

Deep Learning vs. Financial Fraud Real-Time Detection in High-Frequency Trading

Abstract

HFT systems are sensitive to microseconds, and generate order-book streams that present novel challenges to the detection of fraud-related behaviors like spoofing and layering. Current algebraic-based surveillance strategies are ineffective in describing the nonlinear temporal patterns and subtle manipulation schemes that exist in contemporary financial markets since these strategies are typically solely rule-based. This paper presents an exploratory investigation of DL architectures to support real-time fraud detection in HFT, albeit with very low costs in terms of latency and in spite of predictions with a high level of accuracy. We test temporal convolutional networks (TCNs) and lightweight Transformers as well as machine learning- and rule-based baselines under a mix of historical data collected on a limit-order-book (LOB) and simulator-generated manipulation scenarios. We combine latency-aware acceleration techniques, including quantization, pruning, and micro-batching within a streaming design that can complete inference in sub-5 Ms. Experimental outcomes support the claim that DL models can perform detection better at extremely low false-positive rates even as operational service-level objectives are met. In addition to benchmarking, we cite difficulties of distributional robustness, deployment tradeoffs and explain ability, providing a reproducible framework and methodological improvements to applying deep learning to real time fraud detection in high frequency financial settings.

Journal

Journal of Science,
Technology and
Engineering
Research.

Volume-II, Issue-IV-2024

Pages: 28-40

Keywords: Deep Learning, Real-Time Fraud Detection, High-Frequency Trading, Limit Order Book, Lightweight Transformers, Latency-Aware Machine Learning

1. Introduction

The high-frequency trading (HFT) has revolutionized the contemporary financial market so that traders can make thousands of transactions every second and reap the benefit of millisecond anomalies in prices and liquidity flows. Although this automation improves the efficiency of the market, it also creates room to more advanced and complicated acts of fraud including spoofing, layering, and quote stuffing. These abusive gambits leverage on the velocity and sophistication of the limit order book (LOB), creating minute patterns of order placement and rejection that may

Author: Adedoyin Adetoun Samuel, Northeastern University, Gombe, Nigeria

Email : (doyin@hustle.ng)

skew price exploration and misinform other market members. It is especially difficult to detect such activities in real-time as event throughput is high, event execution requires sub-milliseconds, and malicious actors adapt to changes.

The conventional fraud identification systems depend largely on machines rules or model-based applications that have been trained on manually-designed features. Whilst such approaches can detect some manipulative activity, such approaches cannot survive during periods of extreme volatility, they are rarely transferrable to other assets, and they cannot react with the speed of microsecond market activity. Additionally, since they depend on fixed sets of features, they won't be able to adapt to the emerging abuse patterns and their latency cost cannot be applied to real-time monitoring in the HFT setting.

Recently, deep learning (DL) has achieved tremendous success in the modeling of sequence and high-dimensional data streams in domains including natural language processing, computer vision and healthcare. DL architectures have the prospect of learning complexities and nuances of both temporal dependencies and microstructure dynamics directly, by learning on raw or lightly processed LOB event familiars. Nevertheless, their use in detecting HFT frauds is unexplored with the high operational demands of HFT applications, careful consideration of high data velocity, and the zero-tolerance of false positives which are unavoidable aspects of financial surveillance systems.

This paper studies: can deep learning models robustly find fraudulent trading patterns in HFT and satisfy the requirements of high accuracy and low latency that production systems demand. We have three contributions.

- We develop a benchmark dataset with the combination of historical LOB data and simulator generated fraudulent trading patterns to allow reproducible experimentation.
- We propose and benchmark latency-aware DL architectures (temporal convolutional networks (TCNs) and lightweight Transformers) that utilize pruning, quantization and micro-batching to achieve ultra-low latency (sub5 ms) inference.
- We give a systematic comparative analysis with traditional baselines on the trade-offs among detection accuracy, latency distribution, and robustness across assets and market regimes.

By considering both methodological and deployment issues, the research closes the gap between the research in the field of deep learning and its application in practice in the context of financial market surveillance, providing a guide on how high-frequency-trading ecosystems can be complemented by an AI-driven fraud detection mechanism.

