

AI-INTEGRATED FPGA FOR MARKET MAKING IN VOLATILE ENVIRONMENTS

by

Shreejit Verma

A THESIS

Submitted to the Faculty of the Stevens Institute of Technology
in partial fulfillment of the requirements for the degree of

MASTER OF SCIENCE - FINANCIAL ENGINEERING

Shreejit Verma, Candidate

ADVISORY COMMITTEE

Dragos Bozdog, Chairperson

Date

Ionut Florescu, Reader

Date

STEVENS INSTITUTE OF TECHNOLOGY
Castle Point on Hudson
Hoboken, NJ 07030
2026

©2026, Shreejit Verma. All rights reserved.

AI-INTEGRATED FPGA FOR MARKET MAKING IN VOLATILE ENVIRONMENTS

ABSTRACT

This Master's thesis investigates the integration of artificial intelligence (AI) with field-programmable gate arrays (FPGAs) to develop adaptive market-making strategies tailored for volatile financial environments in a high-frequency trading (HFT) environment. Traditional market-making models often struggle with latency and adaptability during periods of high volatility, leading to increased inventory risk and suboptimal performance. We propose a hybrid system that embeds reinforcement learning (RL) algorithms for dynamic bid-ask spread optimization directly onto FPGA hardware, enabling sub-microsecond inference times while responding to real-time market signals such as order flow imbalances and volatility clusters. Using historical tick-level data from major exchanges and simulated volatile scenarios, the framework is evaluated through metrics including Sharpe ratio, adverse selection mitigation, and latency benchmarks. The expectation is to achieve significant improvements in profitability and risk management compared to software-based RL or standalone FPGA implementations. This work will contribute to a scalable AI-FPGA architecture, highlighting practical deployment challenges, and offering insights for future quant engineering in dynamic markets, with potential applications in NYC's competitive HFT landscape.

Keywords: High-Frequency Trading, FPGA Acceleration, Reinforcement Learning, Market Making, Low Latency

Author: Shreejit Verma

Advisor: Dragos Bozdog

Date: November 2025

Program: Financial Engineering

Degree: MASTER OF SCIENCE - FINANCIAL ENGINEERING

Acknowledgments

To be completed

Table of Contents

Abstract	iii
Acknowledgments	v
List of Tables	ix
List of Figures	x
Chapter 1 Introduction and System Overview	1
Chapter 2 Literature Review	16
2.1 Introduction	16
2.2 Theoretical Background	17
2.3 Overview of Existing Research	18
2.3.1 From Classical Market Making to Learning-Based Control	18
2.3.2 Hardware Acceleration and FPGA Design in HFT	19
2.3.3 AI on FPGAs and Edge Inference	19
2.3.4 Hybrid Systems Integrating Learning and Hardware	20
2.4 Critical Evaluation	20
2.5 Identification of Gaps	21
2.6 Relevance to the Present Research	21
2.7 Conceptual Framework	22

2.8	Summary and Transition	22
Chapter 3	High-Level Architecture of an FPGA Matching Engine	24
3.1	Description of Key Modules	24
Chapter 4	FPGA System Architecture and Methodology	27
4.1	The Ultra-Low-Latency T2T Data Path (FPGA)	27
4.2	Core Integration: The RL-Inference Module	28
4.3	The Hybrid Control Plane for Model Management	29
Chapter 5	Ultra-Low-Latency Software Baseline Implementation	31
5.1	Target Performance Metrics	31
5.2	Software System Architecture	31
5.3	Core Infrastructure Components	32
5.4	Network and Market Data Processing	33
5.5	Strategy and Risk Management	33
5.6	Performance Tuning and Benchmarking	34
Chapter A	Appendix	1
A.1	Setup Guide: Verilog on Apple Silicon	1
A.2	FPGA Code Explanation and Module Overview	2
A.2.1	Network Ingress Modules	3
A.2.2	Feed Handler Modules	3
A.2.3	Strategy and Decision Logic	3
A.2.4	Order Transmission Modules	3
A.2.5	Control Plane and Integration	4
A.3	Appendix: Complete Verilog and VHDL Code Listings	4
A.3.1	1. Network Ingress: Parser	4

A.3.2	2. Network Ingress: Multicast Gate	6
A.3.3	3. Feed Handler: Frame Extractor	7
A.3.4	4. Feed Handler: Order Book Update	8
A.3.5	5. Strategy Logic: Decision Block	10
A.3.6	6. Hardware Pre-Trade Risk Gate	10
A.3.7	7. Order Encoder	12
A.3.8	8. Order Egress: TX Bridge	13
A.3.9	9. Control Plane: AXI-Lite Registers	14
A.3.10	10. Top-Level Integration	15
A.4	Appendix: C++ ULL Software Baseline Code Listings	18
	Bibliography	43

List of Tables

List of Figures

1.1	High Level Design - Part 1	3
1.2	Order Book Management and Event-Driven Propagation	4
1.3	Event-Driven Pipeline and Nanosecond-Precision Timing	6
1.4	FPGA-Accelerated Tick-to-Trade Pipeline	7
1.5	Internal Architecture of the FPGA Acceleration Module	8
1.6	Order Processing and Post-Trade Analytics	10
1.7	System-Wide Monitoring and Metrics Infrastructure	12
1.8	Unified High-Frequency Trading System Architecture	13
3.1	High-Level Data Flow in an FPGA Matching Engine	24

Chapter 1

Introduction and System Overview

Market making is central to modern financial markets, as liquidity providers continuously quote bid and ask prices to facilitate trading while managing inventory risk and seeking profit from the spread. The advent of high-frequency trading (HFT) has transformed this role, demanding ultra-low-latency decision-making under dynamic and volatile conditions. Traditional stochastic frameworks, such as the Avellaneda–Stoikov model, offer valuable theoretical foundations but often falter in environments characterized by sudden price swings, order book imbalances, and liquidity shortages. Their reliance on fixed assumptions limits adaptability in real-time, high-volatility settings [14, 20]. This has prompted significant interest in more flexible approaches capable of balancing adverse selection, inventory control, and profitability at microsecond scales.

Recent advances in artificial intelligence (AI), particularly reinforcement learning (RL), have provided promising alternatives. RL agents learn adaptive quoting strategies through iterative interactions with either simulated markets or historical order book data, enabling dynamic policies that outperform rule-based methods [19, 7]. Deep RL models, including recurrent architectures, have demonstrated enhanced performance in inventory management and order placement, particularly in stressed environments [22, 5]. Moreover, multi-objective RL frameworks have applied Pareto optimization to manage trade-offs between latency, volatility resilience, and slippage reduction [10]. Despite these advances, most implementations remain software-based, and the computational overhead of deep learning methods introduces latency that undermines their applicability in production-grade HFT environments.

Parallel to developments in AI, field-programmable gate arrays (FPGAs) have emerged as critical enablers of ultra-low-latency trading. Unlike CPUs and GPUs, FPGAs can process market data feeds, order matching, and risk checks at nanosecond scales, leveraging hardware-level parallelism [4?]. Studies highlight their advantages in power efficiency, deterministic execution, and real-time data handling, which are essential for HFT infrastructures [11, 6]. Applications range from pre-built IP libraries for networking and protocol parsing [11] to FPGA-based system-on-chip frameworks for algorithmic execution [1]. Furthermore, integration of FPGAs with technologies like RDMA has enabled near-zero-copy communication pipelines, further reducing latency across trading networks [6]. Nevertheless, traditional FPGA deployments have typically been limited to deterministic, pre-specified functions, lacking the adaptability required for volatile and evolving market conditions.

The convergence of AI and FPGA technology offers a promising path forward. Recent studies have explored FPGA-accelerated AI inference in trading contexts, enabling models such as XGBoost or neural networks to run at sub-microsecond scales [? 21]. Dynamic FPGA reconfiguration has been proposed as a means to accommodate evolving RL or deep learning models in real time [2]. Early prototypes of AI-augmented FPGA trading systems suggest significant reductions in end-to-end latencies for market-making strategies [9, 8]. Yet, gaps remain: few studies systematically evaluate AI–FPGA integration under extreme volatility, flash-crash conditions, or across multi-asset contexts. Moreover, while industry reports emphasize the growing adoption of FPGA–AI platforms for finance [? 17], rigorous academic investigations into resilience, scalability, and security trade-offs remain limited.

This thesis aims to bridge these gaps by developing an integrated reinforcement learning–FPGA framework for market making under high-volatility scenarios. By combining RL’s adaptive decision-making capabilities with FPGA’s hardware-level

acceleration, the proposed system seeks to reduce latency bottlenecks, improve inventory risk management, and enhance robustness in stressed market conditions. In doing so, it builds upon the strengths of prior RL and FPGA research while addressing limitations in real-world deployment contexts.

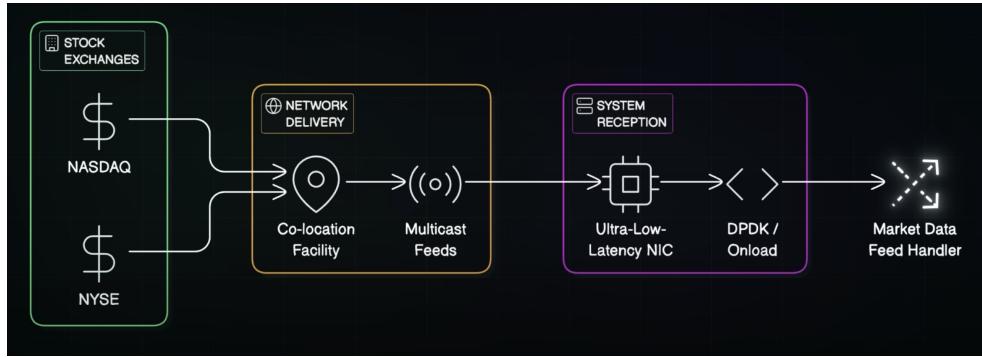


Figure 1.1: High Level Design - Part 1

The first stage of the proposed system begins with the **ingestion of raw market data from stock exchanges such as NASDAQ and NYSE**. These exchanges publish continuous streams of tick-level information, including order book updates and trade events, using **multicast feeds**. To minimize latency, trading firms typically deploy their infrastructure within **co-location facilities** physically situated near the exchange servers, ensuring that market updates traverse the shortest possible physical distance before reaching the system.

Within this co-located environment, incoming data is captured through an **ultra-low-latency network interface card (NIC)**. Unlike conventional NICs, these specialized devices are optimized for deterministic packet capture at microsecond and even nanosecond granularity. To further reduce overhead, the system employs **kernel-bypass mechanisms** such as DPDK or Solarflare Onload, allowing direct user-space access to packet streams and eliminating delays introduced by the operating system's networking stack.

The processed feed is then passed to the **market data feed handler**, which performs protocol decoding, normalization, and transformation into an internal format suitable for downstream components. This module acts as the critical bridge between raw exchange data and the trading logic of the system, ensuring that millions of messages per second can be ingested and translated without loss.

This stage, illustrated in Figure 1.1, provides the **foundational data layer for the AI-integrated FPGA framework**. By guaranteeing ultra-low-latency and reliable data delivery, it enables the reinforcement learning agents and FPGA-accelerated decision modules in later stages of the architecture to operate on timely and accurate market signals—an essential requirement for market making in volatile, high-frequency environments.

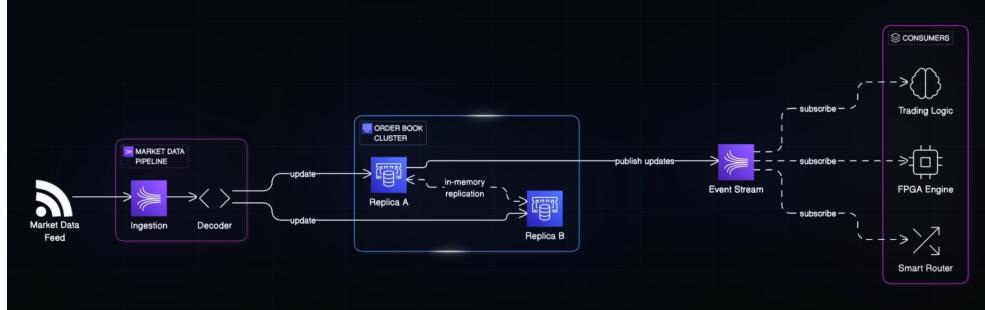


Figure 1.2: Order Book Management and Event-Driven Propagation

Following the initial ingestion and decoding phase, Figure 1.2 illustrates the core processing architecture responsible for maintaining the market state and disseminating updates to decision-making components. This event-driven design is critical for achieving nanosecond-level determinism. The decoded data from the **Market Data Pipeline** is first used to update the **Order Book Cluster**. To eliminate disk I/O latency and ensure high availability, the system maintains a complete, live snapshot of the order book entirely in-memory. The architecture employs a fault-tolerant design, featuring two synchronized instances, **Replica A** and **Replica B**, which are

maintained through continuous **in-memory replication**. Should one instance fail, the system can seamlessly failover to the other without interruption. Upon each update to the order book, a state change event is published to a lock-free, multi-consumer **Event Stream**. This stream serves as the central backbone of the system, broadcasting timestamped market events to all downstream modules. This publish-subscribe model decouples the order book from the logic engines, allowing for parallel, independent processing. The primary **Consumers** of this event stream are:

- **Trading Logic:** A software-based strategy engine that subscribes to the stream to evaluate market conditions, manage inventory risk, and execute algorithmic strategies. This component allows for complex, nuanced decision-making that may be difficult to implement directly in hardware.
- **FPGA Engine:** The core of the proposed system. This hardware component also subscribes directly to the event stream, enabling "tick-to-trade" execution. The embedded reinforcement learning model on the FPGA can react to market events in sub-microsecond timeframes, bypassing the overhead associated with the CPU and operating system entirely.
- **Smart Router:** This module consumes market data to make optimal routing decisions, determining the best venue and method for order execution based on factors like liquidity, fees, and latency.

