



Latency in High Performance Trading Systems

Feb 2010

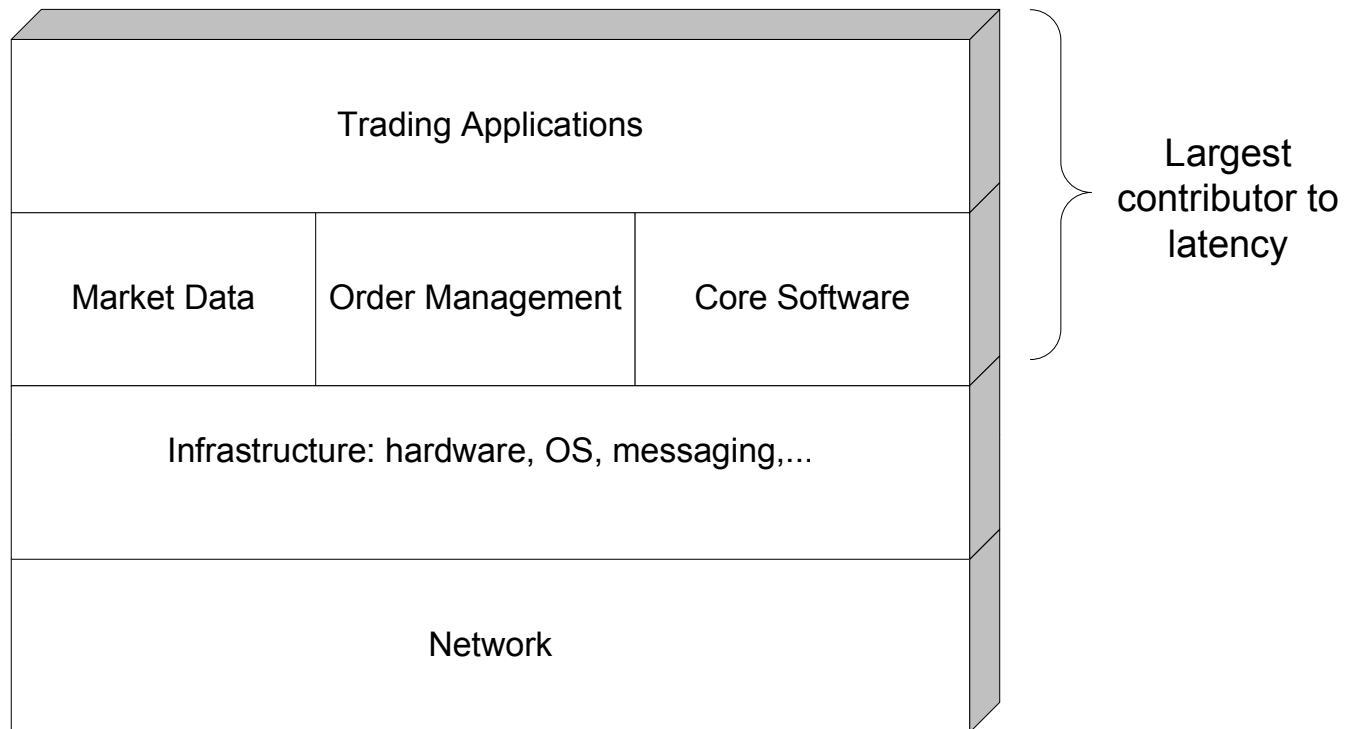
Stephen Gibbs
Automated Trading Group

 **Bank Financial Group**

Overview

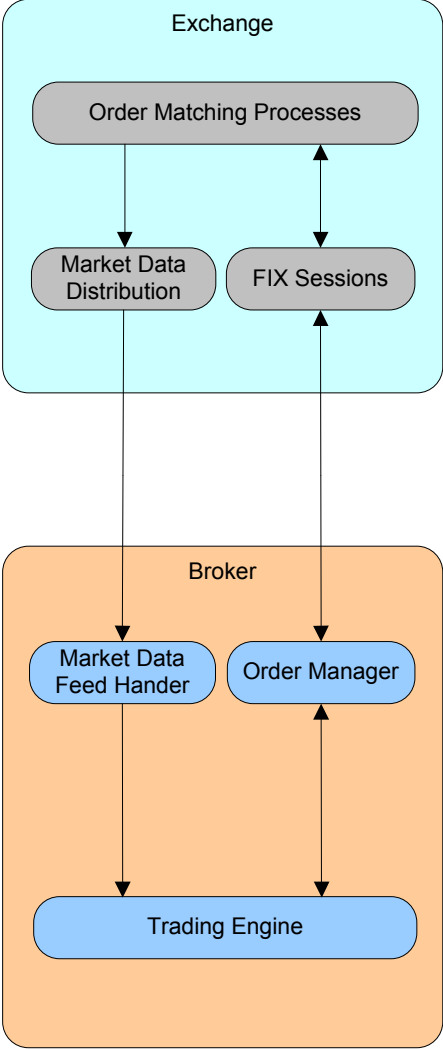
- Review the architecture of a typical automated trading system
- Review the major sources of latency, many at the infrastructure level you should be aware of already and we won't discuss low level networking issues, but focus big picture of the overall architecture from exchange to TD's trading engine and back again.
- Introduce some market jargon
- Lots of examples of latency in real order and market data
- Show a little video using an order book visualizer
- Wrap up with a brief discuss about some advanced technologies to mitigate latency issues (some of which I have used
- The goal is to give provide detail so that if you choose you can dig into some of the more technical aspects of latency as it relates to trading