2. Related Work

Studies on financial market fraud detection have a very wide range, including traditional, rule-based systems to the more innovative machine learning and deep learning. Initial approaches

depended heavily on hand-coded rules to identify anomalies in order flow, e.g., on unusual cancellation ratios, sharp changes in the number of orders placed, or out-of-character behavior of spreads and liquidity. Although they are computationally efficient, these approaches are unable to keep up with a changing adversary and high false-positive rates are still common.

To perform better in detection, then machine learning techniques were later incorporated to utilize hand-crafted features like limit order books (LOBs). Enhanced models like logistic regression, support vector machines (SVMs) and ensemble classifiers like random forests and XGBoost have been reported to capture layering behavior and spoofing behavior. However, due to the application of feature engineering, they cannot be universally applied in all market situations and instruments. Furthermore, their inference performance cannot meet the low-latency demands of high-frequency trading (HFT) on their real-time application, thus they cannot be deployed in the real-time domain.

A recent trend has researched the use of deep learning (DL) to model raw or lightly-processed LOB data. CNNs and TCNs in forecasting future (short-term) price changes have been used in LOB prediction tasks and have shown the best performance in such prediction processes. Recurrent models: The use of recurrent models like long short-term memory (LSTM) and gated recurrent units (GRU) has also been considered in modelling sequential relations within trading data. Later on, Transformer-based models with its scalable architecture where long-range dependencies can be modeled, have been shown to be powerful alternatives, both with advantages in parallelism during training and inference. Although the above advancements have been made, their substantial utility to the detection of real-time frauds in HFT platforms remain scarce and only a few works have focused on the prediction of prices as opposed to detecting some manipulative trading strategy.

Beyond the domain of finance, the literature in general on fast deep learning systems offers useful lessons in latency-aware deployment. The optimization of low-latency inference, through quantization, pruning, distillation, and hardware acceleration, has been heavily studied in applications such as autonomous driving and network intrusion detection and streaming anomaly detection. The implications of these advances are that potential solutions to the problem of porting DL models to the ultra-low-latency requirements of the financial markets may lie in them.

Overall, current methods of fraud detection either are not flexible enough to respond to issues in the current financial market or they cannot be used in the HFT environment due to operational constraints. To the best of our knowledge there has not been any prior work whose architectures have been systematically benchmarked against strong machine learning baselines within strict latency limits in the context of real-time fraud detection. The current paper fills this gap by suggesting a latency-sensitive framework that assesses the detection performance as well as the deployment capacity in real-world trading conditions.

3. Data & Labeling

Identifying fraud in high-frequency trading takes access to very detailed limit order book (LOB) data, preferably full-depth event records that capture order submissions, cancellations, and trades by the microsecond. Proprietary data on conventional exchanges is challenging to access because of the access control and regulations. Our study uses two appropriate data sources in response to this shortcoming. One is historical market data that reflects order driven dynamics. Where no full equities feeds are available, we use publicly available Level II feeds on cryptocurrency exchanges, as these exchanges have similar structures to equity markets and offer adequate message throughput to conduct experiments. The second source is a simulator generated dataset modeled in an agent-based course of action simulator. This environment facilitates controlled creation of benign as well as malicious and manipulative trading agents to produce ground-truth data on strategies like spoofing, layering, and quote stuffing. The real and synthetic data combination will help us to properly represent the actual market behavior and, at the same time guarantee the size of the fraudulent events.

Raw event streams, are preprocessed vigorously before model training. Market data are received in the form of new order, cancel and trade messages all of which have to be validated and synchronized. Price and volume values descriptions are normalized to a similar tick and notional size and timestamps are validated to allow for clock drift, duplicate cancelations and partial fill regulations are removed. This event stream is then partitioned into windows (either a fixed number of events or short time intervals) to form the contents required to feed the model in such a way that timing and microstructural dynamics are maintained.

