

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™

Solace
Systems™

Distributed Data Fabrics and Hardware WAN Optimization

Achieving 10X WAN Efficiency in Globally
Distributed Applications

Hans Jespersen

Systems Engineer

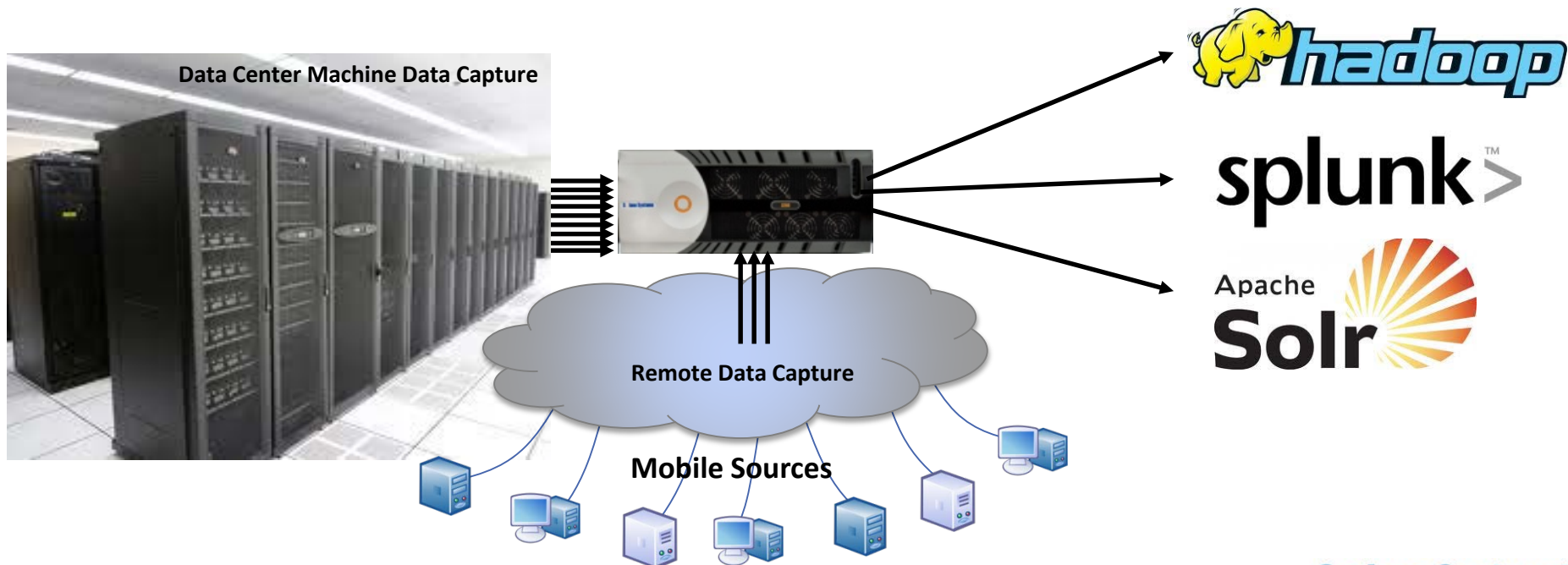
hans.jespersen@solacesystems.com

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**

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.***



Old Faithful

TCP Header

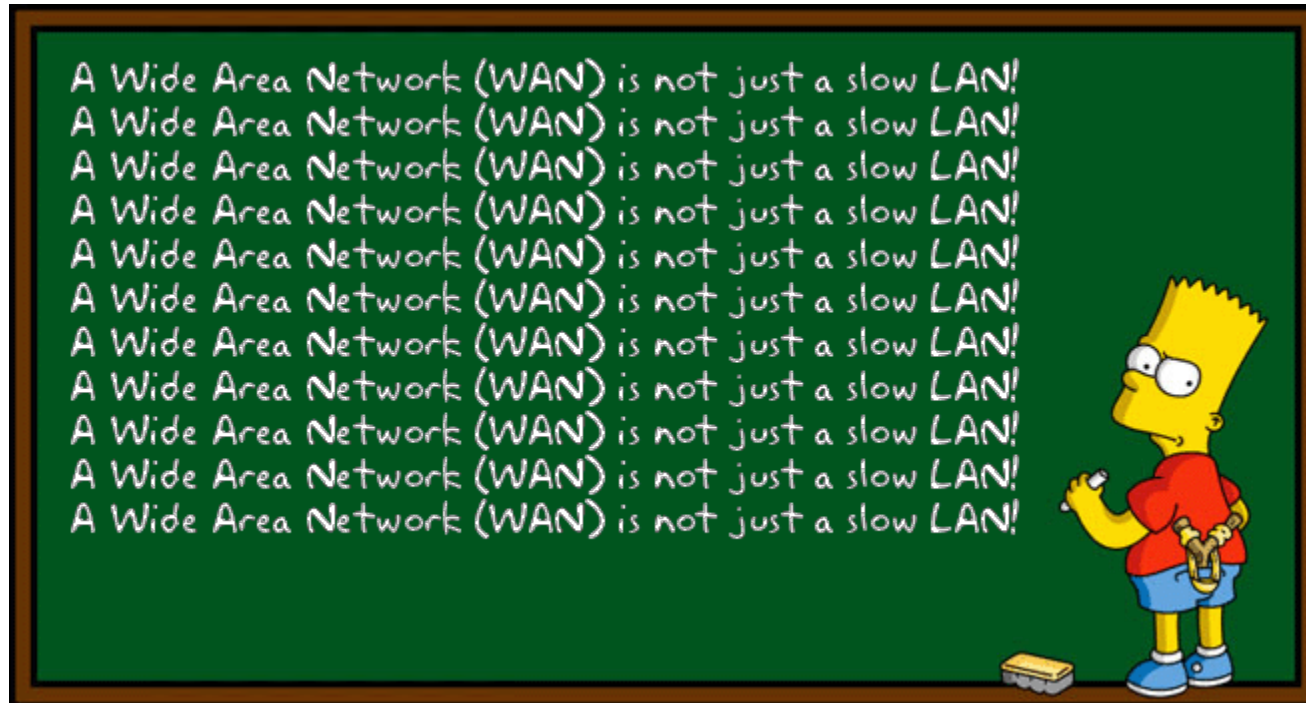
Offsets	Octet	0								1								2								3							
Octet	Bit	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31
0	0	Source port																Destination port															
4	32	Sequence number																															
8	64	Acknowledgment number (if ACK set)																															
12	96	Data offset	Reserved 0 0 0			N S	C W R E	E C R E	U R G K	A C K H	P S S T	R S Y N	F I N N	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)																															
...																															



