

Issues in The Design of Exascale Networks

HiPINEB

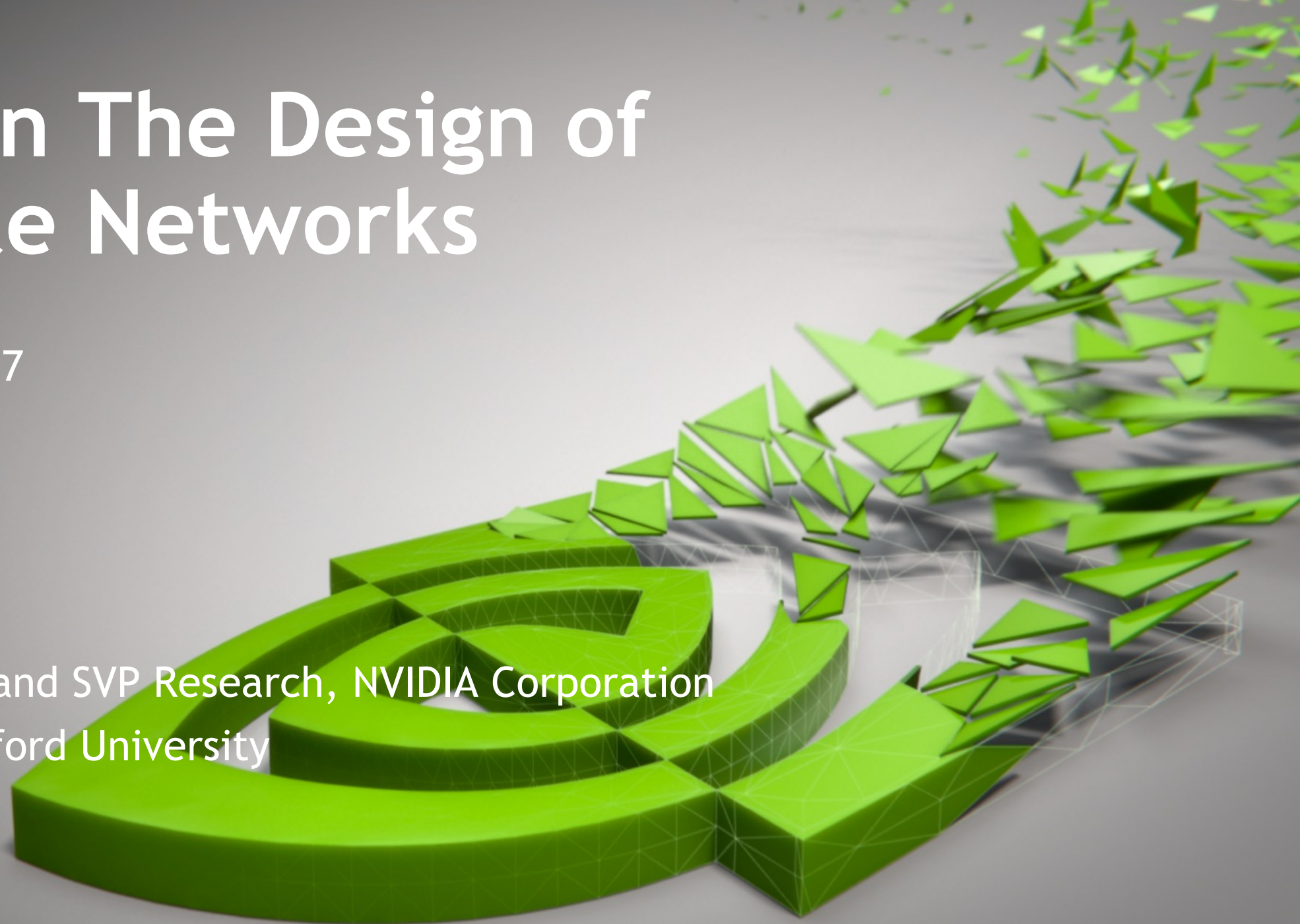
February 5, 2017



Bill Dally

Chief Scientist and SVP Research, NVIDIA Corporation

Professor, Stanford University



Acknowledgement

This talk presents the work of the NVIDIA Network Research Group

Larry Dennison (Director)

Nir Arad

Matt Blumrich

Hans Eberle

Ted Jiang

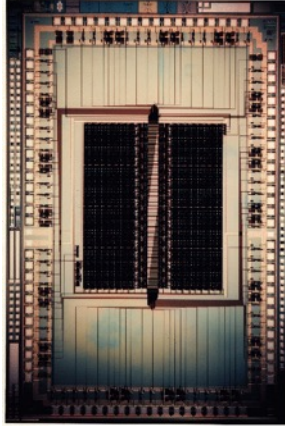
Alex Ishii (DGX Product Group)

Outline

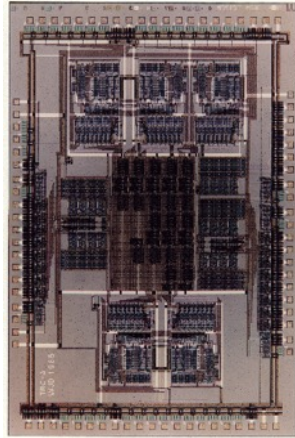
- Desiderata - Exascale network requirements
 - Efficient fine-grain communication, cost effective bandwidth, resilience
- Topology - Engineering to optimize available technology
- Routing
- Flow control and congestion avoidance
- Error control
- Ordering
- The role of photonics
- System sketch

Some History

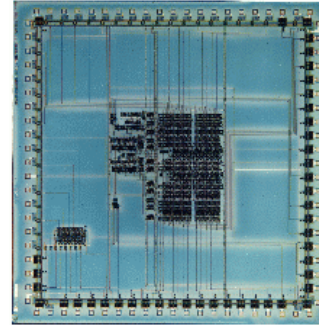
MARS Router
1984



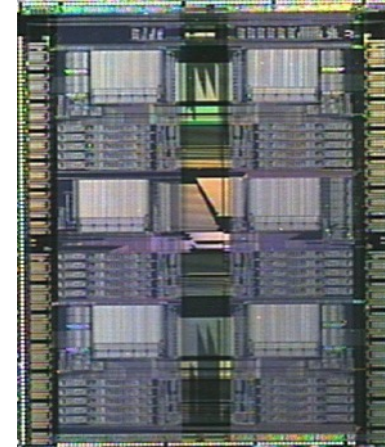
Torus Routing Chip
1985



Network Design Frame
1988



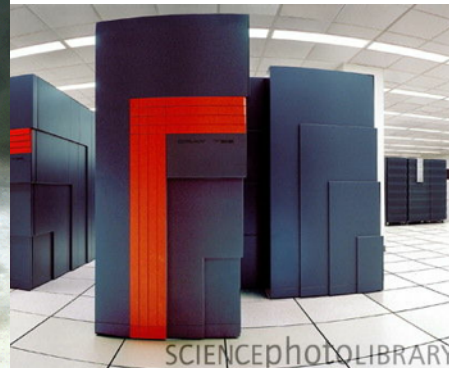
Reliable Router
1994



J-Machine
1992



Cray T3D
1992

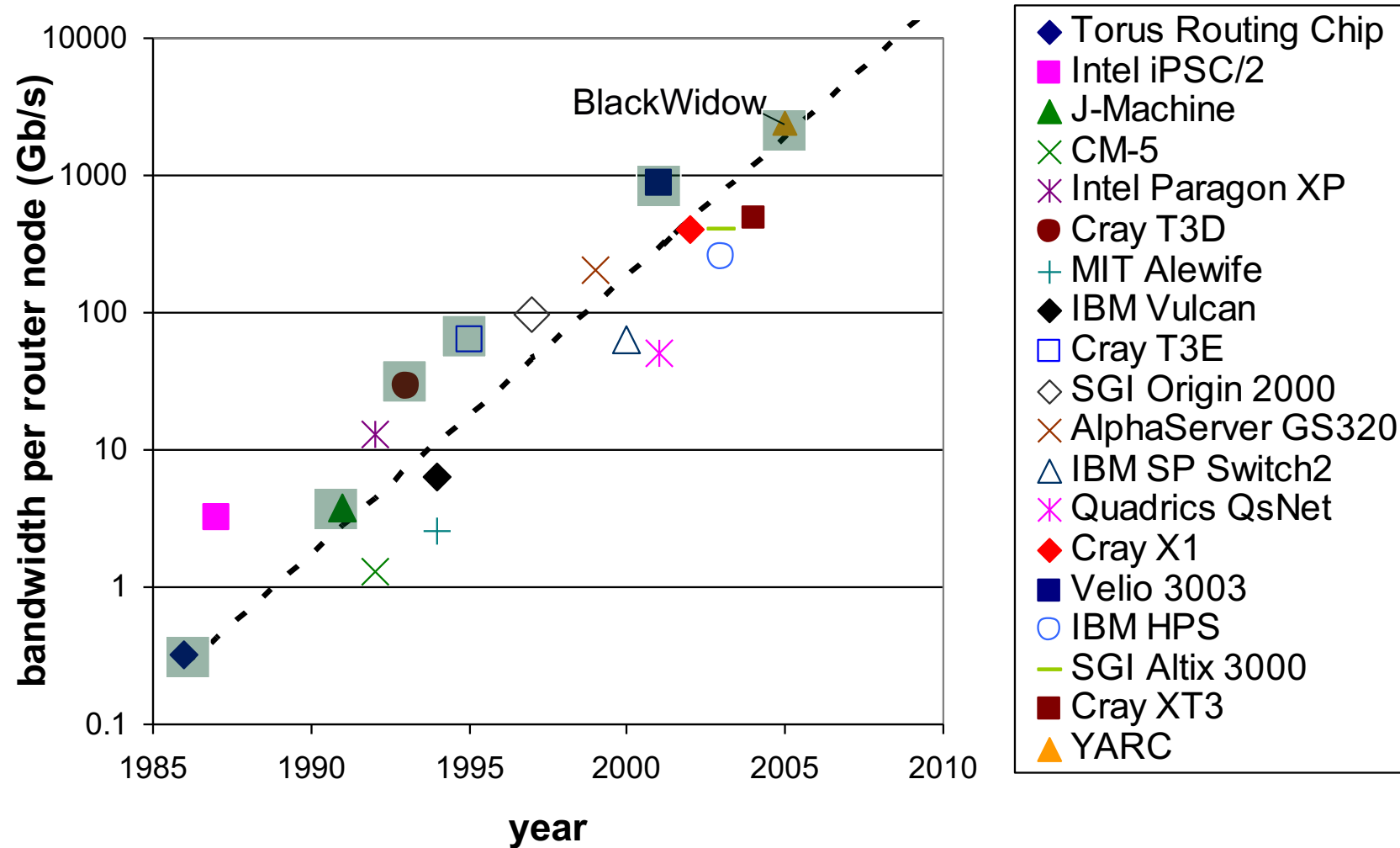


Cray T3E 1995



Cray Black Widow
2006

Trend Line



NVIDIA DGX-1

WORLD'S FIRST DEEP LEARNING SUPERCOMPUTER



170 TFLOPS

8x Tesla P100 16GB

NVLink Hybrid Cube Mesh

Optimized Deep Learning Software

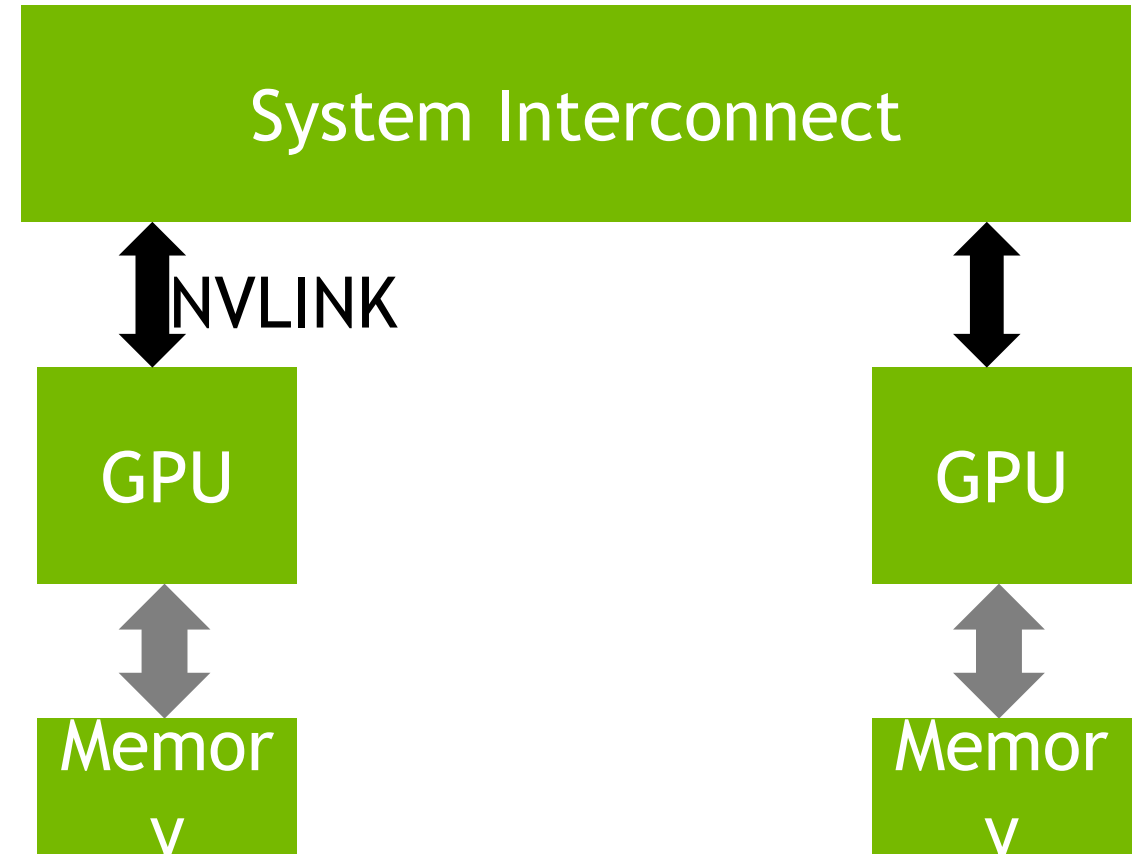
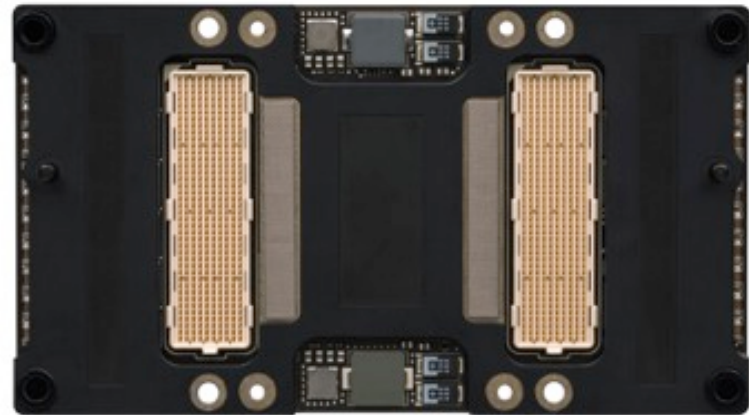
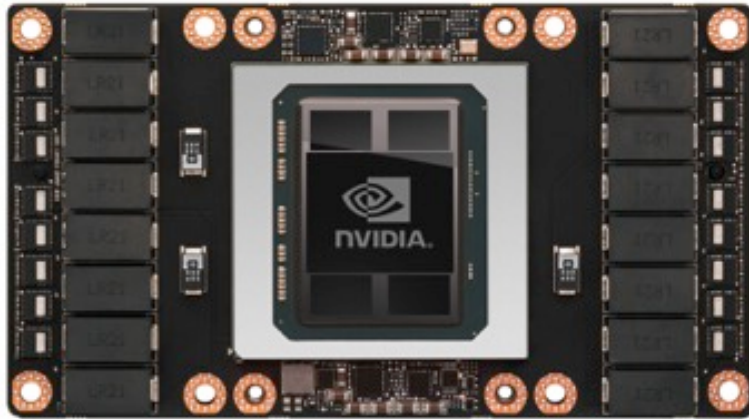
Dual Xeon