This architecture ensures that the AI-driven FPGA engine, alongside other critical components, receives a synchronized and near-instantaneous view of the market, which is essential for the efficacy of any high-frequency market-making strategy.

Figure 1.3 details the architecture's event-driven core, which is responsible for propagating market state changes with deterministic, nanosecond-level precision.

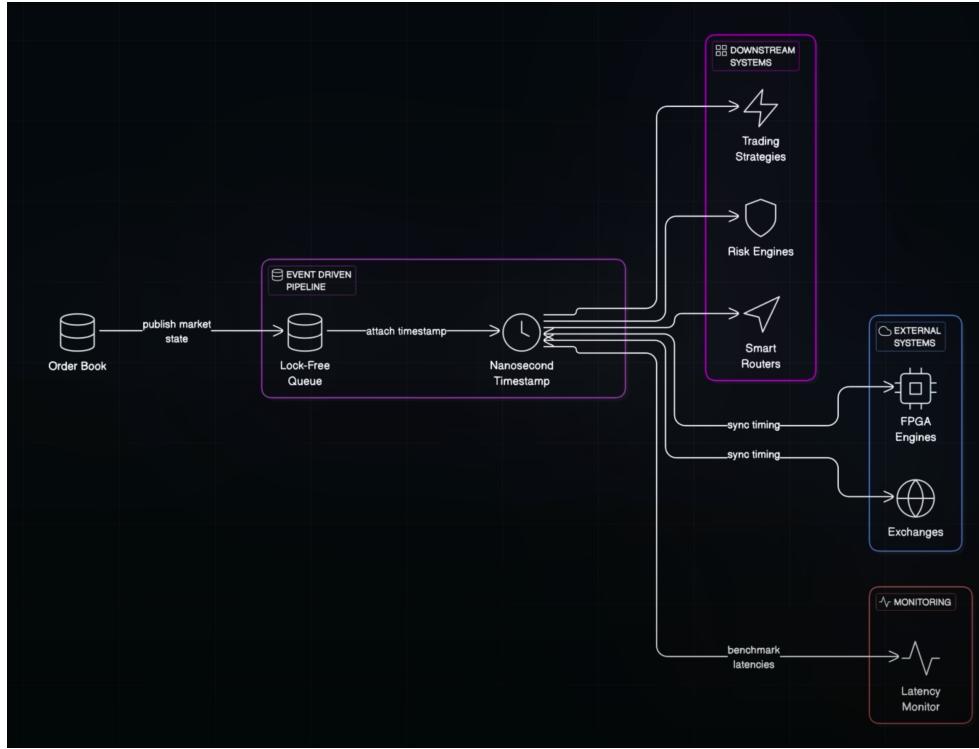


Figure 1.3: Event-Driven Pipeline and Nanosecond-Precision Timing

This pipeline is the central nervous system of the trading platform, ensuring all components operate on a synchronized and chronologically exact sequence of market events. The process begins when an update to the in-memory **Order Book** is published into the **Event Driven Pipeline**. To manage high-throughput, concurrent data access without introducing latency from thread contention, the event is first placed into a **Lock-Free Queue**. As the event is dequeued, it is immediately stamped with a **Nanosecond Timestamp**. This high-resolution timestamp is fundamentally important for several reasons: it establishes an indisputable sequence of events, enables precise latency benchmarking across different system components, and provides a synchronization clock for time-sensitive external systems. The timestamped event is then multicast to a variety of consumers:

- **Downstream Systems:** These software-based components consume the event

stream to perform their respective functions. This includes the **Trading Strategies** engine, the pre-trade **Risk Engines** that enforce safety checks, and the **Smart Routers** that determine optimal execution venues.

- **External Systems:** These systems require precise time synchronization to function correctly. The **FPGA Engines**, which execute the hardware-accelerated RL models, rely on this timing to align their actions perfectly with the market data tick they are processing. Likewise, order messages sent to the **Exchanges** must be correctly sequenced and timestamped for compliance and clearing.
- **Monitoring:** A dedicated **Latency Monitor** subscribes to the event stream to benchmark the performance of the entire "tick-to-trade" pipeline. By comparing timestamps at various stages, the system can be continuously optimized to eliminate bottlenecks.

This architecture guarantees that the AI-integrated FPGA, along with all other decision-making modules, operates on a coherent and precisely timed representation of the market, which is a non-negotiable requirement for competitive high-frequency trading.



Figure 1.4: FPGA-Accelerated Tick-to-Trade Pipeline

Figure 1.4 illustrates the most latency-sensitive segment of the architecture: the hardware-accelerated High-Frequency Trading (HFT) pipeline. This diagram demonstrates the end-to-end flow from market data reception to order execution, with the

Field-Programmable Gate Array (FPGA) acting as the primary decision-making engine. The process initiates with the **Market Data Feed**, which is processed by the **Feed Handler**. The resulting normalized data is then timestamped and placed into a lock-free **Event Queue**. This queue streams **direct tick events** to the FPGA, ensuring the hardware receives market data with minimal jitter and the lowest possible latency. The central component is the **FPGA Acceleration** module. Within this module, the custom **FPGA Logic**, which contains the synthesized reinforcement learning (RL) model, receives the tick event. At hardware speed, without the overhead of a CPU or operating system, the logic evaluates the market state and decides on the optimal quoting strategy based on its learned policy. This decision is passed to the hardware **Execution Engine**, which formulates the corresponding order message. The entire process, from receiving the tick to generating a response, occurs in sub-microsecond timeframes. The resulting **sub-microsecond orders** are then forwarded to the **Order Router**. Before being sent to the exchange, these orders would pass through the pre-trade risk checks (as detailed in Figure 1.3) to ensure compliance and safety. Finally, the order is dispatched to the **Exchange**. This "tick-to-trade" pathway represents the system's critical advantage, leveraging the parallelism and deterministic, low-latency nature of FPGAs to react to market opportunities faster than any software-based equivalent could.

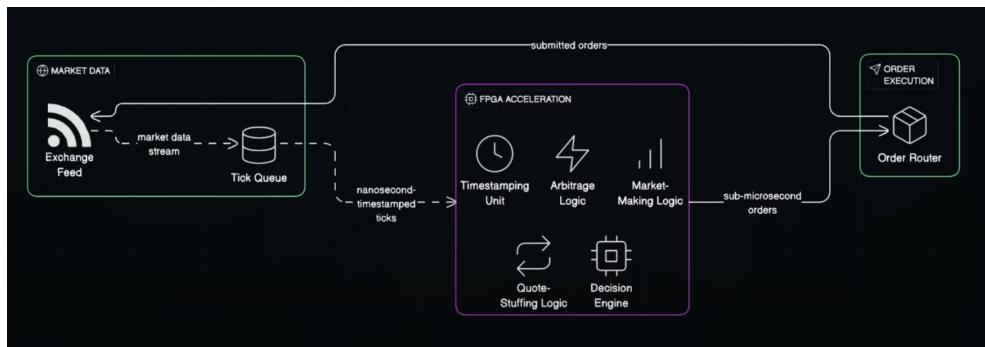


Figure 1.5: Internal Architecture of the FPGA Acceleration Module

Figure 1.5 provides a granular view of the internal architecture of the **FPGA Acceleration** module, detailing the parallel hardware logic blocks responsible for strategy execution. The module operates on two primary data streams. The first is the **Market Data** stream, where the **Exchange Feed** is buffered into a **Tick Queue**, providing a continuous flow of **nanosecond-timestamped ticks**. The second is a critical feedback loop of **submitted orders** from the **Order Router**. This feedback is essential for state-aware strategies, allowing the FPGA to manage its inventory and outstanding orders in real-time. Within the FPGA, several specialized logic blocks operate in parallel:

- **Timestamping Unit:** An internal unit for precise latency measurement and ensuring synchronization of all internal processes relative to the incoming market data.
- **Strategy Logic Blocks:** The FPGA is programmed with multiple, concurrent trading strategies, each implemented as a distinct hardware circuit. The diagram shows examples such as **Arbitrage Logic** and **Quote-Stuffing Logic**. Critically, the **Market-Making Logic** block is where the proposed adaptive reinforcement learning (RL) model is synthesized. This block is responsible for dynamically calculating optimal bid-ask spreads based on the learned policy.
- **Decision Engine:** This is the central processing core of the FPGA. It integrates the signals and outputs from all parallel strategy blocks, considers the current state of submitted orders, and makes the final, unified trading decision. This engine effectively performs the inference step of the embedded RL model.

The output of the Decision Engine is a stream of **sub-microsecond orders**, which are sent to the **Order Router** for execution. This modular, parallel hardware design

enables the system to evaluate multiple complex market conditions simultaneously and react with deterministic, ultra-low latency, achieving performance unattainable by conventional software-based systems.

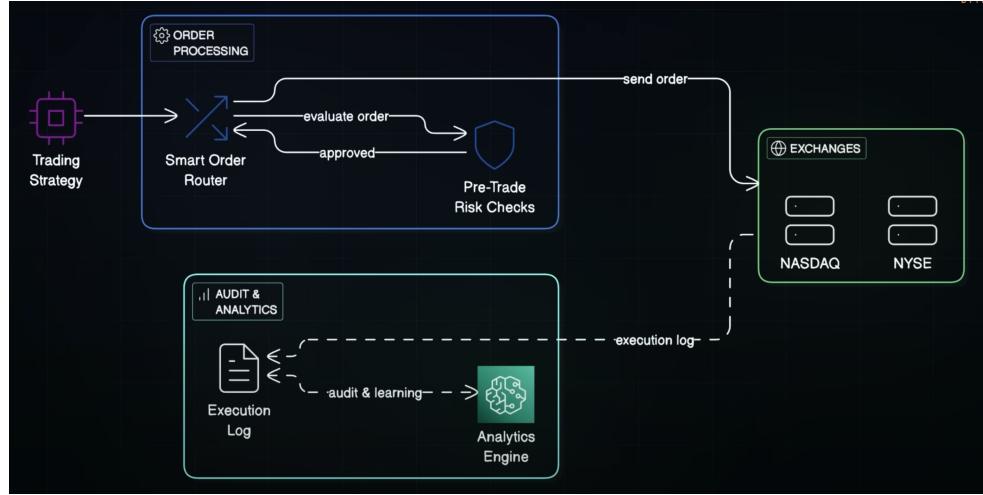


Figure 1.6: Order Processing and Post-Trade Analytics

Figure 1.6 illustrates the final stages of the trade lifecycle: order processing, risk management, and post-trade analysis. These components ensure that trading is executed safely and provide a critical feedback loop for continuous strategy improvement. The process begins when a **Trading Strategy** (originating from either the FPGA or a software-based engine) generates a desire to trade. This signal is sent to the **Order Processing** module.

- **Smart Order Router (SOR):** The first component within this module is the SOR. It receives the trade signal and determines the optimal venue and method for execution based on factors like liquidity, latency, and exchange fee structures.
- **Pre-Trade Risk Checks:** Before an order is sent to an exchange, it is evaluated by a series of mandatory **Pre-Trade Risk Checks**. This critical safety

layer validates the order against predefined limits, such as maximum order size, position limits, and rate checks, to prevent erroneous trades that could lead to significant financial loss. Once the order is approved, it is dispatched to the selected exchange (e.g., NASDAQ, NYSE).

After an order is executed at an exchange, an **execution log** is generated and sent to the **Audit & Analytics** module.

- **Execution Log:** This component is an immutable, timestamped record of all trading activity, including fills, partial fills, and rejections. It serves as the primary source for compliance reporting and post-trade analysis.
- **Analytics Engine:** The logs are consumed by the **Analytics Engine**. This engine is responsible for audit and learning. For the proposed system, this engine would analyze the performance of the RL-based market-making strategy, providing data on profitability, adverse selection, and inventory risk. The insights derived here are crucial for retraining and refining the AI models, thus closing the loop and enabling continuous, data-driven strategy optimization.

This complete feedback architecture, combining ultra-low-latency execution with robust risk management and intelligent analytics, forms a comprehensive framework for deploying adaptive, high-performance trading strategies in volatile market environments.