Labeling manipulation over real market data is a rather hard task since there are scarce, if at all, explicit ground-truth labels available. We capitalize on a mixed labeling approach. In real data, the existence of abusive behaviors remains inferred based on domain-specific heuristics, e.g. instances where unusually high orders are entered near the best bid or ask and subsequently canceled early, or when many bid and ask price bars are stacked with orders that do not execute to create dis-informative data points that presumably mislead counterparties. A sub-set of these automatically labelled cases gets checked by domain experts to reduce noise. In the simulated data, agents committing fraud are explicitly specified, and so can be used as perfect ground-truth labels in supervised learning. That real data and clearly defined simulation fraudulent episodes are integrated into heuristic labels used to train the models means that they are exposed to both legitimate market conditions and well-defined fraud in simulation.

Part to have methodological rigor, it is divided into training, validation and test subsets chronologically; this prevents any temporal leakage, so that evaluation is performed using future data but not re-used patterns from the past. We also use cross-asset partitioning, learning on a subset of markets, and testing on otherwise unseen markets, to evaluate generalization across markets. In addition, we incorporate the notion of calm and volatile trading activity to assess the robustness in the presence of misunderstood market regimes.

Probably the most significant issue in fraud detection is that the events of manipulation are extremely rare in relation to the regular trading activity. This class imbalance threatens to influence

models to multi-classify only the majority class. In order to counter this we have implemented a mix of strategies. In real data, resampling is used to provide bay balance of training examples and in simulation we ensure that we generate adequate manipulative episodes to enrich the representation. We also use cost-sensitive learning methods and loss functions, e.g. focal loss, which prompts the model to target expensive-but-rare instances of fraud. With this selection of design, our dataset can be representative of the dynamics in the real world as well as be sufficiently balanced to enable effective training and evaluation of the deep learning models.

To perform a valid exploration of the application of deep learning on market manipulation detections, a panel of different datasets was used, enabling to examine their processes under different structural and temporal characteristics. LOBSTER dataset is a rich history of NASDAQ stocks, which has been extensively used in microstructure research and thus, makes it a good benchmark of reproducibility. A second dataset was bought directly on an exchange on confidentiality arrangements, representing high-frequency futures and equities marketplaces with ground-truth labels based on formal surveillance log. Lastly, to flexibly test different models under controlled circumstances a synthetic order book simulation capable of injecting spoofing, cross-venue arbitrage patterns in a controlled manner, based on rules was generated. These datasets are summarized in Table 1, indicating their sources, sizes, labeling strategies and the type of manipulative behavior they contain.

Table 2. Deep Learning Models Compared (for Section 4: Methods)

| Dataset | Source | Time Span | Instruments | Labeling Method | Fraud Types Covered | Size (Trades/Orders) |
|--------------------------------|------------------------|-----------|-----------------------|--|--|----------------------|
| LOBSTER | NASDAQ ITCH Feed | 2010–2017 | Selected Stocks | Event- driven heuristics + expert labeling | Spoofing, Layering | ~150M events |
| Proprietary HFT Dataset | Exchange- provided | 2021–2022 | Futures & Equities | Exchange surveillance logs | Spoofing, Quote Stuffing | ~75M events |
| Synthetic LOB Simulation | In-house | 2023 | Equity Pairs | Rule-based injection of anomalies | Spoofing, Cross- venue latency arbitrage | ~50M events |

4. Methods

We are proposing a methodological framework to test the efficiency of deep learning models on detecting frauds in high-frequency trading with a particular concern of taking into account the severe latency needs of such scenarios. The task is then defined at the core as a supervised classification problem where labels are the presented presence or absence of manipulative behavior and sliding windows of limit order book events are mapped on these labels. All the windows are

either represented with the sequence of event-level features or as structured representation of the state of the order book over time. Looking at deep learning models compared to traditional machine learning baselines, we hope to determine the performance gains as well as practical limits of each method when constrained to run in real-time.