7 TB SSD Deep Learning Cache

Dual 10GbE, Quad IB 100Gb

3RU - 3200W

NVLINK - Enables Fast Interconnect, PGAS Memory



DGX SATURNV

World's Most Efficient AI Supercomputer



Fastest AI Supercomputer in TOP500

4.9 Petaflops Peak FP64 Performance
19.6 Petaflops DL FP16 Performance
124 NVIDIA DGX-1 Server Nodes



Most Energy Efficient Supercomputer

#1 on Green500 List
9.5 GFLOPS per Watt
2x More Efficient than Xeon Phi System

13 DGX-1 Servers in Top500

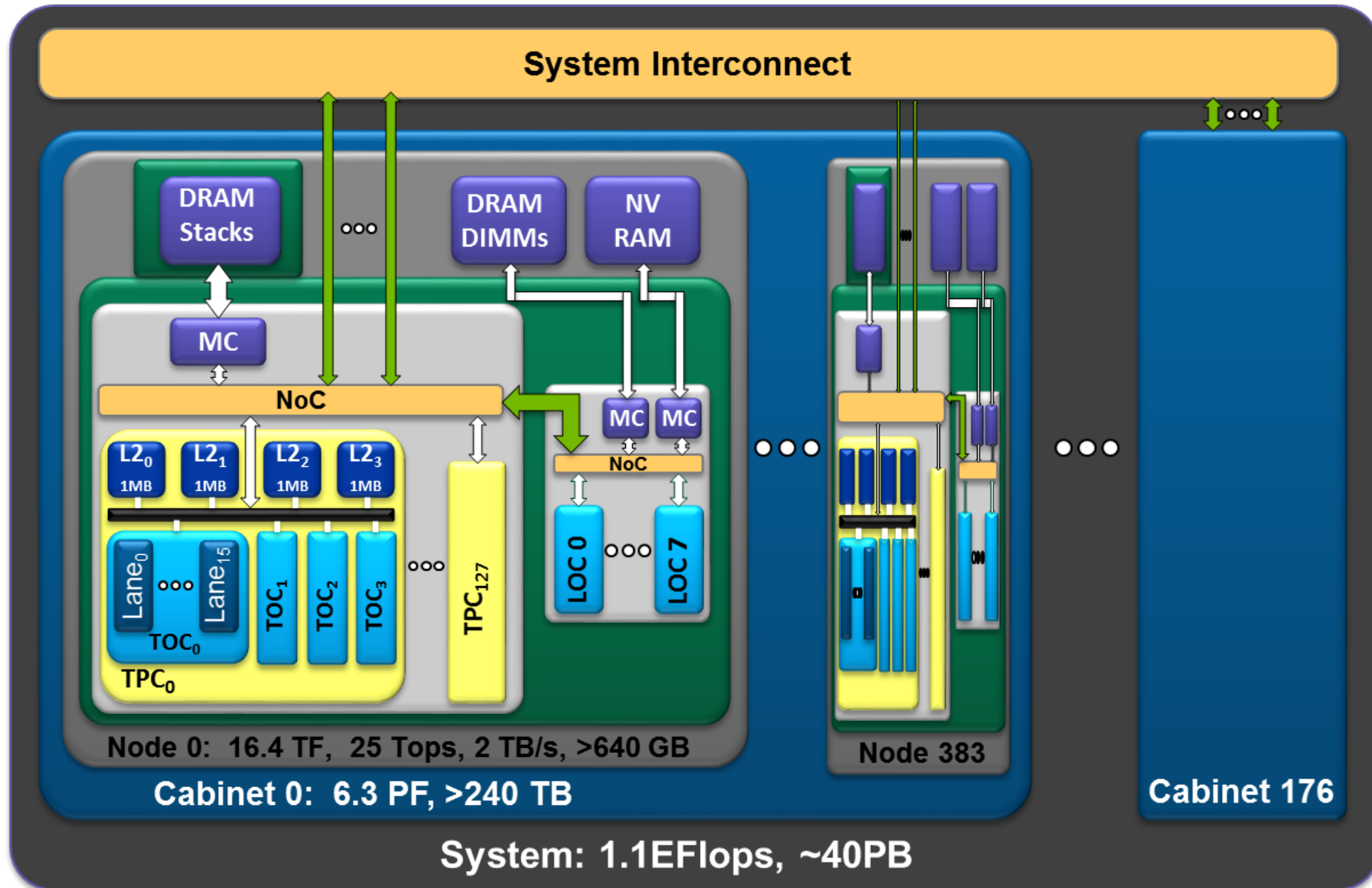
FACTOIDS

38 DGX-1 Servers for Petascale supercomputer

55x less servers, 12x less power vs CPU-only
supercomputer of similar performance

Desiderata

Exascale System Sketch

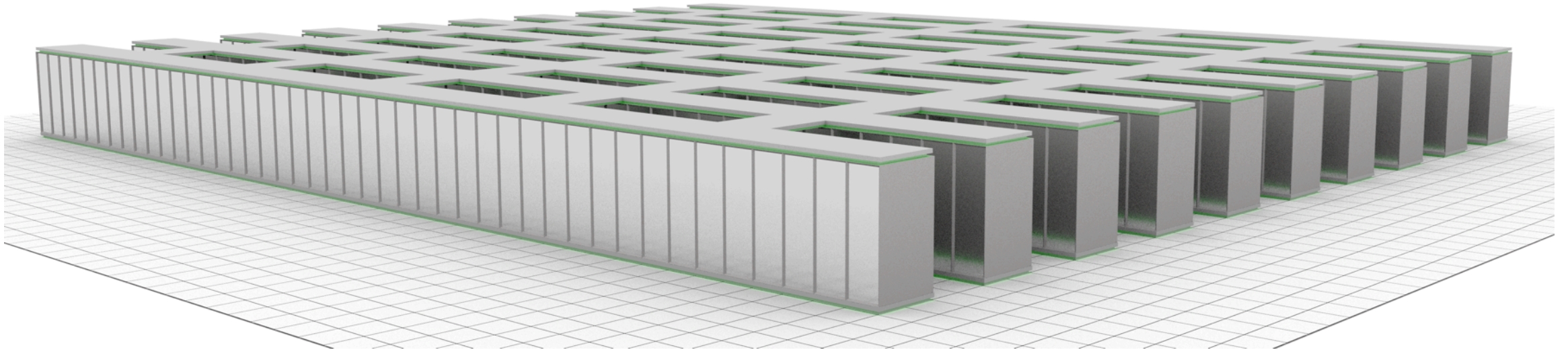


Desiderata

- Scale - 10^5 endpoints
- Cost-effective bandwidth (injection and bisection) - B/s\$
- Low latency (dominated by time-of-flight)
- Reliable exactly-once delivery - ($\text{BER} < 10^{-21}$)
- Low overhead (latency and occupancy)

- Enable strong scaling
- Low-overhead shared-memory operations
- Highly concurrent operation

Scale - $\sim 10^5$ Powerful Endpoints

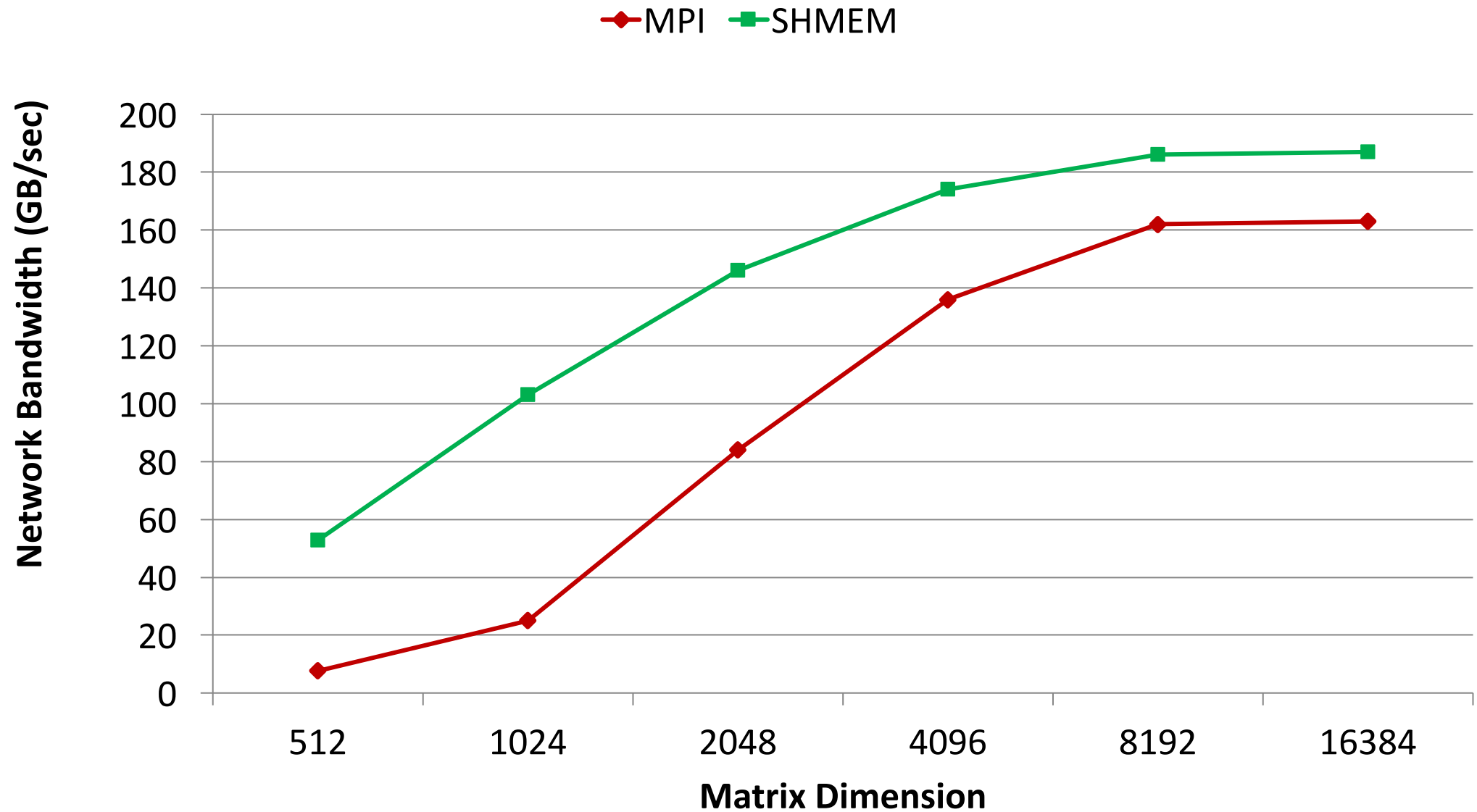


- Each endpoint is a 16.4TFLOPS (DP) GPU with 4TB/s memory bandwidth
- 400GB/s injection bandwidth is 10:1 local memory to neighbor memory
 - Pascal GP100 today is ~5:1 (750GB/s memory: 160GB/s NVLINK)