Figure 1.7 provides a detailed overview of the system's comprehensive monitoring and metrics infrastructure. This architecture operates in parallel to the main trading pipeline and is critical for ensuring operational stability, performance tuning, and regulatory compliance. The core of this infrastructure is the **Order Management System (OMS)**, which acts as the central nervous system for all trade-related information. The OMS tracks critical data points for every order, including the **Routes**

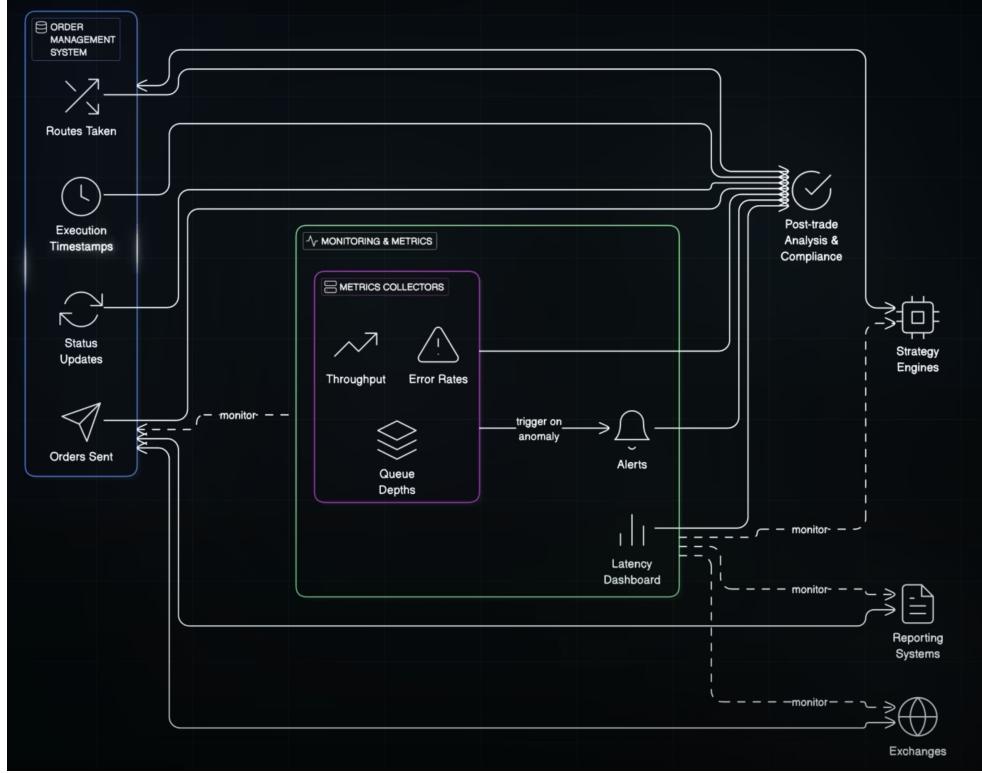


Figure 1.7: System-Wide Monitoring and Metrics Infrastructure

Taken, precise **Execution Timestamps**, real-time **Status Updates** (e.g., filled, partially filled, rejected), and the initial **Orders Sent**. This wealth of data from the OMS feeds into two primary downstream processes:

1. **Monitoring & Metrics:** A dedicated module continuously monitors the system's health and performance.
 - **Metrics Collectors:** These components actively gather key performance indicators (KPIs) in real-time, such as system **Throughput**, **Error Rates**, and the depth of various message queues (**Queue Depths**).
 - **Alerts:** If any of the collected metrics breach predefined thresholds, indicating a potential anomaly (e.g., a sudden drop in throughput or a spike in errors), the system automatically triggers **Alerts** to notify system op-

erators.

- **Latency Dashboard:** This component provides a real-time visualization of critical latency measurements, such as the "tick-to-trade" time, allowing for immediate identification of performance bottlenecks.

2. **Post-Trade Analysis & Compliance:** The data from the OMS is also funneled into this module, which is responsible for regulatory reporting and strategy performance analysis. The collected metrics and logs are used to monitor the performance of **Strategy Engines**, **Reporting Systems**, and the interactions with **Exchanges**.

This robust monitoring framework ensures that every aspect of the trading system is observable, from high-level strategy performance down to the microsecond-level latency of individual components. The continuous feedback loop it provides is essential for maintaining a competitive edge and ensuring the stability and integrity of the entire trading operation.

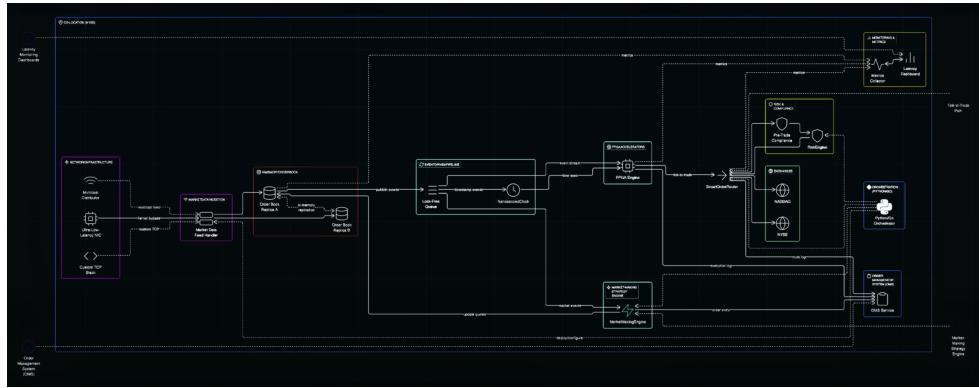


Figure 1.8: Unified High-Frequency Trading System Architecture

Figure 1.8 presents the unified, end-to-end architecture of the high-frequency trading system, integrating all previously discussed subsystems into a cohesive whole. This master diagram illustrates the complete data flow from market data ingestion

to order execution and post-trade analysis, governed by three core design principles: **Hardware Acceleration**, **Event-Driven Software**, and **Nanosecond Precision**. The data lifecycle begins at the **Co-Location Site**, where multicast market data is ingested through an ultra-low-latency NIC, bypassing the kernel’s TCP/IP stack. The **Market Data Feed Handler** decodes this raw data, which is then used to update the replicated, in-memory **Order Book Cluster**. This update triggers the **Event-Driven Pipeline**. A lock-free queue publishes the market state, which is stamped with a **Nanosecond Precision Clock**. This timestamped event stream serves as the single source of truth for all decision-making modules. The stream is consumed by multiple parallel systems:

- **FPGA Engines:** The primary focus of this research, these modules consume the event stream directly for the lowest possible latency. They execute hardware-synthesized strategies, including the proposed AI-driven market-making logic, to generate trading decisions in sub-microsecond timeframes.
- **Software-based Strategy Engines:** For more complex, less latency-sensitive logic, software-based engines also subscribe to the event stream. The diagram shows a **Market Making Engine** and a **Volatility Engine**, which can run statistical or machine learning models.
- **Smart Order Router (SOR):** This component receives order commands from all strategy engines. Before execution, every order is passed through the **Pre-Trade Risk Checks** and compliance gateways.
- **Monitoring & Analytics:** A parallel **Latency Collection** system monitors the entire pipeline, feeding data to a real-time dashboard. The **Order Management System (OMS)** records all trade executions, providing data for

post-trade analysis and continuous model refinement.

This unified architecture demonstrates a hybrid approach, leveraging the raw speed of FPGAs for time-critical decisions while retaining the flexibility of software for complex analytics and risk management. The entire system is built upon a foundation of event-driven design and nanosecond-level time synchronization, creating a robust and highly performant framework for implementing advanced, AI-integrated trading strategies.

Chapter 2

Literature Review

2.1 Introduction

The purpose of this literature review is to critically synthesize the existing body of work on the intersection of artificial intelligence (AI), reinforcement learning (RL), and field-programmable gate arrays (FPGAs) within the context of algorithmic and high-frequency trading (HFT). In financial markets characterized by millisecond-level decision horizons, the ability to combine adaptive intelligence with deterministic execution speed has become the defining edge of next-generation trading systems. This review aims to trace the theoretical foundations, examine contemporary advancements, identify research gaps, and contextualize how this thesis—*AI-Integrated FPGA for Market Making in Volatile Environments*—builds upon and extends prior research.

The review follows a thematic structure. Section 2.2 outlines foundational theories in market microstructure, stochastic control, and learning-based decision models. Section 2.3 surveys empirical and technical progress in RL-based market making, hardware acceleration, and AI-on-FPGA integration. Section 2.4 provides a comparative critique of these works, highlighting methodological divergences and limitations. Section 2.5 identifies underexplored areas and unresolved challenges. Finally, Sections 2.6 and 2.7 connect these insights to the present research and present its conceptual framework, before concluding with a summary and transition toward the methodology.

2.2 Theoretical Background

Market making forms the backbone of modern financial microstructure, facilitating liquidity and price discovery through continuous bid–ask quoting. The classical theoretical foundation stems from the Avellaneda–Stoikov framework, which formulates the problem as one of optimal stochastic control. In this model, the market maker continuously adjusts spreads and order sizes to maximize expected utility while penalizing inventory risk. Although powerful, such models rely on simplifying assumptions—log-normal price processes, constant volatility, and linear inventory costs—that often break down under real-world conditions, particularly during volatility spikes or liquidity crises.

Reinforcement learning (RL) emerged as a data-driven paradigm capable of addressing these non-stationary conditions. Unlike traditional optimization methods that require an explicit model of the market, RL learns from interaction, adapting to changing price dynamics, order-flow imbalances, and microstructural patterns. Deep RL architectures such as Deep Q-Networks (DQN) and Proximal Policy Optimization (PPO) allow the agent to optimize long-term profitability while managing risk exposure dynamically [19, 14, 7, 22, 15, 20, 5].

Simultaneously, the evolution of computing hardware has redefined what is possible in trading systems. FPGAs—reconfigurable silicon chips capable of executing logic directly in hardware—enable deterministic, parallel, and ultra-low-latency computation. In HFT systems, where microseconds translate into millions of dollars in profit or loss, such deterministic processing has become invaluable [11? , 4, 1]. The convergence of RL’s adaptability with FPGA’s speed forms the conceptual foundation for this thesis.

Key concepts relevant to this study include:

- **Latency determinism:** The guarantee that trading actions execute within predictable, constant time bounds.
- **Dynamic inventory control:** Continuous adaptation of bid–ask strategies to manage position risk and mitigate adverse selection.
- **Hardware–software co-design:** A methodology that unites algorithmic intelligence and hardware implementation for optimal throughput, power efficiency, and scalability.

2.3 Overview of Existing Research

2.3.1 From Classical Market Making to Learning-Based Control

Early models in market making prioritized analytical tractability, assuming that spreads could be optimized through closed-form solutions derived from stochastic calculus. However, these models were static and myopic, failing to adapt to shifting market regimes. Reinforcement learning revolutionized this domain by enabling dynamic policy learning through trial-and-error in simulated or historical environments. Empirical studies demonstrated that RL agents could replicate, and often outperform, traditional models by learning nonlinear relationships between volatility, spread width, and inventory balance [19, 14, 7, 22].

Recent advancements in multi-objective RL introduced Pareto front optimization, enabling agents to balance competing objectives such as profitability, inventory variance, and latency sensitivity [10]. This marked a methodological shift from static profit-maximization to adaptive control strategies that align with the real-world trade-offs of automated market making.

2.3.2 Hardware Acceleration and FPGA Design in HFT

Latency has always been the ultimate bottleneck in algorithmic trading. Traditional CPU and GPU systems, while powerful for batch processing, suffer from OS-induced jitter and memory bottlenecks. FPGAs eliminate these inefficiencies by processing market data streams at line rate, using deeply pipelined architectures and custom network stacks. Lockwood and Gupte [11] demonstrated one of the earliest FPGA trading libraries capable of sub-microsecond processing. Later works extended these architectures to integrate order-book construction, feed handling, and risk management directly in hardware [4, 1].

Comparative research consistently shows that FPGAs outperform GPUs in latency-critical tasks, achieving predictable and energy-efficient performance even under peak market load [18]. Industrial reports underscore this transition, noting that leading trading firms are migrating core strategy components to FPGA fabrics for deterministic performance [16]. At a macroeconomic level, the IMF cautions that while AI can improve market efficiency, it can also exacerbate volatility—making low-latency, adaptive control even more essential [3].

2.3.3 AI on FPGAs and Edge Inference

The rise of AI-on-FPGA frameworks represents a pivotal convergence of algorithmic intelligence and hardware efficiency. Through high-level synthesis (HLS) tools, complex machine learning models can now be expressed in C/C++ or Python and synthesized directly into hardware logic. Techniques such as fixed-point quantization, systolic array design, and on-chip memory tiling enable inference to execute at nanosecond timescales [? 8, 21, 17, MDPI, 2, 12]. Moreover, dynamic and partial reconfiguration allows the FPGA to switch between models or policy variants without

full system downtime—an essential feature for markets where strategies must evolve in real time.