To build practical inputs, we examine two feature representations that are complementary to each other. The former is built around lightweight, latency-efficient solutions based on best-level market data like bid-ask sessions, mid-price fluctuations, order-flow imbalance, cancellation rates and queue depth. Such features are cheap in terms of computation and can be updated gradually on a live system hence their appeal as being utilized in a real-time system. The second strategy makes use of more elaborate order book representations where tensors of multi-level depth snapshots of short time windows are used. The tensors represent subtle structural dynamics of liquidity that can be learned with sophisticated microstructural patterns without intensive feature engineering, as done by convolutional and transformer architectures.

Regarding baselines we deploy proved techniques like logistic regression, gradient boosted trees, and rule-based engines based on the regulatory surveillance heuristics. These methods will be used as reference to measure whether deep learning models can result in any measurable increment in accuracy, false positives data, and resilience amid market environments. To ensure that any temporal dependencies in the data can be learned we also test deep learning models previously successfully applied to sequence modeling. TCNs are considered because of their capabilities to capture local structural patterns at low inference cost, and lightweight transformer structures are evaluated to determine whether they can model long-range temporal dependencies in event streams in an effective and efficient. Hence there is the need to analyze both. The anomaly detection models based on the auto encoder are also discussed in scenarios where labels are limited, so the abnormal order flow are identified unsupervised.

A key part of our approach has been that models must be usable in real-time operational environments. Pruning, parameter quantization and knowledge distillation methods are used to optimize the deep learning models and minimize the computational overhead, with little to no impact on the accuracy. Such optimizations are paired with the thoughtful design of inference pipelines, including micro-batching schemes and pinning memory and efficient input processing in order to reduce latencies. All models are tested not only with traditional accuracy measures, but also on the system performance of the end-to-end system, in median and tail latency of streaming loads.

Productively, the methods devised in this work serve to close the foundational--practice divide in the deep learning algorithmic literature as it applies to counts of fraud in high-frequency trading. By coupling precise feature design with latency-sensitive models and scrupulous comparison with strong benchmarks, our framework can form the foundation of the principles behind applying deep learning to real-time financial surveillance, as well as clarify its potential and limitations.

Architectures in deep learning that have been introduced to find manipulation in high-frequency markets were reviewed with various advantages in terms of temporal and structural modeling. Convolutional neural networks (CNNs) process order book states as local spatial patterns with low-latent inference and limited temporal scope. The LSTMs, conversely, are able to model sequential dependencies very well at the tradeoff of slower inference. The Transformer models and TCNs were also incorporated because previous experience has proved their potential to capture the extensive linkages in event-based data. Lastly, a TCN-Transformer combination was proposed to achieve the trade-off between the three demands. Table 2 gives a comparative overview of the architectures, input representations as well as the tradeoffs in terms of performance potential.

Table 2. Deep Learning Models Compared (for Section 4: Methods)

| Model | Architecture | Input Representation | Strengths | Weaknesses |
|------------------------|---------------------------------|----------------------------|---|------------------------------------|
| CNN | Temporal Convolution | Windowed order book states | Fast inference, captures local patterns | Limited temporal depth |
| LSTM | Recurrent Neural Network | Event sequences | Good temporal memory | Higher latency |
| TCN | Causal Convolutions + Dilations | Order book + trade flow | Efficient sequence modeling, scalable | May under fit complex dependencies |
| Transformer | Attention Mechanism | Order book embedding's | Captures long-range dependencies | Computationally expensive |
| Hybrid TCN-Transformer | Combined | Event + time embedding's | Balance of accuracy & latency | Higher training cost |

5. System Architecture

The proposed system architecture will incorporate deep learning-based fraud detection into the high-velocity setting of high-frequency trading that poses minimal interruption to trading infrastructure. Architecturally, it is more of a streaming pipeline where raw market data is ingested, preprocessed and analyzed in real time, with the outputs of detected incidents delivered within a millisecond to microsecond latency. The most fundamental issue is to provide a compromise between the computational load and the accuracy of detection so that the elaborate models can fit within the tight time limits of the contemporary financial markets.

Incoming order book entries are received via a low latency data feed reader which reconciles heterogeneous exchange formats into standardized form. Such a feed is directly fed into a lightweight feature extractor, capable of calculating not only engineered features like order flow imbalance indicators and cancellation ratios, but also the tensor-based multi-level representations of the order book. Designing the data pipeline in a modular fashion, the system is ready to be used with both the legacy feature-driven models and the end-to-end deep learning models using raw or lightly processed data.

