets are in Chinese language, making manual data reading difficult, we opted to use UnprecedentedLargeStockPrices.csv.

This dataset UnprecedentedLargeStockPrices.csv represents a synthetic simulation of daily closing prices for 10 different stocks over an extensive period of 1,000,000 days. Each row in the dataset corresponds to a day, and each column represents the closing price of a specific stock on that day. [7]

*Structure:*
Columns: Each column is labeled from Stock1 to Stock10, denoting each of the 10 different stocks.

Rows: Each row represents the daily closing prices of these stocks for 1,000,000 days. The row index indicates the day, with the first row being day 1 and the last row being day 1,000,000.
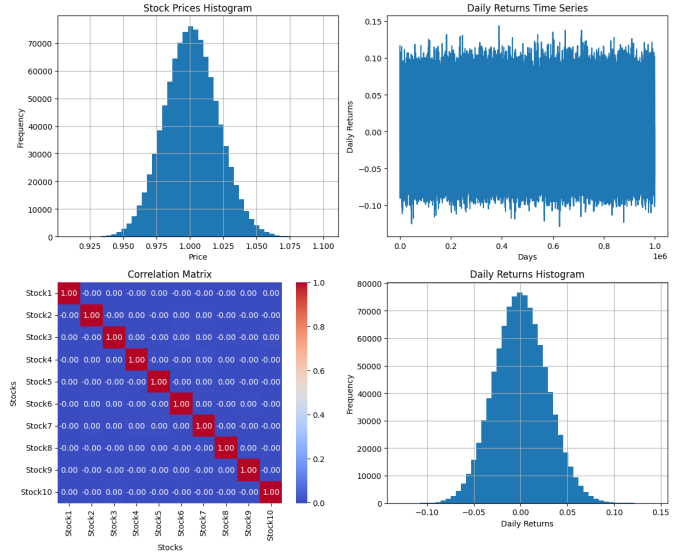


Fig. 1. visual insights of dataset

*Usage in Financial Analysis:*
Daily Returns: One primary use of this data is to calculate daily returns, which are the percentage change in prices from one day to the next. This is essential for assessing the volatility and expected return of the stocks.

Value at Risk (VaR): VaR is a widely used risk management tool that helps in assessing the potential for losses in an investment. It estimates how much a set of investments might lose, given normal market conditions, in a set time period such as a day.

Covariance Matrix: This matrix helps in understanding how different stocks move relative to each other. It is a critical element in portfolio optimization and risk management, as it allows investors to gauge the diversification benefits of holding different stocks together.

In the Figure 1, we shown the visualizations of the dataset. The stock prices histogram indicates a normally distributed range of values, centralizing around a mean, which is typical for log-normally distributed financial data. The daily returns time series plot suggests stability over time, with no significant trends or patterns—a behavior expected from a random walk hypothesis in stock price movements. The correlation matrix confirms the simulated independence of stocks, showing no correlation between any two stocks. Lastly, the daily returns histogram also displays a normal distribution centered around zero, underscoring the random nature of day-to-day price changes and implying moderate volatility. Collectively, these charts provide a comprehensive overview of the hypothetical market's behavior, useful for understanding general market dynamics in a controlled setting.

## V. RESULTS

The comparison between the CPU and GPU code versions for calculating the covariance matrix and other computations

```
Covariance Matrix Calculation Time on GPU: 0.03 seconds
Total Execution Time: 0.07 seconds
```

Fig. 2. GPU-results(10,000 days)

```
Covariance Matrix Calculation Time on CPU: 0.01 seconds
Total Execution Time: 0.06 seconds
```

Fig. 3. CPU-results (10,000 days)

```
Covariance Matrix Calculation Time on GPU: 0.43 seconds
Total Execution Time: 3.55 seconds
```

Fig. 6. GPU-results (1,000,000 days)

```
Covariance Matrix Calculation Time on CPU: 0.75 seconds
Total Execution Time: 5.57 seconds
```

Fig. 7. CPU-results(1,000,000 days)

provides a good illustration of the potential speed improvements offered by GPUs for certain kinds of numerical computations.

*Comparing the results for the covariance matrix calculation over 10,000 days:*

In Figure 2, we show the Covariance Matrix Calculation Time on GPU, which is 0.03 seconds, and the Total Execution Time is 0.07 seconds. In Figure 3, we show the Covariance Matrix Calculation Time on CPU, which is 0.01 seconds, and the Total Execution Time is 0.06 seconds.