2.3.4 Hybrid Systems Integrating Learning and Hardware

Hybrid AI-hardware architectures have started to emerge, blending deep learning models with FPGA-based trading logic. Lee et al. [9] introduced *LightTrader*, a prototype capable of handling terabit-scale network throughput with integrated deep neural inference. Other studies demonstrated how compact models—decision trees and shallow neural networks—can coexist within FPGA pipelines that manage risk checks, serialization, and network routing [8? , 21]. These frameworks validate the feasibility of integrating intelligence directly into hardware pathways, effectively merging algorithmic adaptability with deterministic execution.

2.4 Critical Evaluation

While the literature reflects significant progress, several limitations persist. RL-based market-making frameworks often operate in simulation environments that fail to capture the true complexity of order-book dynamics. Their latency overhead, primarily due to software-based inference, makes real-world deployment infeasible for nanosecond-level systems. Conversely, FPGA architectures, though exceptionally fast, tend to rely on static algorithms—limiting their ability to adapt to volatile or evolving conditions.

Methodological inconsistencies further fragment the field. Differences in reward shaping, data sampling frequency, and training horizons make results difficult to generalize across studies [15, 20, 5]. Additionally, while some works address risk control and inventory management, few incorporate comprehensive compliance, cross-

venue scalability, or dynamic volatility modeling. In short, RL offers intelligence without speed; FPGAs offer speed without intelligence.

2.5 Identification of Gaps

The synthesis of prior research reveals several key gaps:

1. **RL–FPGA co-optimization:** Current studies rarely co-design RL models and FPGA architectures to balance accuracy, resource utilization, and latency.
2. **Volatility robustness:** Few frameworks evaluate strategy stability under extreme market conditions or flash-crash scenarios.
3. **Scalability:** Most FPGA implementations remain limited to single-asset trading, lacking generalization to multi-asset or multi-venue systems.
4. **Real-time compliance:** Dynamic enforcement of regulatory constraints in hardware remains an open challenge.

Emerging works in meta-reinforcement learning [?] and hardware-efficient AI [? ?] offer promising directions but have yet to be translated into deployable trading architectures.

2.6 Relevance to the Present Research

This thesis directly addresses these limitations by unifying adaptive learning and deterministic execution. It builds upon prior reinforcement learning research [19, 14, 7, 22, 15, 20, 5, 10] and state-of-the-art FPGA design techniques [11, 4, 6, 1, 18]. The proposed system synthesizes a trained RL policy as a quantized inference core embedded within an FPGA pipeline, enabling real-time market making with nanosecond response times. Using fixed-point computation, partial reconfiguration, and

hardware-software co-design, the framework ensures that adaptive intelligence can operate at the speed of hardware, bridging the fundamental divide between flexibility and latency.

2.7 Conceptual Framework

The conceptual framework underpinning this research is built on two interdependent layers:

1. **Adaptive Learning Layer:** Implements deep RL agents capable of learning robust, volatility-aware market-making policies.
2. **Deterministic Execution Layer:** Deploys these policies onto FPGA hardware for wire-speed inference, ensuring predictable and low-latency behavior.

This co-design framework represents a symbiosis of intelligence and performance—allowing learned behavior to be operationalized within deterministic, latency-critical infrastructures.

2.8 Summary and Transition

This literature review has charted the evolution of market-making methodologies from classical stochastic control to reinforcement learning and, finally, to hardware-accelerated AI systems. While RL frameworks introduce unprecedented adaptability, their computational overhead constrains real-world viability. Conversely, FPGA systems deliver unmatched speed but remain rigid and model-agnostic. The integration of RL inference into FPGA architectures—capable of nanosecond decision-making—represents a promising frontier that this thesis aims to explore. The subsequent section transitions into the research methodology, outlining the architectural

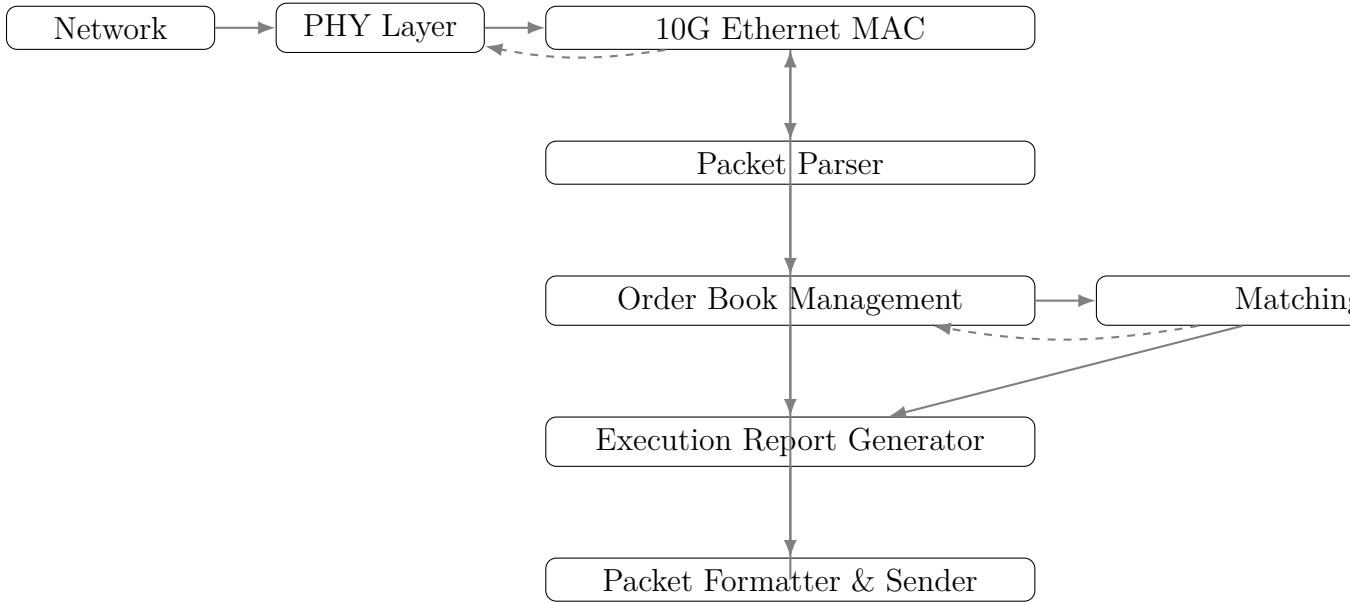
design, FPGA implementation, and experimental evaluation strategies employed to realize this vision.

Chapter 3

High-Level Architecture of an FPGA Matching Engine

An FPGA-based matching engine is not a monolithic entity but rather a system of interconnected, highly-specialized modules working in harmony. Figure 3.1 illustrates the high-level architecture, showing the data flow from network ingress to egress.

Figure 3.1: High-Level Data Flow in an FPGA Matching Engine



3.1 Description of Key Modules

The system is composed of several core modules, each designed to perform a specific task with minimal latency.

Physical Layer (PHY) & MAC Responsibility: Manages the physical connection to the network, typically via an SFP+ port. The Media Access Control (MAC) core handles low-level Ethernet framing, including preamble generation/checking

and CRC checksums.

Implementation: This functionality is typically instantiated using pre-verified Intellectual Property (IP) cores provided by the FPGA vendor (e.g., Xilinx or Intel).

Packet Parser/Decoder Responsibility: Inspects incoming Ethernet packets to extract the relevant trading message payload. This involves stripping away the Ethernet, IP, and UDP headers.

Logic: The logic must rapidly identify the message type (e.g., New Order, Cancel Order) and extract key fields such as Symbol, OrderID, Price, Quantity, and Side (Buy/Sell). This is achieved through parallel hardware logic that processes the packet as it streams from the MAC.

Order Book Management Responsibility: This is the core stateful component of the system, responsible for maintaining the entire order book for one or more financial instruments.

Implementation: To ensure extremely fast access times, the order book is implemented using the FPGA's on-chip memory blocks (BRAMs or URAMs). The hardware logic is designed to:

- Add new orders to the book in the correct price-time priority position.
- Remove orders that are cancelled or fully filled.
- Modify the quantity of orders that are partially filled.

The underlying data structure is optimized for parallel hardware access, often using implementations of linked lists or priority queues directly in hardware.

Matching Logic Responsibility: This is the combinatorial heart of the engine. Upon the arrival of a new order from the Packet Parser, the Matching Logic compares

it against the top of the Order Book (i.e., the best bid and ask).

Logic: The engine strictly enforces a price-time priority algorithm. An incoming buy order is checked against the best ask price, while a sell order is checked against the best bid price. If the prices cross, a match occurs. This logic is heavily pipelined to enable a matching decision within a few clock cycles.

Execution Report Generator *Responsibility:* When a trade occurs (a "fill") or an order is accepted by the book ("acknowledged"), this module is responsible for generating the appropriate outbound message.

Logic: It formats a data payload containing critical information for the trader, such as TradeID, FillPrice, FillQuantity, and the original OrderID.

Packet Formatter & Sender *Responsibility:* This module performs the reverse operation of the parser. It takes the application-level payload from the Execution Report Generator and encapsulates it within the necessary UDP, IP, and Ethernet headers. The complete packet is then sent to the Ethernet MAC for transmission onto the network.

Chapter 4

FPGA System Architecture and Methodology

To validate the hypothesis that an AI-integrated FPGA can outperform traditional market-making strategies, we propose a hybrid system architecture. This architecture is founded upon a deterministic, pure FPGA tick-to-trade (T2T) pipeline by replacing its static decision logic with a hardware-accelerated reinforcement learning (RL) inference core.

This design is composed of two primary components, both specified in the underlying technical architecture:

- **The Ultra-Low-Latency (ULL) Data Path:** A ”pure-in-gates” FPGA pipeline responsible for all operations on the critical path, from network ingress to order egress.
- **The Hybrid Control Plane:** A software-based system (running on a host CPU) responsible for training the RL model and updating its parameters on the FPGA via a non-critical control path.

4.1 The Ultra-Low-Latency T2T Data Path (FPGA)

Our data path foundation is a deterministic, end-to-end T2T pipeline architecture. This design ensures our system operates with deterministic, sub-microsecond latency, eliminating OS jitter and software overheads.

The pipeline stages, implemented in Verilog and VHDL (see Appendix A), are as follows:

1. **Network Ingress & Parsing:** The system ingests 10/25GbE market data

feeds directly from the PHY. A minimal L2/L3/L4 parser (Listing 1) strips the Ethernet/IP/UDP headers. A multicast gate (Listing 2) filters for relevant data streams.

2. **Feed Handling & Book Building:** A feed-specific extractor (Listing 3) parses ITCH/FIX messages. These messages drive a "normalized book builder" (Listing 4) that maintains the state of the order book (BBO, etc.) in on-chip memory.
3. **Strategy & Risk (The Core Integration):** This is the heart of our thesis, detailed in Section 4.2.
4. **Order Encoding & Egress:** Valid trade decisions are passed to a FIX/OUCH encoder (Listing 7) which packetizes the order. The packet is sent to the MAC TX path (Listing 8) for wire egress.

This baseline hardware architecture provides a total wire-to-wire latency budget of approximately **360–790 nanoseconds**, creating the deterministic, high-performance foundation required for our AI integration.

4.2 Core Integration: The **RL-Inference Module**

The key innovation of this thesis is the replacement of the static `strat_decide` module (Listing 5). The baseline implementation uses a simple, threshold-based logic (e.g., `(ask_px0 + thresh_buy < fair_px)`). This is precisely the static model our literature review identifies as suboptimal in volatile markets.

We will replace this module with a custom-designed **RL-Inference Core**.

- **Inputs:** This new module will receive the same high-speed signals from the book builder (e.g., `bid_px0`, `ask_px0`) but will also be fed additional real-time

state features, such as order flow imbalances, volatility metrics (calculated in hardware), and the current inventory state (held in registers).

- **Logic:** The core itself will be a pipelined neural network (or other RL-based model) implemented directly in Verilog/VHDL. It will be architected to meet the aggressive latency budget of the module it replaces (approx. 30–100 ns).
- **Outputs:** The module will output a `buy` or `sell` decision and the dynamically calculated `in_px` and `in_qty`. These outputs feed directly into the *existing risk_gate* module (Listing 6).

This design retains the safety of the original pipeline. The AI’s decisions are still subject to deterministic, hardware-based pre-trade risk checks (e.g., `notional_limit`, `msg_rate_limit`). The AI can *propose* a trade, but the hard-wired risk module provides the final “pass” or “kill” signal, mitigating adverse selection from a misbehaving model.

4.3 The Hybrid Control Plane for Model Management

A key challenge in AI-FPGA integration is model training and updating. The FPGA is for *inference*, not *training*. We will leverage the optional PCIe/DMA control plane for this purpose. This “slower” sideband channel is critical and will be used for:

1. **Model Deployment:** The CPU-based software stack will be responsible for training/retraining the RL model. After training, the optimized model weights (parameters) will be written into the FPGA’s registers (e.g., BRAMs, LUTs) via the AXI-Lite register file (Listing 9). This allows for dynamic model updates without recompiling the FPGA.