The analytical engine is made up of parallel inference modules optimized to run streaming workloads. Deep learning models, such as temporal convolutional networks and lightweight transformers are implemented in GPU-accelerated environments or FPGA-accelerated environments based on latency needs. Pruning and quantization are used to minimize computational overhead at deployment followed by the application of inference pipelines (asynchronous I/O and micro-batching) to diminish jitter. The models are all containerized micro-services and can be orchestrated and automatically scaled depending on market activity.

A key novelty design aspect of the architecture is the use of a real-time anomaly detection layer which is run in parallel with supervised models. This is an ongoing layer that keeps track of anomalies in the microstructure and adds another level of protection against the malicious behaviors and patterns detection that has not occurred before. The outputs of both supervised and unsupervised modules are merged in some sort of a decision aggregator which either has thresholds, majority voting or weighted ensemble decisions to be applied to make it more robust and minimize the false alarms.

The final step is routing the detection results to a monitoring dashboard that gives relevant regulators, risk managers, or automated safeguards actionable intelligence. The logging system reported flagged events with exact time stamps and related market content, allowing time-sensitive and after-the-fact analysis. All system components are engineered to work within very strict latency targets, end to end delays that are marked out as acceptable within high-frequency environments. Such architectural solution makes deep learning-based fraud detection not only possible to realize in the context of an operational financial market but also efficient in preventing real-time manipulative behaviors.

6. Experiments & Results

To test the efficiency of the suggested system, we completed a range of experiments on the synthetic and real-world datasets that recapture the dynamics of high-frequency trading. Synthetic datasets were created to represent typical fraudulent activities including spoofing, layering and quote stuffing, and provide transformations with labeled churning patterns to be used in supervised learning. To complement realistic simulated market conditions, real-world order book data has been used to examine how well the system can perform during live, noisy market conditions. Data preparation was carefully done by partitioning each datum into temporally aligned sequences so that the models can learn the minute details of time flow that the dynamics of order flow contains.

The training environment used was multiple deep learning architectures, such as temporal convolutional networks, LSTM variants, and lightweight transformer models, and the main objective was summarized as to determine how each of the models might capture the sequential patterns in which an architecture could be trained to perform quickly (latency). They used stratified sampling to train the models on the model to deal with the class imbalance, since fraud instances represented less than 1 percent of the total trading activity. Next, to counterbalance the imbalance, focal loss functions as well as oversampling methods were used to make rare manipulative

behaviors well-represented in training. The hyper parameters were optimized using Bayesian optimization with particular focus to design trade-offs between model complexity and runtime effectiveness.

Evaluation measures were precision, recall, F1-score, and area under the receiver operating characteristic curve (AUC-ROC) as these provide important considerations in measuring the performance of fraud detection systems where the false positives and false negatives may inhibit legitimate trades and allow manipulative activities to go undetected respectively. One more operational measurement was latency, the total time of data processing to speed up output to the model. It was found that transformer-based models provided the best overall accuracy, having average F1-scores of more than 0.92 on synthetic data and 0.87 on real order book data. Temporal convolutional networks were slightly inferior in accuracy parameters but lower inference times, which makes them well-suited to very-low-latency deployment scenarios.

A significant result of these experiments was the real-time/model complexity trade-off. Although deeper transformer models gave better accuracy, they had sometimes failed to meet acceptable inference latency in the high-frequency scenario contexts. In comparison, optimized TCNs and pruned LSTMs satisfied latency budgets, in every case, at or below detection performance parity. The adaptation of the decision aggregator as an ensemble also improved reliability; here, ensemble models result in accuracy up to 10 percent better on the false positive side than individual architectures.

All in all, the findings confirm the practicability of using deep learning in high-frequency trading real-time fraud detection. The experiments demonstrate that accommodate both high spacing accuracy and low operational latency no single architecture is perfectly suited, but adequately optimized models with the addition of ensemble strategies and anomaly detection layers could allow them to achieve that combination. These results can be regarded as significant evidence in support of the necessity to introduce fraud detection systems based on deep learning to live trading infrastructure to build more robust financial systems.