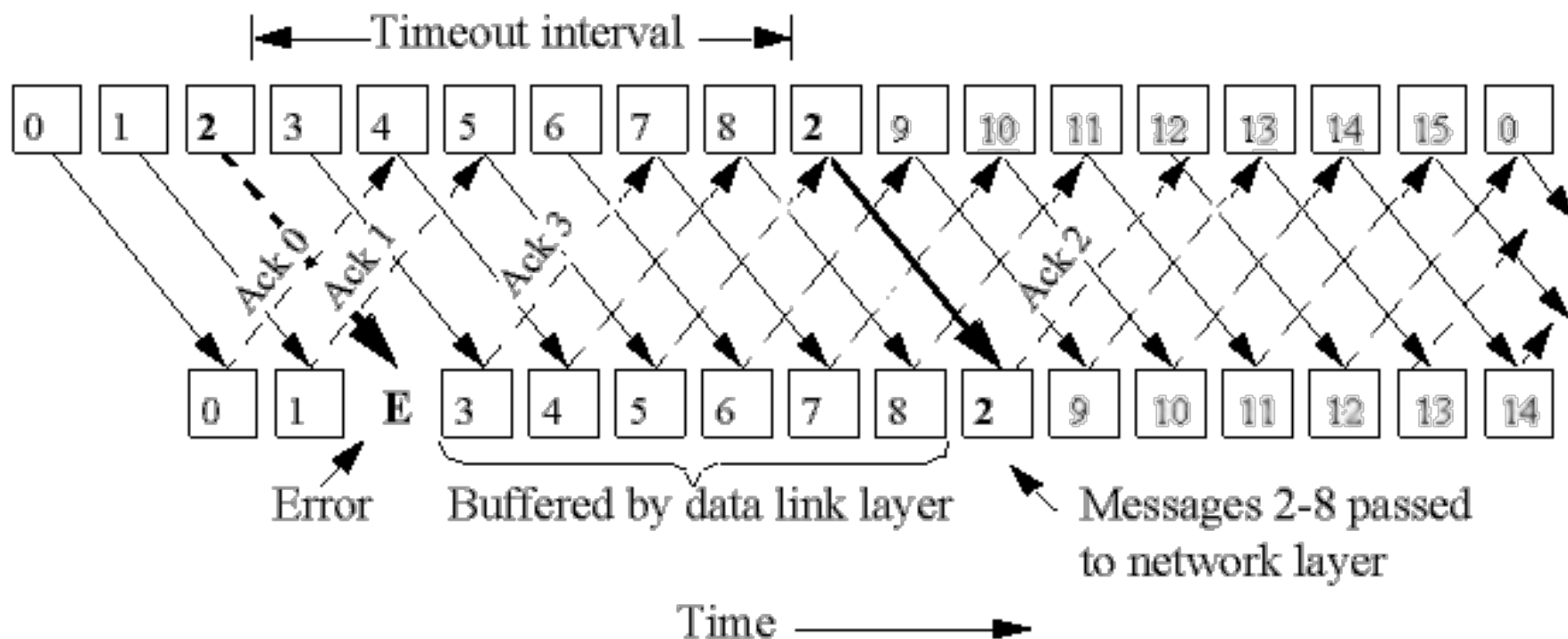
What do we get?

- **Reliable**
- **Ordered delivery (of a stream of octets)**
- **Error-free data transfer**
- **Flow control**
- **Congestion control**

Everything comes with a price



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/>

Why is this problem growing?

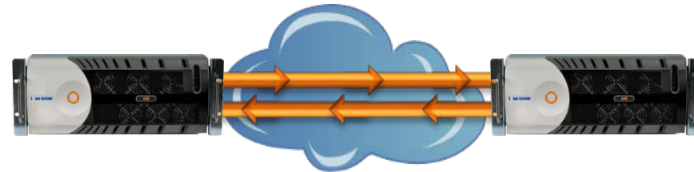
- Globalization
 - Public Internet backbone
 - RDBMS -> NoSQL, IMDG
 - Rich Data Types
 - Mobile Apps
 - Client Side Data
 - DR and BCP
1. Bandwidth
 2. Latency (RTT)
 3. Error Rate (Loss)



