# Abstract

**Specialized hardware is giving architects new high-efficiency options to accelerate the WAN and avoid "long fat network" problems. This session will explore how network processors, FPGAs, flash storage, ultra capacitors, and other exotic silicon is increasing the capabilities and performance of WAN-based applications. Specific use cases include Distributed Message Routing, Web Data Streaming, Sensor Nets, and Active/Active Data Grid Replication.**

**Solace Systems**™

# Distributed Data Fabrics and Hardware WAN Optimization

Achieving 10X WAN Efficiency in Globally Distributed Applications
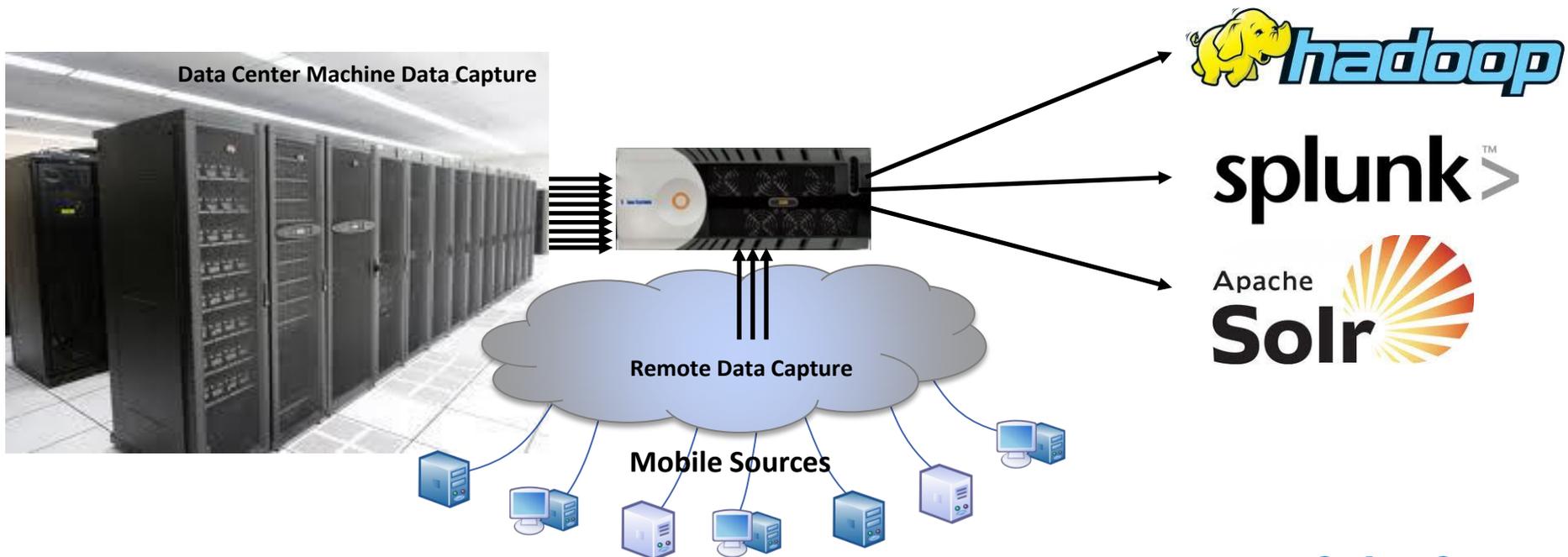
Hans Jespersen

Systems Engineer

hans.jespersen@solacesystems.com

# Agenda

- Introduce the use case
- TCP/IP and the Long Fat Network Problem
- Technology & Industry Trends
- How do traditional WANop solutions help (HW & SW)
- What isn't addressed with network layer WANop
- Message Brokers and application specific WANop
- Advanced Silicon and Exotic Hardware
- Benchmarking Performance
- Q&A

Solace Systems™

# Real-time Streaming Big Data

*Need is for efficiently collecting, aggregating and moving large amounts of streaming machine generated data from **multiple sources** to **multiple data stores** across **multiple locations**.*