Sources of Latency

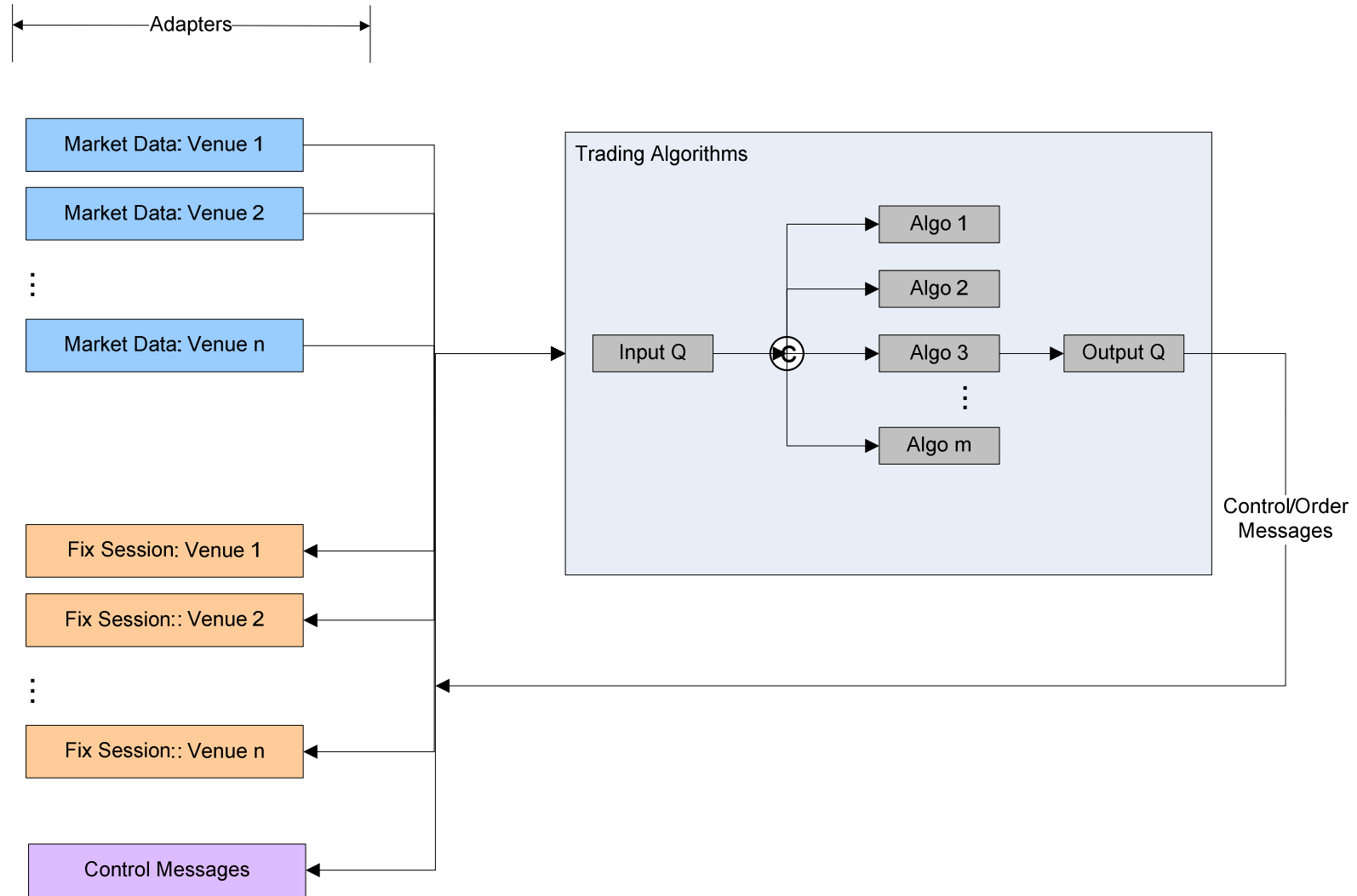


Equity Trading System Architecture

High Level Data Flow



Algo Data Flow



Message Protocols

Financial Information eXchange

- A self describing protocol used for transmitting trade related information.
- ASCII with a header, body and trailer, including a checksum
- FIX Session is layered on TCP
- delimiter SOH (0x01)
- Admin and application messages