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



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**

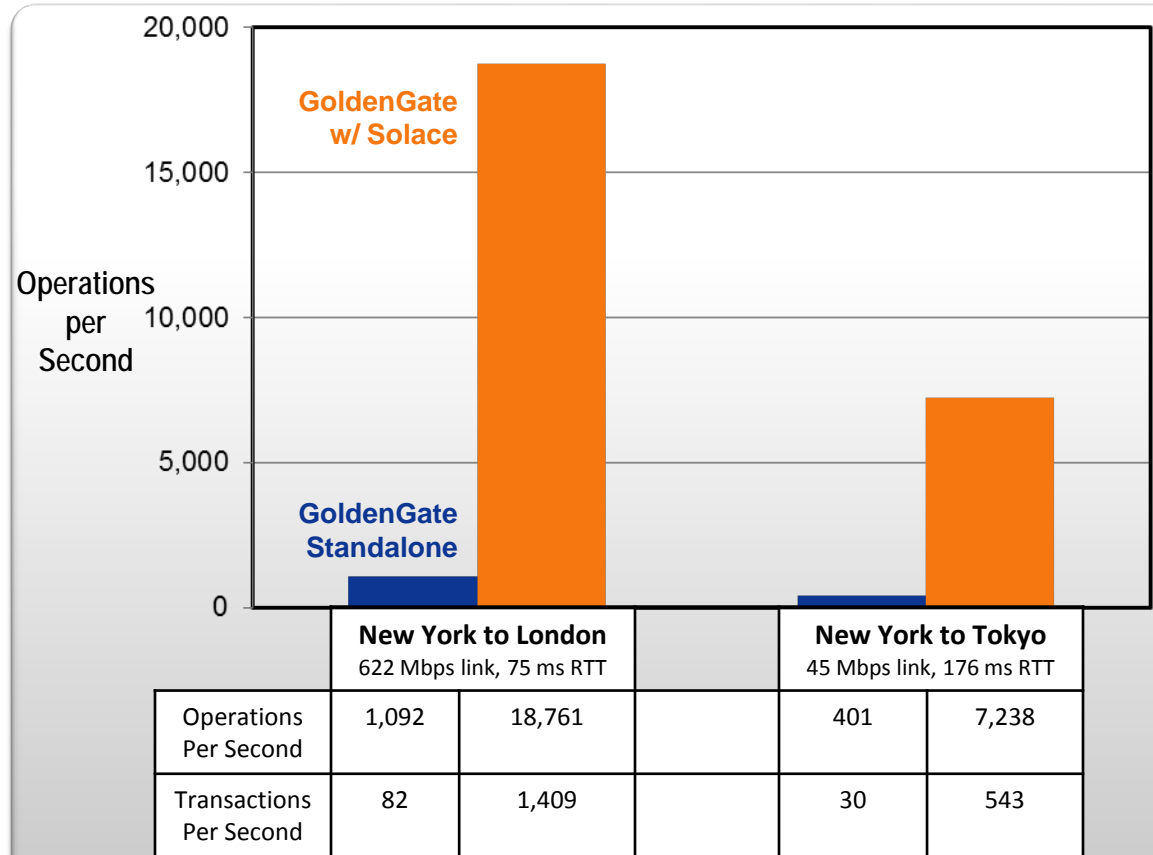
Millions of concurrent connections

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 ms round trip time

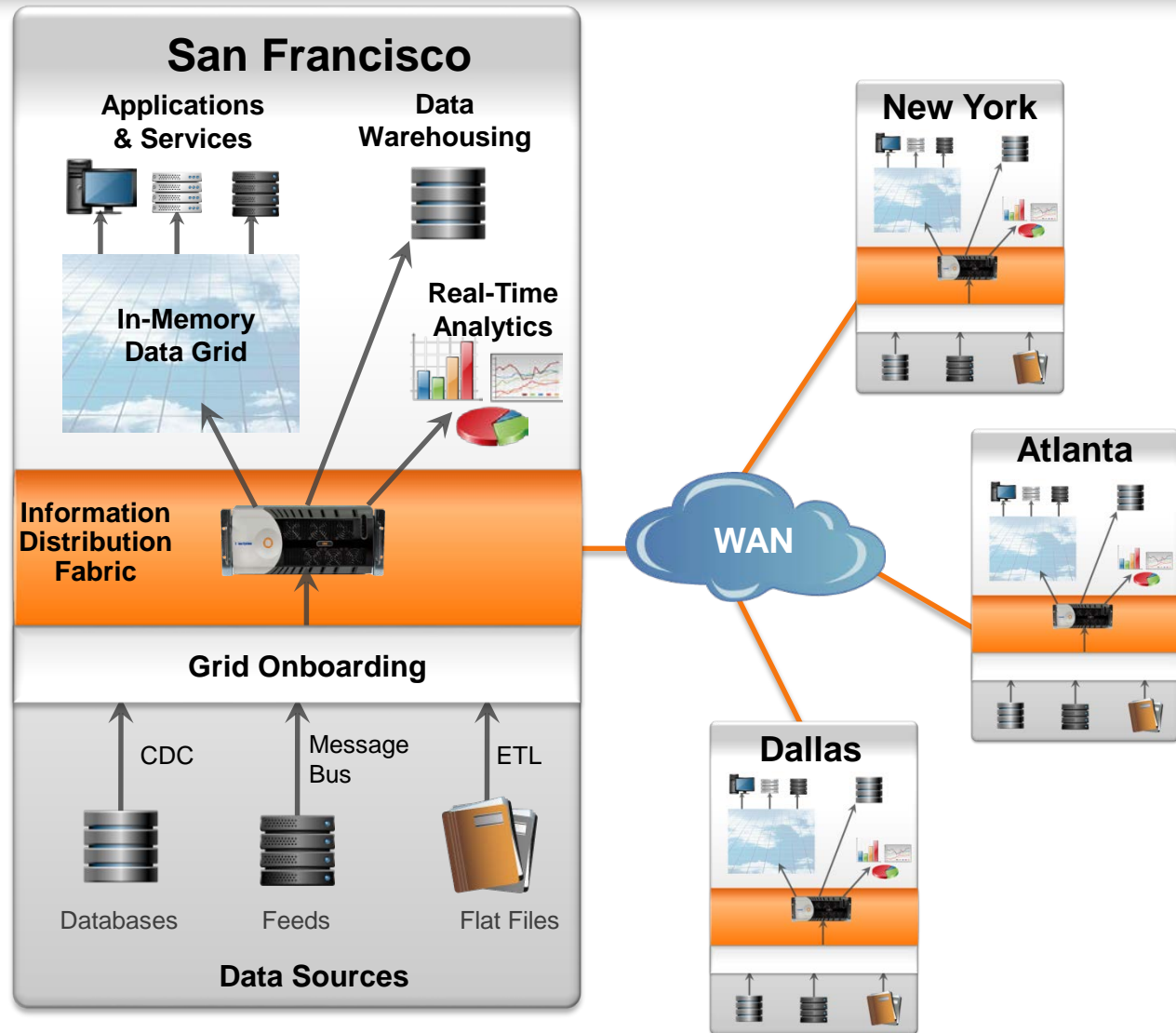
Solace 18x Faster



Back to the use case

Transactions != Packets
Database Records != Packets
Objects != Packets

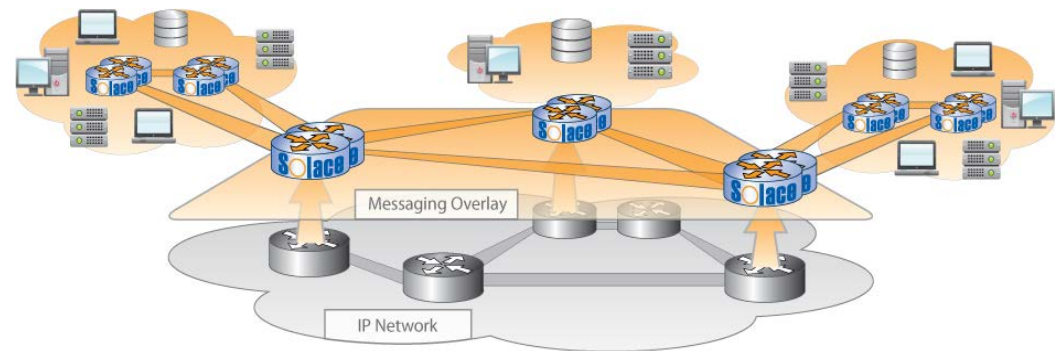
The Modern Information Distribution Fabric



Messaging Middleware Value

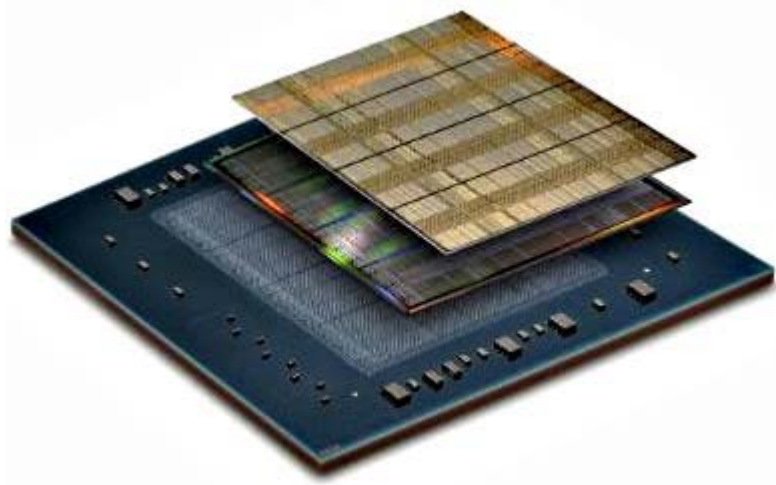


- Producer/Consumer Decoupling
- Disconnected Operation
- Location Independence
- Multipoint Delivery
- Advanced Filtering & Routing
- Message-based Granularity