Data Center Machine Data Capture

Remote Data Capture
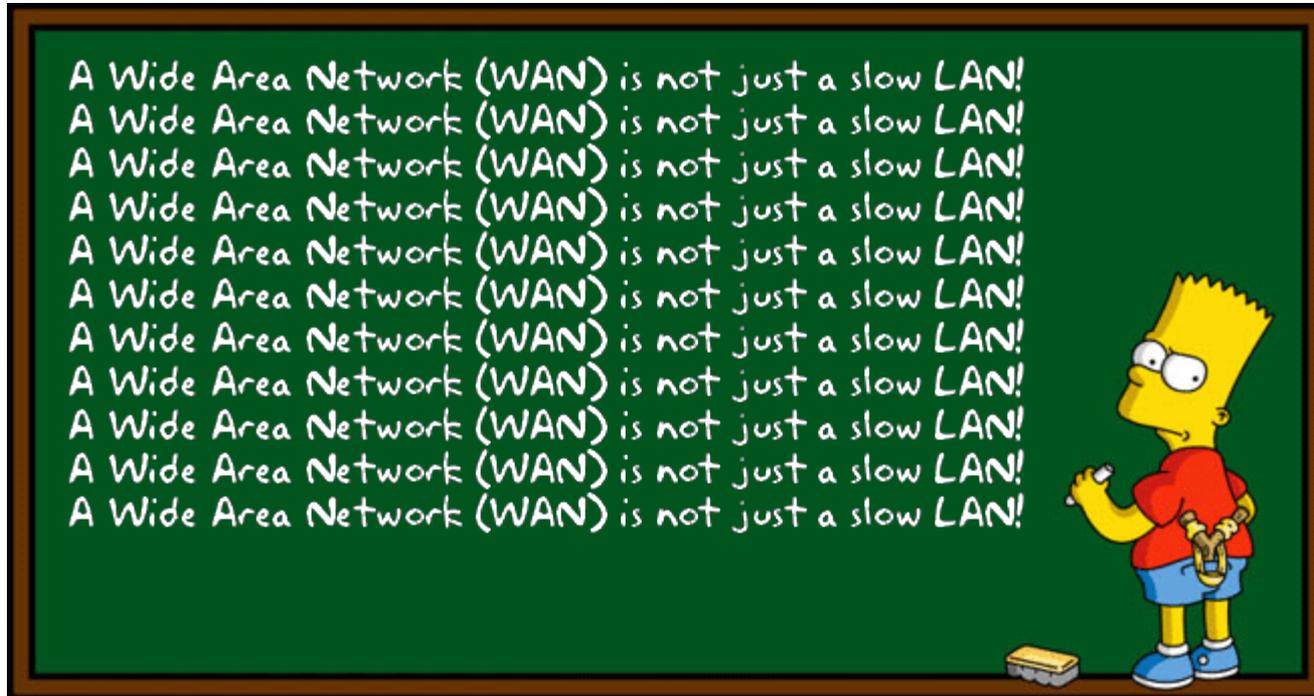
Mobile Sources

**SOlace Systems**™

# Old Faithful

**TCP Header**

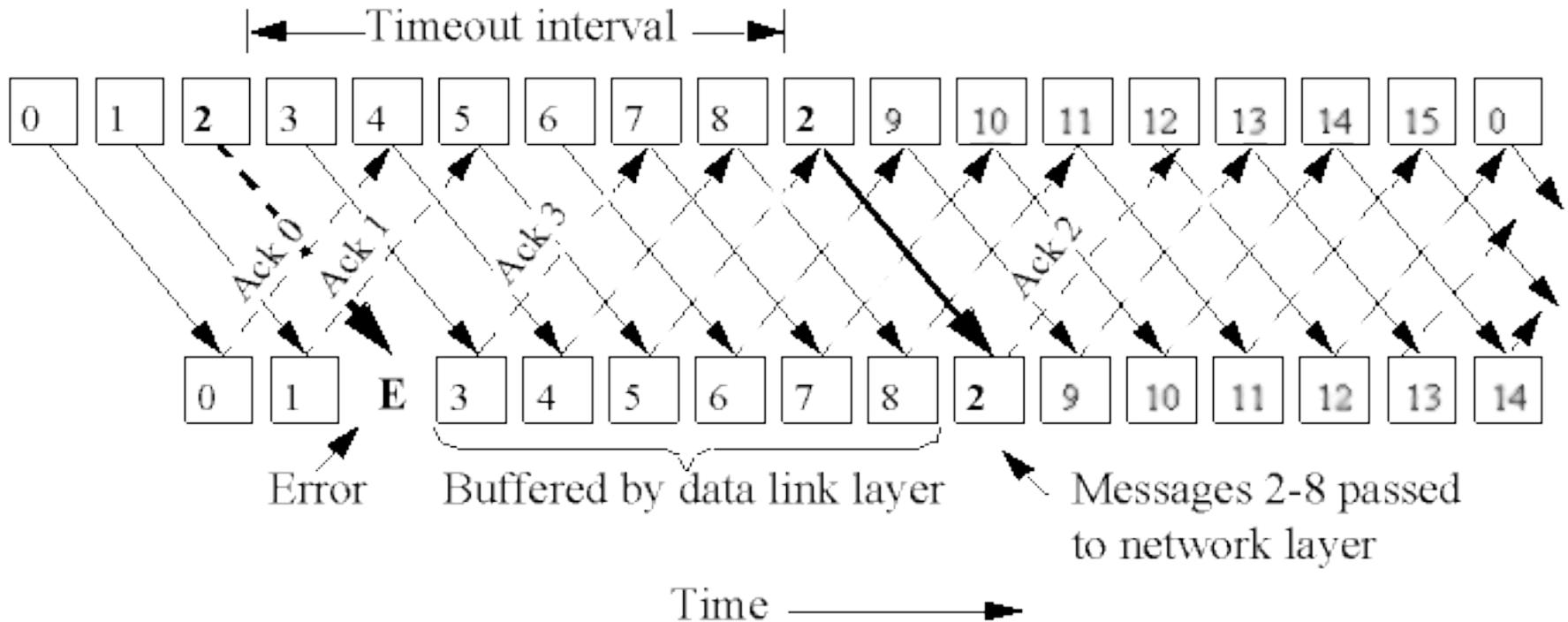| Offsets | Octet | 0 | | | | | | | | 1 | | | | | | | | 2 | | | | | | | | 3 | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Octet | Bit | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
| 0 | 0 | Source port | | | | | | | | | | | | | | | Destination port | | | | | | | | | | | | | | | |
| 4 | 32 | Sequence number | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 8 | 64 | Acknowledgment number (if ACK set) | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 12 | 96 | Data offset | | | | Reserved 0 0 0 | | | NS | CWR | ECE | URG | ACK | PSH | RST | SYN | FIN | Window Size | | | | | | | | | | | | | | | |
| 16 | 128 | Checksum | | | | | | | | | | | | | | | | Urgent pointer (if URG set) | | | | | | | | | | | | | | | |
| 20 ... | 160 ... | Options (if Data Offset > 5, padded at the end with "0" bytes if necessary) ... | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |

**SOlace Systems™**

# What do we get?

- **Reliable**

- **Ordered delivery (of a stream of octets)**

- **Error-free data transfer**

- **Flow control**

- **Congestion control**

**Solace Systems**™

# Everything comes with a price

**Solace Systems™**

# The LFN (Elephant) in the room

# Throughput != Bandwidth

○ **Bandwidth**

○ **Latency (RTT)**

○ **Error Rate (Loss)**

○ **Handy online calculators of effective throughput**
  - http://www.silver-peak.com/calculator/

Solace Systems™

# Why is this problem growing?

- Globalization
- Public Internet backbone
- RDBMS -> NoSQL, IMDG
- Rich Data Types
- Mobile Apps
- Client Side Data
- DR and BCP