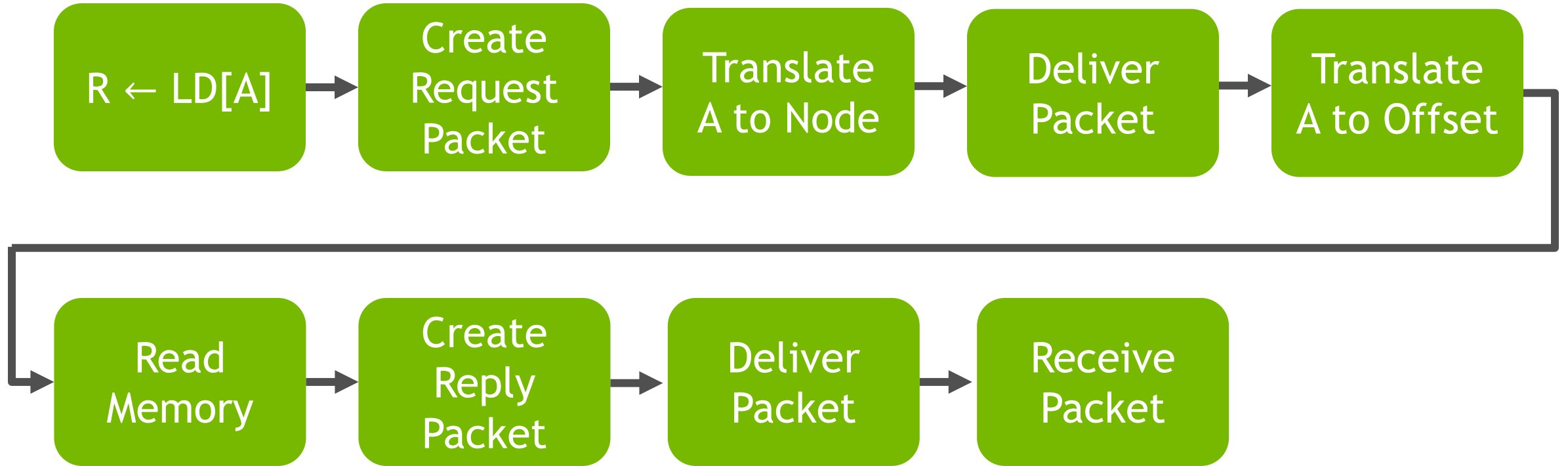
Cost-Efficient Bandwidth (B/s\$)

- Network cost dominated by links (AOCs)
- Use minimum number of expensive links per route (1 dragonfly)
- Operate each link near capacity (flow-control and congestion avoidance)
- Make payload a high fraction of bits on the wire
 - For small payloads (16B) as well as large

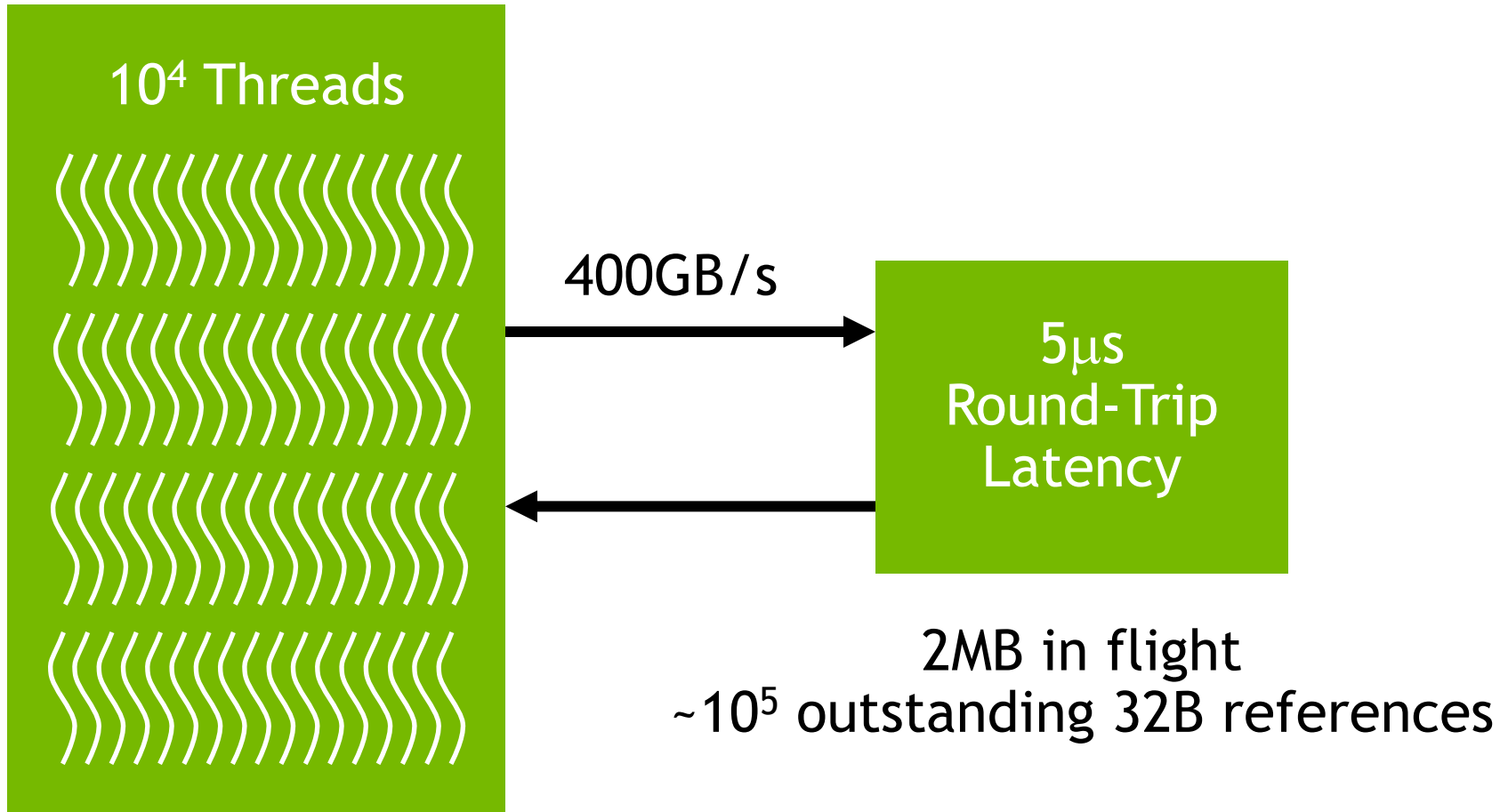
The Need for PGAS



Remote Load/Store



$\sim 10^5$ Outstanding References per Endpoint



Topology

Cost of 100Gb/s



\$500

100m



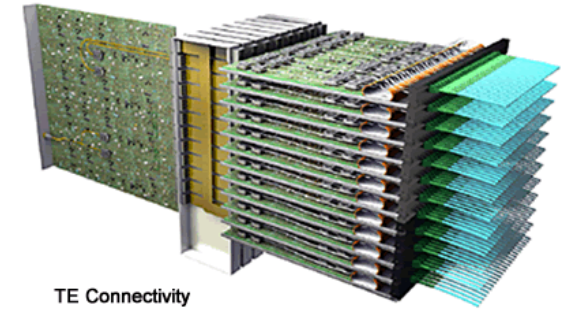
\$50

5m



\$10

1m

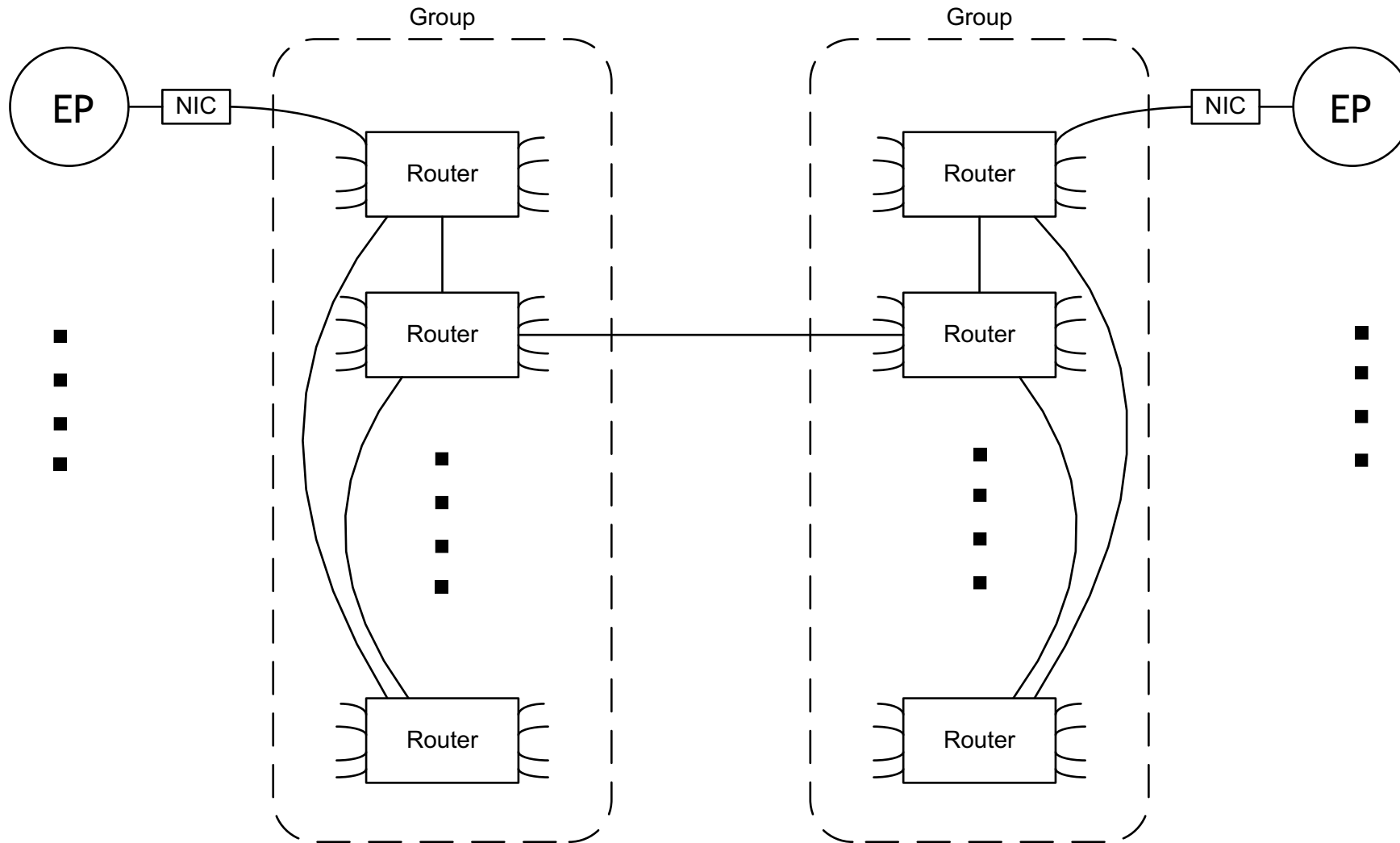


TE Connectivity

\$5

0.3m

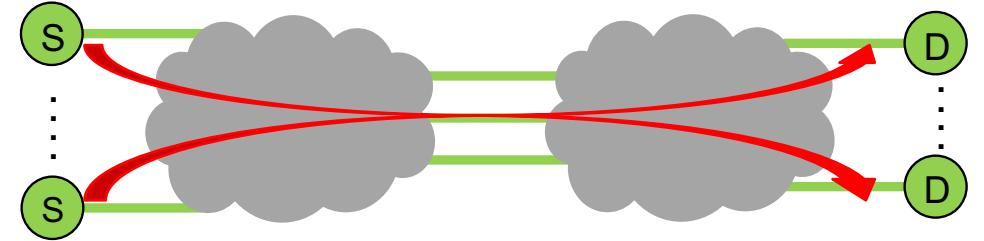
Dragonfly Topology



Adaptive Routing

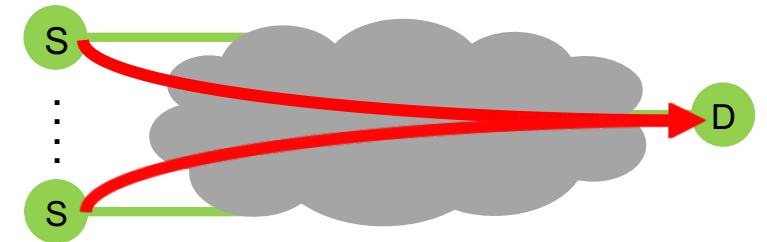
Sources of Congestion

Fabric congestion



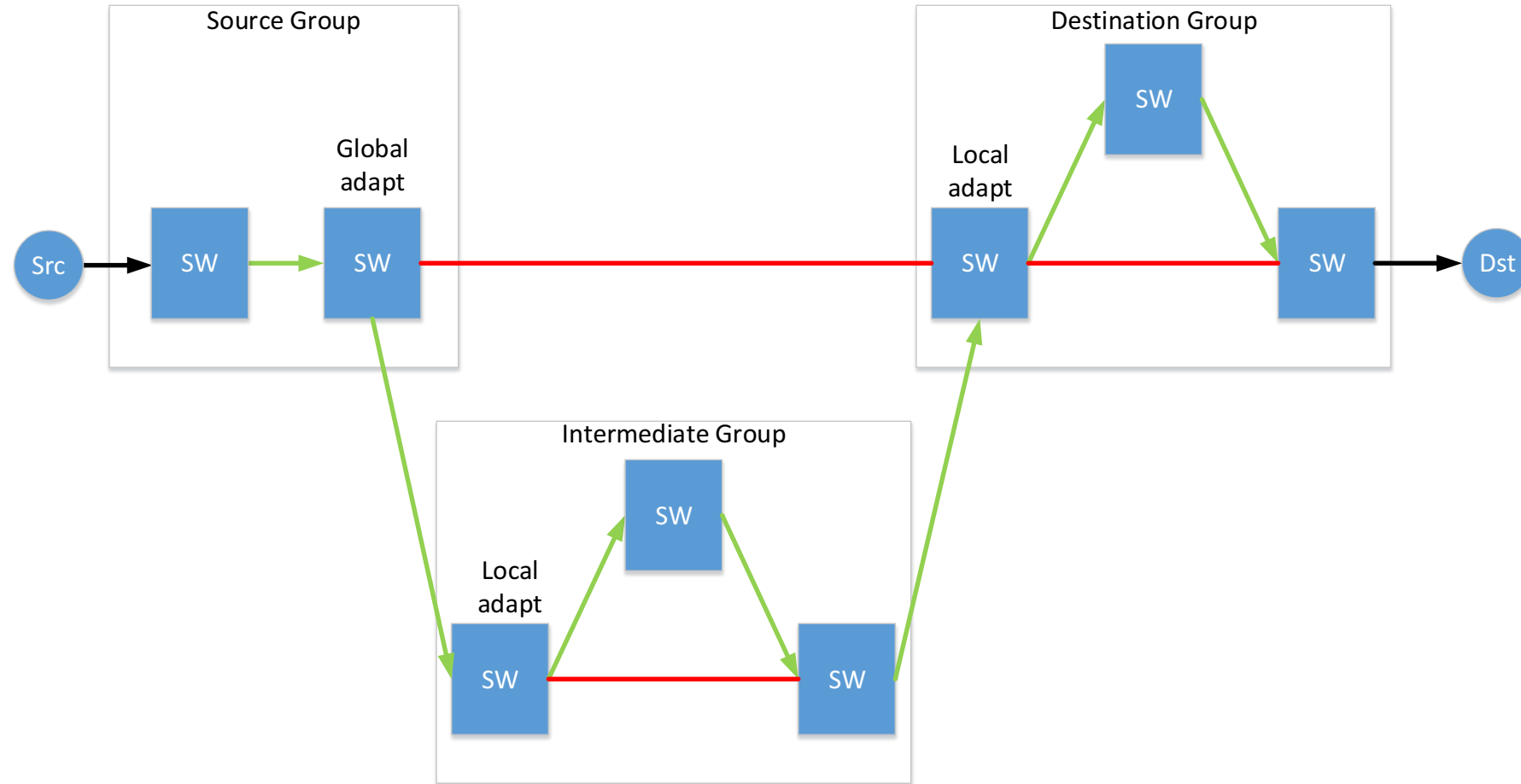
- Cause: low bisection bandwidth or load imbalance
- Solution: add bandwidth, improve load-balance using adaptive routing

