

# How to do 100K+ TPS at less than 1ms latency

**Martin Thompson & Michael Barker**

**QCon SF 2010**

**LMAX**



# Agenda

- **Context Setting**
- **Tips for high performance computing (HPC)**
- **What is possible on a single thread???**
- **New pattern for contended HPC**
- **Q & A**

# **Who/What is LMAX?**

- **The London Multi-Asset Exchange**
- **Spin-off from Betfair into retail finance**
- **Access the wholesale financial markets on equal terms for retail traders**
- **We aim to build the highest performance financial exchange in the world**

# What is Extreme Transaction Processing (XTP)?



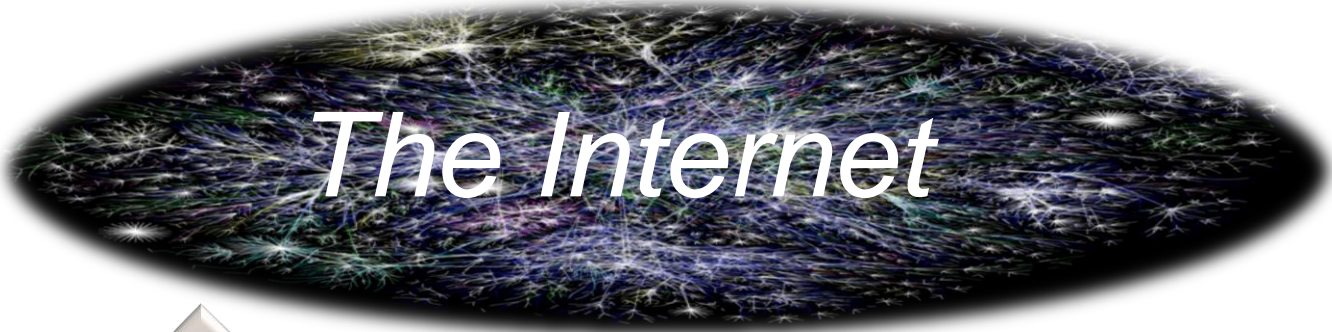
The Betfair Experience



v

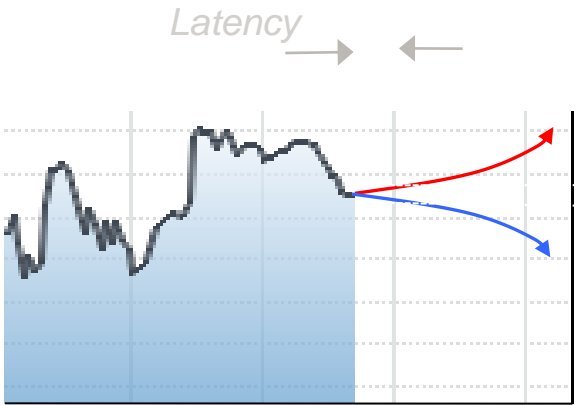


# What is Extreme Transaction Processing (XTP)?

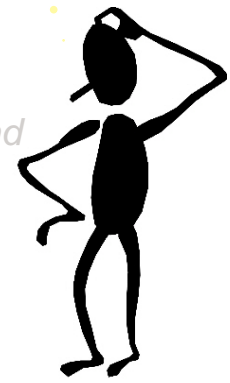


GBP / USD					
172					
171					
170	1.6028	- 1.33%	10.201	N/A	10.201
169	1.6011	+ 1.55%	13.293	N/A	20.169
168	0.8602	- 0.81%	N/A	20.169	N/A
167	0.8605	- 0.11%	20.169	N/A	N/A
166	1.238	+ 0.11%	N/A	N/A	1.662
165	1.6577	+ 1.12%	N/A	1.662	10.201
164	0.873	+ 3.23%	1.662	10.201	0.873
163	0.1150	- 2.14%	10.201	0.873	
162	0.1123	+ 2.18%	0.873		

The LMAX Model



Risky!