1. Bandwidth
2. Latency (RTT)
3. Error Rate (Loss)

# Traditional WAN optimization techniques

- **Deduplication**
- **Compression**
- **Latency optimization**
- **Caching/proxy**
- **Forward error correction**
- **Protocol spoofing**
- **Traffic shaping**
- **Equalizing / Prioritizing**
- **Connection limiting**
- **Rate limiting**

**Bi-directional Message Streaming**

**Hardware Compression**

**Multiple Parallel Connections**

**Solace Systems**™

# Offloading the IP Stack to Hardware



**Cavium Octeon II**
- **32 core MIPS64 Processor**

**Pre-built application acceleration engines**
- **Packet Processing**
- **Encryption/Decryption**
- **Deep Packet Inspection (RegEx)**
- **Compression/decompression**
- **De-duplication**
- **RAID**

**Millions of concurrent connections**

**Solace Systems**™

# Improving the Speed of GoldenGate Synchronization

- **Real customer results**
  - Tested synchronization over a 622 Mbps link between New York and London with 75 ms round trip time
  - And over a 45 Mbps link between New York and Tokyo with 175 round trip time

- **Solace 18x Faster**



| | New York to London 622 Mbps link, 75 ms RTT | | | New York to Tokyo 45 Mbps link, 176 ms RTT | |
|---|---|---|---|---|---|
| Operations Per Second | 1,092 | 18,761 | | 401 | 7,238 |
| Transactions Per Second | 82 | 1,409 | | 30 | 543 |

**Solace Systems™**

# Back to the use case

**Transactions != Packets**

**Database Records != Packets**

**Objects != Packets**

**Solace Systems**™

# The Modern Information Distribution Fabric

# Messaging Middleware Value

- Producer/Consumer Decoupling

- Disconnected Operation

- Location Independence

- Multipoint Delivery

- Advanced Filtering & Routing

- Message-based Granularity



*Messaging layer on top of
your IP network, so you can make
messaging a shared optimized service.*

**S⦾lace Systems™**

# FPGA



**Xilinx Virtex-7 2000T FPGA**
- **More than twice the capacity and bandwidth offered by the largest monolithic devices**
- **2 million logic cells (equivalent to 20 million ASIC gates)**
- **6.8 billion transistors**

# Horizontal scalability on a single chip



**Intel Westmere-EX**
- **2.6 billion transistors**
- **10 64-bit cores @ 2.4 GHz**
- **7,200 MIPS**
- **130 watts TDP**

**Xilinx Virtex-7 2000T FPGA**
- **6.8 billion transistors**
- **3,600 8-bit processors @ 100 MHz**
- **180,000 MIPS**
- **20 watts TDP**

**SOlace Systems**™

# Reliable Messaging

**Pure hardware solution**

- No operating system
- No context switching
- No interrupts
- No data copies

**10 million messages/second**

- Can be any combination, e.g. 5M in & 5M out, 2M in & 8M out



| Bulk Message Rate | Message Size (bytes) | Message Rate (msgs/sec) | User Payload Bandwidth (Mbps) | |
|---|---|---|---|---|
| | 100 | 5,930,000 | 4,744 | |
| | 500 | 2,080,000 | 8,320 | *10GigE Line Rate the is Limit* |
| | 1,000 | 1,080,000 | 8,640 | |
| | 12,000 | 92,000 | 8,832 | |
| | 30,000 | 34,000 | 8,160 | |

**Solace Systems**™

# In Memory Message Caching

## Non-persistent last-value cache can handle any payload

o Cache by number or timeframe

o Can run on appliance,
or as a 64-bit app on a Linux server

o Centralized management of all caches.

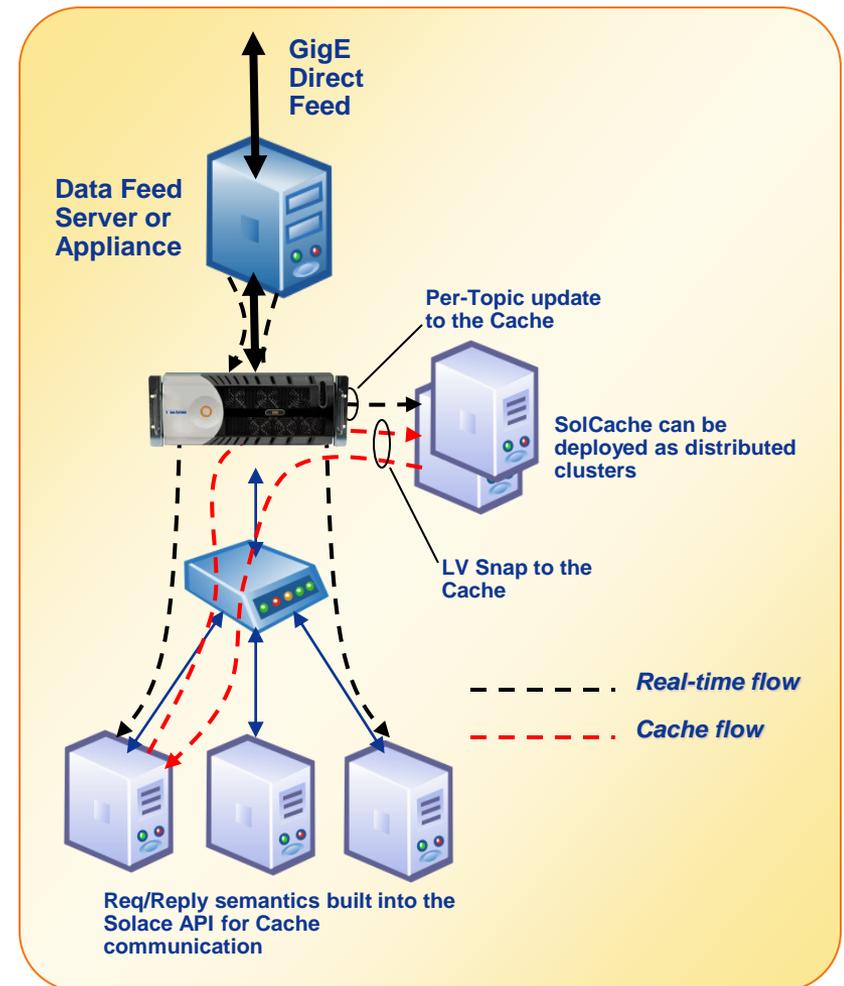o Clustering for load balancing and redundancy