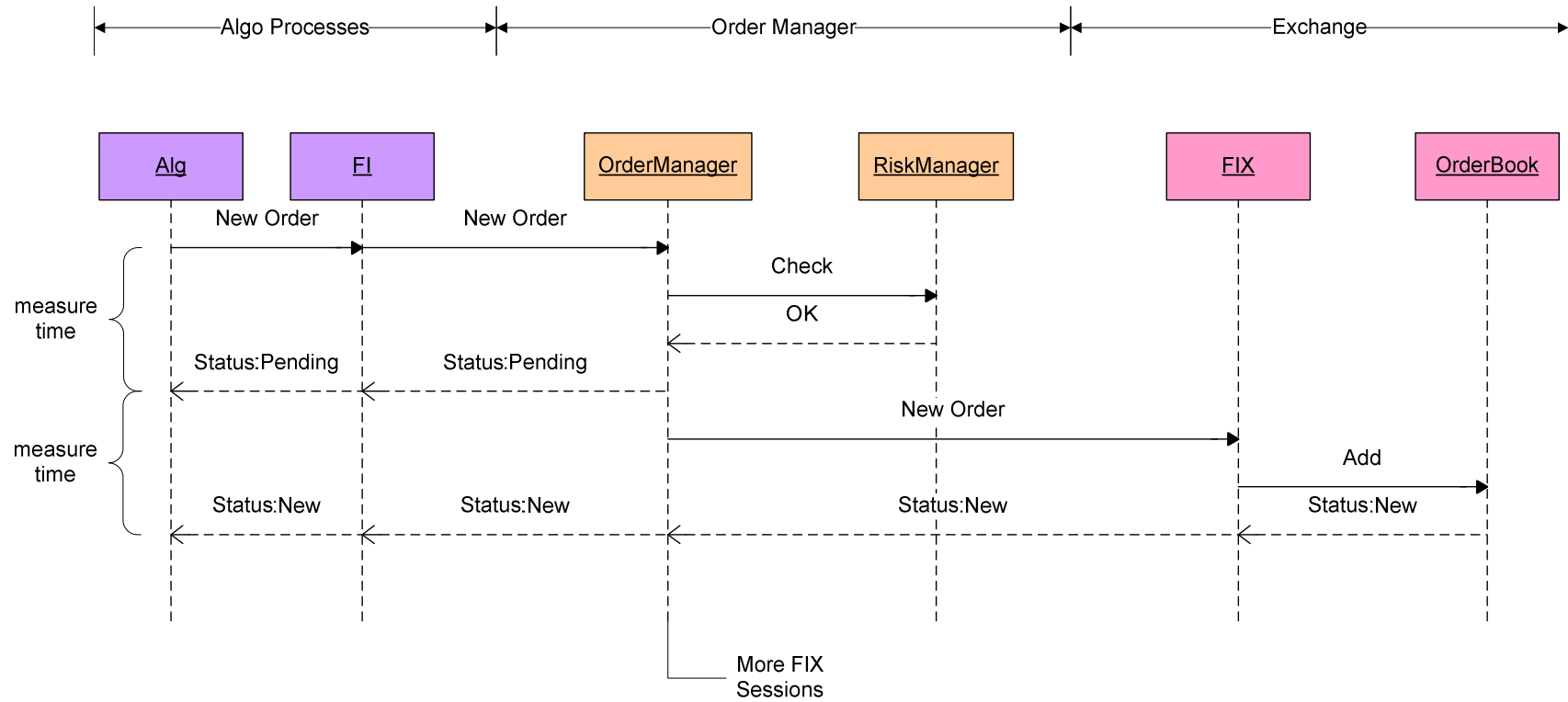
Sample Execution Report

- 8=FIX.4.2 | 9=67 | 35=8 | 49=NASDAQ | 56=ABC | 11=CLIENT_ORDER_ID | 52=20091123-18:30:00.000 | 20=3 | 150=E | 39=E | 55=AAPL | 167=CS | 54=1 | 38=15 | 40=2 | 44=15 | 58=NASDAQ SAMPLE | 59=0 | 47=C | 32=0 | 31=0 | 151=15 | 14=0 | 6=0 | 10=102 |

Market Data

- No market standards so feed handlers need to be written for every trading venue.
- Generally variable length ASCII bytes delivered on a protocol that guarantees delivery (SoupTCP, MouldUDP)