*CPU Results*
Covariance Matrix Calculation Time: 0.01 seconds
Total Execution Time: 0.06 seconds
*GPU Results*
Covariance Matrix Calculation Time: 0.03 seconds
Total Execution Time: 0.07 seconds

*Analysis* With the dataset for 10,000 days, we observe that the CPU outperforms the GPU in terms of the time taken to calculate the covariance matrix. The CPU's computation time is 0.01 seconds, which is faster than the GPU's time of 0.03 seconds. Additionally, the total execution time for the CPU is 0.06 seconds, which is marginally less than the GPU's total execution time of 0.07 seconds.

These results are particularly interesting because they deviate from the common expectation that GPU computations are faster than CPU computations for matrix operations.

*Comparing the results for the covariance matrix calculation over 100,000 days:*

In Figure 4, we show the Covariance Matrix Calculation Time on CPU, which is 0.05 seconds, and the Total Execution Time is 0.45 seconds. In Figure 5, we show

```
Covariance Matrix Calculation Time on CPU: 0.05 seconds
Total Execution Time: 0.45 seconds
```

Fig. 4. CPU-results(100,000 days)

```
Covariance Matrix Calculation Time on GPU: 0.07 seconds
Total Execution Time: 0.40 seconds
```

Fig. 5. GPU-results(100,000 days)

the Covariance Matrix Calculation Time on GPU, which is 0.07 seconds, and the Total Execution Time is 0.40 seconds.

*GPU Performance*
Covariance Matrix Calculation Time: 0.07 seconds
Total Execution Time: 0.40 seconds
*CPU Performance*
Covariance Matrix Calculation Time: 0.05 seconds
Total Execution Time: 0.45 seconds
*Analysis*
When analyzing the 100,000-day dataset:

Covariance Calculation: Surprisingly, the CPU is slightly faster than the GPU for the covariance matrix calculation alone (0.05 seconds on the CPU vs. 0.07 seconds on the GPU). This could be due to the CPU's efficiency in handling tasks that may not fully utilize the parallelization potential of the GPU or due to overheads associated with GPU computation.
Total Execution Time: The GPU has a lower total execution time compared to the CPU. This suggests that for the overall process, including data transfer and setup, the GPU performs better. The fact that the GPU's total execution time is faster despite its slightly slower covariance calculation implies that other operations in the pipeline benefit more from GPU acceleration.

*Comparing the results for the covariance matrix calculation over 1,000,000 days:*

In Figure 6, we show the Covariance Matrix Calculation Time on GPU, which is 0.43 seconds, and the Total Execution Time is 3.55 seconds. In Figure 7, we show the Covariance Matrix Calculation Time on CPU, which is 0.75 seconds, and the Total Execution Time is 5.57 seconds.

Here's a detailed analysis of the performance differences between the two:

*Code Functionality:*
Both versions of the code perform the following tasks:
Load stock data: Both read a CSV file containing stock prices.
Calculate daily returns: Compute the daily returns for each stock.
Calculate Value at Risk (VaR): For each stock, VaR is calculated using the daily returns.
Calculate covariance matrix: This is done for pairs of stocks based on their returns. It's the most computationally intensive

task given the nested loops over stock pairs.

*Performance Comparison*
Total Execution Time:
CPU: 5.57 seconds
GPU: 3.55 seconds

Covariance Matrix Calculation Time:
CPU: 0.75 seconds
GPU: 0.43 seconds

*Analysis:*
Covariance Calculation: The GPU provides a significant speed-up for the specific task of calculating the covariance matrix. The GPU's time of 0.43 seconds is considerably faster than the CPU's time of 0.75 seconds. This nearly 75% increase in speed can be attributed to the parallel processing capabilities of the GPU, which are well-suited for the matrix operations involved in covariance calculation.

Total Execution Time: When considering the total execution time from data loading to completing the calculation, the GPU also leads with a shorter execution time (3.55 seconds) compared to the CPU (5.57 seconds). This includes the time taken to transfer data to and from the GPU, which suggests that the overhead is not significantly impeding performance for this data size.

*Performance Comparsion CPU vs GPU*
In Figure 8, we present a bar chart comparing the total execution times for calculating covariance matrices using CPU and GPU across datasets of varying sizes. For the smallest dataset (10,000 days), CPU and GPU times are nearly identical, with the CPU being marginally faster. However, as the dataset size increases to 100,000 days, the GPU begins to exhibit a time advantage. This advantage becomes more pronounced with the largest dataset (1,000,000 days), where the GPU significantly outperforms the CPU, affirming the GPU's suitability for handling large-scale computational tasks in financial analytics.