o Support wildcard requests

o Request/reply with the cache, and control of synchronization of cache requests

o Topic names partitioned amongst instances to scale storage



**GigE Direct Feed**

**Data Feed Server or Appliance**

**Per-Topic update to the Cache**

**SolCache can be deployed as distributed clusters**

**LV Snap to the Cache**

**Real-time flow**

**Cache flow**

**Req/Reply semantics built into the Solace API for Cache communication**

**Solace Systems**™

# Cascading Cache



**New York**

**London**

Publisher of
lon/fx/gbp

Cache for
ny/fx/usd

Cache for
lon/fx/gbp

Cache for
lon/fx/gbp

**Solace Systems™**

# Incremental Updates

## Cache Contents for NY/EQ/JNPR

SYMBOL: JNPR

VENUE: NYSE

**LAST: 19.19**

**VOLUME: 31,870**

DAY LOW: 19.03

DAY HIGH: 19.21

52-WEEK LOW: 15.13

52-WEEK HIGH: 23.98

## Updated Cache Contents

SYMBOL: JNPR

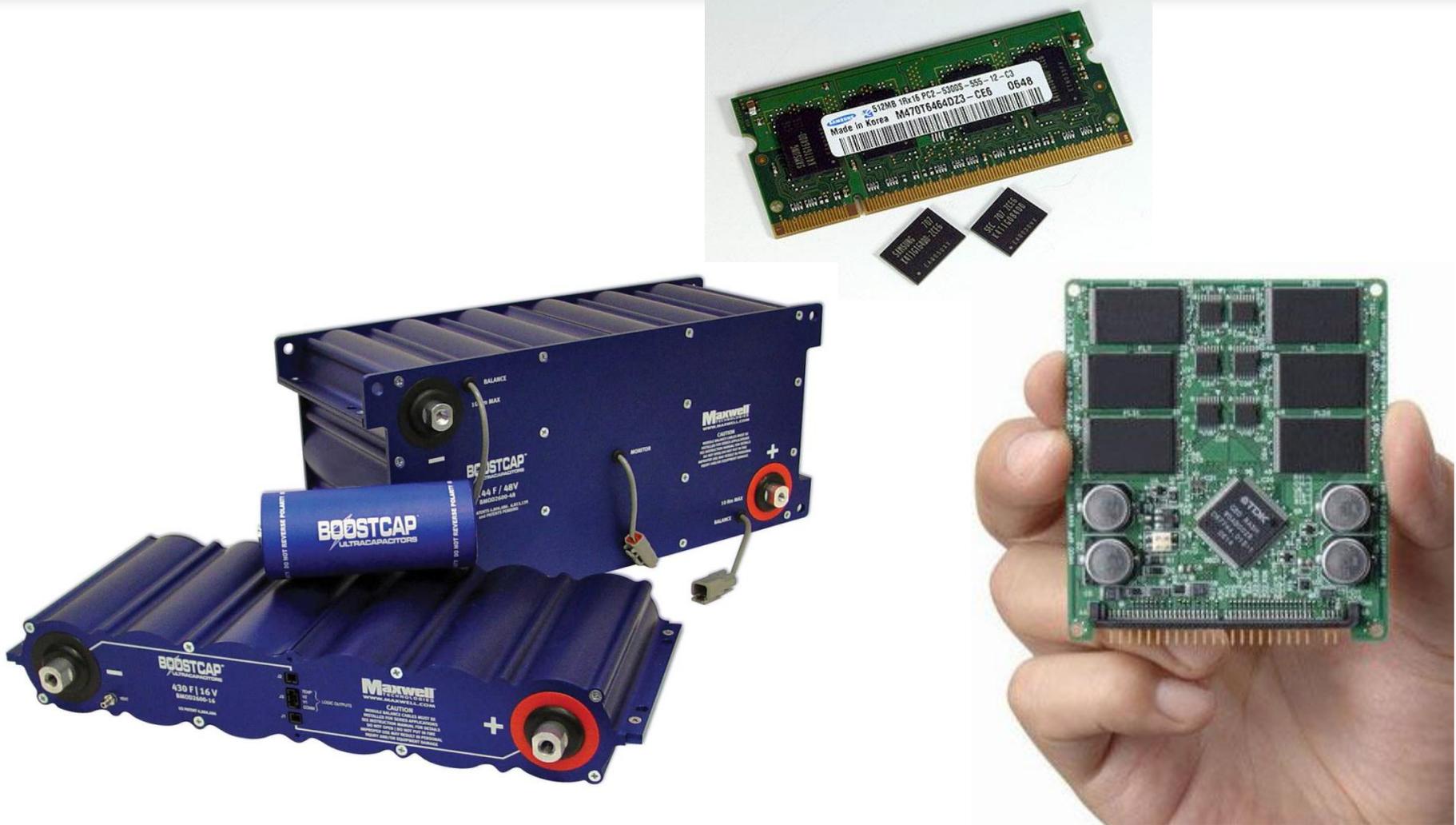VENUE: NYSE

**LAST: 19.19**

**VOLUME: 31,870**

DAY LOW: 19.03

DAY HIGH: 19.21

52-WEEK LOW: 15.13

52-WEEK HIGH: 23.98

**Solace Systems™**

# Hybrid Storage

# Fault Tolerant Clustering of Messaging Nodes

**Publisher**

**Subscriber**

**3a** Receipt acknowledged since message is guaranteed

**3b** If subscriber available, message is delivered immediately

**1**

Message persisted in on-board RAM

**2**

Message and state replicated to mate, which confirms receipt

**Redundant Mate**

- Connected to primary via two fiber links
- Same connectivity to L2 and storage as primary

**4**

If subscriber is slow or disconnected, their backlog is spooled to SAN, delivered as soon as client is able to receive