Endpoint congestion



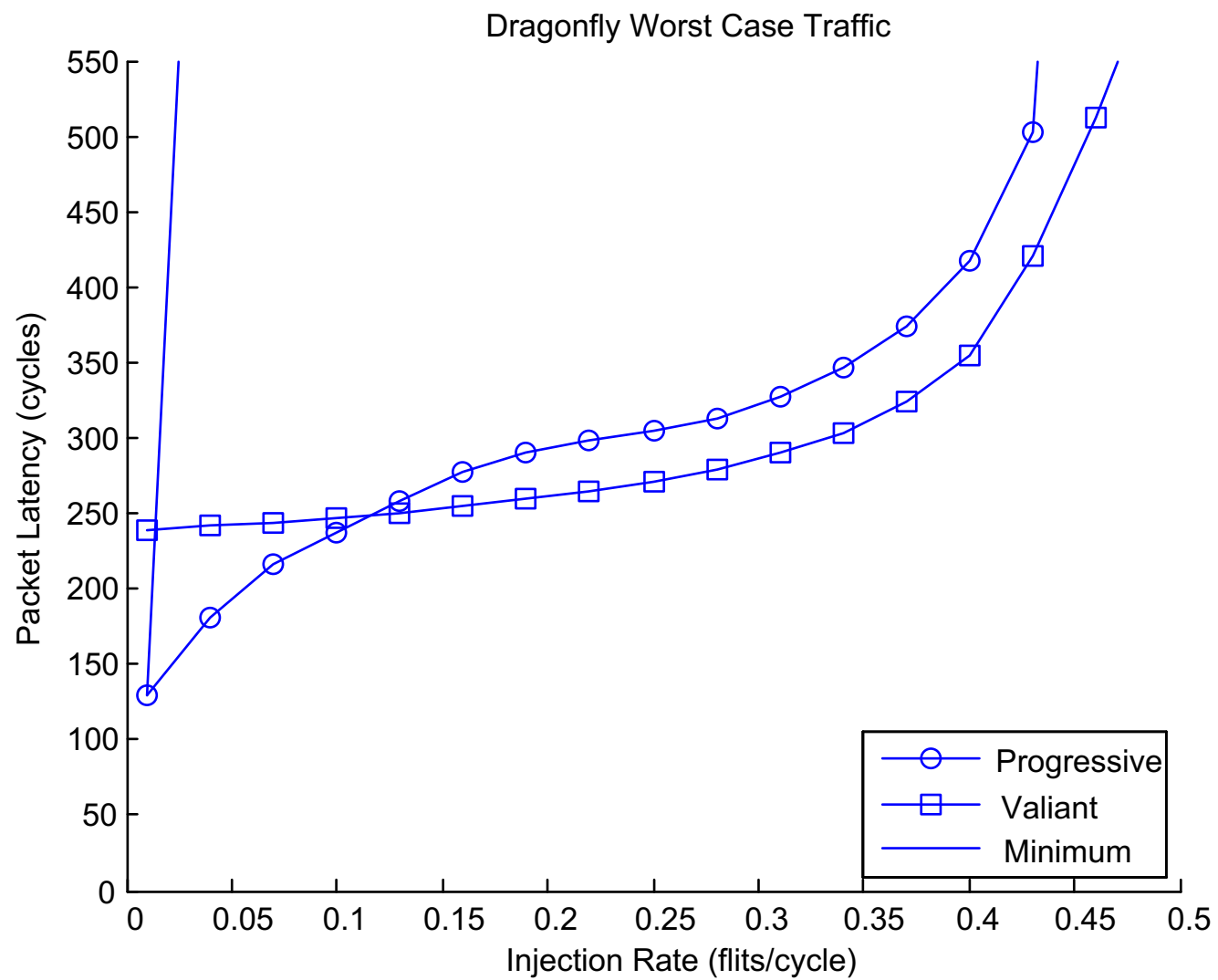
- Cause: endpoint bandwidth over-subscription
- Solution: reduce bandwidth demand by throttling traffic sources

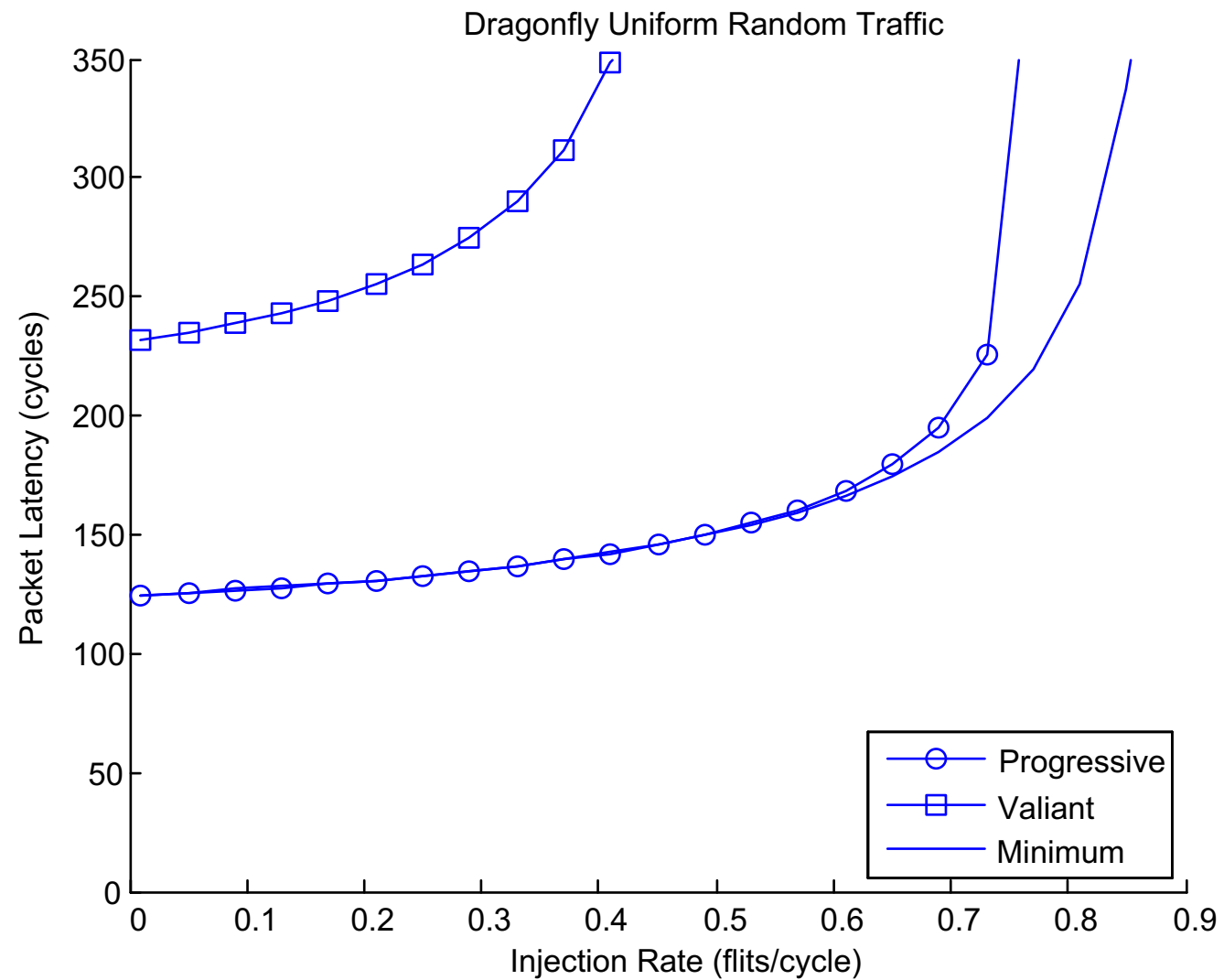
Progressive Adaptive Routing with Local misroute



Singh, A., 2005. *Load-balanced routing in interconnection networks* (Doctoral dissertation, Stanford University).

Jiang, Nan, William J. Dally, and John Kim. "Indirect Adaptive Routing on Large Scale Interconnection Networks," *ISCA 2009*





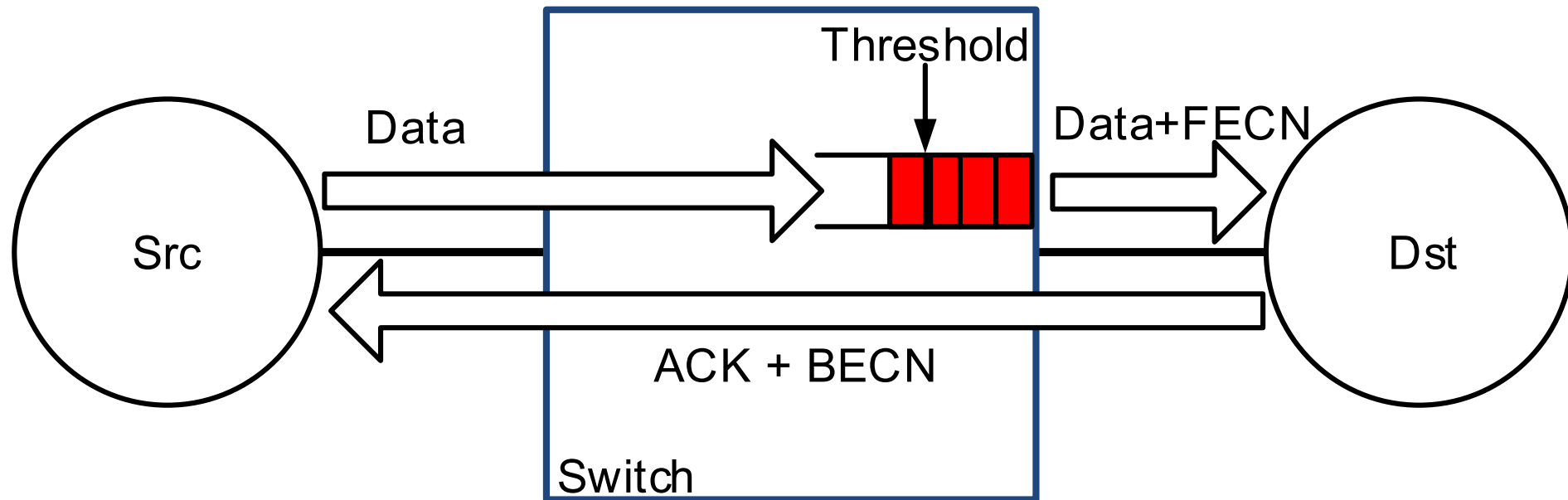
Congestion Avoidance

LHRP

Congestion Notification

$150\text{GB/s} * 3\mu\text{s RTT} = 450\text{ KB}$ inflight before first notification

