

HPC has transformed financial services

Yihe (Jordan) Zhang

2025 Spring – CS455

4/17/2025



Contents

Hardware

Scalable FPGA systems accelerate High-Frequency-Trading (HFT)

Software

Parallel algorithms help speedup the financial models and algorithms

Future Directions

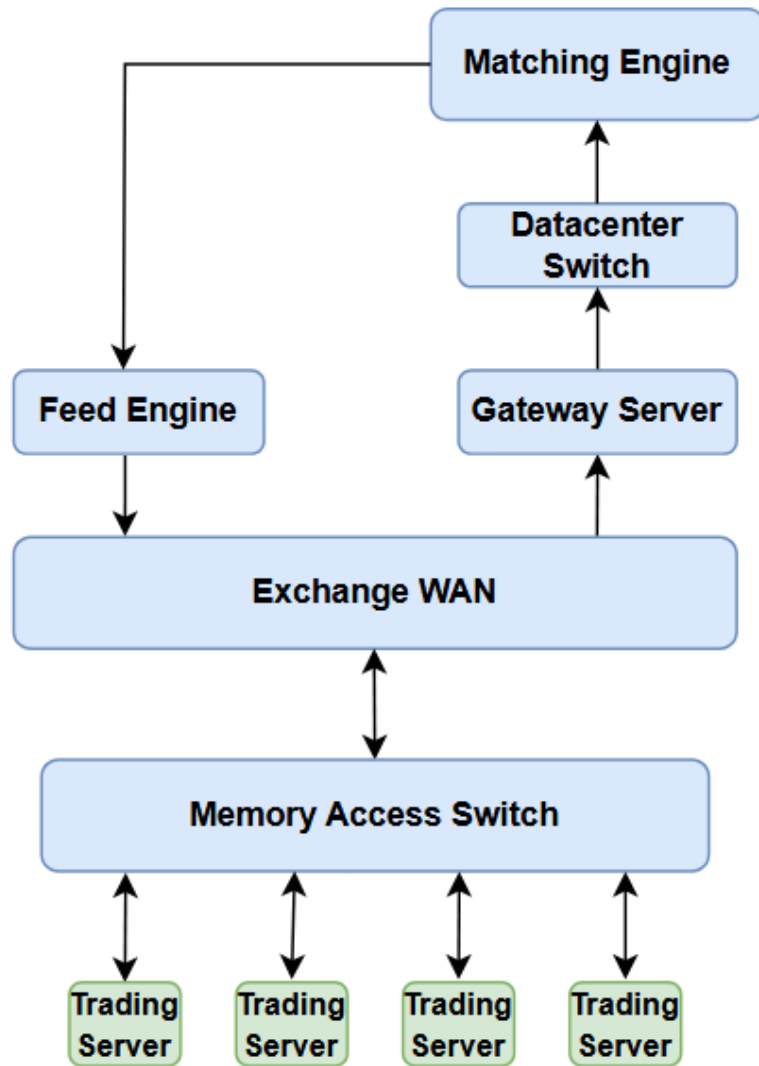
Hardware



What is High-Frequency Trading (HFT)?

- A trading strategy that executes thousands to millions of orders in fractions of a second
- Relies on ultra-low-latency data processing and decision making
- Success is determined by **speed** – microseconds can make or break a trade
- Used for arbitrage, market making, and short-term opportunities

Hardware



⚡ How FPGA Accelerates HFT ?

- **Bypasses OS latency:** Uses UDP offloading to process packets directly in hardware
- **Real-time protocol decoding:** Supports FAST message parsing with pipelined microcode architecture
- **Parallel data paths:** Simultaneous decoding of multiple market data streams
- **Low and deterministic latency:** Ideal for sub-microsecond trading decisions
- **Custom hardware logic:** Tailored for specific trading strategies and faster execution

Software

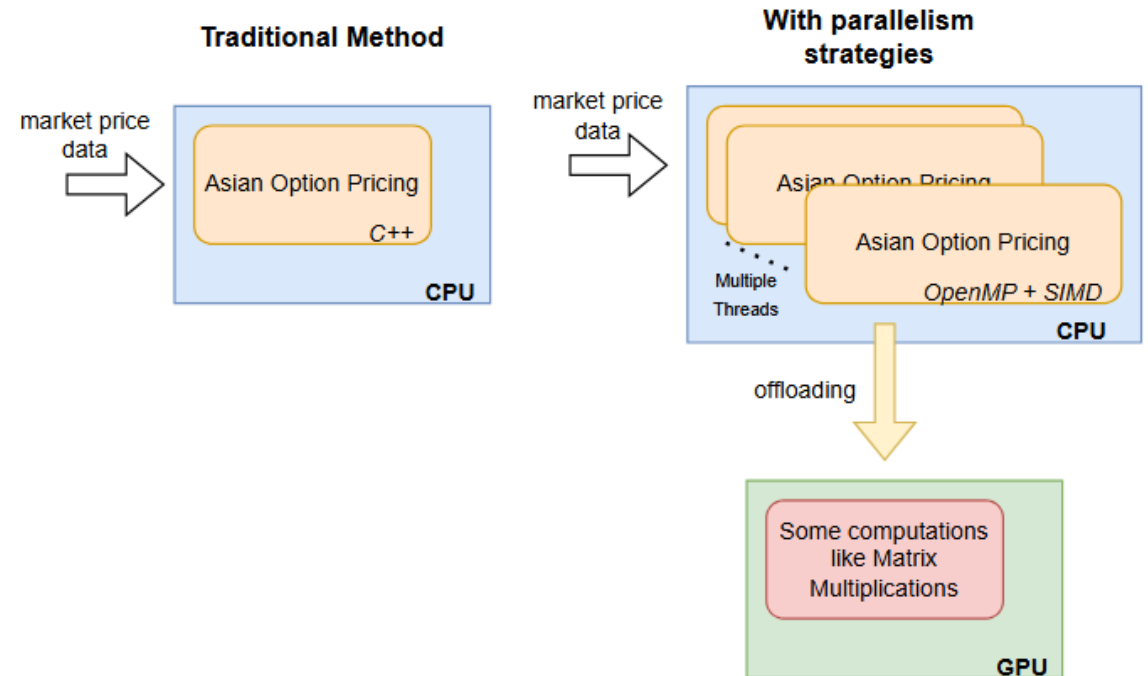
How HPC Software works in Modern Financial Systems