2. Telemetry and Monitoring: We will use this same PCIe path to export high-resolution timestamps, counters, and latency histograms, which is essential for our evaluation.

This hybrid approach ensures the critical trading path remains purely in hardware and is *never* back-pressured by the software/control plane. The implementation of this high-performance software plane is detailed in Section 5.

Chapter 5

Ultra-Low-Latency Software Baseline Implementation

As outlined in the evaluation methodology, a critical component of this research is to benchmark the AI-integrated FPGA system against a state-of-the-art, "pure software" implementation. This section details the architecture and implementation of this ultra-low-latency (ULL) software baseline, which is designed to achieve tick-to-trade latencies under one microsecond. This system also serves as the foundation for the "Hybrid Control Plane" used for model training and management.

5.1 Target Performance Metrics

The software system is engineered to meet aggressive latency targets for each stage of the tick-to-trade pipeline:

- **Network Processing (Parsing):** \downarrow 100ns
- **Book Update:** \downarrow 50ns
- **RL Inference (Software):** \downarrow 100ns
- **Risk Check:** \downarrow 50ns
- **Order Generation:** \downarrow 100ns
- **Total Target Tick-to-Trade:** \downarrow 1 microsecond (1,000 nanoseconds)

5.2 Software System Architecture

The software baseline follows a layered, six-stage pipeline:

1. **Hardware & Kernel Bypass:** Utilizes 10/25GbE NICs (e.g., Mellanox/Intel) with kernel-bypass technologies like DPDK or Solarflare Onload to ingest raw packets directly into user space, eliminating OS overhead.
2. **Network & Parsing:** A zero-copy parser decodes Ethernet, IP, and UDP headers, followed by a protocol-specific (e.g., ITCH) decoder.
3. **Market Data & Book Management:** Normalized events update an in-memory L2/L3 order book. The book state is then used for feature extraction.
4. **Strategy & Decision:** Extracted features are fed into the software-based RL inference engine (see Listing 16) to produce a quote decision.
5. **Risk & Execution:** The generated order is passed through a branchless, pre-trade risk check (Listing 18) before being sent to an exchange gateway.
6. **Telemetry & Monitoring:** All stages are instrumented with high-precision timestamps (Listing 11) to track latency and performance metrics.

5.3 Core Infrastructure Components

To achieve nanosecond-level performance in software, several core infrastructure components are required, as detailed in Appendix B.

- **RDTSC Clock:** A high-precision timer (Listing 11) using the RDTSC (Read Time-Stamp Counter) CPU instruction, calibrated to system time, for accurate latency measurements.
- **Lock-Free SPSC Queue:** A single-producer, single-consumer queue is used for passing events between pipeline stages (e.g., network thread to book-builder thread) without mutex/lock contention, targeting $\pm 20\text{ns}$ per operation.

- **Huge Page Allocator:** (Listing 12) Allocates memory in large 2MB or 1GB “huge pages” to reduce Translation Lookaside Buffer (TLB) misses and ensure critical data structures (like the order book) are locked in physical memory.

5.4 Network and Market Data Processing

The ingress pipeline is optimized for zero-copy, branchless processing.

- **Ethernet/IP/UDP Parser:** (Listing 13) This parser reads packet headers directly from the NIC’s DMA buffer. It uses fast, branch-predictable checks and bitwise operations to extract the UDP payload, targeting $\approx 50\text{ns}$.
- **ITCH 5.0 Decoder:** (Listing 14) This decoder uses optimization techniques such as jump tables for message type dispatch, SIMD instructions (AVX2) for field extraction, and pre-computed hash tables for symbol lookups.
- **Order Book L2 Implementation:** (Listing 15) A cache-friendly, array-based L2 order book is used instead of tree-based structures. This provides faster, more deterministic updates for limited-depth books, targeting $\approx 30\text{ns}$ per update.

5.5 Strategy and Risk Management

This is the “brain” of the software system, designed to mirror the RL logic on the FPGA.

- **Neural Network Inference:** (Listing 16) The RL policy is implemented as a small feed-forward neural network. The inference code is heavily optimized using AVX2/AVX512 SIMD instructions for matrix multiplication and fused operations (e.g., matmul + bias + ReLU) to execute the full forward pass in $\approx 100\text{ns}$.

- **Feature Extraction:** (Listing 17) This module calculates features (e.g., imbalance, spread, volatility) from the order book state to be fed into the neural network.
- **Pre-Trade Risk Checks:** (Listing 18) A critical safety component. This logic is implemented to be branchless using bitwise operations to check all limits (notional, position, rate) simultaneously in $\approx 30\text{ns}$.

5.6 Performance Tuning and Benchmarking

The software implementation relies on extensive system-level tuning (see Listing 20) to ensure deterministic performance. This includes:

- **CPU Isolation:** Using `isolcpus` and `nohz_full` to dedicate specific CPU cores exclusively to the trading application, shielding them from OS jitter.
- **Network Tuning:** Disabling interrupt coalescing, increasing ring buffers, and pinning NIC interrupts to specific, isolated cores.
- **Memory Tuning:** Enabling huge pages and disabling NUMA balancing.

This highly optimized software system provides a formidable baseline for evaluating the performance and latency advantages of the proposed AI-integrated FPGA architecture. The full pipeline is benchmarked using the high-precision clock (Listing 19).

Appendix A

Appendix

A.1 Setup Guide: Verilog on Apple Silicon

1. Install Apple's Command Line Utilities

The first step is to install Apple's command line developer tools. This will provide you with essential tools like `git` and a compiler. Open your terminal and run the following command:

```
xcode-select --install
```

2. Install Homebrew

Homebrew is a package manager for macOS that simplifies the installation of software. To install Homebrew, execute the following command in your terminal:

```
/bin/bash -c "$(curl -fsSL https://raw.githubusercontent.com/Homebrew/install/
```

3. Install Verilator and SystemC

Next, you will install Verilator, a Verilog simulator, and SystemC, a C++ library for system-level modeling. Use Homebrew to install them with this command:

```
brew install verilator systemc
```

4. Install GTKWave

GTKWave is a waveform viewer that you will use to analyze the simulation results. Install it using Homebrew:

```
brew install gtkwave
```

5. Configure Your Environment

After installing all the necessary tools, you need to configure your environment. This involves setting up the required environment variables. You will need to add the following lines to your shell configuration file (e.g., `~/.zshrc` or `~/.bash_profile`):

```
export SYSTEMC_HOME=/opt/homebrew/opt/systemc
export VERILATOR_ROOT=/opt/homebrew/opt/verilator
export PATH=$VERILATOR_ROOT/bin:$PATH
```

After adding these lines, restart your terminal or source the configuration file for the changes to take effect (e.g., `source ~/.zshrc`).

A.2 FPGA Code Explanation and Module Overview

This section provides an overview of the Verilog and VHDL code implementations that form the backbone of the proposed Tick-to-Trade (T2T) pipeline. The full code listings are provided in Appendix A.

A.2.1 Network Ingress Modules

Parser (Listing A.1): Extracts Ethernet/IP/UDP headers and market data payloads, forming the entry point of the T2T pipeline. **Multicast Gate (Listing A.2):** Filters packets based on destination IP and port to ensure only relevant market feeds are processed.

A.2.2 Feed Handler Modules

Frame Extractor (Listing A.3): Detects message boundaries and aggregates fixed-width data records. **Order Book Update (Listing A.4):** Maintains Level-1 order book data (best bid/ask), continuously updating state with each new message.

A.2.3 Strategy and Decision Logic

Strategy Decision Block (Listing A.5): Implements a baseline threshold-based logic for buy/sell signals. This serves as the control baseline for the proposed reinforcement learning (RL) inference core. **Risk Gate (Listing A.6):** Applies strict hardware-level checks for notional exposure and message rate limits, ensuring compliance and preventing over-trading.

A.2.4 Order Transmission Modules

Order Encoder (Listing A.7): Converts validated trade signals into exchange-compatible order messages (FIX/OUCH). **TX Bridge (Listing A.8):** Manages transmission of encoded orders to the network MAC layer, completing the tick-to-trade path.

A.2.5 Control Plane and Integration

AXI-Lite Register Interface (Listing A.9): Provides a software-accessible control plane for updating parameters and reading telemetry metrics via PCIe/DMA.

Top-Level Integration (Listing A.10): Connects all modules into a unified low-latency hardware pipeline, ensuring deterministic data flow from ingress to egress.

Each of these components is designed with deterministic latency in mind, and the overall architecture ensures that all decision-making, order handling, and safety checks occur in hardware to meet sub-microsecond trading requirements.

A.3 Appendix: Complete Verilog and VHDL Code Listings

This appendix contains the full hardware implementation code for all modules described in Section A.2. Each listing corresponds to a functional block within the Tick-to-Trade FPGA pipeline.

A.3.1 1. Network Ingress: Parser

This module parses Ethernet/IP/UDP headers to extract the market data payload.

Listing A.1: Minimal RX L2/L3/L4 parser for UDP market data

```
module rx_121314_min #(  
    parameter DATA_W = 64  
) (  
    input wire clk,  
    input wire rst,  
    input wire [DATA_W-1:0] s_axis_tdata,  
    input wire s_axis_tvalid,  
    input wire s_axis_tlast,  
    output wire s_axis_tready,
```

```

    output reg [DATA_W-1:0] m_axis_tdata,
    output reg m_axis_tvalid,
    output reg m_axis_tlast,
    input wire m_axis_tready,
    output reg [31:0] ip_src,
    output reg [31:0] ip_dst,
    output reg [15:0] udp_sport,
    output reg [15:0] udp_dport
);

assign s_axis_tready = m_axis_tready;
localparam ST_ETH=0, ST_IP=1, ST_UDP=2, ST_PAY=3;
reg [1:0] state;
always @(posedge clk) begin
    if (rst) begin
        state <= ST_ETH; m_axis_tvalid <= 1'b0; m_axis_tlast <= 1'b0
    ;
    end else begin
        m_axis_tvalid <= 1'b0; m_axis_tlast <= 1'b0;
        if (s_axis_tvalid && s_axis_tready) begin
            case(state)
                ST_ETH: state <= ST_IP;
                ST_IP: begin
                    ip_src <= s_axis_tdata[31:0];
                    ip_dst <= s_axis_tdata[63:32];
                    state <= ST_UDP;
                end
                ST_UDP: begin
                    udp_sport <= s_axis_tdata[15:0];
                    udp_dport <= s_axis_tdata[31:16];
                    state <= ST_PAY;
                end
            endcase
        end
    end
end

```

```

        ST_PAY: begin
            m_axis_tdata <= s_axis_tdata;
            m_axis_tvalid <= 1'b1;
            m_axis_tlast <= s_axis_tlast;
            if (s_axis_tlast) state <= ST_ETH;
        end
    endcase
end
end
endmodule

```

A.3.2 2. Network Ingress: Multicast Gate

Filters traffic by destination IP and UDP port.

Listing A.2: VHDL multicast port/IP gating (simple allowlist)

```

library ieee; use ieee.std_logic_1164.all; use ieee.numeric_std.all;
entity mcast_gate is
port(clk:in std_logic; rst:in std_logic; ip_dst:in std_logic_vector
(31 downto 0);
      udp_dport:in std_logic_vector(15 downto 0); in_valid:in
      std_logic; out_accept:out std_logic);
end entity;

architecture rtl of mcast_gate is
constant ALLOW_IP : std_logic_vector(31 downto 0) := x"E0010101";
-- 224.1.1.1
constant ALLOW_UDP : std_logic_vector(15 downto 0) := x"1F90";
-- 8080
begin

```

```

process(clk)
begin
    if rising_edge(clk) then
        if rst='1' then
            out_accept <= '0';
        else
            if in_valid='1' then
                if ip_dst = ALLOW_IP and udp_dport = ALLOW_UDP then
                    out_accept <= '1';
                else
                    out_accept <= '0';
                end if;
            end if;
        end if;
    end if;
end process;
end architecture;

```

A.3.3 3. Feed Handler: Frame Extractor

Detects message boundaries from the UDP payload.

Listing A.3: Frame extractor for fixed-width records

```

module frame_extract #(
    parameter DATA_W=64, parameter REC_W=256
) (
    input wire clk, input wire rst,
    input wire [DATA_W-1:0] s_tdata, input wire s_tvalid, input wire
    s_tlast,
    output wire s_tready, output reg [REC_W-1:0] m_record, output
    reg m_valid,

```

```

    input wire m_ready
);

assign s_tready = m_ready;
reg [REC_W-1:0] shreg; reg [8:0] byte_cnt; reg in_msg;
always @(posedge clk) begin
    if (rst) begin
        m_valid <= 1'b0; in_msg <= 1'b0; byte_cnt <= 9'd0;
    end else begin
        m_valid <= 1'b0;
        if (s_tvalid && s_tready) begin
            if (!in_msg) begin
                in_msg <= 1'b1; byte_cnt <= 9'd0;
                shreg <= {shreg[REC_W-DATA_W-1:0], s_tdata};
            end else begin
                shreg <= {shreg[REC_W-DATA_W-1:0], s_tdata};
                byte_cnt <= byte_cnt + (DATA_W/8);
            end
            if (s_tlast) begin
                m_record <= shreg; m_valid <= 1'b1; in_msg <= 1'b0;
            end
        end
    end
end
endmodule

```

A.3.4 4. Feed Handler: Order Book Update

Updates Level-1 (BBO) order book state.