Reaction is too late, slow response time, large transient



Last-Hop Reservation Protocol

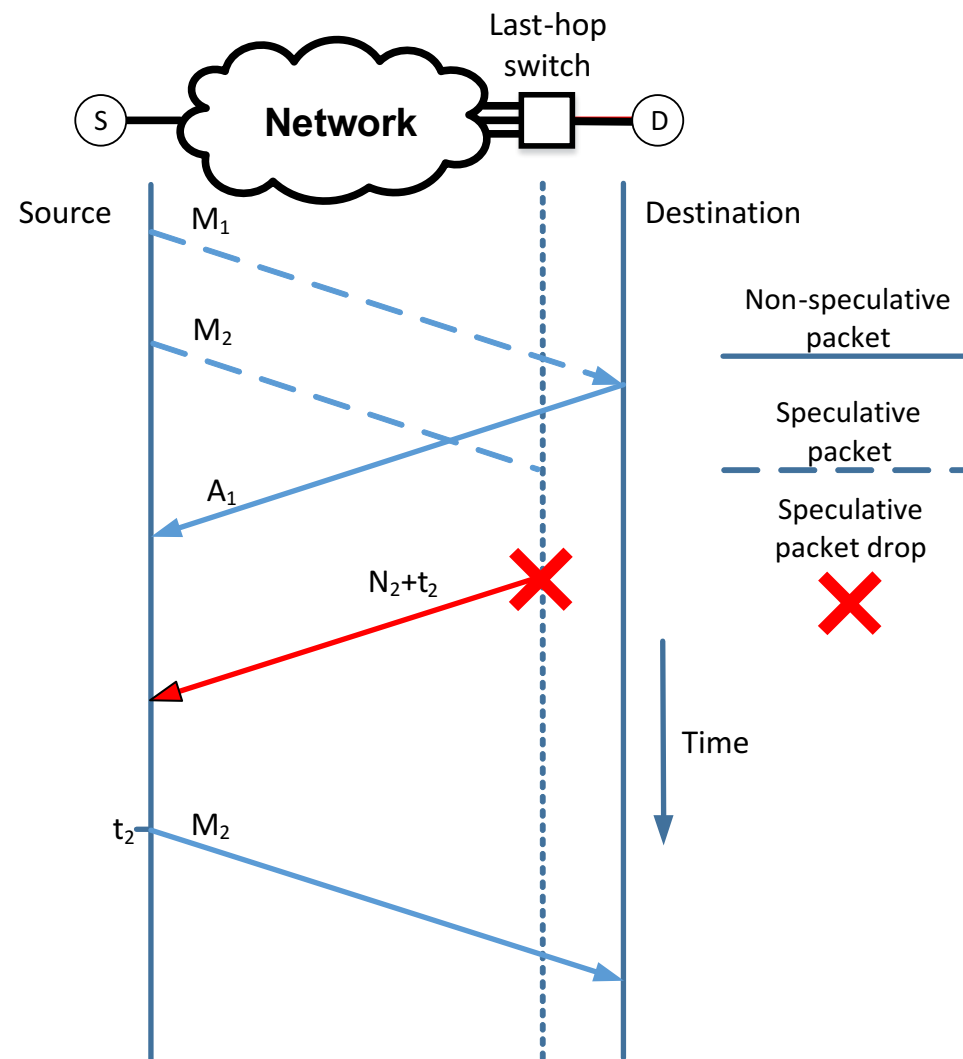
Observation: preserve ejection channel bandwidth for data packets

Move the endpoint reservation scheduler to the last-hop switch

Messages are first transmitted speculatively

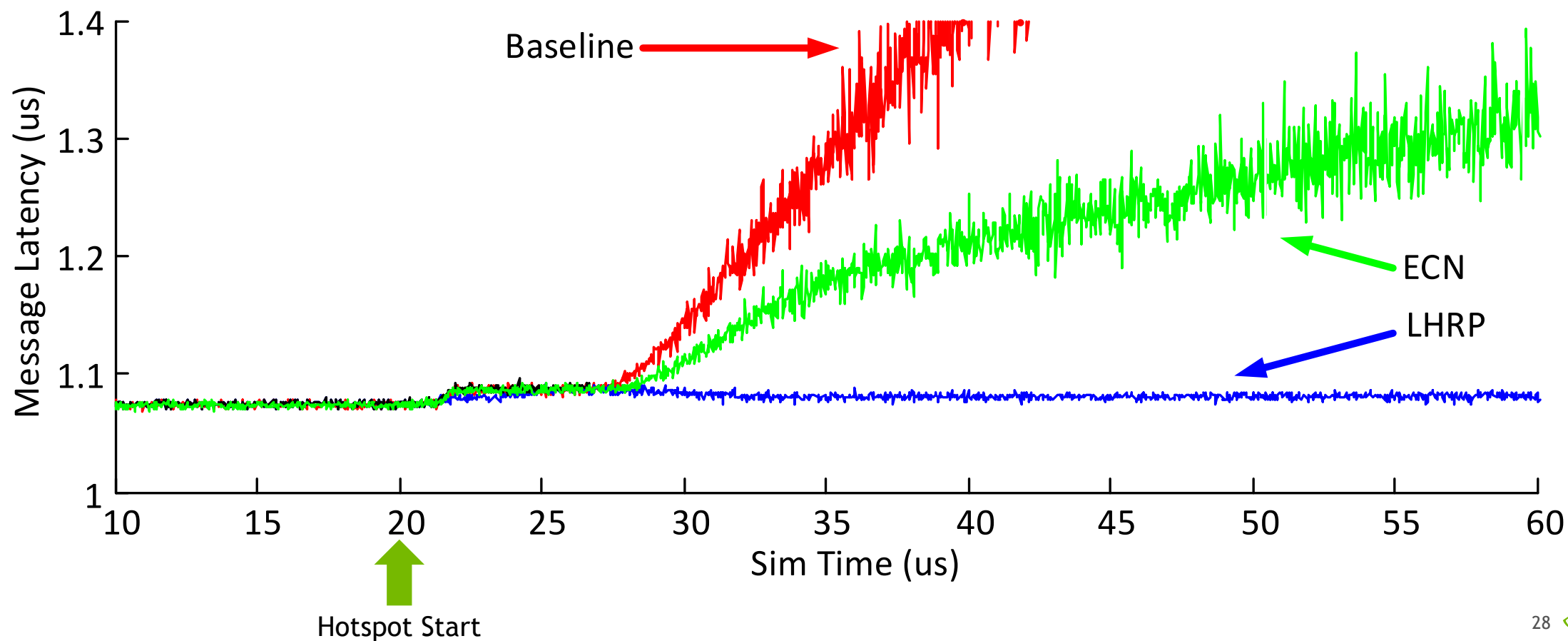
No congestion: speculative messages will arrive successfully

With congestion: speculative message is dropped by the last-hop switch, reservation is sent back with the nack



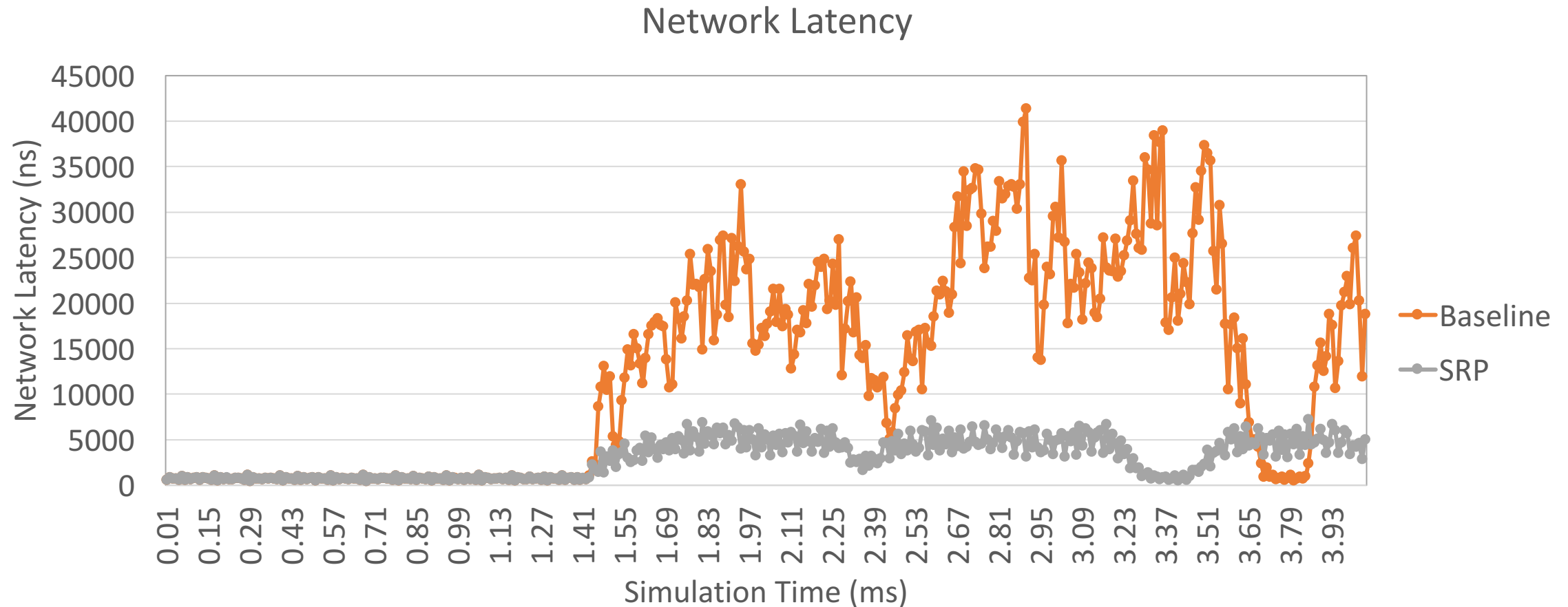
Initial Congestion Response

40% uniform random + 60:4 Hotspot @ 20us - 4 flits/message



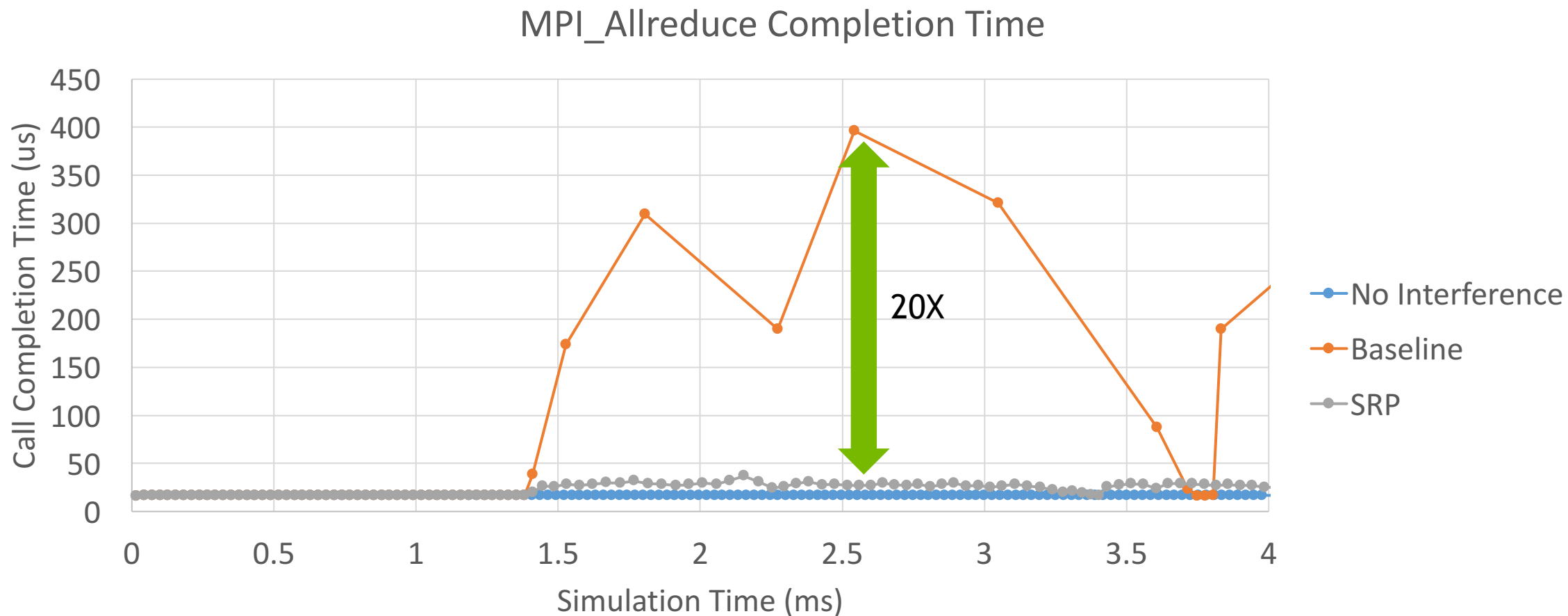
Impact of Congestion on Network Performance