**SAN**

**Solace Systems™**

# Multiple Datacenters and Disaster Recovery

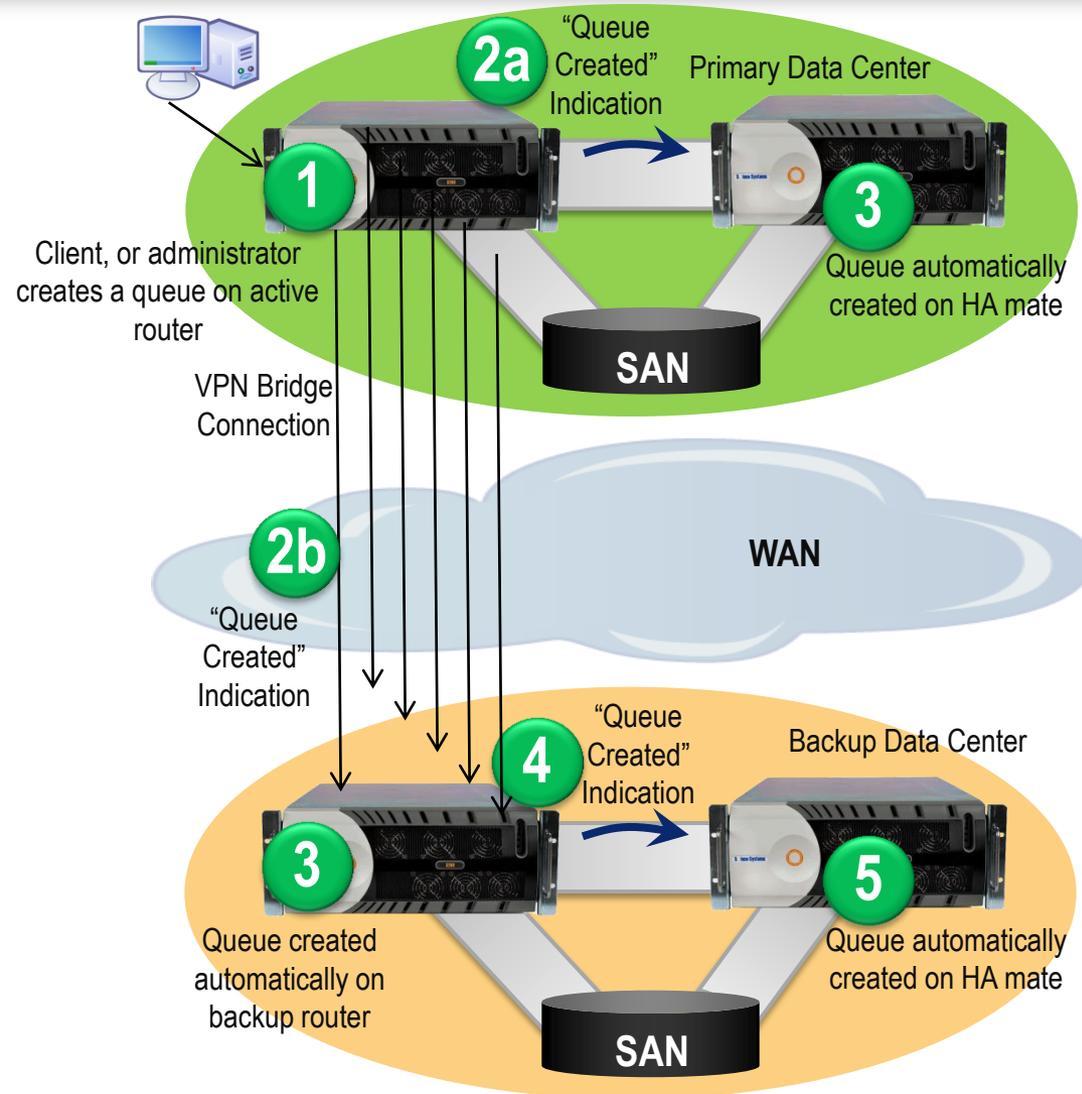- **Automatic replication of:**
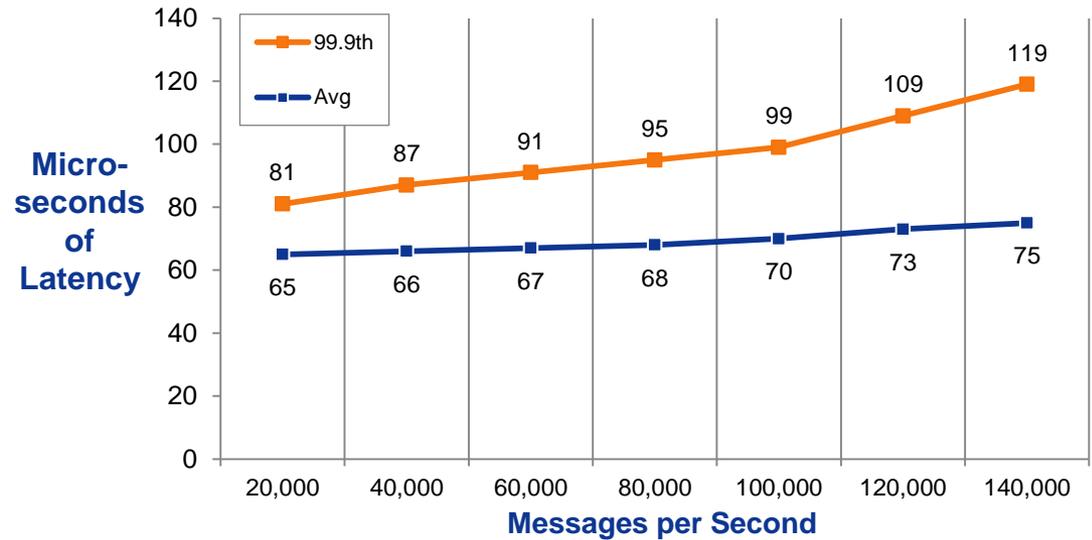  - Client-created endpoints
  - Configuration data
  - Transactional state