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

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**

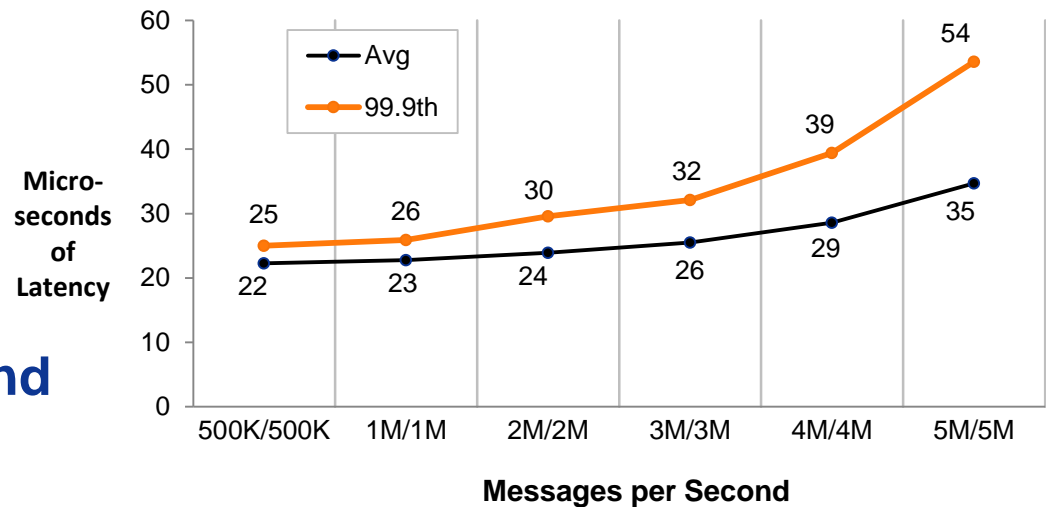
Reliable Messaging

○ 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

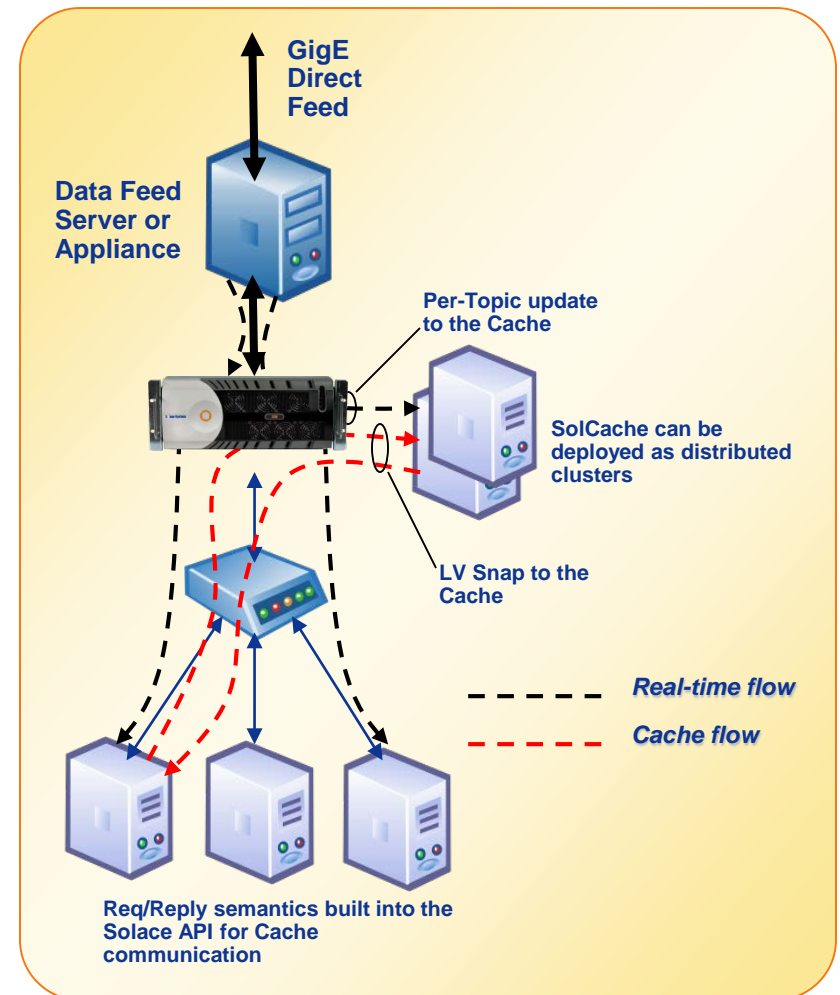


Bulk Message Rate	Message Size (bytes)	Message Rate (msgs/sec)	User Payload Bandwidth (Mbps)	10GigE Line Rate the is Limit
	100	5,930,000	4,744	
	500	2,080,000	8,320	
	1,000	1,080,000	8,640	
	12,000	92,000	8,832	
	30,000	34,000	8,160	

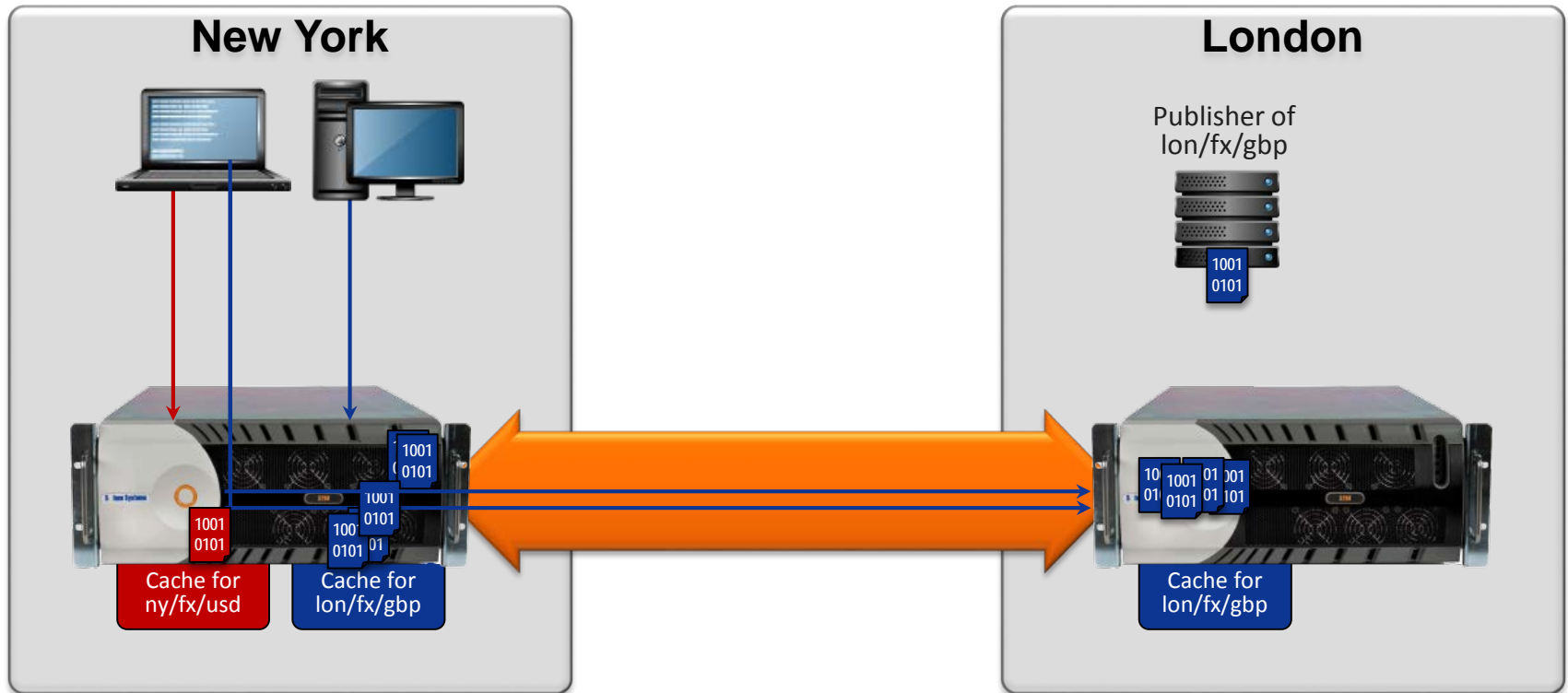
In Memory Message Caching

Non-persistent last-value cache can handle any payload

- Cache by number or timeframe
- Can run on appliance, or as a 64-bit app on a Linux server
- Centralized management of all caches.
- Clustering for load balancing and redundancy
- Support wildcard requests
- Request/reply with the cache, and control of synchronization of cache requests
- Topic names partitioned amongst instances to scale storage