# How not to solve this problem



# Phasers or Disruptors?



# Tips for high performance computing

1. Show good “Mechanical Sympathy”
2. Keep the working set In-Memory
3. Write cache friendly code
4. Write clean compact code
5. Invest in modelling your domain
6. Take the right approach to concurrency



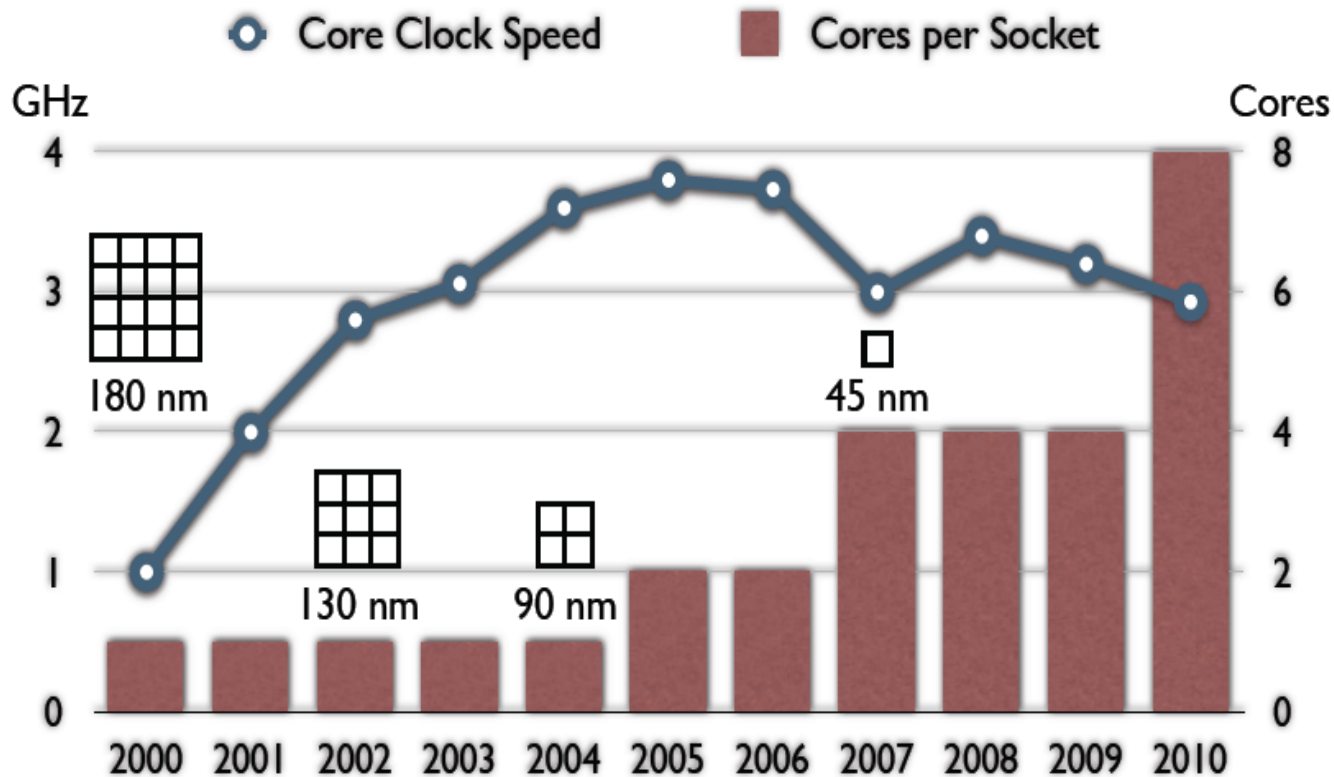
# 1. Mechanical Sympathy – 1 of 2

## Memory

- Latency not significantly changed
- Massive bandwidth increase
- 144GB in a commodity machine

## CPUs

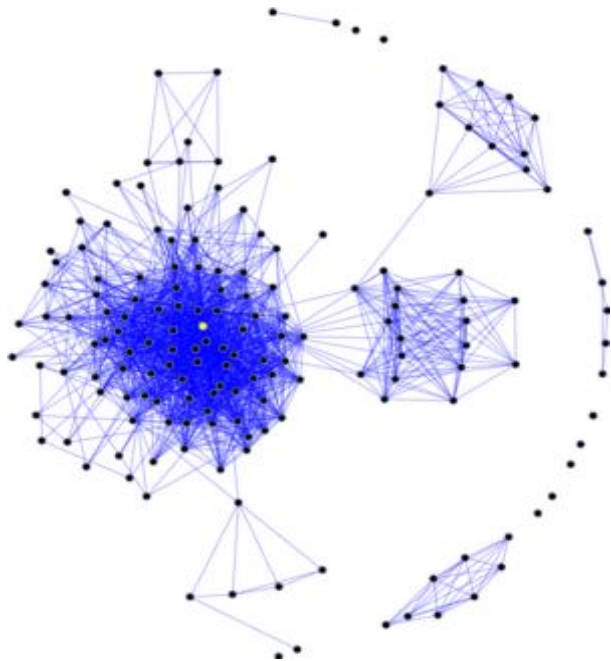
- The GHz race is over
- Multi core
- Bigger smarter caches