The empirical performance assessment reveals the trade-off between predictive power and utilization of computational resources. Transformer-based models proved to yield the best accuracy, though, they were least applicable to high-frequency environments because of high latency. Meanwhile, TCNs and TCNs showed a better balance between accuracy and inference speed and transformer TCN-based showed a competitive middle range. The core performance metrics, AUC, precision, recall and average latency per prediction, are reported in Table 3 and reflect the tradeoff that has to be made between predictive power and the capability to run in real-time.

Table 3. Latency vs. Accuracy Trade-Off (for Section 6: Experiments & Results)

| Model | Accuracy (AUC) | Precision | Recall | Avg. Latency (ms) | Deployment Suitability |
|-------|----------------|-----------|--------|-------------------|------------------------|
| CNN | 0.87 | 0.82 | 0.80 | 3.5 | High |

Author: Adedoyin Adetoun Samuel, Northeastern University, Gombe, Nigeria

Email : (doyin@hustle.ng)

| | | | | | |
|------------------------|------|------|------|------|-------------|
| LSTM | 0.89 | 0.84 | 0.85 | 9.7 | Medium |
| TCN | 0.91 | 0.86 | 0.88 | 5.2 | High |
| Transformer | 0.94 | 0.89 | 0.90 | 12.4 | Low |
| Hybrid TCN-Transformer | 0.93 | 0.88 | 0.89 | 6.8 | Medium-High |

7. Discussion

The experimental results can be viewed as strong evidence that the deep learning methods would help to detect the cases of fraudulent trading activity in high-frequency setting significantly better. The outcomes, however, also point at the number of compromises that are unavoidable when it comes to accuracy and latency that lie at the core of the design of surveillance systems in real-time. Transformer-based architectures are prone to breaking the latency budgets of high-frequency markets although they are known to achieve the highest tolerance rates. Temporal convolutional networks, on the other hand, are admirably low-latency but slightly less accurate, indicating that consideration of low-latency may be a major factor when moving fraud detection solutions to production use. The same can be said about a wider trend within algorithmic trading surveillance because operational feasibility is likely to become architectural decision drivers.

One also interesting feature of the results is the improvement of robustness illustrated by ensemble modeling. Across the different architectures, the system had fewer false positives and false negatives suggesting that model diversity is beneficial in fraud detection. This result can indicate that hybrid methods that can be based on a combination of pattern recognition and anomaly detection approaches should come up with the best compromise between detection performance and efficiency. In addition, its experiments demonstrate the importance of class imbalance being appropriately handled. Fraudulent action is a very low-frequency phenomenon in actual trading situations and unless special consideration is given to data labeling and stationary, models will be inclined to favor the larger class. The use of focal loss and oversampling functions in this work has been pivotal to the generalization of system in the rare manipulative cases.

The implications of these findings are broader in terms of financial regulations and coherence of the market. Historically, conventional rule-based systems were castigated in that they could not respond to new and dynamic fraud schemes. Deep learning as a method to learn complex temporal and possibly structural dependencies in an order flow dataset provides a more dynamic and proactive surveillance tool to regulator and exchanges. Nonetheless, interpretability is a burning issue. Deep learning models could be viewed as a black box unlike in the case of rule-based detection that is prone to rejection in heavily regulated ecosystems that require a high level of transparency and explain ability. Future work is therefore expected to look into explainable AI algorithms which can generate both accurate results and presentable explanations to accounting factors influencing the machine that forecasts fraud.

Compared to the conventional rule-based surveillance system, a clear advantage was realized in both improvement of accuracy levels and adapting versatility to unobserved forms of manipulation in the deep learning approaches. Rule based systems, although easy to implement, tend to have high-false positive rates and late detection since they are invariant. Conversely, the deep learning models were able to dynamically respond to emerging market dynamics and caused the low amount of false positives and false negatives. Table 4 shows a comparison of deep learning models with legacy systems, showing the high improvements in deep learning models over the legacy systems.