Cascading Cache



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

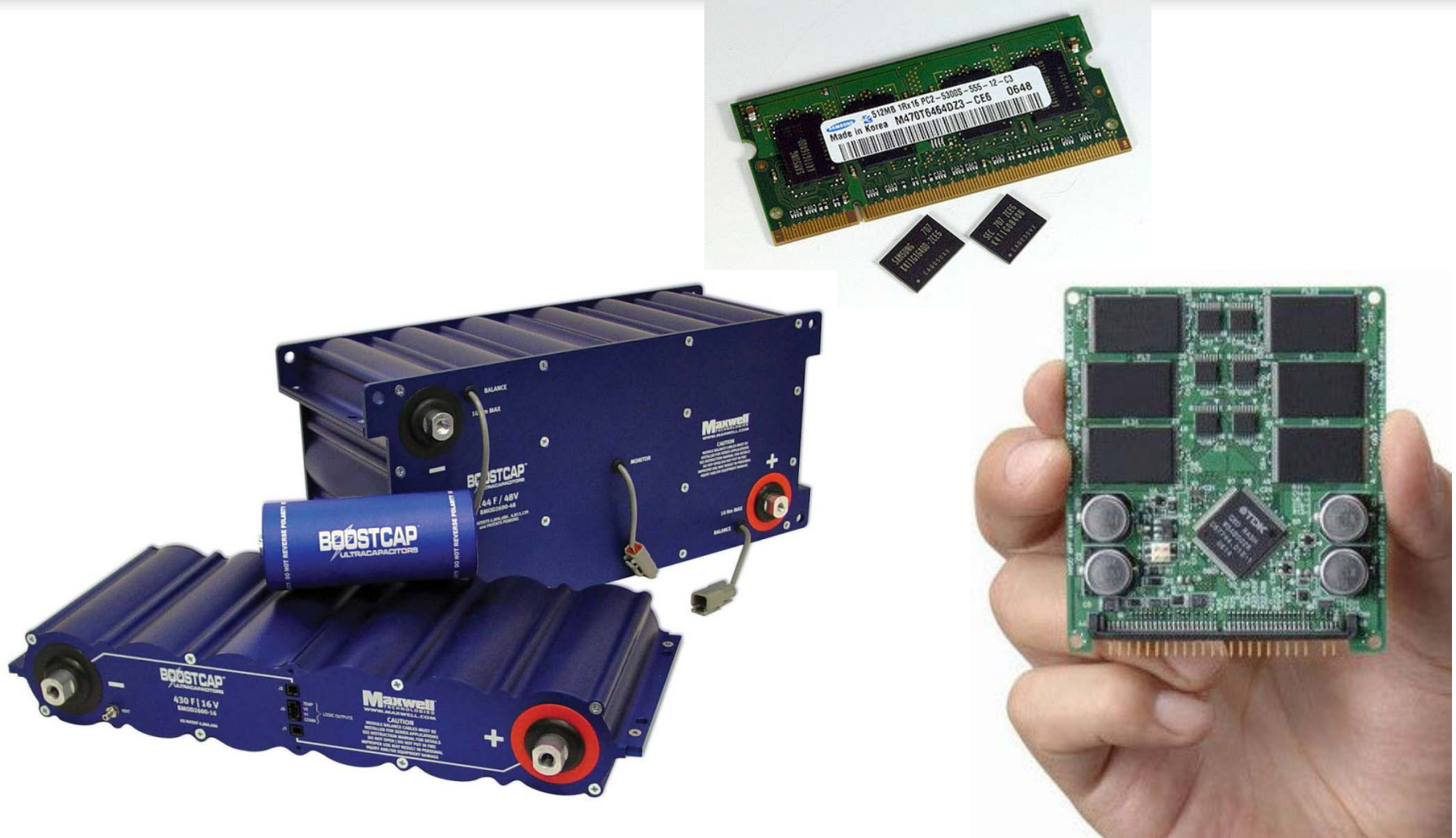
DAY LOW: 19.03

DAY HIGH: 19.21

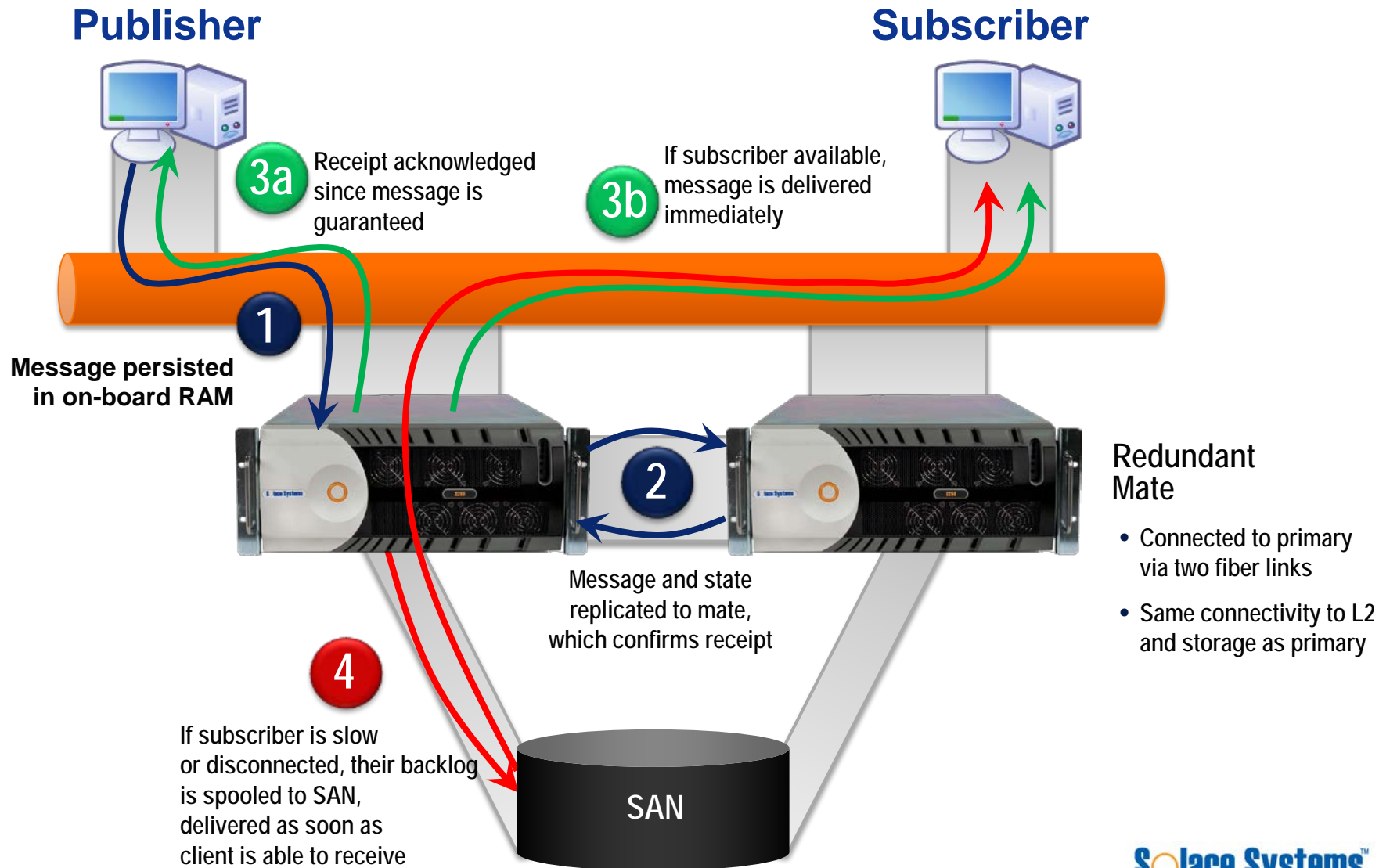
52-WEEK LOW: 15.13

52-WEEK HIGH: 23.98

Hybrid Storage