# 1. Mechanical Sympathy – 2 of 2

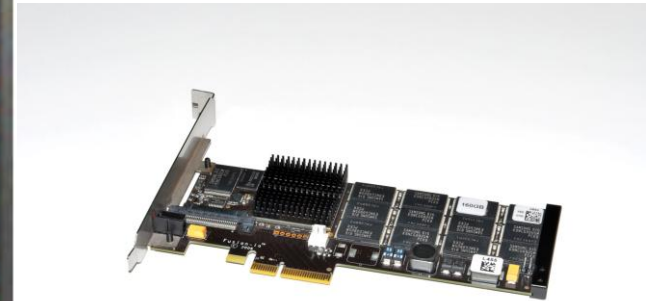
## Networks

- Sub 10 microseconds for local hop
- Wide area bandwidth is cheap
- 10GigE is now a commodity
- Multi-cast is getting traction



## Storage

- Disk is the new tape! Fast for sequential access
- SSDs for random threaded access
- PCI-e connected storage



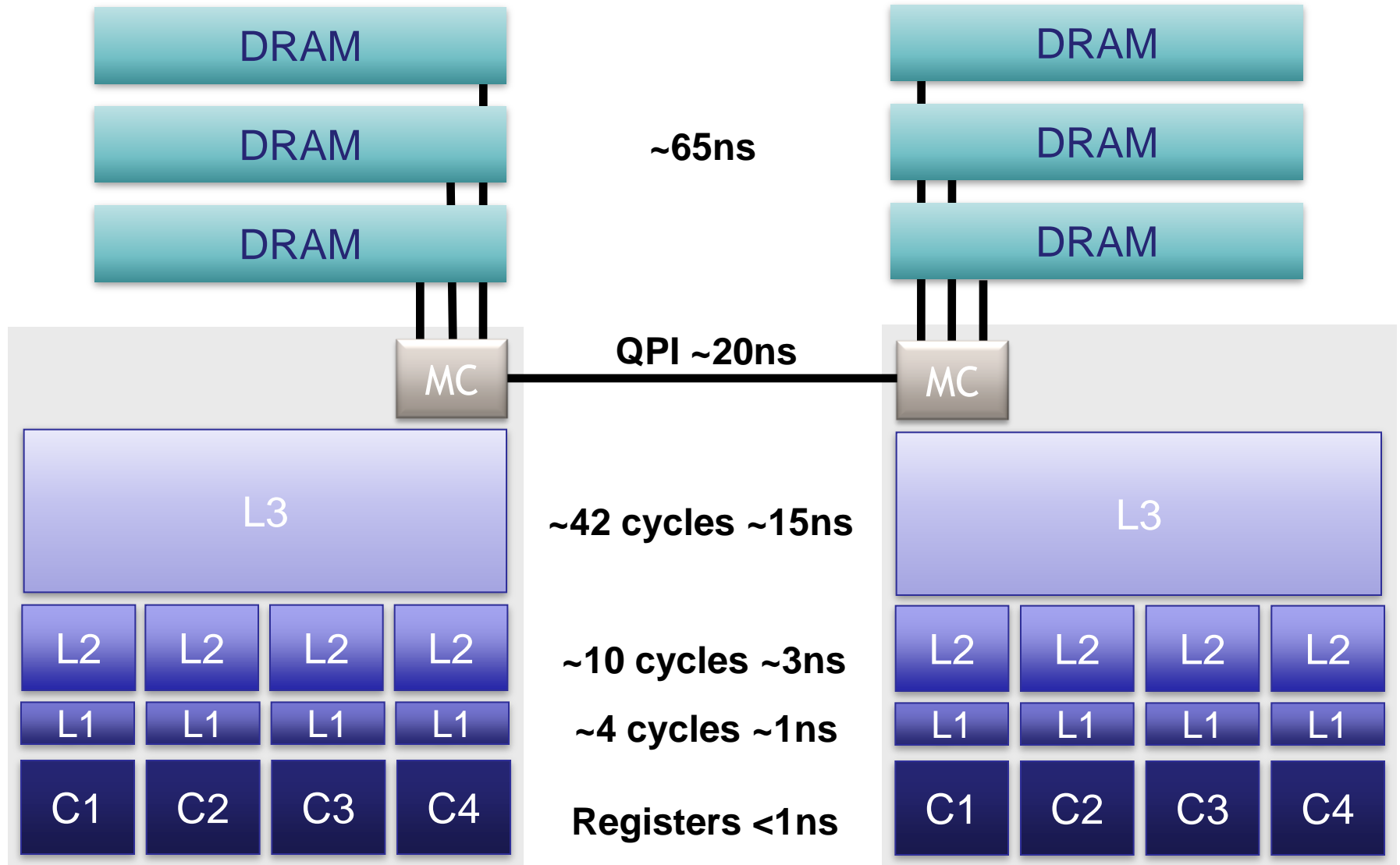
## 2. Keep the working set In-Memory

Does it feel awkward working with data remote from your address space?

- Keep data and behaviour co-located
- Affords rich interaction at low latency
- Enabled by 64-bit addressing