1024-rank BIGFFT MPI trace on 2.5K-node dragonfly



Impact of Congestion Interference on Network Operations

1024-rank MPI_Allreduce + 1024-rank BIGFFT on 2.5K-node dragonfly



Interference, An Open Problem

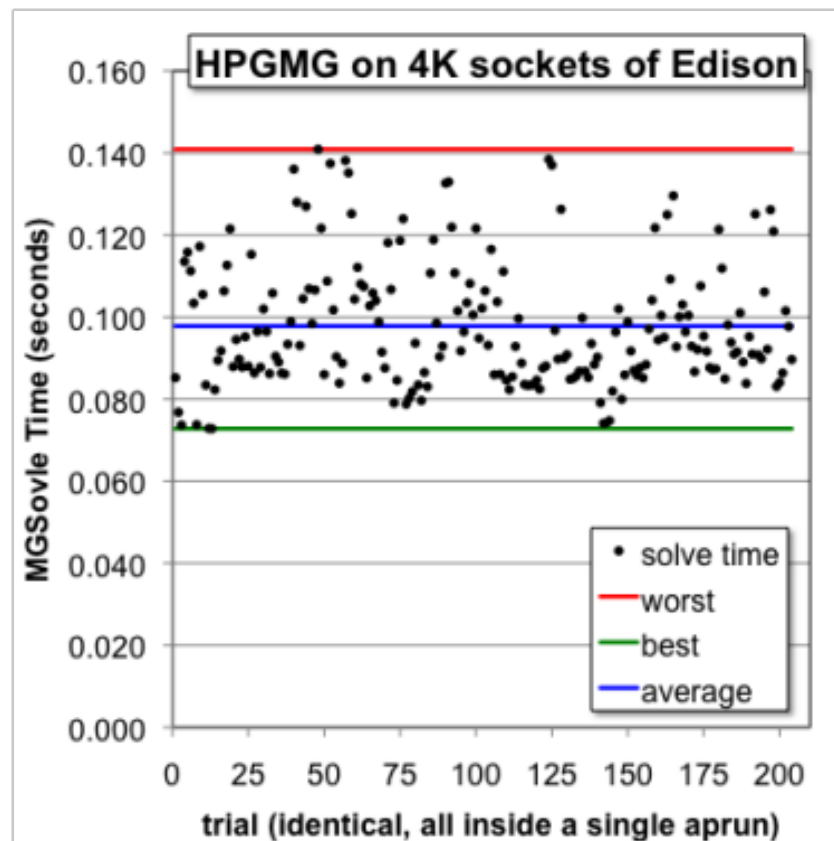
Interference



Performance Variability

F U T U R E T E C H N O L O G I E S G R O U P

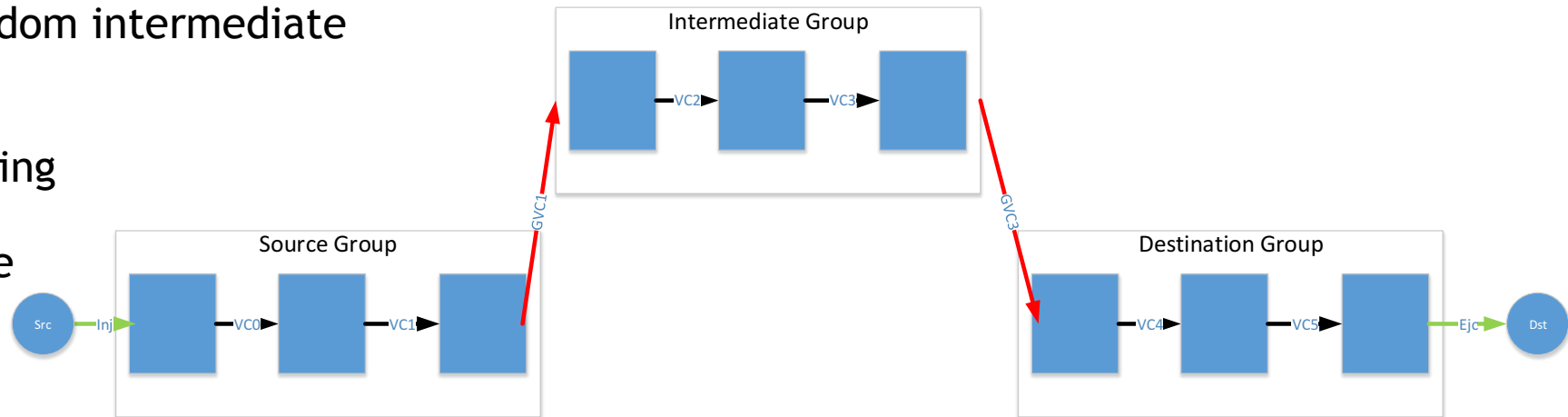
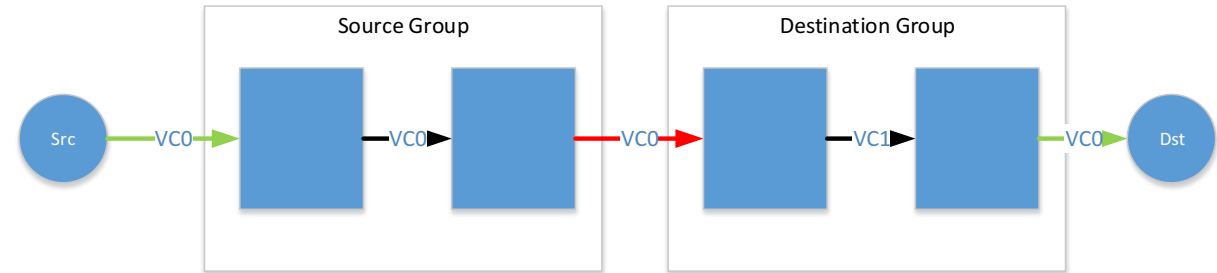
- ❖ Performance is highly variable at scale on Edison (Dragonfly)
- ❖ Actively investigating the cause with NERSC/Cray.
- ❖ Figure to the right shows the individual solves times obtained when solving
 - the exact same problem
 - always using 4K sockets (and 32K cores)
 - with the same decomposition
 - within a single aprun (while loop in single execution)
- ❖ **Average performance is 33% lower than best**



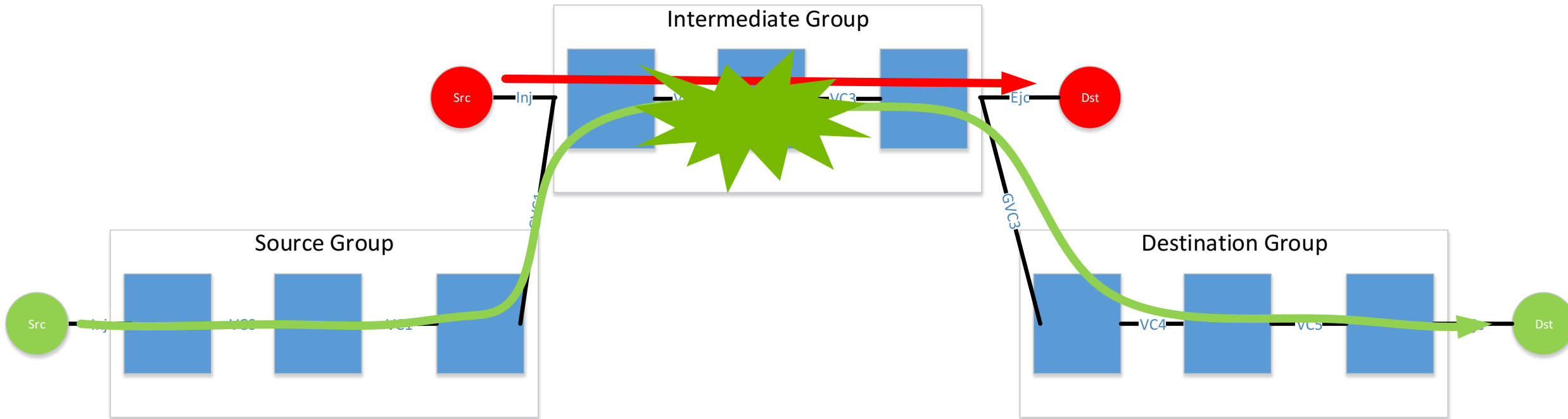
Sam Williams
Designforward Tech Talk July 2014

ADAPTIVE ROUTING

- ▶ Minimal routing is insufficient for all traffic patterns
 - ▶ Two switches shares a single local channel
 - ▶ Two groups shares a single global channel
- ▶ Utilize non-minimal network paths
 - ▶ “Bounce” off of a random intermediate switch/group
 - ▶ Creates resource sharing
 - ▶ Source of interference

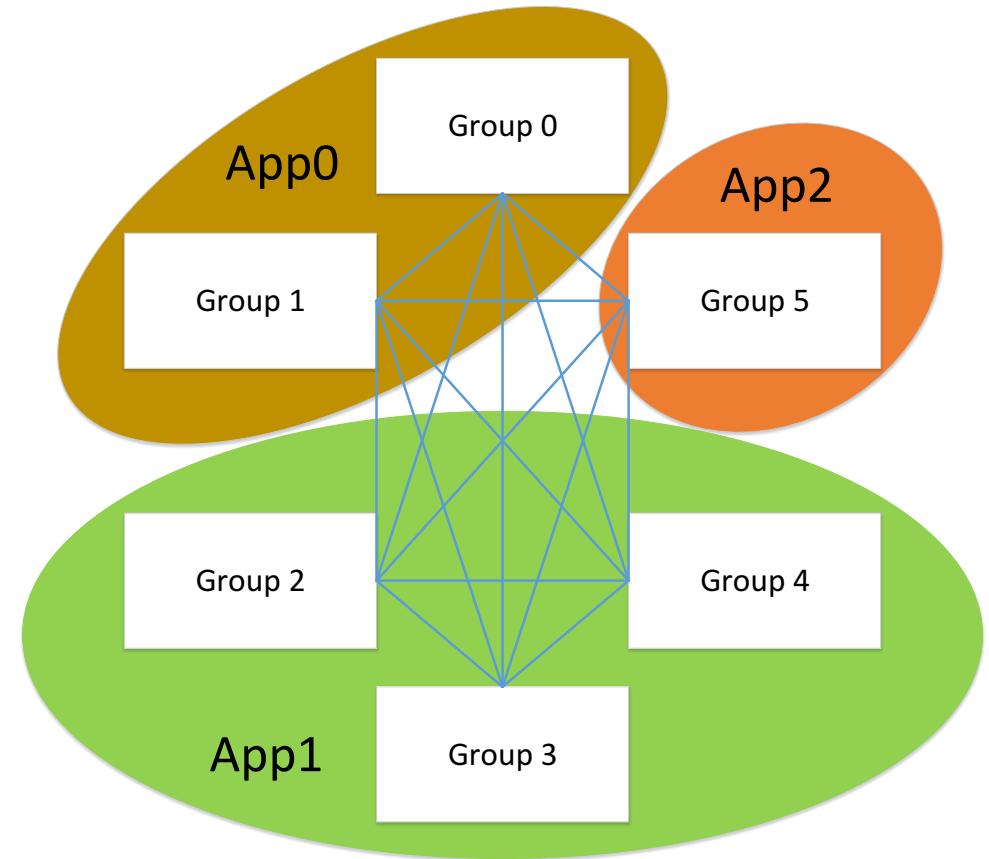


Adaptive Routing Interference



Bandwidth Partitioning

- ▶ Multiple applications running on a network each in a partition
- ▶ Each partition
 - ▶ Monitor the fraction of intra-partition adaptive traffic on the global links
 - ▶ Adjust the adaptive routing bias to maintain 50% fraction
 - ▶ Simpler alternative to physically partitioning network
- ▶ Partitions are not perfectly isolated
 - ▶ Subject to transient traffic variations
 - ▶ Adaptive routing reaction time