Fault Tolerant Clustering of Messaging Nodes

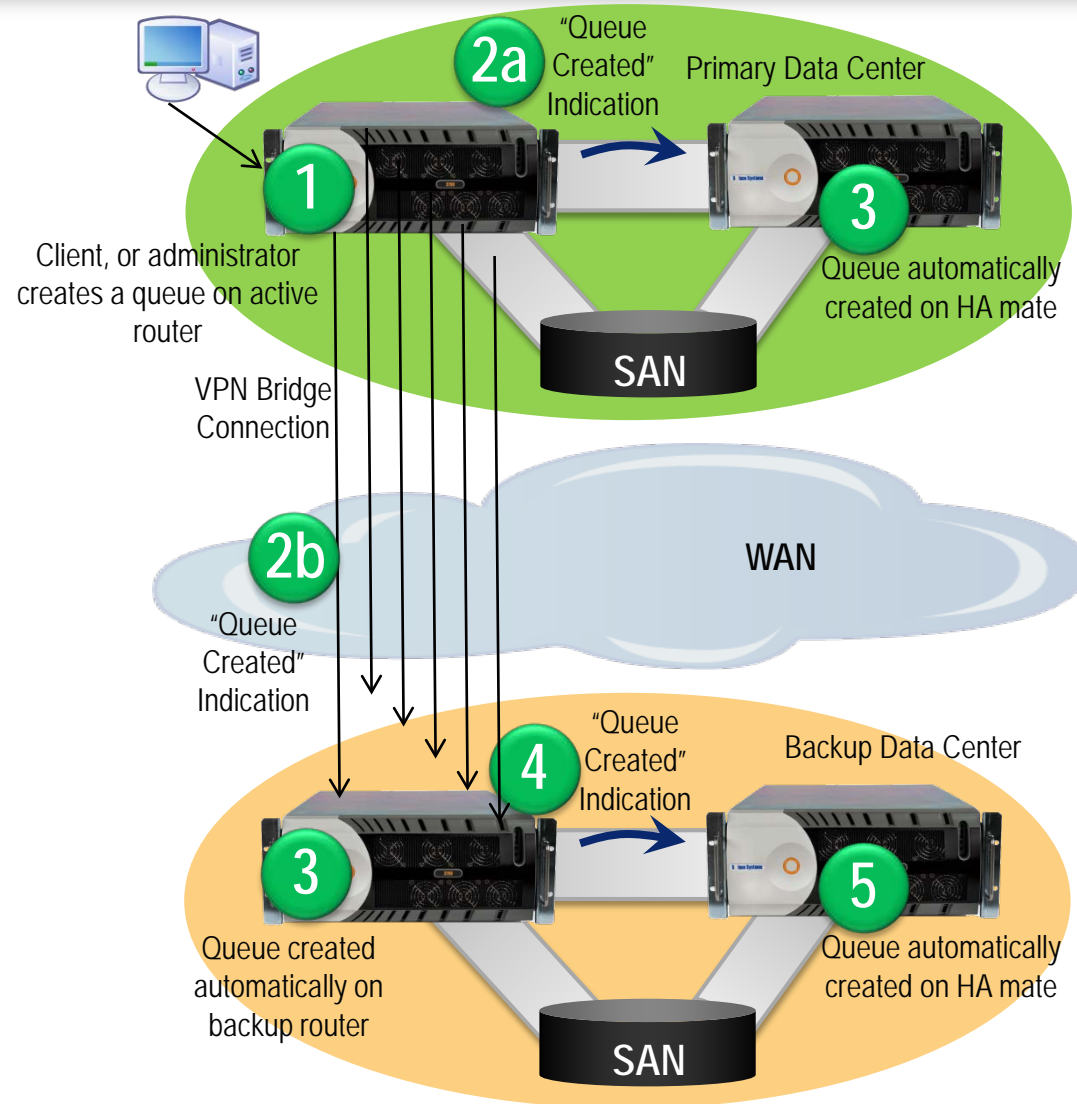


Multiple Datacenters and Disaster Recovery

- Automatic replication of:

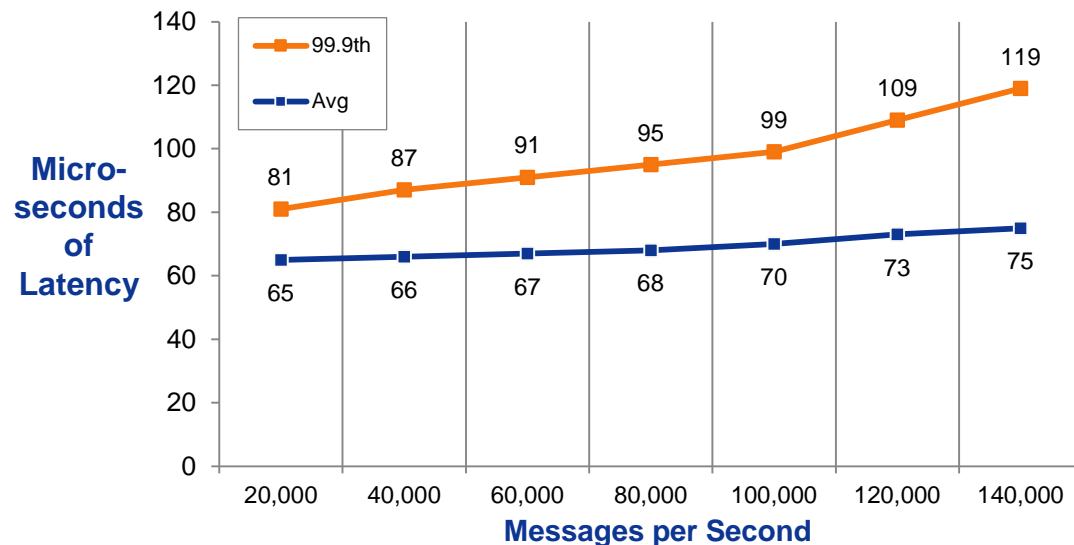
- Client-created endpoints
- Configuration data
- Transactional state