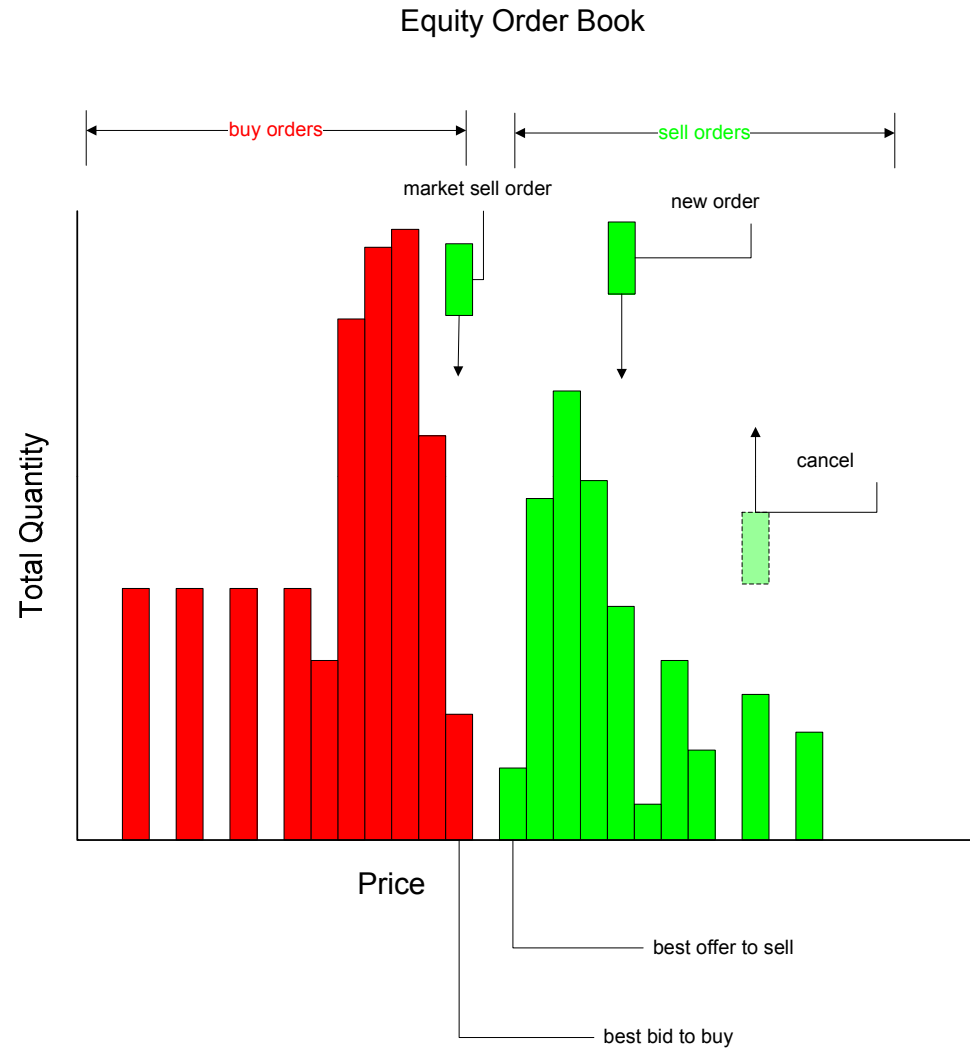
New Order Sequence



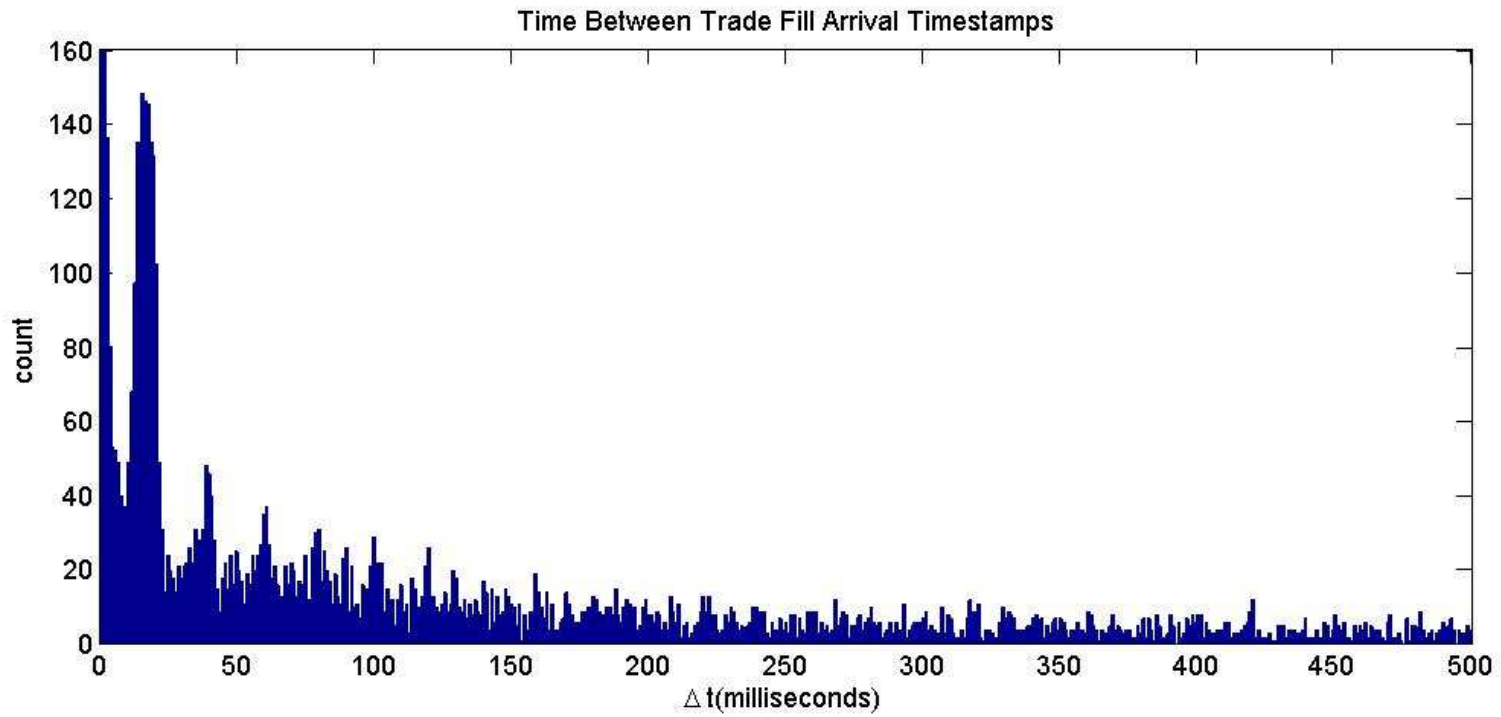
Market Data Latency

- Anatomy of the Order Book
- Order Book Visualizer
- Examples

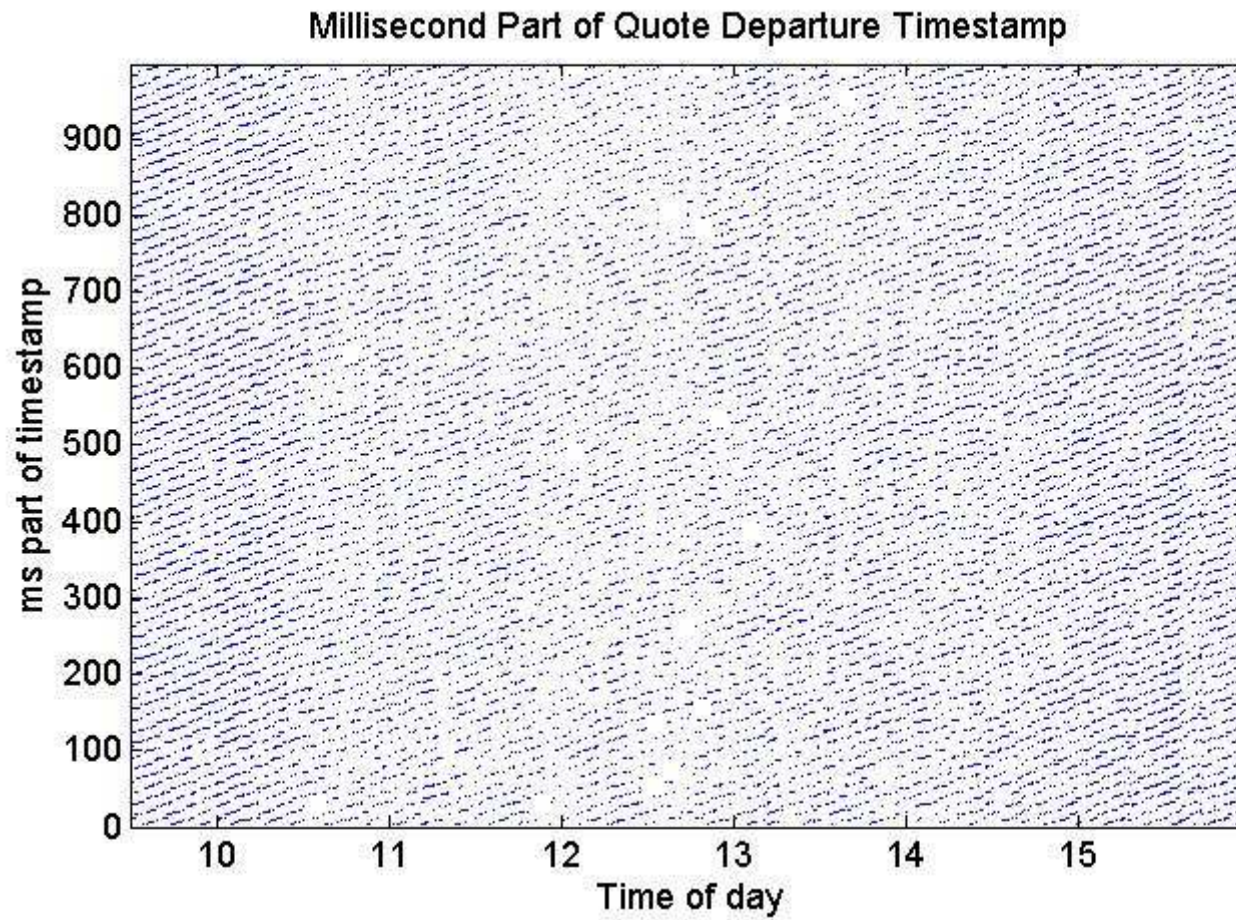
Order Book



Banding



Banding



Advanced Technologies

Field Programmable Gate Arrays (FPGA)

- Ideal for feed handlers and bulk analytics
- high throughput, low latency
- Advertised as being capable of processing millions of messages per second with 10-25 microsecond latency
- must be programmed specifically for each exchange
- can be purchased for many venues ([Celoxica](#) , [Exegy](#))

Graphical Processing Units (GPU)

- 100s of processing cores yield a massively parallel programming platform
- Growing support for computational sciences including finance.