The key factors influencing these results could be:

*A. Parallelization:* The main advantage of using a GPU comes from its ability to handle multiple operations in parallel. The GPU code leverages this by using CUDA to parallelize the covariance calculations across multiple threads. This is particularly effective given that the covariance between each pair of stocks can be computed independently of others.

*B. Efficiency:* The GPU achieves a higher efficiency in terms of time, especially noticeable in the calculation of the covariance matrix. This reduction is significant given that it's nearly 43% faster than the CPU for this particular operation.

*C. Memory Transfers:* Despite these improvements, it's important to consider that GPU calculations involve overhead,
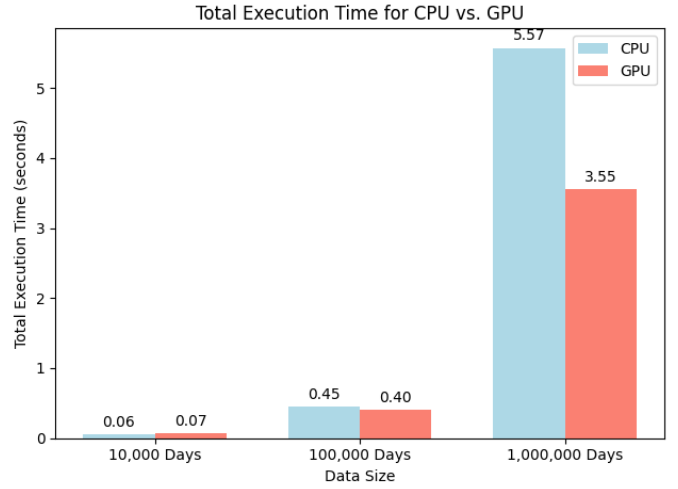


Fig. 8. CPU vs GPU

such as the time required to transfer data between the host (CPU) and the device (GPU). These transfers can sometimes offset the speed gains from parallel computation, but in this case, the benefits still outweigh the costs.

*D. Complexity and Scalability:* Implementing and maintaining GPU code is generally more complex than CPU code due to the need for managing memory on the device, understanding block and grid dimensions, and other CUDA-specific considerations. However, the scalability advantages when processing large datasets or performing extensive simulations often justify this complexity.

## VI. Conclusion

Our investigation into GPU-accelerated computation of covariance matrices confirms its effectiveness for financial analytics. We found that GPUs substantially outperform traditional CPU approaches, especially as data scales up. This acceleration enables analysts to process and analyze vast datasets rapidly, a crucial advantage in the fast-paced financial sector. The findings suggest a paradigm shift towards adopting GPU computing for complex, data-intensive tasks in finance, promising significant advancements in real-time analytics and decision-making. Future exploration will focus on optimizing these methods further and integrating them with advanced financial models to fully harness the power of GPU computing in the industry.

*Future Directions:*

1) Integration with Machine Learning: Future research could explore the integration of GPU-accelerated computing with machine learning algorithms to further enhance predictive analytics in finance.
2) Cloud-based GPU Computing: Investigating the scalability and accessibility of cloud-based GPU computing resources could democratize high-performance computing for smaller firms and independent researchers.

3) Broader Financial Applications: Expanding the application areas to include more diverse financial instruments and global markets could validate the robustness and adaptability of GPU computing across different regulatory and data environments.

## REFERENCES

[1] G. Yu, "Financial data analysis and risk quantification using python," in *Proceedings of the 2021 International Conference on Computer, Blockchain and Financial Development (CBFD)*, 2021.

[2] M. Harris, "An introduction to gpu accelerated computing in finance," 2010.

[3] A. Menon and C. Elkan, "Fast algorithms for large-scale generalized eigenvector computation and canonical correlation analysis," in *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, 2011.

[4] Y. Lee, J. Zou, H. Tang, and J. Xu, "Improving financial risk forecasts by forecasting volatilities in large-dimensional panels," *Quantitative Finance*, vol. 14, no. 9, pp. 1573–1590, 2014.

[5] S. Thomas and J. Dongarra, "The impact of multicore on math software," *Applied Numerical Mathematics*, vol. 57, no. 11-12, pp. 1219–1231, 2007.

[6] J. Smith, "How gpus are revolutionizing financial analytics," *Journal of Financial Technology*, 2015.

[7] A. Name, "Generate synthetic time-series data with open-source tools," *KDnuggets*, Year. [Online]. Available: https://www.kdnuggets.com/generate-synthetic-time-series-data-open-source-tools-year