- Needs Configurability for Sync & Asynchronous replication



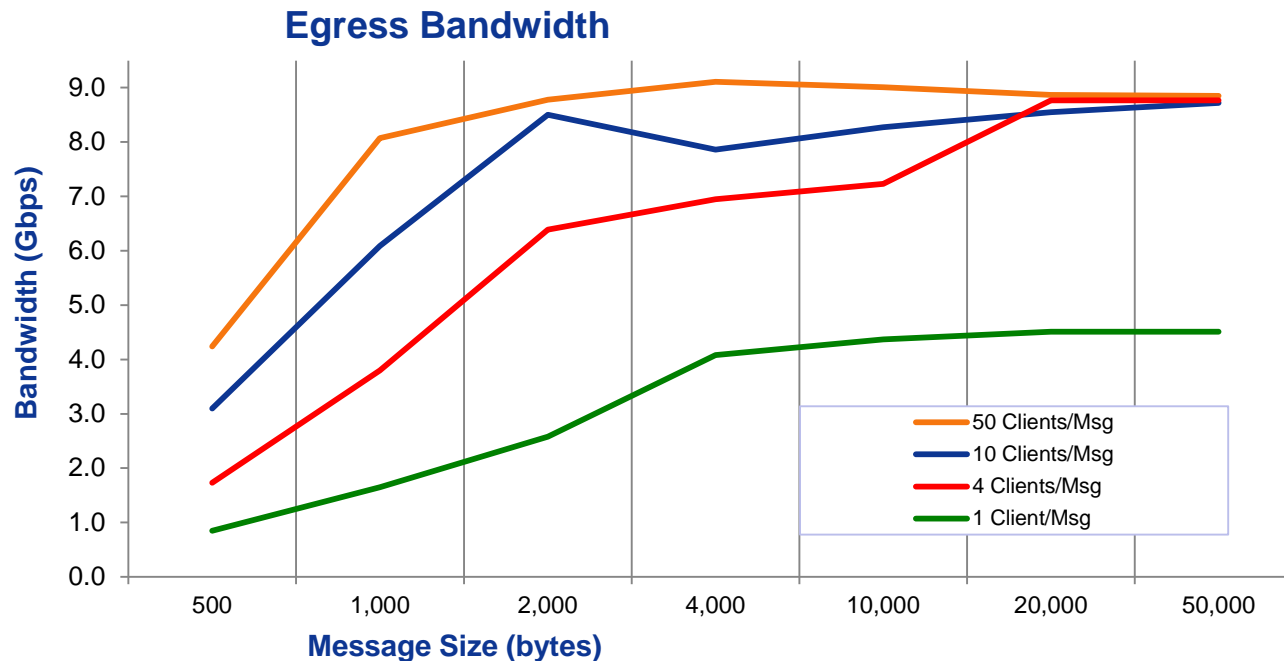
Guaranteed Messaging; Store & Forward Performance

- **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**



Bulk Message Rate	Message Size (bytes)	Message Rate (msgs/sec)	User Payload Bandwidth (Mbps)
	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

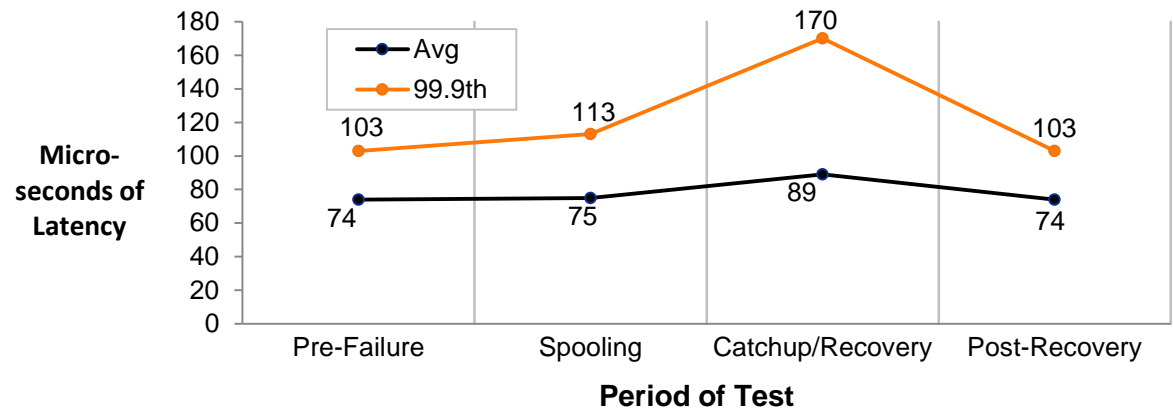
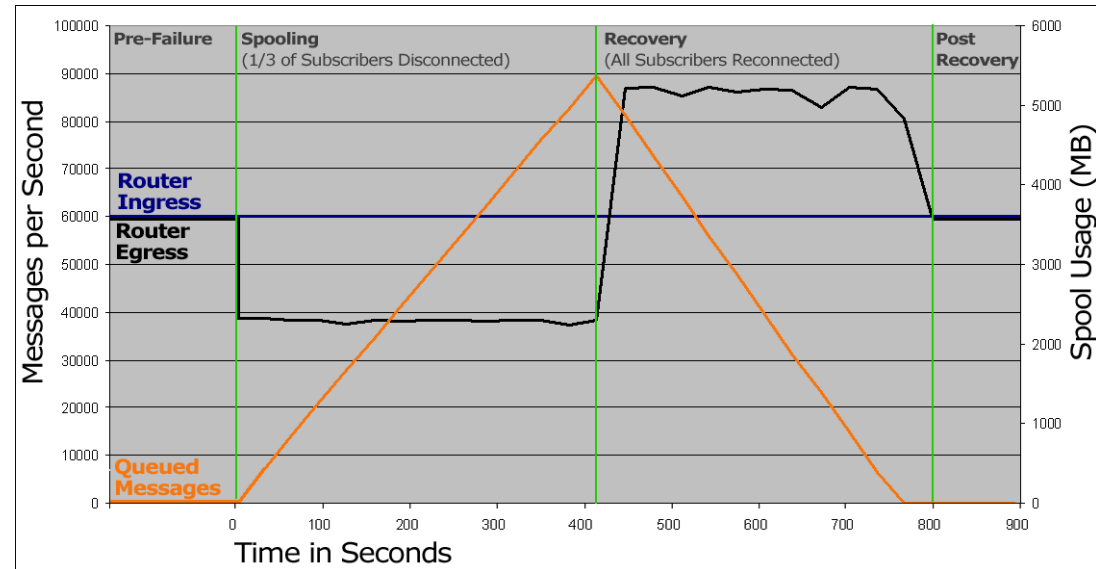
Guaranteed Messaging; Fan-out performance