- **Needs Configurability for Sync & Asynchronous replication**



Primary Data Center

2a "Queue Created" Indication

1 Client, or administrator creates a queue on active router

3 Queue automatically created on HA mate

SAN

VPN Bridge Connection

2b "Queue Created" Indication

WAN

Backup Data Center

4 "Queue Created" Indication

3 Queue created automatically on backup router

5 Queue automatically created on HA mate
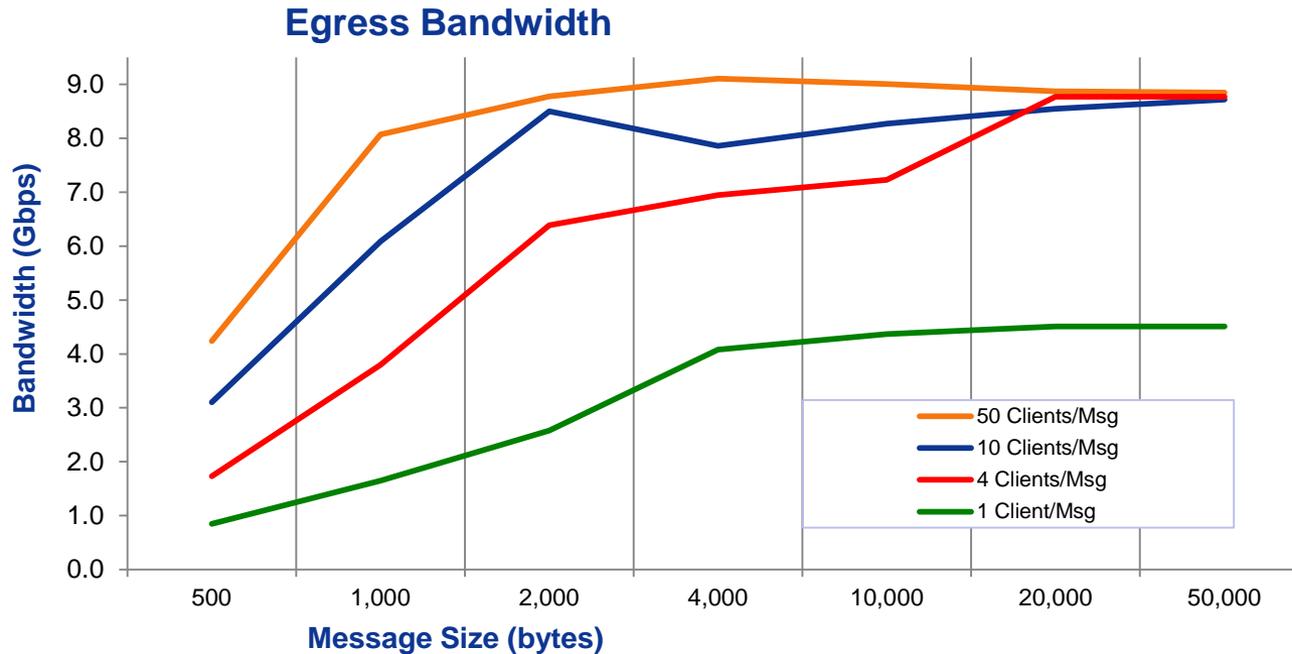
SAN

**Solace Systems**™

# Guaranteed Messaging;
# Store & Forward Performance

- **Failsafe w/o overhead of persisting every message to disk**

- **205K msgs/sec ingress and 205K msgs/sec egress**

- **Up to 4.5 Gbps of guaranteed messaging bandwidth**

- **Consistent latency even when servicing slow or recovering subscribers**

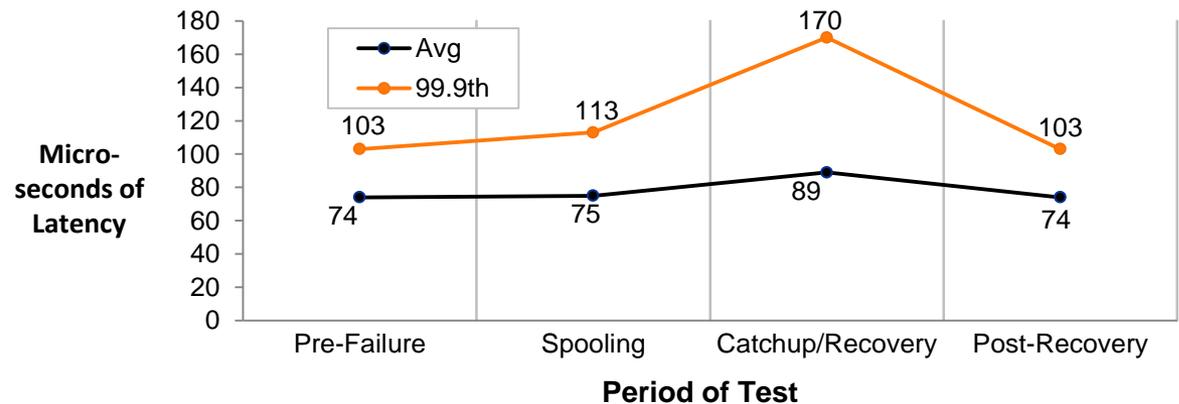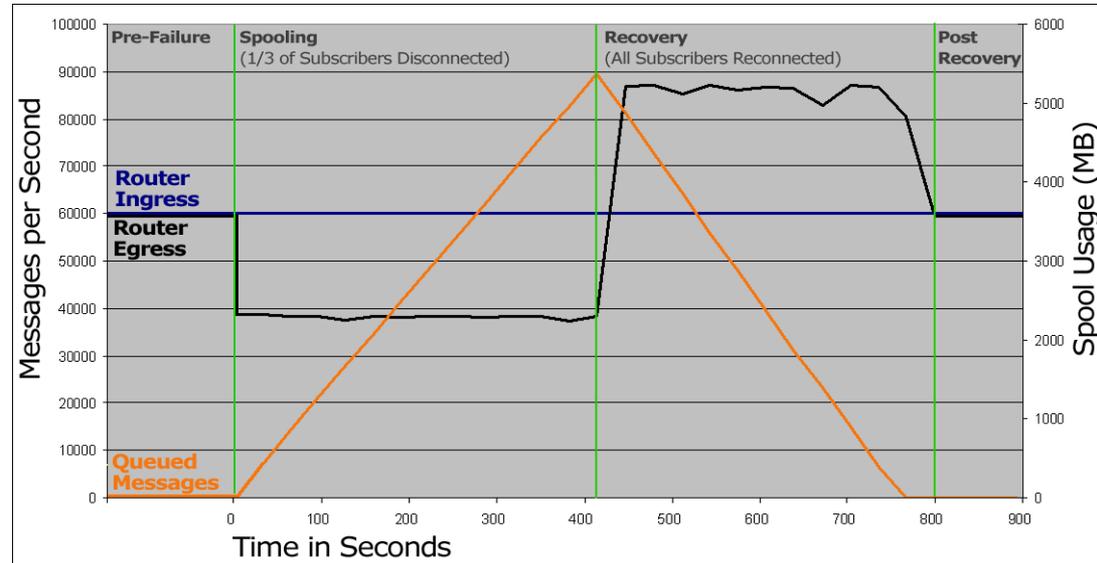| | Message Size (bytes) | Message Rate (msgs/sec) | User Payload Bandwidth (Mbps) |
|---|---|---|---|
| **Bulk Message Rate** | 100 | 206,400 | 165 |
| | 512 | 206,400 | 845 |
| | 1,024 | 202,000 | 1,655 |
| | 2,048 | 157,500 | 2,580 |
| | 4,096 | 124,400 | 4,076 |
| | 10,240 | 53,400 | 4,375 |
| | 20,480 | 27,500 | 4,506 |
| | 51,200 | 11,000 | 4,506 |