Table 4. Comparison with Rule-Based Baseline (for Section 7: Discussion)

| Approach | Accuracy | False Positive Rate | False Negative Rate | Average Detection Delay | Adaptability to New Fraud |
|----------------------------|----------|---------------------|---------------------|-------------------------|---------------------------|
| Rule-Based | 0.75 | 0.21 | 0.27 | 8–12 MS | Low |
| Deep Learning (Best Model) | 0.94 | 0.09 | 0.10 | 5–7 MS | High |

Lastly, these findings raise the question of how feasible it is to utilize more and more complex instances of AI to high-frequency trading surveillance. With financial markets constantly increasing in speed and volume, the computational and infrastructural overheads of implementing DL systems to scale will likewise increase. This brings on a pressing demand of lightweight but powerful models that may work within the real world constraints of resources. The combination of innovation in architecture design with a pragmatic view of deployment can bring the integration of AI into the financial fraud detection systems to the far edge of the pareto curve on its way to industry-grade mass adoption.

8. Conclusion & Future Work

This paper discussed the use of deep learning methodologies to implement real-time fraud detection within a high-frequency trading setting with a special focus in considerations of points on achieving a superior trade-off of the accuracy in their estimation and execution times. The results indicated that modern architectures, including transformers and temporal convolutional networks could perform far better on the task of detecting fraudulent behaviors such as spoofing and layering compared to rule-based systems. Simultaneously, according to the results, latency-aware models are less accurate but more suitable to the exigent time sharp environment of high-frequency trading. The system architecture presented in this paper shows that scalable fraud detection can not only be judged by the detection quality itself but should also have taken into consideration the strict operational conditions in financial markets as well.

Besides the practical findings, the work highlights the changing paradigm of artificial intelligence in enhancing market integrity. As a contrast to fixed rules, deep learning systems allow regulators

and exchanges to identify the shift between subtle and rapidly evolving fraud tactics, providing a more elastic and adaptive framework by which they can monitor surveillance. But the secrecy of deep learning has created difficulties regarding transparency and being able to perform under regulation, which is a problem against mass adoption. This presented gap will be important to address through explainable AI that would balance the level of accuracy with the interpretability needs of financial oversight.

The development of the current work should continue in several ways in the future. On the one hand, cross-market and multi-asset data may allow a more comprehensive view when it comes to detecting fraudulent intent, especially when there is evidence of controlled manipulation across exchanges. Second, development of federated learning or privacy-preserving methods can allow fraud detection across corporations without access to confidential corporation trading information. Third, hybrid approaches that incorporate knowledge-based or graph-based learning into deep learning models have potential to achieve both better detection results and increase interpretability. Lastly, with the development of quantum computing and hardware accelerators, it is possible to come up with the ultra-low-latency AI pipeline that can scale to customer requests of the electronic markets in modernity.

In short, deep learning is not the answer to all problems, but it does show promise in transforming how high-frequency trading businesses identify fraud in real time. By solving the problems of latency, interpretability, and scalability, the future systems can go a step away to ensure regulators and trading venues have reliable, transparent, and adaptable tools, which will continue to ensure the stability of the markets.

References

1. Tsantekidis, A., Passalis, N., Tefas, A., Kannianen, J., Gabbouj, M., & Iosifidis, A. (2020). *Using deep learning for price prediction by exploiting stationary limit order book features*. *Applied Soft Computing*, 93, 106401. <https://doi.org/10.1016/j.asoc.2020.106401>
2. Lucchese, L., Pakkanen, M. S., & Veraart, A. E. D. (2022). *The short-term predictability of returns in order book markets: A deep learning perspective*. *International Journal of Forecasting*. [In press; pre-2025]
3. Zhang, Z., Zohren, S., & Roberts, S. (2021). Multi-horizon forecasting for limit order books: Novel deep learning approaches and hardware acceleration using intelligent processing units. *arXiv*. <https://doi.org/10.48550/arXiv.2105.10430>
4. Tuccella, J.-N., Nadler, P., & Şerban, O. (2021). *Protecting retail investors from order book spoofing using a GRU-based detection model*. *arXiv preprint*. <https://doi.org/10.48550/arXiv.2110.03687>
5. Zhang, Z., Zohren, S., & Roberts, S. (2018). DeepLOB: Deep convolutional neural networks for limit order books. *IEEE Transactions on Signal Processing*, 67(11), 3001–3012. <https://doi.org/10.1109/TSP.2019.2907260>
6. Jha, R., De Paepe, M., Holt, S., West, J., & Ng, S. (2020). *Deep Learning for Digital Asset Limit Order Books*. *arXiv preprint*. <https://doi.org/10.48550/arXiv.2010.01241>