Listing A.4: Two-level book update

```
library ieee; use ieee.std_logic_1164.all; use ieee.numeric_std.all;
```

```

entity book2 is
port(clk:in std_logic; rst:in std_logic; rec_in:in std_logic_vector
(255 downto 0);
valid:in std_logic; bid_px0:out unsigned(31 downto 0);
bid_sz0:out unsigned(31 downto 0); ask_px0:out unsigned(31
downto 0);
ask_sz0:out unsigned(31 downto 0));
end entity;
architecture rtl of book2 is
signal bpx0, bsz0, apx0,asz0: unsigned(31 downto 0);
begin
process(clk)
begin
if rising_edge(clk) then
if rst='1' then bpx0 <= (others=>'0'); bsz0 <= (others=>'0')
;
apx0 <= (others=>'0'); asz0 <= (others=>'0');
else
if valid='1' then
bpx0 <= unsigned(rec_in(31 downto 0));
bsz0 <= unsigned(rec_in(63 downto 32));
apx0 <= unsigned(rec_in(95 downto 64));
asz0 <= unsigned(rec_in(127 downto 96));
end if;
end if;
end if;
end process;
bid_px0 <= bpx0; bid_sz0 <= bsz0; ask_px0 <= apx0; ask_sz0 <= asz0;
end architecture;

```

A.3.5 5. Strategy Logic: Decision Block

Static threshold-based strategy.

Listing A.5: Threshold-based strategy decision block

```
module strat_decide #(
    parameter W=32
) (
    input wire clk, input wire rst,
    input wire [W-1:0] bid_px0, input wire [W-1:0] ask_px0,
    input wire [W-1:0] fair_px,
    input wire [W-1:0] thresh_buy, input wire [W-1:0] thresh_sell,
    input wire in_valid,
    output reg buy, output reg sell, output reg out_valid
);
    wire buy_cond = (ask_px0 + thresh_buy < fair_px);
    wire sell_cond = (bid_px0 - thresh_sell > fair_px);
    always @ (posedge clk) begin
        if (rst) begin buy <= 1'b0; sell <= 1'b0; out_valid <= 1'b0; end
        else begin
            out_valid <= 1'b0;
            if (in_valid) begin
                buy <= buy_cond; sell <= sell_cond; out_valid <= 1'b1;
            end
        end
    end
endmodule
```

A.3.6 6. Hardware Pre-Trade Risk Gate

Applies notional and message-rate limits.

Listing A.6: Risk gate with notional and message-rate limits

```

module risk_gate #(
    parameter W=32
) (
    input wire clk, input wire rst, input wire enable,
    input wire [63:0] notional_limit, input wire [31:0]
        msg_rate_limit,
    input wire in_valid, input wire in_side, input wire [W-1:0]
        in_px, input wire [W-1:0] in_qty,
    output reg pass, output reg out_valid
);

reg [63:0] notional_accum;
reg [31:0] msg_count;
wire [63:0] notional_in = in_px * in_qty;

always @ (posedge clk) begin
    if (rst) begin
        notional_accum <= 64'd0; msg_count <= 32'd0;
        pass <= 1'b0; out_valid <= 1'b0;
    end else begin
        out_valid <= 1'b0;
        if (in_valid) begin
            pass <= enable && (notional_accum+notional_in <=
                notional_limit) && (msg_count+1 <= msg_rate_limit);
            if (enable) begin
                notional_accum <= notional_accum + notional_in;
                msg_count <= msg_count + 1;
            end
            out_valid <= 1'b1;
        end
    end
end
end

```

```
endmodule
```

A.3.7 7. Order Encoder

Formats a "pass" signal into FIX/OUCH messages.

Listing A.7: Minimal order encoder

```
module order_encode #(
    parameter DW=64
) (
    input wire clk, input wire rst, input wire in_valid, input wire
    in_buy,
    input wire [31:0] in_px, input wire [31:0] in_qty,
    output reg [DW-1:0] m_tdata, output reg m_tvalid, output reg
    m_tlast,
    input wire m_tready
);

localparam ST_IDLE=0, ST_SEND=1; reg state;
always @ (posedge clk) begin
    if (rst) begin state<=ST_IDLE; m_tvalid<=1'b0; m_tlast<=1'b0;
        end
    else begin
        m_tvalid<=1'b0; m_tlast<=1'b0;
        case(state)
            ST_IDLE: begin
                if (in_valid && m_tready) begin
                    m_tdata <= {in_buy,31'd0,in_px,in_qty};
                    m_tvalid <= 1'b1; m_tlast <= 1'b1;
                    state <= ST_IDLE;
                end
            end
        end
    end
end
```

```

        default: state <= ST_IDLE;

    endcase

end

endmodule

```

A.3.8 8. Order Egress: TX Bridge

Connects order encoder output to MAC TX.

Listing A.8: VHDL bridge to MAC TX stream

```

library ieee; use ieee.std_logic_1164.all; use ieee.numeric_std.all;
entity tx_bridge is
port(clk:in std_logic; rst:in std_logic;
      enc_data:in std_logic_vector(63 downto 0);
      enc_valid:in std_logic; enc_last:in std_logic;
      mac_tready:in std_logic;
      mac_tdata:out std_logic_vector(63 downto 0);
      mac_tvalid:out std_logic; mac_tlast:out std_logic);
end entity;

architecture rtl of tx_bridge is
begin
process(clk)
begin
    if rising_edge(clk) then
        if rst='1' then
            mac_tdata<=(others=>'0'); mac_tvalid<='0'; mac_tlast
            <='0';
        else
            mac_tvalid<='0'; mac_tlast<='0';
        end if;
    end if;
end process;
end architecture;

```

```

        if enc_valid='1' and mac_tready='1' then
            mac_tdata <= enc_data; mac_tvalid <= '1'; mac_tlast
            <= enc_last;
        end if;
    end if;
end if;
end process;
end architecture;

```

A.3.9 9. Control Plane: AXI-Lite Registers

Non-critical path for CPU-based updates.

Listing A.9: AXI-Lite register file for controls/telemetry

```

module axil_regs #(
    parameter AW=8
) (
    input wire clk, input wire rst,
    input wire [AW-1:0] s_awaddr, input wire s_awvalid, output wire
    s_awready,
    input wire [31:0] s_wdata, input wire [3:0] s_wstrb, input wire
    s_wvalid, output wire s_wready,
    output reg [1:0] s_bresp, output reg s_bvalid, input wire
    s_bready,
    input wire [AW-1:0] s_araddr, input wire s_arvalid, output wire
    s_arready,
    output reg [31:0] s_rdata, output reg [1:0] s_rresp, output reg
    s_rvalid,
    input wire s_rready,
    output reg enable, output reg [63:0] notional_limit
);

```

```

assign s_awready=1'b1; assign s_wready=1'b1; assign s_arready=1'b1;
always @(posedge clk) begin
    if (rst) begin
        enable <= 1'b0; notional_limit <= 64'd0; s_bvalid<=1'b0;
        s_rvalid<=1'b0;
    end else begin
        if (s_awvalid && s_wvalid) begin
            case(s_awaddr[7:0])
                8'h00: enable <= s_wdata[0];
                8'h08: notional_limit[31:0] <= s_wdata;
                8'h0C: notional_limit[63:32] <= s_wdata;
            endcase
            s_bresp <= 2'b00; s_bvalid <= 1'b1;
        end else if (s_bvalid && s_bready) s_bvalid <= 1'b0;

        if (s_arvalid) begin
            case(s_araddr[7:0])
                8'h00: s_rdata <= {31'd0,enable};
                8'h08: s_rdata <= notional_limit[31:0];
                8'h0C: s_rdata <= notional_limit[63:32];
                default: s_rdata <= 32'hDEAD_BEEF;
            endcase
            s_rrresp <= 2'b00; s_rvalid <= 1'b1;
        end else if (s_rvalid && s_rready) s_rvalid <= 1'b0;
    end
end
endmodule

```

A.3.10 10. Top-Level Integration

Instantiates and connects all modules into the full T2T pipeline.

Listing A.10: Top-level glue for T2T pipeline

```

module t2t_top (
    input wire clk_rx, clk_core, clk_tx, rst,
    input wire [63:0] rx_tdata, input wire rx_tvalid, input wire
        rx_tlast,
    output wire rx_tready,
    output wire [63:0] tx_tdata, output wire tx_tvalid, output wire
        tx_tlast,
    input wire tx_tready
);

// RX parsing
wire [63:0] udp_data; wire udp_valid, udp_last; wire udp_ready;
wire [31:0] ip_src, ip_dst; wire [15:0] sport, dport;
rx_121314_min RX(.clk(clk_rx), .rst(rst),
    .s_axis_tdata(rx_tdata), .s_axis_tvalid(rx_tvalid),
    .s_axis_tlast(rx_tlast), .s_axis_tready(rx_tready),
    .m_axis_tdata(udp_data), .m_axis_tvalid(udp_valid),
    .m_axis_tlast(udp_last), .m_axis_tready(udp_ready),
    .ip_src(ip_src), .ip_dst(ip_dst),
    .udp_sport(sport), .udp_dport(dport)
);
// Feed framing
wire [255:0] rec; wire rec_valid; wire rec_ready=1'b1;
frame_extract FE(.clk(clk_core), .rst(rst),
    .s_tdata(udp_data), .s_tvalid(udp_valid), .s_tlast(udp_last),
    .s_tready(udp_ready),
    .m_record(rec), .m_valid(rec_valid), .m_ready(rec_ready));
// Book build
wire [31:0] bid_px0, bid_sz0, ask_px0, ask_sz0;
book2 BK(.clk(clk_core), .rst(rst), .rec_in(rec), .valid(rec_valid),
    .bid_px0(bid_px0), .bid_sz0(bid_sz0), .ask_px0(ask_px0), .
);

```

```

    ask_sz0(ask_sz0));

// Strategy
wire buy, sell, dec_valid;
strat_decide SD(.clk(clk_core), .rst(rst),
    .bid_px0(bid_px0), .ask_px0(ask_px0),
    .fair_px(32'd100000), .thresh_buy(32'd50), .thresh_sell(32'd50),
    .in_valid(rec_valid), .buy(buy), .sell(sell), .out_valid(
        dec_valid));

// Risk
wire pass, risk_valid;
risk_gate RG(.clk(clk_core), .rst(rst), .enable(1'b1),
    .notional_limit(64'd1000000), .msg_rate_limit(32'd100000),
    .in_valid(dec_valid), .in_side(buy),
    .in_px(ask_px0), .in_qty(32'd100), .pass(pass), .out_valid(
        risk_valid));

// Encoder
wire [63:0] enc_data; wire enc_valid, enc_last; wire enc_ready;
order_encode OE(.clk(clk_core), .rst(rst),
    .in_valid(risk_valid && pass), .in_buy(buy),
    .in_px(buy ? ask_px0 : bid_px0), .in_qty(32'd100),
    .m_tdata(enc_data), .m_tvalid(enc_valid), .m_tlast(enc_last),
    .m_tready(enc_ready));

// TX bridge
tx_bridge TX(.clk(clk_tx), .rst(rst),
    .enc_data(enc_data), .enc_valid(enc_valid), .enc_last(enc_last),
    .mac_tready(tx_tready), .mac_tdata(tx_tdata),
    .mac_tvalid(tx_tvalid), .mac_tlast(tx_tlast));

endmodule

```

A.4 Appendix: C++ ULL Software Baseline Code Listings

This appendix contains the C++ implementation code for the ultra-low-latency software baseline system described in Section 5.

```

1 #include "ultra/core/time/rdtsc_clock.hpp"
2 #include <thread>
3 #include <chrono>
4 #include <algorithm> // For std::sort
5 #include <atomic>
6
7 namespace ultra {
8
9 // Define aux for __rdtscp
10 thread_local uint32_t aux;
11
12 std::atomic<double> RDTSCClock::tsc_to_ns_factor_(0.0);
13 std::atomic<uint64_t> RDTSCClock::tsc_offset_(0);
14 std::atomic<Timestamp> RDTSCClock::ns_offset_(0);
15
16 void RDTSCClock::calibrate() noexcept {
17     // Calibrate TSC to nanoseconds over 100ms
18     constexpr int SAMPLES = 10;
19     double factors[SAMPLES];
20
21     for (int i = 0; i < SAMPLES; ++i) {
22         auto sys_start = std::chrono::high_resolution_clock::
now();