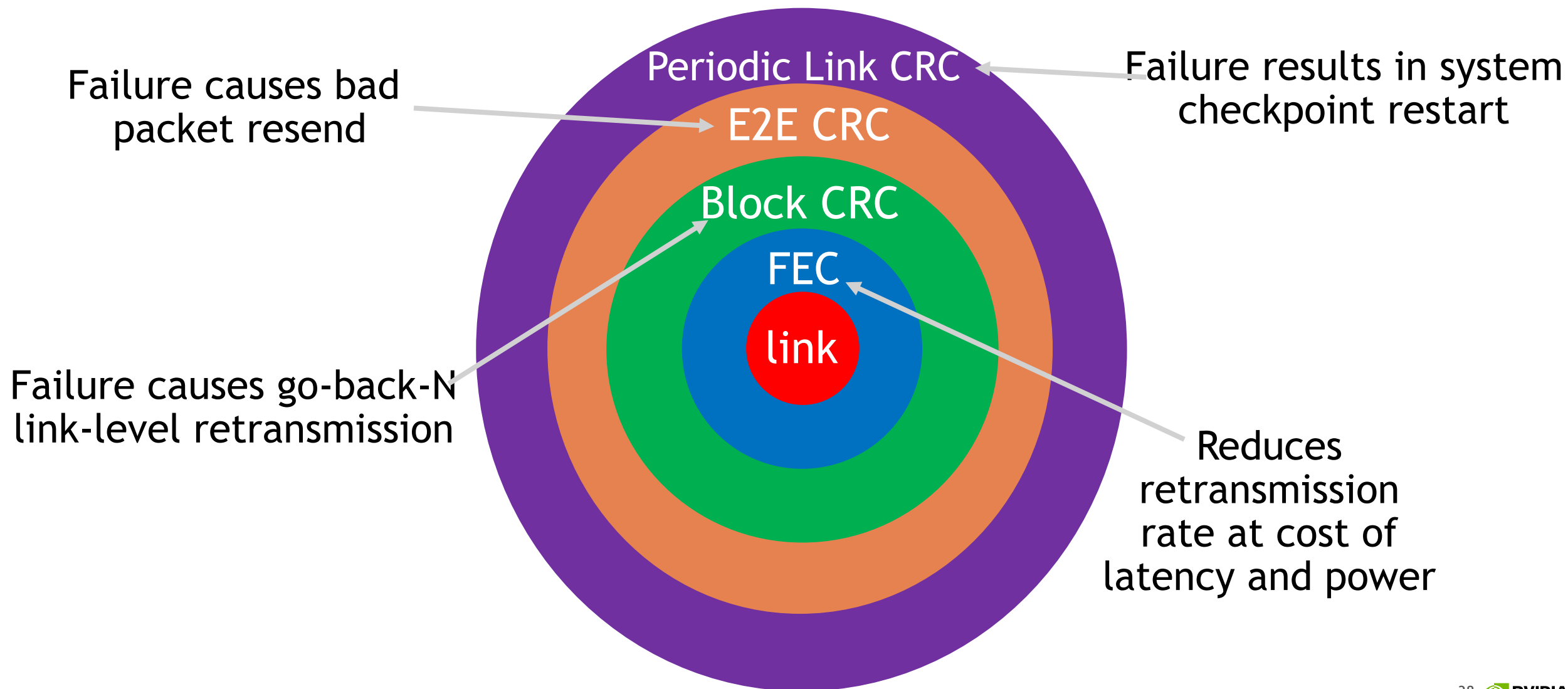


Error Control

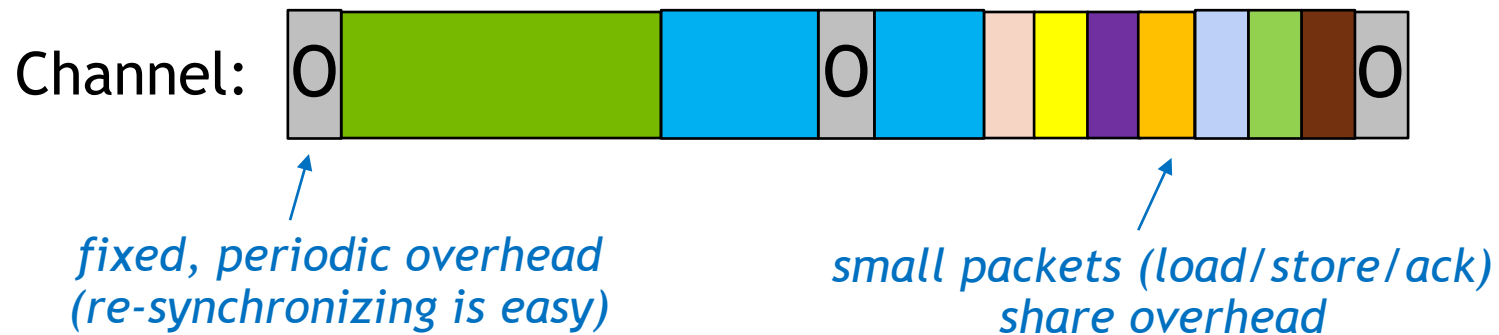
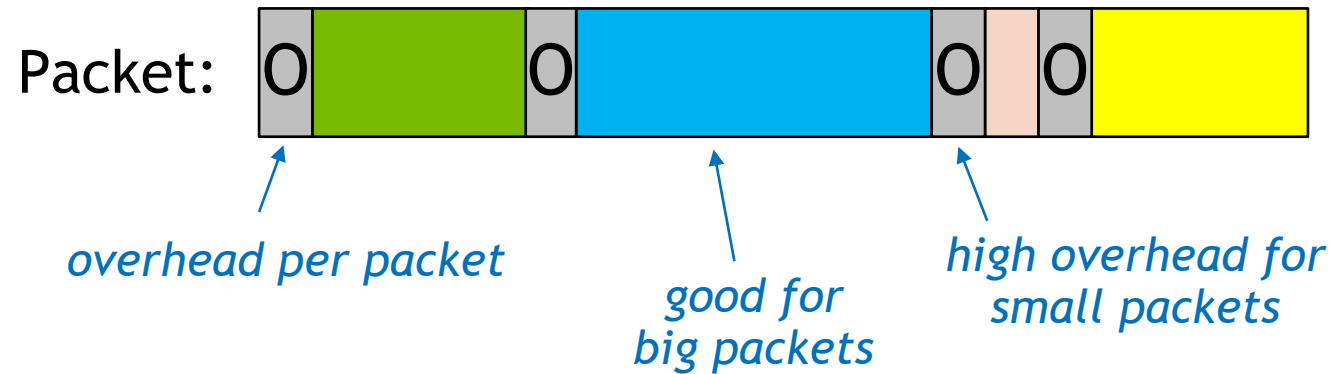
Error Control Problem

- $\sim 10^6$ links in an Exascale system
- Bandwidth of 2×10^{11} b/s each (total bandwidth of $\sim 10^{17}$ b/s)
- Bit error rate of 10^{-4} to 10^{-12} (total error rate of 10^5 - 10^{13} errors/s)
- System wide error rate of 10^{-4} errors/s (1 week MTBF)
- Spec network error rate at 10^{-5} errors/s (order of magnitude less)
- Requires bit error rate of 10^{-22}