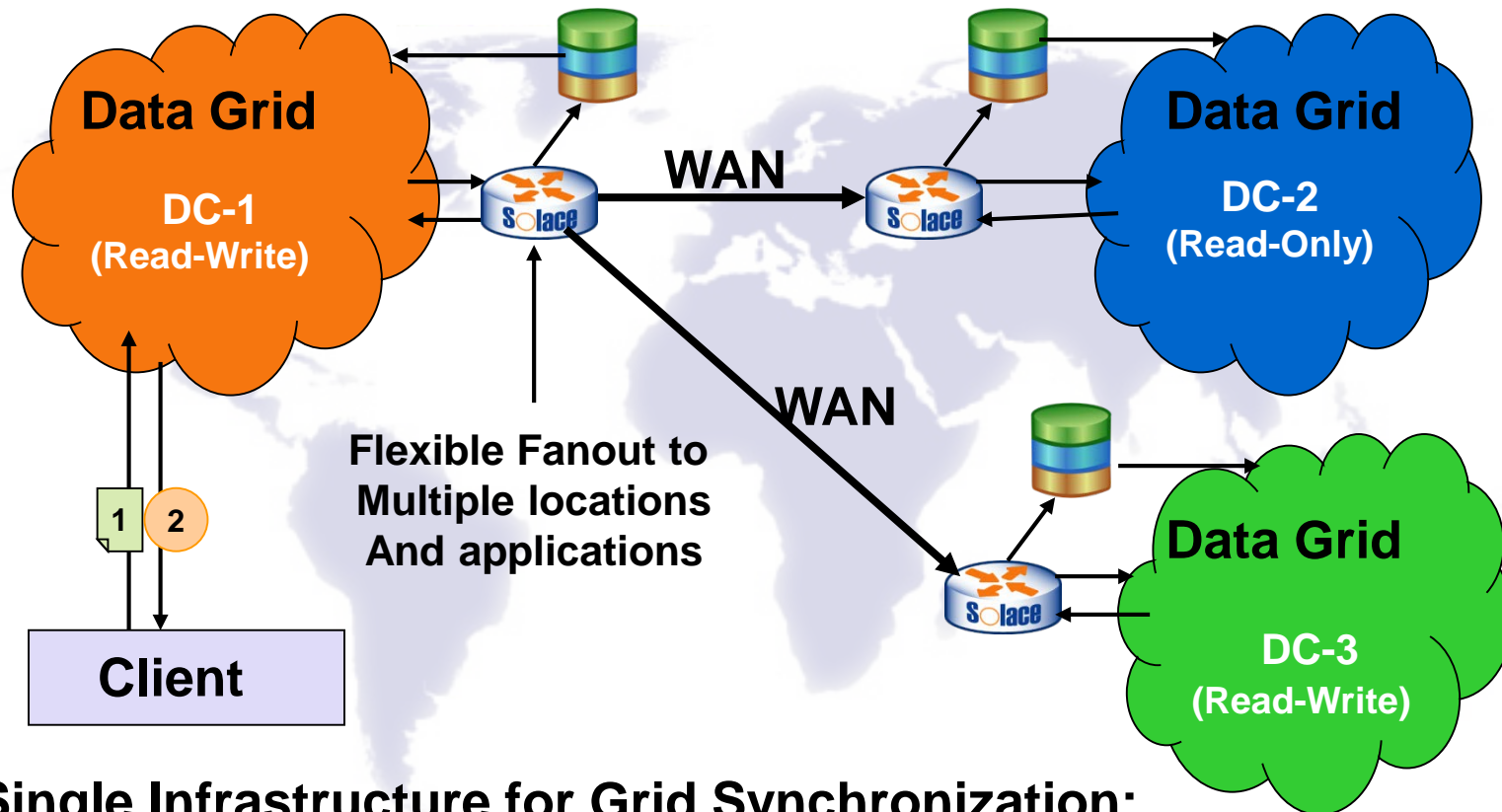
Message Size (bytes)	Egress Rate, 1 client/msg (msgs/sec)	Egress Rate, 4 clients/msg (msgs/sec)	Egress Rate, 10 clients/msg (msgs/sec)	Egress Rate, 50 clients/msg (msgs/sec)
512	206,400	422,000	756,000	1,035,000
1,024	202,000	464,000	744,000	985,000
2,048	157,500	390,000	519,000	536,000
4,096	124,400	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

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



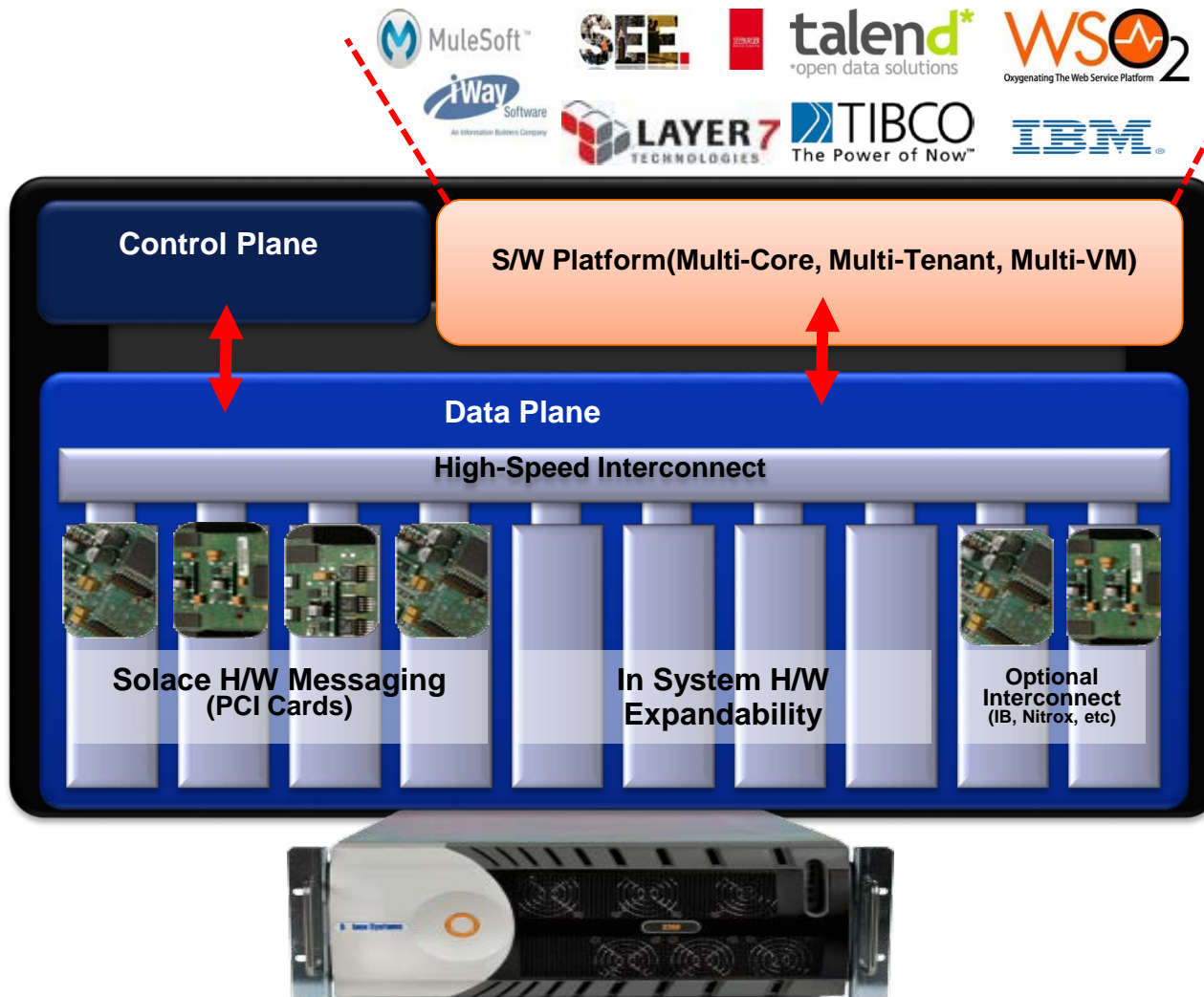
Optimized In-memory Grid Replication



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 as an Appliance Platform



- 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.

The Modern Information Distribution Fabric

High Volume
Onboarding

Fast, Efficient
WAN Sync