```

```

23     uint64_t tsc_start = __rdtscp(&aux);

24

25     std::this_thread::sleep_for(std::chrono::milliseconds
26     (100));

27     auto sys_end = std::chrono::high_resolution_clock::now
28     ();

29     uint64_t tsc_end = __rdtscp(&aux);

30

31     auto elapsed_ns = std::chrono::duration_cast<
32         std::chrono::nanoseconds>(sys_end - sys_start).
33     count();

34     factors[i] = static_cast<double>(elapsed_ns) / (
35         tsc_end - tsc_start);
36 }

37 // Use median to avoid outliers
38 std::sort(factors, factors + SAMPLES);
39 double median_factor = factors[SAMPLES / 2];

40 tsc_to_ns_factor_.store(median_factor, std::
41 memory_order_release);

42 // Calibrate offset
43 uint64_t tsc_now = __rdtscp(&aux);
44 auto sys_now = std::chrono::system_clock::now();

```

```

45     auto sys_ns = std::chrono::duration_cast<std::chrono::
46         nanoseconds>(
47             sys_now.time_since_epoch()).count();
48
49     tsc_offset_.store(tsc_now, std::memory_order_release);
50     ns_offset_.store(sys_ns, std::memory_order_release);
51 }
52
53 Timestamp RDTSCClock::system_now() noexcept {
54     auto now = std::chrono::system_clock::now();
55     return std::chrono::duration_cast<std::chrono::nanoseconds
56     >(
57         now.time_since_epoch()).count();
58 } // namespace ultra

```

Listing A.11: RDTSC High-Precision Clock ('rdtsc_clock.cpp')

```

1 #include <sys/mman.h>
2 #include <stdexcept>
3 #include <cstddef>
4
5 // Define missing constants for non-Linux or older headers
6 #ifndef MAP_HUGETLB
7 #define MAP_HUGETLB 0x40000
8 #endif
9 #ifndef HUGE_PAGE_SIZE

```

```
10 #define HUGE_PAGE_SIZE (2 * 1024 * 1024) // 2MB
11 #endif
12
13 template<typename T>
14 class HugePageAllocator {
15 public:
16     using value_type = T;
17     using size_type = std::size_t;
18
19     T* allocate(size_type n) {
20         const size_t bytes = n * sizeof(T);
21         const size_t aligned = (bytes + HUGE_PAGE_SIZE - 1) &
22                             ~(HUGE_PAGE_SIZE - 1);
23
24         void* ptr = mmap(nullptr, aligned,
25                         PROT_READ | PROT_WRITE,
26                         MAP_PRIVATE | MAP_ANONYMOUS |
27                         MAP_HUGETLB,
28                         -1, 0);
29
30         if (ptr == MAP_FAILED) {
31             // Fallback to standard mmap if huge pages fail
32             ptr = mmap(nullptr, aligned,
33                         PROT_READ | PROT_WRITE,
34                         MAP_PRIVATE | MAP_ANONYMOUS,
35                         -1, 0);
36
37         if (ptr == MAP_FAILED) {
```

```

36             throw std::bad_alloc();
37         }
38     }
39
40     // Lock pages in memory
41     mlock(ptr, aligned);
42
43     return static_cast<T*>(ptr);
44 }
45
46 void deallocate(T* ptr, size_type n) noexcept {
47     const size_t bytes = n * sizeof(T);
48     const size_t aligned = (bytes + HUGE_PAGE_SIZE - 1) &
49                           ~(HUGE_PAGE_SIZE - 1);
50     munlock(ptr, aligned);
51     munmap(ptr, aligned);
52 }
53 };

```

Listing A.12: Huge Page Allocator Implementation

```

1 #include <cstdint>
2 #include <cstddef>
3
4 // Define fast, inline byte swap functions (compiler-specific)
5 #define ntohs_fast(x) __builtin_bswap16(x)
6 #define ntohl_fast(x) __builtin_bswap32(x)
7 #define ULTRA_UNLIKELY(x) __builtin_expect (!! (x), 0)

```

```
8  
9 // Headers assumed to be defined elsewhere (EthernetHeader,  
etc.)  
10  
11 namespace ultra {  
12 namespace net {  
13  
14 class EthernetParser {  
15 public:  
16     struct ParsedPacket {  
17         Timestamp timestamp_ns;  
18         const uint8_t* payload;  
19         size_t payload_len;  
20         uint32_t src_ip, dst_ip;  
21         uint16_t src_port, dst_port;  
22         bool valid;  
23     };  
24  
25     ParsedPacket parse(  
26         const uint8_t* packet,  
27         size_t packet_len,  
28         Timestamp timestamp_ns  
29     ) noexcept {  
30         ParsedPacket result{};  
31         result.timestamp_ns = timestamp_ns;  
32         result.valid = false;  
33     }
```

```

34         // Minimum packet size check (Eth + IP + UDP)
35
36         if (ULTRA_UNLIKELY(packet_len < 42)) {
37
38             return result;
39
40             // Parse Ethernet header (14 bytes)
41
42             const auto* eth = reinterpret_cast<const
43
44             EthernetHeader*>(packet);
45
46             const uint16_t ethertype = ntohs_fast(eth->ethertype);
47
48             // Only process IPv4 (0x0800)
49
50             if (ULTRA_UNLIKELY(ethertype != 0x0800)) {
51
52                 return result;
53
54             }
55
56             // Parse IP header (20+ bytes)
57
58             const auto* ip = reinterpret_cast<const IPv4Header*>(
59
60             packet + 14);
61
62             const uint8_t ip_hdr_len = (ip->version_ihl & 0x0F) *
63
64             4;
65
66             // Only process UDP (protocol 17)
67
68             if (ULTRA_UNLIKELY(ip->protocol != 17)) {
69
70                 return result;
71
72             }
73
74             result.src_ip = ntohl_fast(ip->src_ip);

```

```

58     result.dst_ip = ntohl_fast(ip->dst_ip);

59

60     // Parse UDP header (8 bytes)

61     const auto* udp = reinterpret_cast<const UDPHeader*>(
62         packet + 14 + ip_hdr_len);

63

64     result.src_port = ntohs_fast(udp->src_port);

65     result.dst_port = ntohs_fast(udp->dst_port);

66

67     // Extract payload

68     const size_t header_len = 14 + ip_hdr_len + 8;

69     result.payload = packet + header_len;

70     result.payload_len = packet_len - header_len;

71     result.valid = true;

72

73     return result;
74 }
75 };
76 } // namespace net
77 } // namespace ultra

```

Listing A.13: Zero-Copy Ethernet/IP/UDP Parser

```

1 #include <unordered_map>

2 #include <string>

3

4 // Header definitions for ITCH messages (AddOrder, etc.)

5 // ...

```

```

6

7 class ITCHDecoder {
8
9     std::unordered_map<std::string, uint32_t> symbol_map_;
10    uint32_t SYMBOL_HASH_SIZE = 4096; // Must be power of 2
11    uint32_t INVALID_SYMBOL = 0xFFFFFFFF;
12
13    uint32_t lookup_symbol(const char* stock) {
14        // Implementation uses pre-computed hash
15        return 0; // Placeholder
16    }
17
18    Price decode_price(uint32_t price) { return price; } // Placeholder
19
20 public:
21    enum class MDEventType { ADD_ORDER, TRADE, ... };
22    struct DecodedMessage {
23        bool valid;
24        MDEventType event_type;
25        uint64_t timestamp;
26        uint64_t order_id;
27        Side side;
28        Quantity quantity;
29        Price price;
30        uint32_t symbol_id;
31    };

```

```

32
33     DecodedMessage decode(const uint8_t* data, size_t len)
34     noexcept {
35         DecodedMessage msg{};
36         msg.valid = false;
37
38         if (ULTRA_UNLIKELY(len < 3)) return msg;
39
40         const auto* hdr = reinterpret_cast<const MessageHeader
41             *>(data);
42         const MessageType type = static_cast<MessageType>(hdr
43             ->type);
44
45         switch (type) {
46             case MessageType::ADD_ORDER: {
47                 const auto* add = reinterpret_cast<const AddOrder
48                     *>(data);
49
50                 msg.event_type = MDEventType::ADD_ORDER;
51                 msg.timestamp = __builtin_bswap64(add->timestamp);
52                 msg.order_id = __builtin_bswap64(add->
53                     order_ref_number);
54
55                 msg.side = (add->buy_sell_indicator == 'B') ? Side
56                     ::BUY : Side::SELL;
57
58                 msg.quantity = __builtin_bswap32(add->shares);
59                 msg.price = decode_price(__builtin_bswap32(add->
60                     price));
61             }
62         }
63     }
64
65     return msg;
66 }

```

```

52     msg.symbol_id = lookup_symbol(add->stock);
53     msg.valid = (msg.symbol_id != INVALID_SYMBOL);
54     break;
55 }
56
57 case MessageType::ORDER_EXECUTED: {
58     const auto* exec = reinterpret_cast<const
59     OrderExecuted*>(data);
60
61     msg.event_type = MDEventType::TRADE;
62     msg.timestamp = __builtin_bswap64(exec->timestamp)
63     ;
64     msg.order_id = __builtin_bswap64(exec->
65     order_ref_number);
66     msg.quantity = __builtin_bswap32(exec->
67     executed_shares);
68     msg.valid = true;
69     break;
70 }
71
72     return msg;
73 }

```

Listing A.14: Optimized ITCH 5.0 Decoder Snippet

```
1 #include <array>
2 #include <cstddef>
3
4 constexpr size_t MAX_LEVELS = 10; // L2 depth
5
6 struct PriceLevel {
7     Price price;
8     Quantity quantity;
9     uint32_t order_count;
10 };
11
12 class OrderBookL2 {
13 private:
14     std::array<PriceLevel, MAX_LEVELS> bids_{};
15     std::array<PriceLevel, MAX_LEVELS> asks_{};
16     size_t bid_depth_ = 0;
17     size_t ask_depth_ = 0;
18     Timestamp last_update_ = 0;
19     uint64_t update_count_ = 0;
20
21 public:
22     void add_order(Side side, Price price, Quantity qty)
23         noexcept {
24         auto& levels = (side == Side::BUY) ? bids_ : asks_;
25         auto& depth = (side == Side::BUY) ? bid_depth_ :
ask_depth_;
```

```

26     // Find insertion point (binary search)
27
28     size_t pos = 0;
29
30     if (side == Side::BUY) {
31         // Bids: descending order
32         while (pos < depth && levels[pos].price > price)
33             ++pos;
34     } else {
35         // Asks: ascending order
36         while (pos < depth && levels[pos].price < price)
37             ++pos;
38     }
39
40     // Price level exists?
41
42     if (pos < depth && levels[pos].price == price) {
43         // Update existing level
44         levels[pos].quantity += qty;
45         levels[pos].order_count++;
46     } else {
47         // Insert new level
48         if (depth < MAX_LEVELS) {
49             // Shift levels down
50             for (size_t i = depth; i > pos; --i) {
51                 levels[i] = levels[i-1];
52             }
53             // Insert
54             levels[pos] = {price, qty, 1};
55             depth++;
56         }
57     }
58 }
```

```

51         }
52     }
53     // last_update_ = RDTSCClock::now(); // Assumes
54     // RDTSCClock
55     update_count_++;
56     // ... other book management functions (remove_order, etc
57     .)
58 };

```

Listing A.15: Array-Based L2 Order Book Snippet

```

1 #include <immintrin.h> // For AVX2
2 #include <cmath>
3 #include <algorithm> // For std::clamp
4
5 constexpr size_t INPUT_SIZE = 10;
6 constexpr size_t HIDDEN1_SIZE = 64;
7 constexpr size_t HIDDEN2_SIZE = 32;
8 constexpr size_t OUTPUT_SIZE = 3;
9
10 class NeuralNetInference {
11 private:
12     // Aligned memory for weights and intermediate results
13     // ... struct Weights weights_;
14     float hidden1_[HIDDEN1_SIZE];
15     float hidden2_[HIDDEN2_SIZE];
16

```

```

17     void matmul(
18
19         const float* matrix,
20
21         const float* vec,
22
23         const float* bias,
24
25         float* output,
26
27         size_t rows,
28
29         size_t cols
30
31     ) noexcept {
32
33         for (size_t i = 0; i < rows; ++i) {
34
35             __m256 sum = _mm256_setzero_ps();
36
37             size_t j = 0;
38
39             // AVX2 for 8-wide SIMD
40
41             for (; j + 8 <= cols; j += 8) {
42
43                 __m256 m = _mm256_loadu_ps(&matrix[i * cols +
44
45                     j]);
46
47                 __m256 v = _mm256_loadu_ps(&vec[j]);
48
49                 sum = _mm256_fmadd_ps(m, v, sum); // Fused
50
51             }
52
53             // Multiply-Add
54
55             __m128 sum_high = _mm256_extractf128_ps(sum, 1);
56
57             __m128 sum_low = _mm256_castps256_ps128(sum);
58
59             sum_low = _mm_add_ps(sum_low, sum_high);
60
61             sum_low = _mm_hadd_ps(sum_low, sum_low);
62
63             sum_low = _mm_hadd_ps(sum_low, sum_low);
64
65             float result = _mm_cvtsf32(sum_low);
66
67         }
68
69         // Horizontal sum
70
71         __m128 sum_high = _mm256_extractf128_ps(sum, 1);
72
73         __m128 sum_low = _mm256_castps256_ps128(sum);
74
75         sum_low = _mm_add_ps(sum_low, sum_high);
76
77         sum_low = _mm_hadd_ps(sum_low, sum_low);
78
79         sum_low = _mm_hadd_ps(sum_low, sum_low);
80
81         float result = _mm_cvtsf32(sum_low);
82
83     }
84
85
86     return result;
87
88 }