Layers of hardware protection



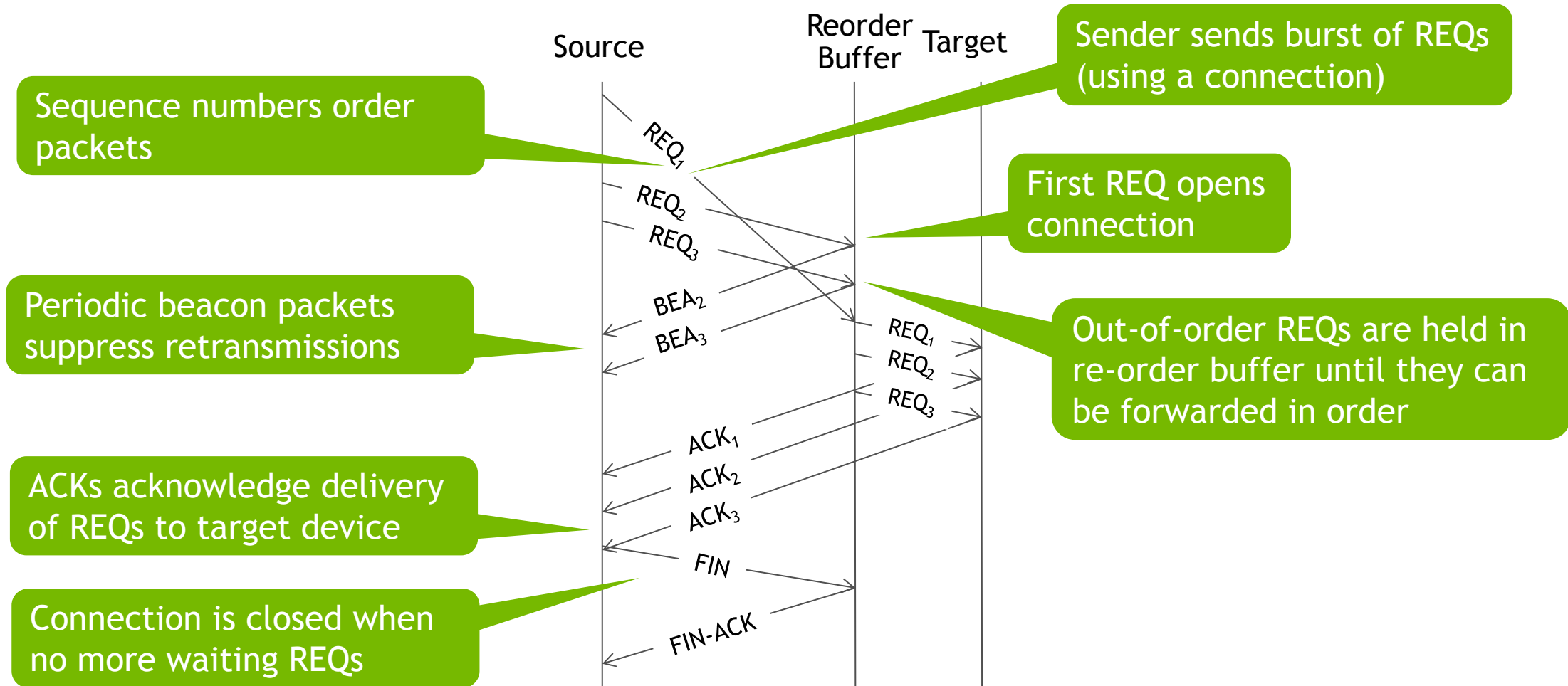
Channel vs packet protocol



Ordering

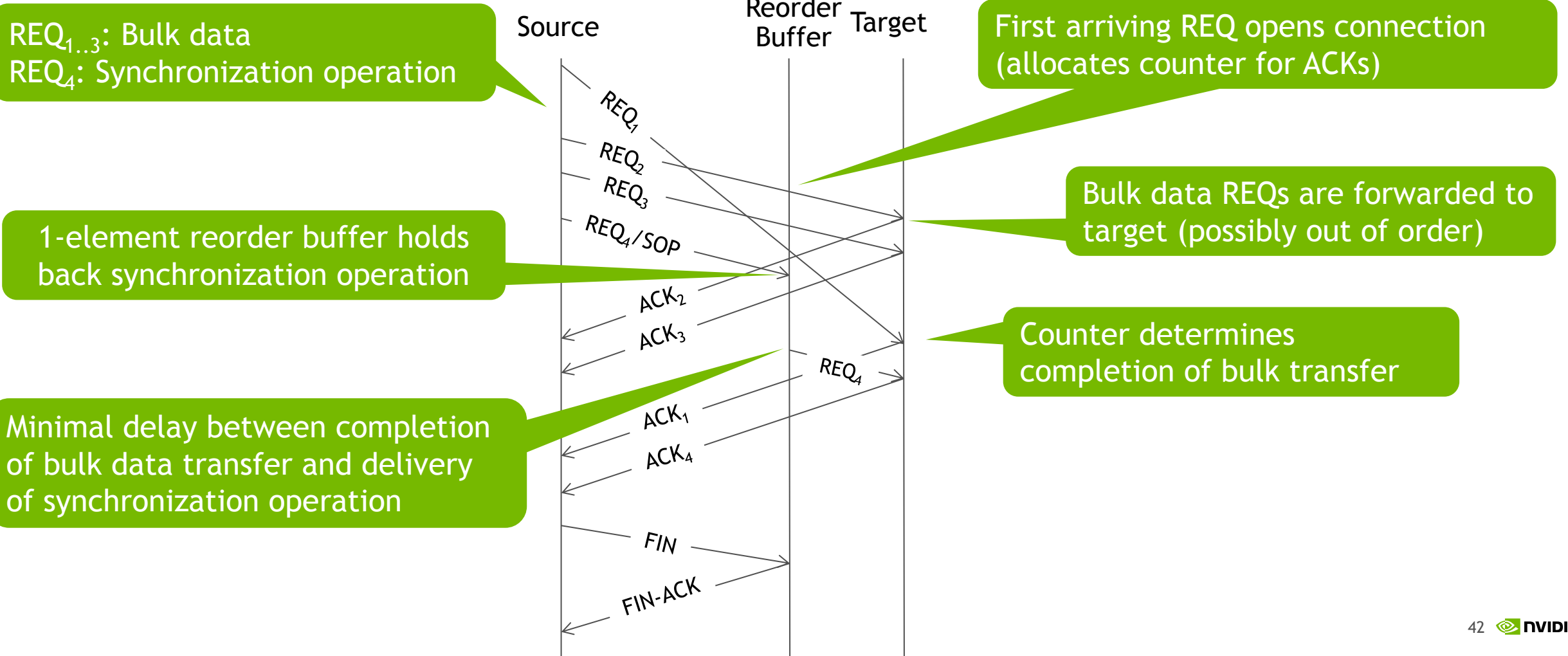
Ordered Transfer Protocol

Lightweight Connection with Low-Overhead Setup and Teardown



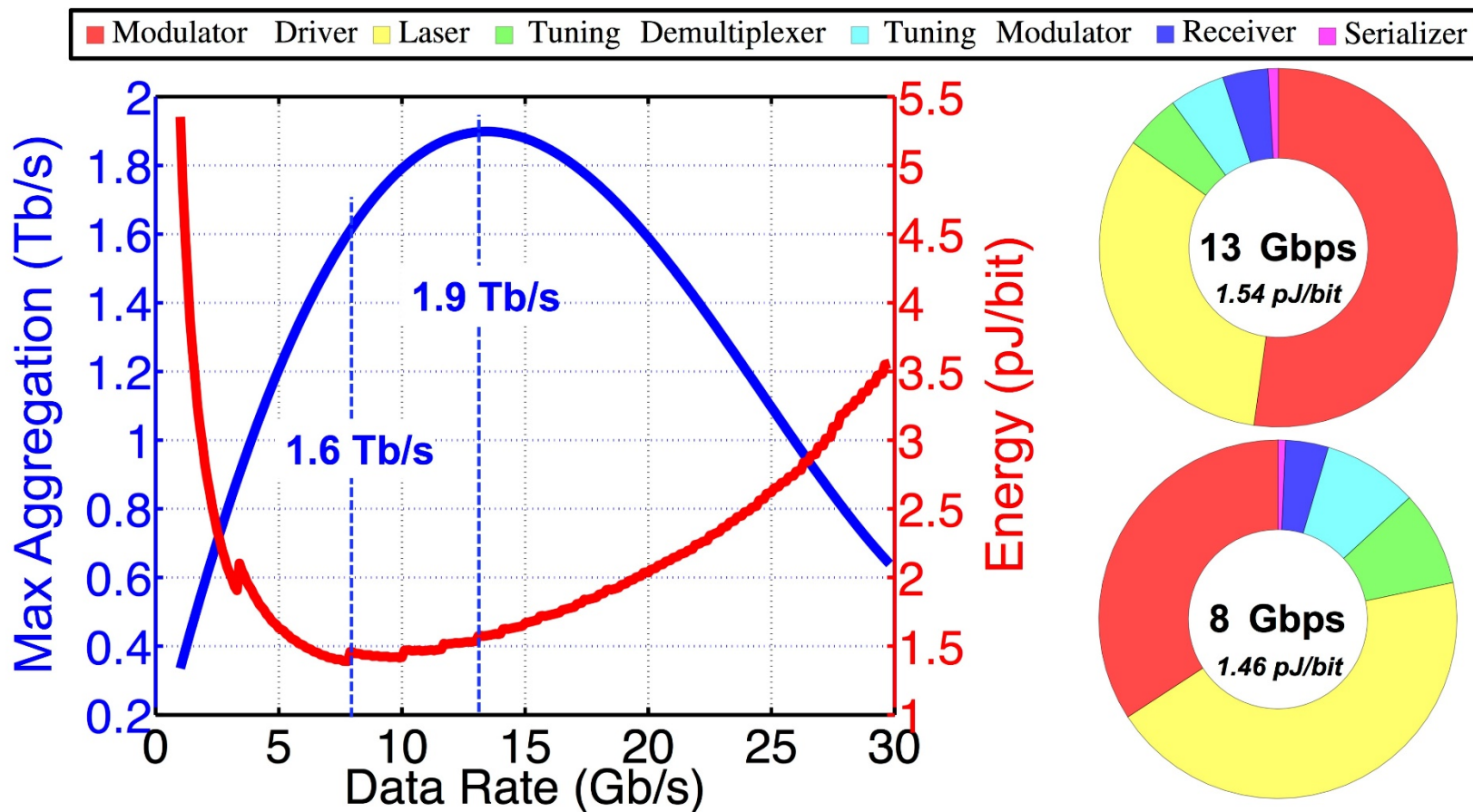
Synchronized Bulk Transfer

Strict Ordering is not Required for many Producer/Consumer Exchanges



The Role of Photonics

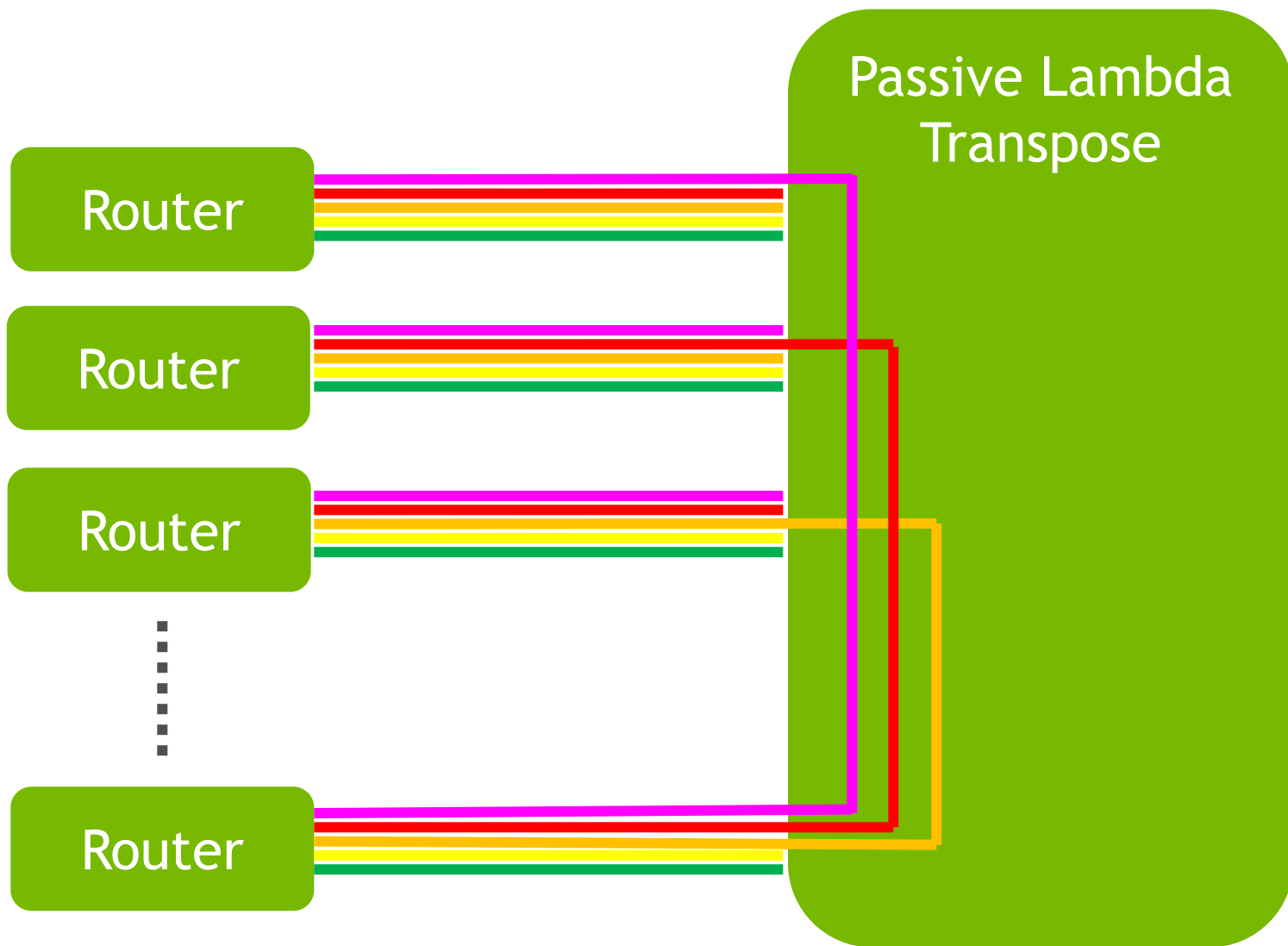
Power/Bandwidth Density



200 lambdas @ 8Gb/s

146 lambdas @ 13Gb/s

Photonic Dragonfly Concept



- 256 Groups
- 16 Fiber Bundles per group
- $16n$ wavelengths per fiber
- 16 Central AWGRs
- Much simpler cable management
- Technology not sufficiently mature for Exascale (2021)
 - Maybe by 2025

Overall System Sketch

Recall Costs



\$500

100m



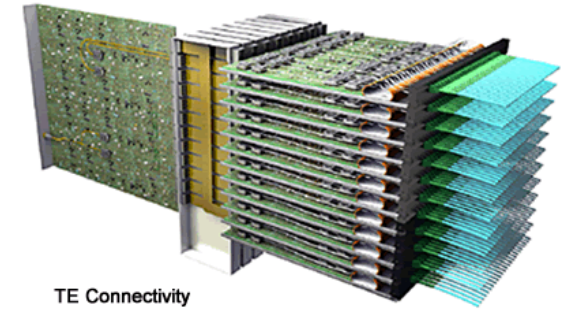
\$50

5m



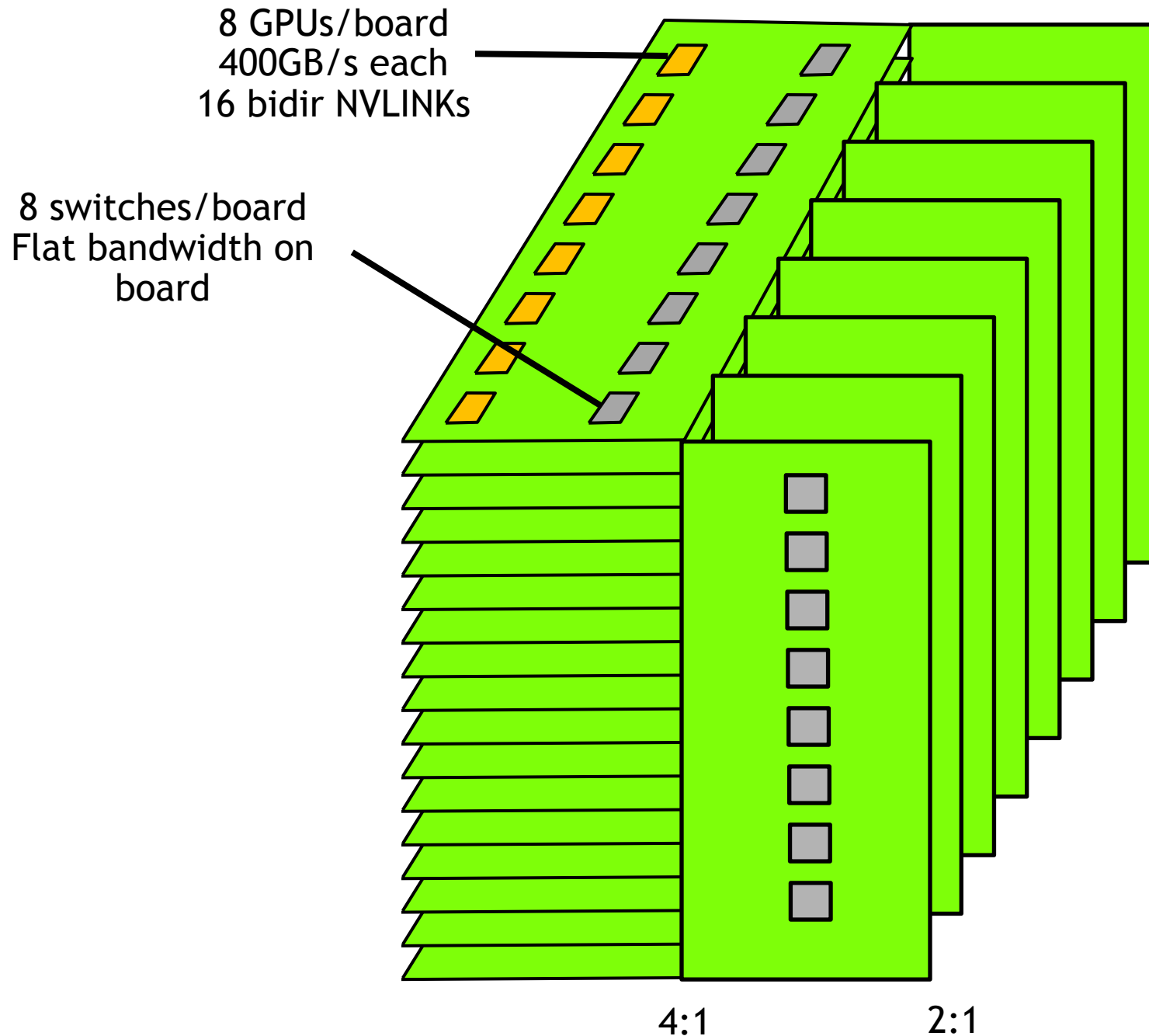
\$10

1m



\$5

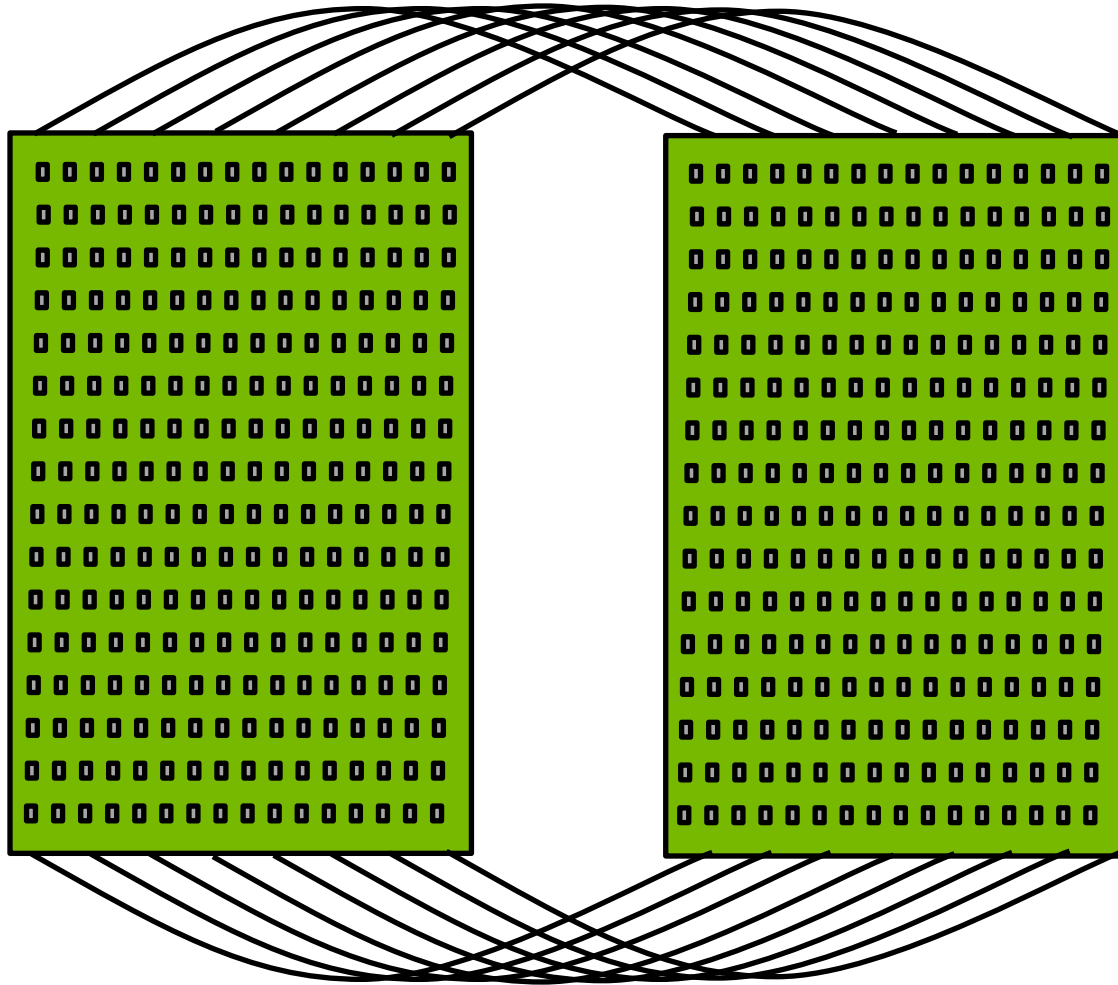
0.3m



Cabinet

- 128 GPUs, 50kW
- 16 GPU x 8 switch boards
- 400GB/s bidir per GPU
 - 3.2TB/s per board
- 192 switches
 - 8 per board
- Flat bandwidth on board
- 100GB/s at crosspoints
 - 32 pairs
 - 12.8TB/s aggregate
- 100GB/s per GPU within cabinet
- All connections electrical

Group



- 1 or 2 cabinets
- Electrical Flex cables between cabinets
- 12.8TB/s between cabinets
 - 512 NVLinks
 - 4096 pairs
 - 64 cables, 64 pairs each
- 256 NVLinks out back of each cabinet - 2 per GPU
 - 50GB/s per GPU global bandwidth
 - 6.4TB/s per cabinet
- Up to 513 groups
 - 131,328 GPUs

System Sketch

- Cost dominated by AOCs - 50GB/s per GPU \$2K per endpoint
- Taper by leaving half the cables out - 25GB/s global bandwidth per GPU
 - Limits maximum system size to 64k GPUs
- Routing - progressive adaptive routing with local misroute (6VCs per class)
- Flow control - flit-level flow control with LHRP
- Error control
 - Channel-level CRC for link
 - Packet-level CRC for ETE
 - FEC for optical cables only (where needed)
- Ordering - per packet or bulk sync
- PGAS support, with two-stage address translation

Conclusion

Conclusion

- An **Exascale** network is not business as usual
- Need fine-grain communication (**PGAS**) for strong scaling
 - **Two-stage** address translation
 - **10^5 outstanding references** per endpoint
 - **Ordering** as needed
- Cost-efficient bandwidth
 - **Topology** driven by communication cost - Dragonfly
 - High **payload efficiency** for small packets (32B)
 - **Congestion avoidance** & **adaptive routing** allows links to operate near capacity
- Error control
 - **Channel-level**, not packet-level CRC
- System sketch
 - 8:2:1 - Board:Cabinet:Global bandwidth taper
 - Cost is dominated by AOCs

Backup - Not in main Talk