- Initially programmed using C with some extensions but the support for other languages is growing.
- Ideal for things like options market making where you have to price a matrix of option values for every change in price

InfiniBand Architecture

- Switched fabric
- Interconnect between processing nodes and I/O devices
- Direct memory access configuration can reduce latency to ~ 1 microsecond

Further Reading

Online references

NASDAQ:Level II Specification

- <http://www.nasdaqtrader.com/content/technicalsupport/specifications/dataproducts/level2spec.pdf>

SoupTCP

- <http://www.nasdaqtrader.com/content/technicalsupport/specifications/TradingProducts/soup.pdf>

MoldUDP

- <http://www.nasdaqtrader.com/content/technicalsupport/specifications/dataproducts/moldudp.pdf>

FIX Protocol

- <http://www.onixs.biz/tools/fixdictionary/index.php>
- http://en.wikipedia.org/wiki/Financial_Information_eXchange

Low latency technologies (<http://www.solacesystems.com/tag/fpgas>)

- GPU assist for algo and Monte Carlo simulations ([NVIDIA](#))
- FPGA-based feed handlers ([Celoxica](#) , [Exegy](#), [Red Line](#))
- Network processor & FPGA-driven messaging ([Solace Systems](#), [TIBCO](#))
- Network acceleration technologies ([Cisco](#), [NetEffect](#), [Arastra](#))
- Analytics ([XtremeData](#))

Nagel's Algorithm

- <http://tools.ietf.org/rfc/rfc896.txt>
- http://en.wikipedia.org/wiki/Nagle%27s_algorithm

Exchange Technology

- http://www.cinnober.com/files/A_Cinnober_whitepaper_on_latency_Oct2009_0.pdf