# Chapter 1

# Tensor-Based Optimization for Next-Gen AI Trading Systems

[1] Summary. This white paper proposes a next-generation Operating System (OS) tailored for small AI models in financial trading systems, leveraging tensor-based optimization to minimize energy consumption, latency, and future costs while maximizing AI model value and usability. The system integrates consensus mechanisms, and Zero-Knowledge Proofs (ZKProof) to ensure optimal performance and security across distributed nodes and validators. By focusing on the unique requirements of AI-driven trading, this approach ensures real-time decision-making, enhanced security, and improved resource management for financial institutions.

## 1.1 Introduction

AI models have become integral to trading, with the growing demand for low-latency, high-performance systems that adapt to volatile markets. While large models dominate many domains, small, efficient AI models are better suited for real-time decision-making in high-frequency trading (HFT). This white paper focuses on the development of an OS optimized for small AI models, using tensor algebra to optimize system efficiency, incorporating martingale theory for decision-making, and embedding ZKProof for security. We aim to build an AI-first OS that scales seamlessly across multiple computing nodes and validators, ensuring efficient resource allocation and secure execution.

## 1.2 Tensor-Based System Architecture

1. Computational Tensor **C** Represents the computation performed by each of the n nodes over time t and resources r. The energy consumption is:

$$E_{comp} = \sum_{i=1}^{n} \sum_{t=1}^{T} \sum_{r=1}^{R} \mathbf{C}(i, t, r) \cdot P_{h,i}(t, r)$$

---

[1]If you're interested in discussing the implementation, feel free to reach out to me at nvh0@yahoo.com

2. Validation Tensor **V** Represents the validation effort required by m validators. The total energy consumed by validation is:

$$E_{val} = \sum_{j=1}^{m} \sum_{t=1}^{T} \sum_{v=1}^{V} \mathbf{V}(j, t, v) \cdot P_{v,j}(t, v)$$

3. Resource Allocation Tensor **R** Tracks resource usage across nodes and validators. The OS ensures that total resource consumption at any time t does not exceed the maximum available resources:

$$\sum_{x=1}^{n+m} \mathbf{R}(x, t, k) \leq R_{max}(t, k)$$

4. AI Value Tensor {**A**} Measures the value of each AI model a across nodes and time steps, capturing its contribution to profit, risk reduction, and execution speed. The total AI value at time t is:

$$V_{total}(t) = \sum_{i=1}^{n} \sum_{a=1}^{A} \mathbf{A}(i, t, a)$$

5. AI Usability Tensor **U** Captures the usability of AI models, including ease of integration and resource requirements. The OS dynamically prioritizes models with higher usability.

## 1.3 Total System Latency

The total system latency $L_{total}$ is a critical factor in high-frequency trading and is explicitly calculated using a latency tensor:

$$\mathbf{L}(x, t, a) \in \mathbb{R}^{(n+m) \times T \times A}$$

The total latency at time t is the sum of the maximum latencies across all nodes and validators:

$$L_{total}(t) = \max_{x \in \{1,2,\ldots,n+m\}} \sum_{a=1}^{A} \mathbf{L}(x, t, a)$$

## 1.4 Consensus and Zero-Knowledge Proof (ZKProof)

1. Consensus Mechanism A consensus tensor **Cns** tracks the agreement among validators for each trade or AI model decision:

$$\mathbf{Cns}(j, t, a) \in \mathbb{R}^{m \times T \times A}$$

The system requires a majority of validators to agree on a decision, with the consensus level calculated as:

$$C_{level}(t, a) = \frac{1}{m} \sum_{j=1}^{m} \mathbf{Cns}(j, t, a)$$

2. ZKProof for Security The ZKProof tensor **ZKP** captures the validation status of AI models using Zero-Knowledge Proofs, ensuring the correctness of decisions without revealing sensitive data:

$$\mathbf{ZKP}(j, t, a) \in \mathbb{R}^{m \times T \times A}$$

The system verifies model outputs securely, with the ZKProof validation level:

$$ZKP_{level}(t, a) = \frac{1}{m} \sum_{j=1}^{m} \mathbf{ZKP}(j, t, a)$$

## 1.5  Tensor-Based Future Cost Optimization

The total future cost of the system, balancing energy consumption, latency, AI model value, usability, consensus, and ZKProof, is given by:

$$C_{future} = \sum_{t=1}^{T} \left( E_{comp}(t) + E_{val}(t) + \lambda \cdot L_{total}(t) \right)$$

$$- \sum_{t=1}^{T} \left( V_{total}(t) + U_{total}(t) + \gamma \cdot C_{level}(t) + \delta \cdot ZKP_{level}(t) \right)$$

## 1.6  Conclusion

The proposed tensor-based OS efficiently integrates AI value, usability, consensus, and Zero-Knowledge Proofs (ZKProof) to enhance system performance, security, and efficiency. By explicitly modeling total system latency and future costs, the system is optimized for real-time AI-driven trading, providing a scalable, low-latency solution for next-generation financial markets. This approach ensures maximum AI model performance, reduced energy consumption, and secure trading in a decentralized, high-frequency trading environment.

**References**

Föllmer, H., & Schied, A. (2011). Stochastic Finance: An Introduction in Discrete Time. Walter de Gruyter.

Boyd, S., & Vandenberghe, L. (2004). Convex Optimization. Cambridge University Press.

Kolda, T. G., & Bader, B. W. (2009). Tensor Decompositions and Applications. SIAM Review, 51(3), 455-500.

Goldreich, O., Micali, S., & Wigderson, A. (1986). Proofs that Yield Nothing but Their Validity and a Methodology of Cryptographic Protocol Design. In Proceedings of the 27th Annual IEEE Symposium on Foundations of Computer Science (pp. 174-187).

Lamport, L., Shostak, R.,  Pease, M. (1982). The Byzantine Generals Problem. ACM Transactions on Programming Languages and Systems (TOPLAS), 4(3), 382-401.

Nesterov, Y. (2018). Lectures on Convex Optimization. Springer.

De Prado, M. L. (2018). Advances in Financial Machine Learning. John Wiley & Sons.

Tanenbaum, A. S., & Van Steen, M. (2007). Distributed Systems: Principles and Paradigms. Prentice Hall.

Papalexakis, E. E., Sidiropoulos, N. D., & Bro, R. (2016). From K-Means to Higher-Way Co-Clustering: Tensor Decomposition with Applications to Big Data. Data Mining and Knowledge Discovery, 31(3), 424-456.

Zhang, Y., & Gupta, V. (2016). Energy Efficient Computing for HFT and Real-Time Financial Analytics. IEEE Transactions on Computers, 65(5), 1562-1574.