- Financial models are becoming more **complex** and **data-intensive**
- Real-time pricing and risk assessment require **high-speed computation**
- Modern hardware (CPUs, GPUs) demands software that can fully utilize **parallelism**
- HPC frameworks help improve **performance**, **scalability**, and **responsiveness**

Software

✚ Exploiting Parallelism in Financial Workloads

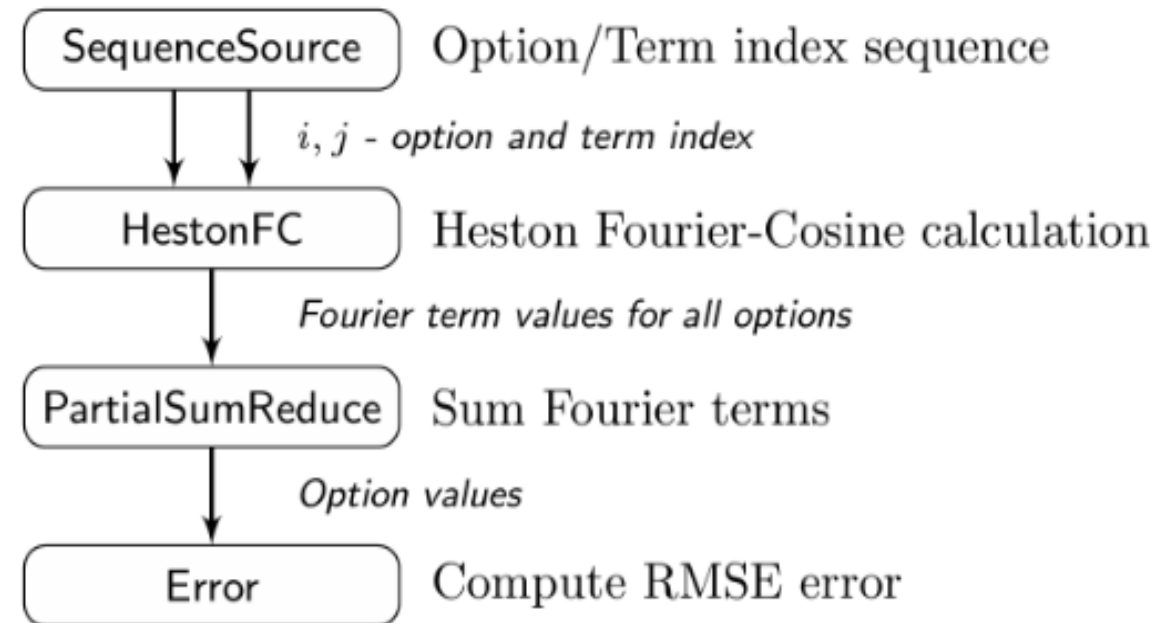
- Financial applications like **option pricing** and **market simulation** are inherently parallel
- Techniques:
 - SIMD vectorization** (data-level parallelism)
 - OpenMP** for thread-level CPU parallelism
 - CUDA** for GPU acceleration



Software

HPC Framework for Heston Model

- Calibration of stochastic models (e.g., **Heston model**) is **computationally expensive**
- **Xcelerit** framework for parallel execution with high-level C++ code
- Modular pipeline design supports efficient computation:
 - Input generation → HestonFC → Aggregation → Error evaluation
- Supports **GPU/CPU** backend without low-level coding



M. Dixon, J. Lotze, and M. Zubair, "A portable and fast stochastic volatility model calibration using multi and many-core processors," in 2014 Seventh Workshop on HighPerformance Computational Finance, pp. 23–28, 2014.

Future HPC&Finance

Future Directions

- Many existing studies apply HPC to finance, but exploration is still limited.
- Leverage the full power of supercomputer to operate large-scale simulations (e.g., stress testing, what-if analysis).
- Hardware accelerators (GPUs, FPGAs) are underutilized for: Real-time risk management during live trading.

Q & A

Thanks for listening !