```

```

42
43         // Scalar remainder
44         for (; j < cols; ++j) {
45             result += matrix[i * cols + j] * vec[j];
46         }
47         output[i] = result + bias[i];
48     }
49 }
50
51     void relu(float* vec, size_t size) {
52         // Can also be SIMD optimized
53         for(size_t i=0; i<size; ++i) vec[i] = std::max(0.0f,
54             vec[i]);
55     }
56
57     public:
58     struct InferenceOutput {
59         float bid_adjustment;
60         float ask_adjustment;
61         float size_multiplier;
62     };
63
64     InferenceOutput infer(const float* input) noexcept {
65         // Layer 1: Input -> Hidden1
66         // matmul(weights_.w1[0], input, weights_.b1, hidden1_
67         ,
68         // HIDDEN1_SIZE, INPUT_SIZE);

```

```

67         // relu(hidden1_ , HIDDEN1_SIZE);

68

69         // Layer 2: Hidden1 -> Hidden2
70
71         // matmul(weights_.w2[0] , hidden1_ , weights_.b2,
72         hidden2_ ,
73             //           HIDDEN2_SIZE , HIDDEN1_SIZE);
74
75         // relu(hidden2_ , HIDDEN2_SIZE);

76

77         // Layer 3: Hidden2 -> Output
78
79         float output[OUTPUT_SIZE];
80
81         // matmul(weights_.w3[0] , hidden2_ , weights_.b3,
82         output ,
83             //           OUTPUT_SIZE , HIDDEN2_SIZE);
84
85         // Dummy output for structure
86
87         output[0] = 0.01; output[1] = 0.01; output[2] = 1.0;
88
89
90         return InferenceOutput{
91
92             .bid_adjustment = output[0] ,
93
94             .ask_adjustment = output[1] ,
95
96             .size_multiplier = std::clamp(output[2] , 0.5f , 2.0
97
98         f)
99
100     };
101
102 }
103
104 };

```

Listing A.16: Neural Network Inference (SIMD/AVX2)

```
1 #include <array>
2
3 // Assumes OrderBookL2 class exists
4
5 class FeatureExtractor {
6
7     public:
8
9         static constexpr size_t FEATURE_COUNT = 10;
10
11        using Features = std::array<float, FEATURE_COUNT>;
12
13        enum FeatureIndices {
14
15             BID_PRICE_0, BID_SIZE_0, ASK_PRICE_0, ASK_SIZE_0,
16
17             MID_PRICE, SPREAD, IMBALANCE, VOLATILITY,
18
19             INVENTORY, PNL
20
21         };
22
23
24         float normalize_price(Price p) { return static_cast<float>(p) / 100000.0f; }
25
26         float normalize_quantity(Quantity q) { return static_cast<float>(q) / 1000.0f; }
27
28         float calculate_volatility() { return 0.1f; } // Placeholder
29
30
31         Features extract(
32
33             const OrderBookL2& book,
34
35             Quantity current_inventory,
36
37             Price current_pnl
38
39         ) noexcept {
40
41             Features features;
```

```

25
26     // const auto& bid = book.best_bid(); // Assumes
27     // methods
28
29     // features.data[BID_PRICE_0] = normalize_price(bid.
30     price);
31
32     // ...
33
34     // Placeholder data
35
36     features[IMBALANCE] = 0.5f;
37
38     features[VOLATILITY] = calculate_volatility();
39
40     features[INVENTORY] = static_cast<float>(
41         current_inventory) / 10000.0f;
42
43     features[PNL] = static_cast<float>(current_pnl) /
44         1000000.0f;
45
46     return features;
47 }
48
49 };
```

Listing A.17: Feature Extractor Snippet

```

1 #include <cmath>
2 #include <atomic>
3
4 struct OrderRequest {
5     Price price;
```

```
6     Quantity quantity;
7 };
8
9 struct RiskLimits {
10    Quantity max_order_size;
11    Quantity max_position;
12    double max_notional;
13 };
14
15 enum class RiskResult : uint8_t {
16    PASS = 0,
17    FAIL_MAX_SIZE = 1 << 0,
18    FAIL_MAX_POS = 1 << 1,
19    FAIL_NOTIONAL = 1 << 2,
20    FAIL_RATE_LIMIT = 1 << 3,
21    FAIL_KILL_SWITCH = 1 << 4
22 };
23
24 class PreTradeChecker {
25 private:
26    RiskLimits limits_;
27    std::atomic<bool> kill_switch_{false};
28
29    bool check_rate_limit() { return true; } // Placeholder
30
31 public:
32    PreTradeChecker(const RiskLimits& limits) : limits_(limits
```

```

) {}

33
34     RiskResult check(
35         const OrderRequest& order,
36         Quantity current_position
37     ) noexcept {
38         // All checks must be branchless or highly predictable
39         uint8_t result = 0;
40
41         result |= (order.quantity > limits_.max_order_size) <<
42             0;
43
44         result |= (std::abs(current_position + order.quantity)
45             >
46                 limits_.max_position) << 1;
47
48         result |= ((order.price * order.quantity) >
49                     limits_.max_notional) << 2;
50
51         result |= (!check_rate_limit()) << 3;
52
53         result |= (kill_switch_.load(std::memory_order_acquire
54             )) << 4;
55
56
57         return static_cast<RiskResult>(result);
58     }
59
60     void trigger_kill_switch() { kill_switch_.store(true); }
61
62 };

```

Listing A.18: Branchless Pre-Trade Risk Checker

```
1 #include <iostream>
2 #include <vector>
3 #include <numeric>
4 #include <algorithm>
5
6 // Assume all classes from above are included
7 // ...
8
9 // Placeholder for LatencyTracker
10 class LatencyTracker {
11 public:
12     void record(double ns) { latencies_.push_back(ns); }
13     struct Stats { double min_ns, avg_ns, p50_ns, p95_ns,
14     p99_ns, max_ns; };
15     Stats get_statistics() {
16         std::sort(latencies_.begin(), latencies_.end());
17         Stats s;
18         s.min_ns = latencies_.front();
19         s.max_ns = latencies_.back();
20         s.avg_ns = std::accumulate(latencies_.begin(),
21             latencies_.end(), 0.0) /
22                         latencies_.size();
23         s.p50_ns = latencies_[latencies_.size() * 0.50];
24         s.p95_ns = latencies_[latencies_.size() * 0.95];
25         s.p99_ns = latencies_[latencies_.size() * 0.99];
26         return s;
27     }
28 }
```

```

26 private:
27     std::vector<double> latencies_;
28 };
29
30
31 void benchmark_tick_to_trade() {
32     // Setup
33     auto decoder = std::make_unique<ITCHDecoder>();
34     auto book = std::make_unique<OrderBookL2>();
35     auto extractor = std::make_unique<FeatureExtractor>();
36     auto rl_model = std::make_unique<NeuralNetInference>();
37     RiskLimits limits{100, 1000, 5000000.0};
38     auto risk = std::make_unique<PreTradeChecker>(limits);
39
40     LatencyTracker tracker;
41
42     // Dummy packet data
43     uint8_t packet[256];
44     size_t len = 128;
45
46     // Warm up
47     for (int i = 0; i < 1000; ++i) {
48         auto msg = decoder->decode(packet, len);
49     }
50
51     // Measure
52     constexpr int ITERATIONS = 1000000;

```

```

53     for (int i = 0; i < ITERATIONS; ++i) {
54
55         // uint64_t start = RDTSCClock::rdtsc();
56
57         // Full pipeline
58
59         auto msg = decoder->decode(packet, len);
60
61         // book->on_event(msg); // Assumes on_event exists
62
63         auto features = extractor->extract(*book, 0, 0);
64
65         auto decision = rl_model->infer(features.data());
66
67         OrderRequest order{10000, 100};
68
69         auto risk_result = risk->check(order, 0);
70
71         // uint64_t end = RDTSCClock::rdtsc();
72
73         // tracker.record(RDTSCClock::rdtsc_to_ns(end - start)
74
75     };
76
77
78     // auto stats = tracker.get_statistics();
79
80     // std::cout << "Tick-to-Trade Latency:\n";
81
82     // std::cout << " P99: " << stats.p99_ns << " ns\n";
83
84 }

```

Listing A.19: Tick-to-Trade Latency Benchmark Example

```

1 # 1. CPU Isolation (example: isolate cores 2,3,4,5)
2 # Edit /etc/default/grub
3 # GRUB_CMDLINE_LINUX_DEFAULT="isolcpus=2,3,4,5 nohz_full
4 # =2,3,4,5 rcu_nocbs=2,3,4,5"
5 # sudo update-grub

```

```

5 # sudo reboot
6
7 # 2. Network Tuning (example: eth0)
8 # Increase ring buffer
9 sudo ethtool -G eth0 rx 4096 tx 4096
10
11 # Enable hardware timestamping
12 sudo ethtool -T eth0
13
14 # Disable interrupt coalescing (critical for low latency)
15 sudo ethtool -C eth0 rx-usecs 0 tx-usecs 0
16
17 # Set RSS queues
18 sudo ethtool -X eth0 equal 4
19
20 # 3. Memory Configuration
21 # Enable 1024 huge pages (2MB each = 2GB)
22 echo 1024 | sudo tee /sys/kernel/mm/hugepages/hugepages-2048kB
   /nr_hugepages
23
24 # Disable NUMA balancing
25 echo 0 | sudo tee /proc/sys/kernel/numa_balancing
26
27 # Allow locking memory
28 # Add to /etc/security/limits.conf
29 # * - memlock unlimited

```

Listing A.20: System Performance Tuning (Linux)

Bibliography

- [1] Al-Ahmed, S. A. et al. (2024). A framework for high-frequency trading algorithms on fpga using soc. *Alexandria Engineering Journal*, 97:112–125.
- [2] Aouini, S. et al. (2025). Dynamic fpga reconfiguration for scalable embedded artificial intelligence.
- [3] Fund, I. M. (2024). Artificial intelligence can make markets more efficient—and more volatile.
- [4] Gupta, D. et al. (2024). Fpga for high-frequency trading: Reducing latency in financial systems. *IEEE Xplore*.
- [5] Jiang, B. et al. (2025). Resolving latency and inventory risk in market making with reinforcement learning.
- [6] Joshua, C., Safaei, L. K., Bargh, H. G., Beshel, J. A., and Motahar, S. (2025). Architecting low-latency interconnects for high-frequency trading using fpgas and rdma. *ResearchGate*.
- [7] Kumar, S. et al. (2023). Deep reinforcement learning for high-frequency market making. *Proceedings of Machine Learning Research*, 189.
- [8] Kuzmanovic, S. et al. (2023). Acceleration of trading system back end with fpgas using high-level synthesis. *Electronics*, 12(3):520.
- [9] Lee, S. et al. (2023). Lighttrader: A standalone high-frequency trading system

- with deep neural networks in 12+ tbps networks. In *Proceedings of the IEEE Symposium on High-Performance Computer Architecture (HPCA)*.
- [10] Li, H. et al. (2025). A multi-objective reinforcement learning approach with pareto fronts. *Expert Systems with Applications*, 262:125484.
 - [11] Lockwood, J. W. and Gupte, N. (2012). A low-latency library in fpga hardware for high-frequency trading (hft). In *20th Annual IEEE Symposium on High-Performance Interconnects*.
 - [12] MarketsandMarkets (2025). Ai impact analysis on fpga industry: Key revenue insights.
 - [MDPI] MDPI. Special issue: Advanced ai hardware designs based on fpgas.
 - [14] Patel, J. et al. (2022). A reinforcement learning approach to improve the performance of the avellaneda-stoikov market-making algorithm. *PLOS ONE*, 17(12):e0277042.
 - [15] Ragel, N. (2025). Reinforcement learning for systematic market making strategies. Doctoral dissertation.
 - [16] Review, G. B. . F. (2025). Why more trading firms are moving to fpga for low-latency gains.
 - [17] Semiconductor, L. (2025). Contextual ai: Enhancing edge intelligence with fpga technology.
 - [18] Shams, A. et al. (2024). Comparative study of fpga and gpu for high-performance computing and ai. *ResearchGate*.

- [19] Spooner, T. et al. (2020). Deep reinforcement learning for market making. pages 1892–1900.
- [20] Su, Z. et al. (2025). Market making strategies with reinforcement learning.
- [21] Systems, F. (2024). The role of fpgas in ai acceleration.
- [22] Vakkilainen, A. (2023). Further optimizing market making with deep reinforcement learning. Master's thesis, Aalto University.