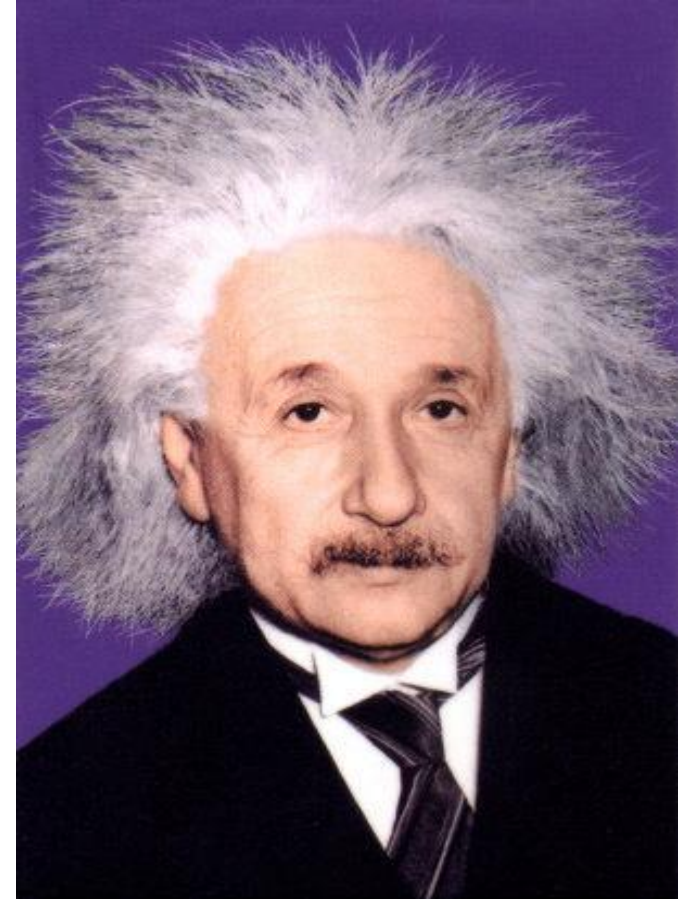
### 3. Write cache friendly code



## 4. Write clean compact code

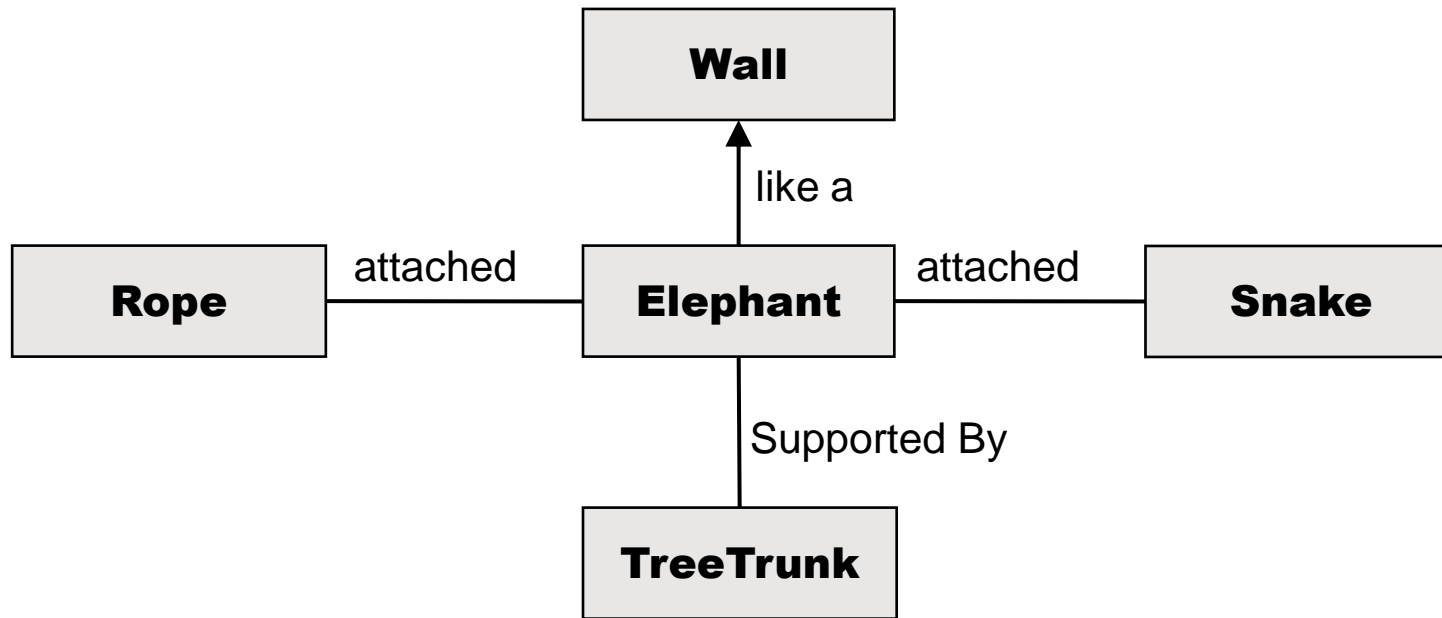
*"Any intelligent fool can make things bigger, more complex, and more violent. It takes a touch of genius -- and a lot of courage -- to move in the opposite direction."*

- Hotspot likes small compact methods
- CPU pipelines stall if they cannot predict branches
- If your code is complex you do not properly understand the problem domain
