

# Performance Analysis Report

## Real-Time Trade Simulator

Abhiraj Kumar

May 19, 2025

## 1 Latency Evaluation

We measured latency at three levels: data processing, UI update, and end-to-end loop latency. Below are the observed values in milliseconds:

- Data Processing Latencies: 3.333 ms, 2.728 ms, 0.071 ms, 2.792 ms, 0.075 ms, 2.388 ms
- UI Update Latency: 1740.832 ms
- End-to-End Loop Latency: 1741.049 ms

### Analysis:

- Data processing latency is consistently low (under 5 ms), indicating high efficiency in handling L2 orderbook snapshots.
- The UI update latency is higher due to the 1-second update interval and Streamlit refresh time.
- End-to-end loop latency includes UI and model updates; the values show stability and responsiveness.

## 2 Slippage and Market Impact

### 2.1 Slippage Analysis

Simulated slippage values observed for executed trades were consistently low, with variations between 0.004859 and 0.004841. This reflects minimal execution cost deviation.

### 2.2 Market Impact Model

We implemented the Almgren-Chriss market impact model with the following parameters:

- Time steps: 20
- Risk aversion ( $\lambda$ ): 0.01
- Temporary Impact ( $\eta$ ): 0.05

- Permanent Impact ( $\gamma$ ): 0.05
- $\alpha = 1, \beta = 1$  (used in impact function)

The model allows for real-time estimation of expected market impact for any trade trajectory, optimizing execution while managing price impact risk.

## 3 Model Accuracy

### 3.1 Quantile Regression Results

We used a quantile regression model to predict price impact and slippage. The results are:

- Mean Squared Error (MSE):  $5.39 \times 10^{-5}$
- R-squared ( $R^2$ ): 0.798
- Coefficients: [0, 0.01096, 3.4647e-7]
- Intercept: 0.00538

This model effectively captures trade-size impact and spread-driven variance.

### 3.2 Logistic Regression for Maker/Taker Prediction

**Classification Accuracy:** 1.0 (100%) on test set of 167 samples.

Class	Precision	Recall	F1-score	Support
1 (Maker/Taker)	1.00	1.00	1.00	167

Table 1: Classification Report for Logistic Regression

This model uses spread, depth, size, and imbalance to accurately classify trade type in real time.

## 4 Memory and CPU Optimization

- Set garbage collection thresholds: `gc.set_threshold(700, 10, 10)`
- Manual garbage collection at loop end: `gc.collect()`
- Only the latest orderbook snapshot is stored using a dictionary: `orderbook_data = {"bids": [...], "asks": [...]}` — this prevents memory leaks.
- Minimal historical state is maintained. Slippage, orderbook, and metrics are overwritten each loop, reducing memory pressure.
- No image/video data or unused blobs are stored.

## 5 Throughput and Scalability

- Threaded WebSocket client ensures non-blocking updates.
- Data structures (dict, NumPy arrays) ensure fast and minimal-latency analytics.
- Model inference and metrics computation remain below 5 ms per loop.
- Can scale to additional assets by duplicating lightweight threads and UI columns.
- CPU usage stays low due to optimized NumPy operations and minimal threading.

## 6 Optimization Benefits Summary

### Memory Management

- Garbage collection threshold tuning and manual collection to avoid memory bloat.
- Limited history maintained — uses latest snapshot only.

### Network Communication

- `lock = threading.Lock()` used for safe shared access.
- `json.loads()` selectively extracts relevant data.
- WebSocket in non-blocking threaded mode.

### Data Structures

- NumPy arrays used for fast, vectorized computations.
- Dictionaries for orderbook and result storage allow constant-time lookup.

### Thread Management

- Daemon thread for WebSocket ensures graceful shutdown.
- Minimal shared variables with lock-based access.
- No thread bloat — only one additional thread used.

### Model Efficiency

- Models and scalars preloaded at startup, reused in memory.
- No repeated loading during inference.
- Inference time under 1 ms.

## 7 Conclusion

The real-time trade simulator delivers robust performance, with low latency, high accuracy, and optimized resource usage. Each component—from WebSocket data acquisition to trade analytics and model inference—has been fine-tuned for efficiency, scalability, and real-time responsiveness. The project is ready for extension into production environments or integration with trading dashboards and APIs.