**Solace Systems**™

# Guaranteed Messaging;
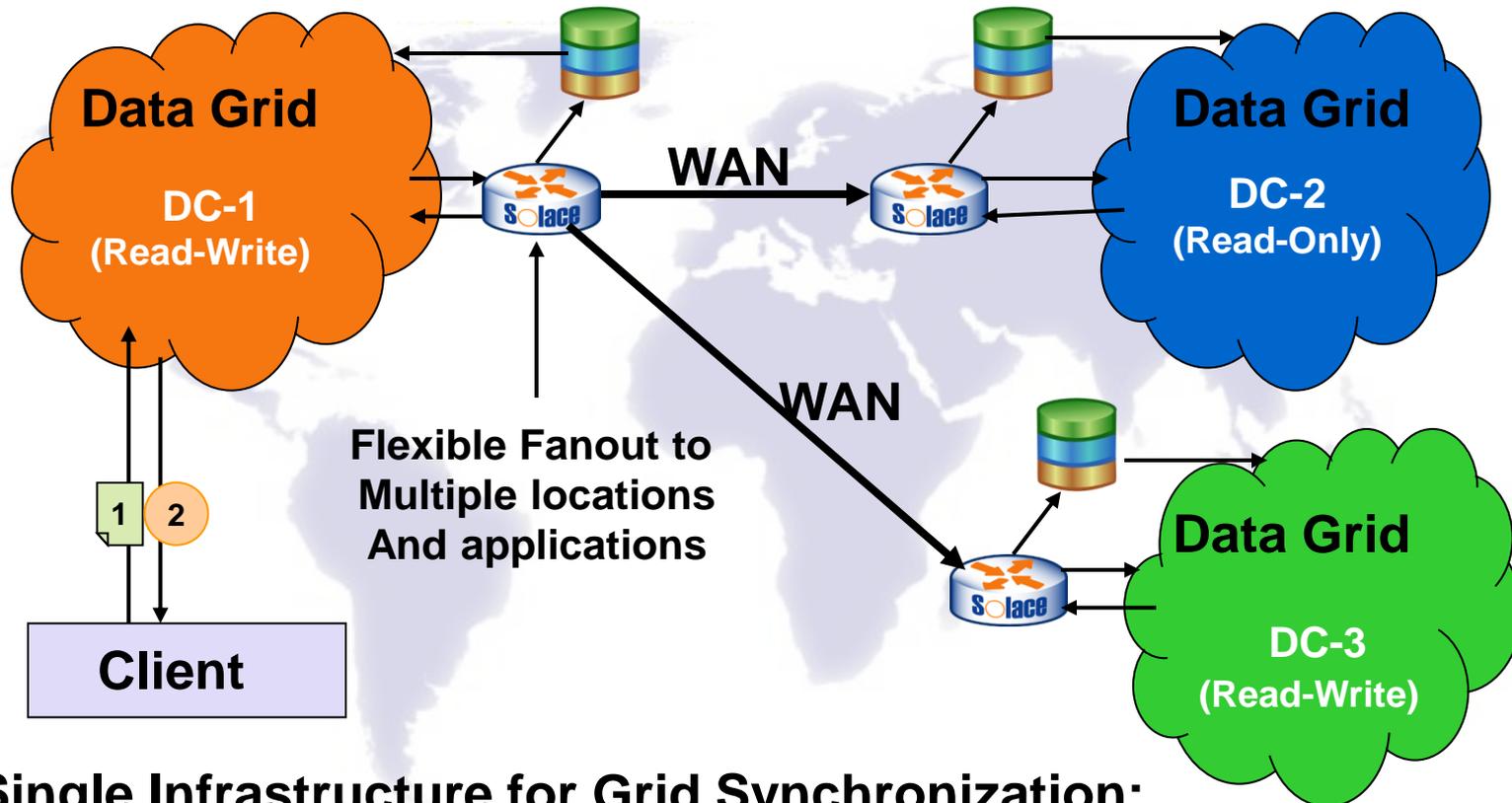# Fan-out performance

**Egress Bandwidth**



| Message Size (bytes) | Egress Rate, 1 client/msg (msgs/sec) | Egress Rate, 4 clients/msg (msgs/sec) | Egress Rate, 10 clients/msg (msgs/sec) | Egress Rate, 50 clients/msg (msgs/sec) |
|---|---|---|---|---|
| 512 | 206,400 | 422,000 | 756,000 | 1,035,000 |
| 1,024 | 202,000 | 464,000 | 744,000 | 985,000 |
| 2,048 | 157,500 | 390,000 | 519,000 | 536,000 |
| 4,096 | 124,400 | 212,000 | 250,000 | 278,000 |
| 10,240 | 53,400 | 88,300 | 101,000 | 110,000 |
| 20,480 | 27,500 | 53,500 | 52,200 | 54,150 |
| 51,200 | 11,000 | 21,400 | 21,300 | 21,600 |