- Nothing in the world is truly complex other than Tax Law



## 5. Invest in modelling your domain

*Model of an elephant based on blind men touching one part each*



- Single responsibility – One class one thing, one method one thing, etc.
- Know your data structures and cardinality of relationships
- Let the relationships do the work

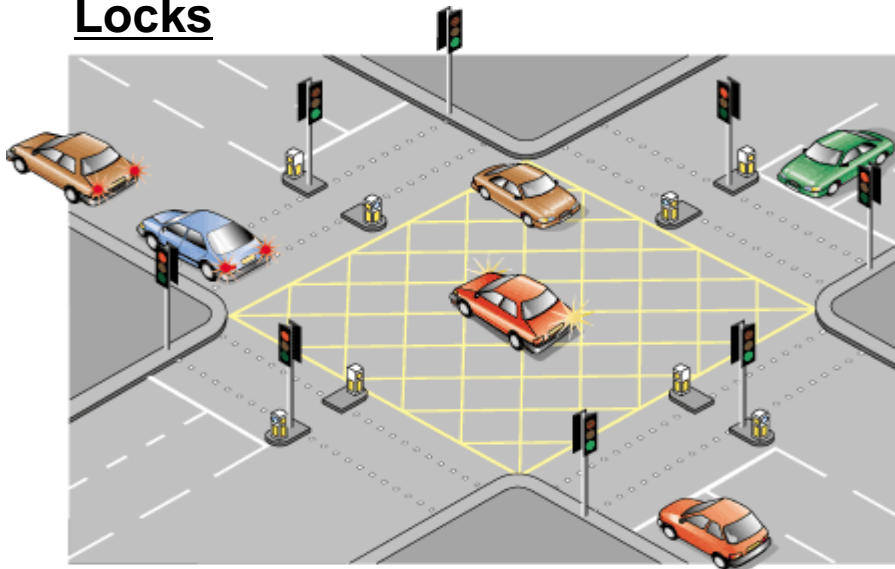
## 6. Take the right approach to concurrency