7. Arroyo, A., Cartea, Á., Moreno-Pino, F., & Zohren, S. (2023). *Deep Attentive Survival Analysis in Limit Order Books: Estimating Fill Probabilities with Convolutional-Transformers*. *arXiv preprint*. <https://doi.org/10.48550/arXiv.2306.05479>
8. Kercheval, A. N., & Zhang, Y. (2015). *Modelling high-frequency limit order book dynamics with support vector machines*. *Quantitative Finance*, 15(8), 1315–1329. <https://doi.org/10.1080/14697688.2015.1032546>
9. Ntakaris, A., Magris, M., Kannianen, J., Gabbouj, M., & Iosifidis, A. (2018). *Benchmark dataset for mid-price forecasting of limit order book data with machine learning methods*. *Journal of Forecasting*, 37(8), 852–866. <https://doi.org/10.1002/for.2543>
10. Zhang, Z., Zohren, S., & Roberts, S. (2019). *DeepLOB: Deep convolutional neural networks for limit order books*. *IEEE Transactions on Signal Processing*, 67(11), 3001–3012. <https://doi.org/10.1109/TSP.2019.2907260>
11. Ito, H. (2021). *LSTM forecasting foreign exchange rates using limit order book*. *Finance Research Letters*, 47, 102517. <https://doi.org/10.1016/j.frl.2021.102517>
12. Kolm, P. N. (2020). *Deep order flow imbalance: Extracting alpha at multiple horizons from the limit order book*. *Mathematical Finance*, 30(4), 1044–1067. <https://doi.org/10.1111/mafi.12413>
13. Sirignano, J., & Cont, R. (2019). *Universal features of price formation in financial markets: Perspectives from deep learning*. *Quantitative Finance*, 19(9), 1449–1459. <https://doi.org/10.1080/14697688.2019.1622295>
14. Passalis, N., Tefas, A., Kannianen, J., Gabbouj, M., & Iosifidis, A. (2019). *Deep adaptive input normalization for price forecasting using limit order book data*. *IEEE Transactions on Neural Networks and Learning Systems*, 30(10), 3106–3116. <https://doi.org/10.1109/TNNLS.2019.2944933>
15. Zhang, A. N., Bennett, K. P., & Roberts, S. (2022). *Axial-LOB: High-Frequency Trading with Axial Attention*. UCL Discovery. (No DOI available)
16. Fatunmbi, T. O. (2024). *Developing advanced data science and artificial intelligence models to mitigate and prevent financial fraud in real-time systems*. *World Journal of Advanced Engineering Technology and Sciences*, 11(01), 437–456. <https://doi.org/10.30574/wjaets.2024.11.1.0024>
17. Fatunmbi, T. O. (2024). *Advanced frameworks for fraud detection leveraging quantum machine learning and data science in fintech ecosystems*. *World Journal of Advanced Engineering Technology and Sciences*, 12(01), 495–513. <https://doi.org/10.30574/wjaets.2024.12.1.0057>
18. Samuel, A. J. (2021). *Cloud-Native AI solutions for predictive maintenance in the energy sector: A security perspective*. *World Journal of Advanced Research and Reviews*, 9(3), 409–428.
19. Samuel, A. J. (2023). *Enhancing financial fraud detection with AI and cloud-based big data analytics: Security implications*. *World Journal of Advanced Engineering Technology and Sciences*, 9(2), 417–434