**Solace Systems™**

# Offline or Slow Consumer Handling

- **Publisher rates not affected by slow/offline consumers**

- **Fast consumers not affected in rate or latency by slow/offline consumers**

- **Re-connected subscribers "catch up" without impacting other clients**

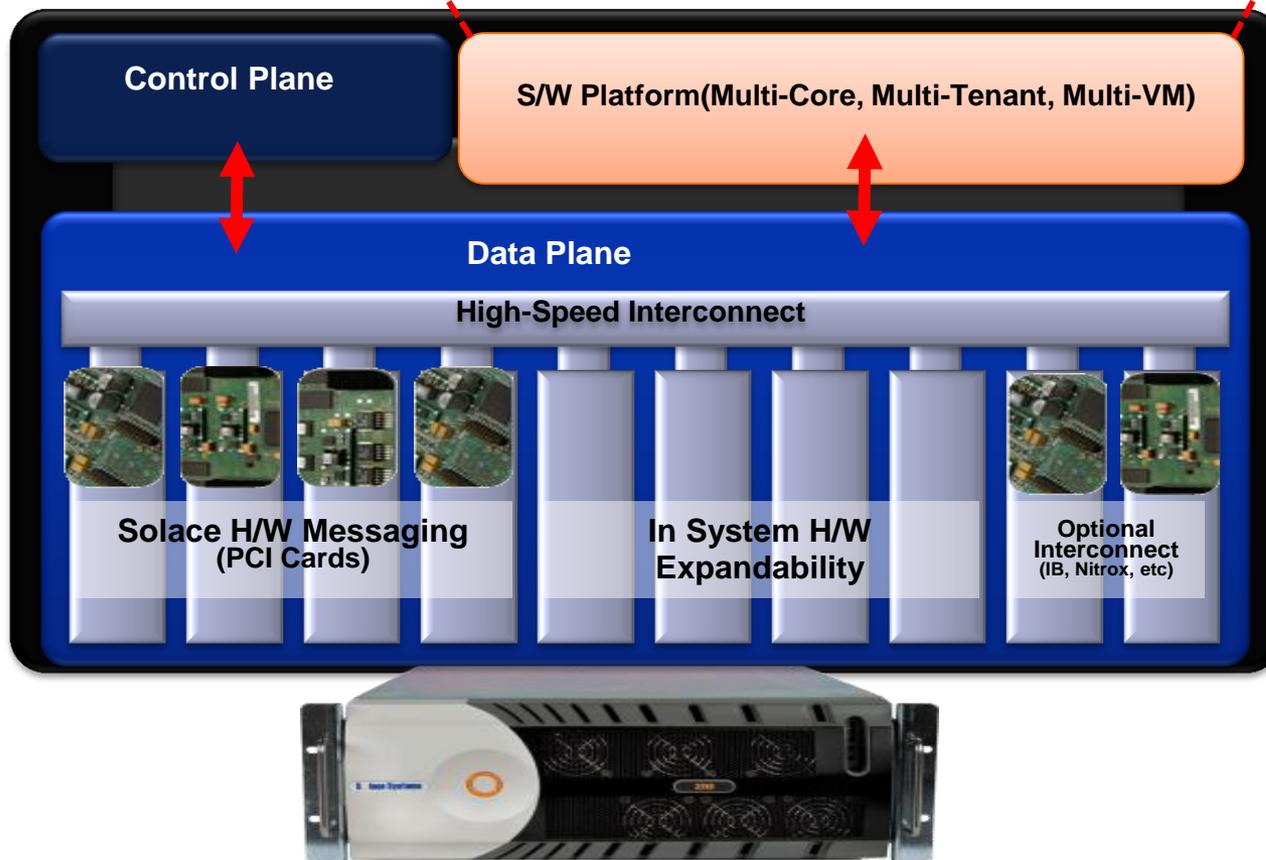- **Behavior & performance cannot be matched by software due to patented technology**

**SOlace Systems™**

# Optimized In-memory Grid Replication



**Data Grid**
DC-1
(Read-Write)

**WAN**

**Data Grid**
DC-2
(Read-Only)

**WAN**

**Data Grid**
DC-3
(Read-Write)

**Flexible Fanout to
Multiple locations
And applications**

**1** **2**

**Client**

## Single Infrastructure for Grid Synchronization:

- **inherently one-to-many, so can propagate
  to many other sites/instances – either locally or over the WAN**

- **Supports DR, Active/Passive, or Active/Active architectures**

**Solace Systems**™

# Solace as an Appliance Platform



**Control Plane**

**S/W Platform(Multi-Core, Multi-Tenant, Multi-VM)**

**Data Plane**

**High-Speed Interconnect**

**Solace H/W Messaging (PCI Cards)**

**In System H/W Expandability**

**Optional Interconnect (IB, Nitrox, etc)**

- Messaging all in hardware
- General purpose processors used to run 3rd party software that interacts seamlessly with hardware messaging internally
- Integration is easy with JMS
- Enables flexible solution options within an appliance

**APIs: JMS,C, .Net, Java, JavaScript, Flash, Silverlight, iOS, Node.js, Ruby, Python, etc.**

**Solace Systems™**

# The Modern Information Distribution Fabric



**High Volume Onboarding**

**Fast, Efficient WAN Sync**