Concurrent programming is about 2 things:

**Mutual Exclusion:** Protect access to contended resources

**Visibility of Changes:** Make the result public in the correct order

### Locks



- Context switch to the kernel
- Can always make progress
- Difficult to get right

### Atomic/CAS Instructions



- Atomic read-modify-write primitives
- Happen in user space
- Very difficult to get right!

# What is possible when you get this stuff right?

On a single thread you have ~3 billion instructions per second to play with:

## 10K+ TPS

- If you don't do anything too stupid

## 100K+ TPS

- With well organised clean code and standard libraries

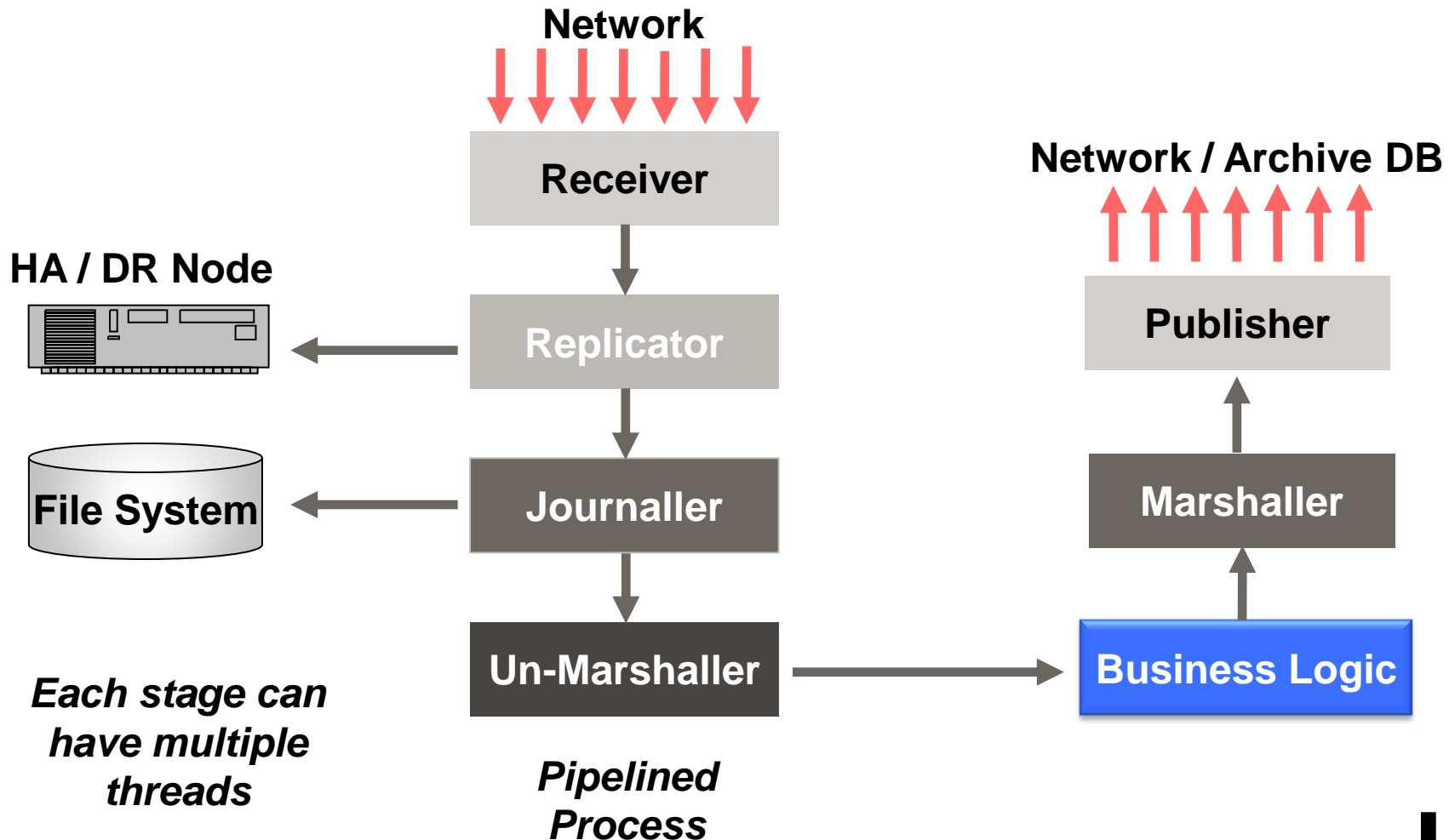
## 1m+ TPS

- With custom cache friendly collections
  - Good performance tests
  - Controlled garbage creation
  - Very well modelled domain
- BTW writing good performance tests is often harder than the target code!!!



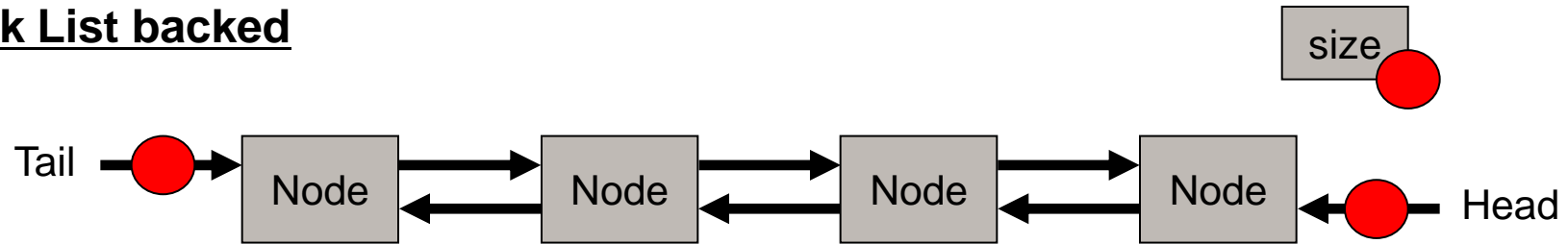
# How to address the other non-functional concerns?

- With a very fast business logic thread we need to feed it reliably
  - > Did we trick you into thinking we can avoid concurrent programming?



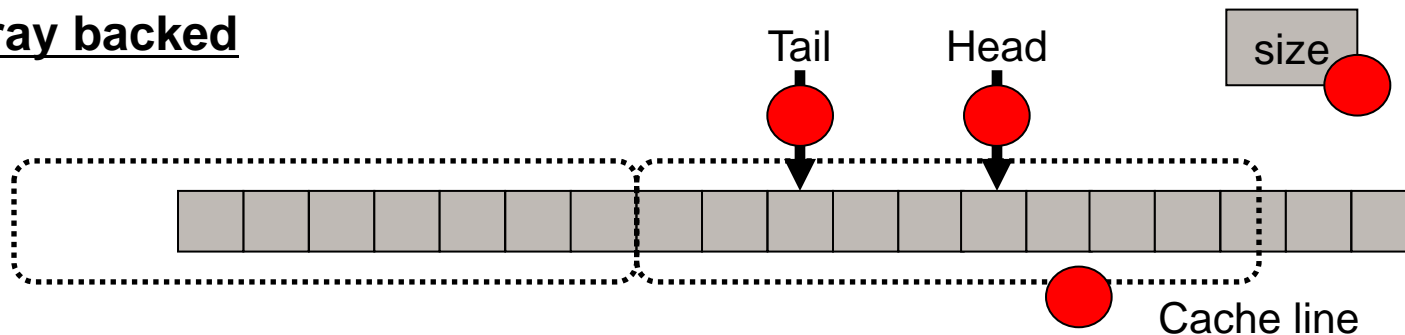
# Concurrent access to Queues – The Issues

## Link List backed



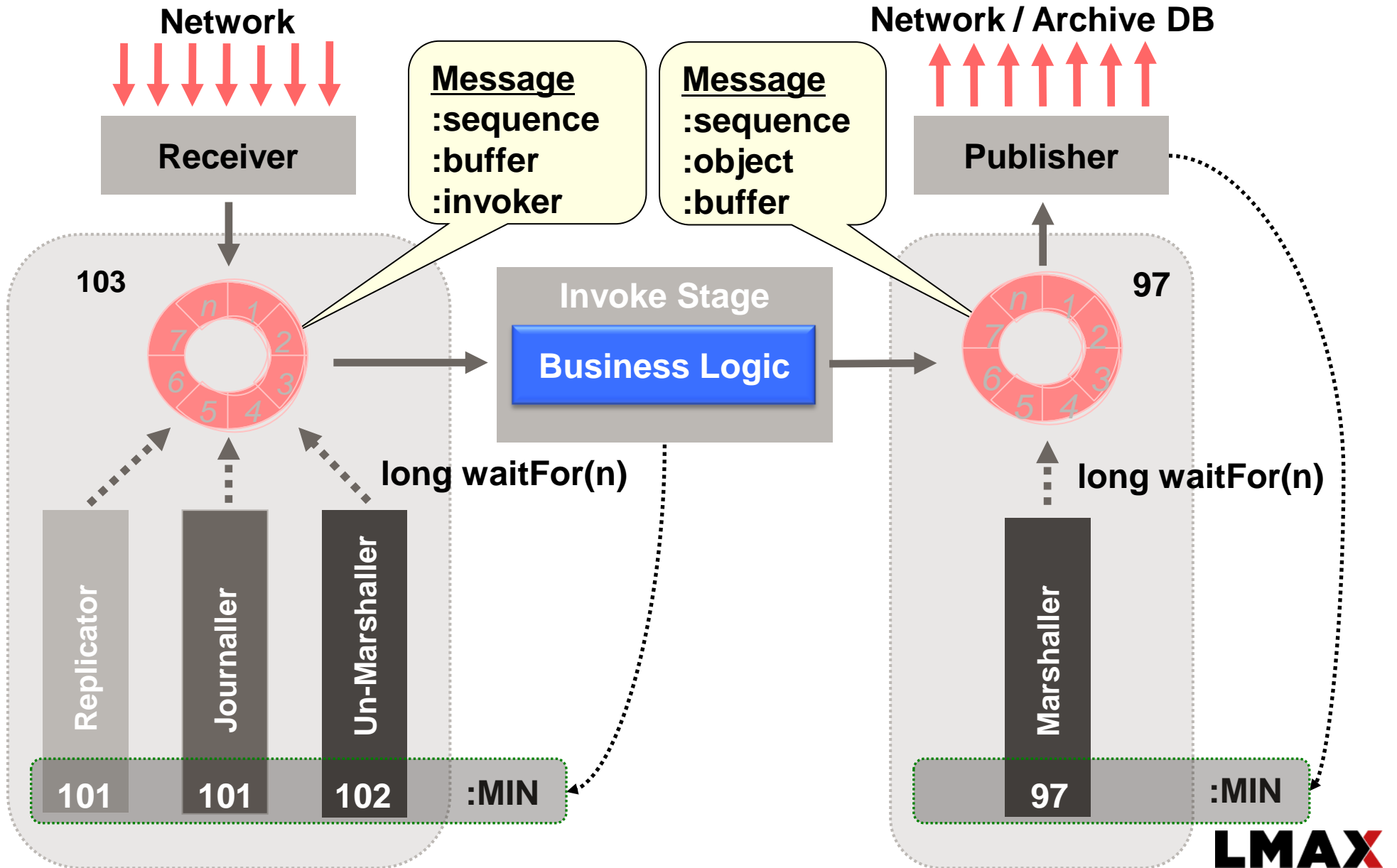
- Hard to limit size
- $O(n)$  access times if not head or tail
- Generates garbage which can be significant

## Array backed



- Cannot resize easily
- Difficult to get \*P \*C correct
- $O(1)$  access times for any slot and cache friendly

# Disruptor



# Quick Recap

- **Most developers have an incorrect view of hardware and what can be achieved on a single thread**
- **On modern processors a cache miss is your biggest cost**
- **Push concurrency into the infrastructure, and make it REALLY fast**
- **Once you have this, you have the world that OO programmers dream of:**
  - > Single threaded
  - > All in-memory
  - > Elegant model
  - > Testable code
  - > No infrastructure or integration worries

**Wrap up**

**Q & A**

**jobs@imax.com**