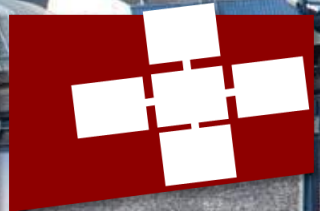


T. HOEFLER

# The three L's in modern high-performance networking: low latency, low cost, low processing load

with support of M. Besta, S. Di Girolamo, K. Taranov @ SPCL -- R. Grant, R. Birghtwell @ Sandia Natl. Labs  
keynote at HiPINEB @ HPCA, Vienna, Austria in February 2017





Low Latency

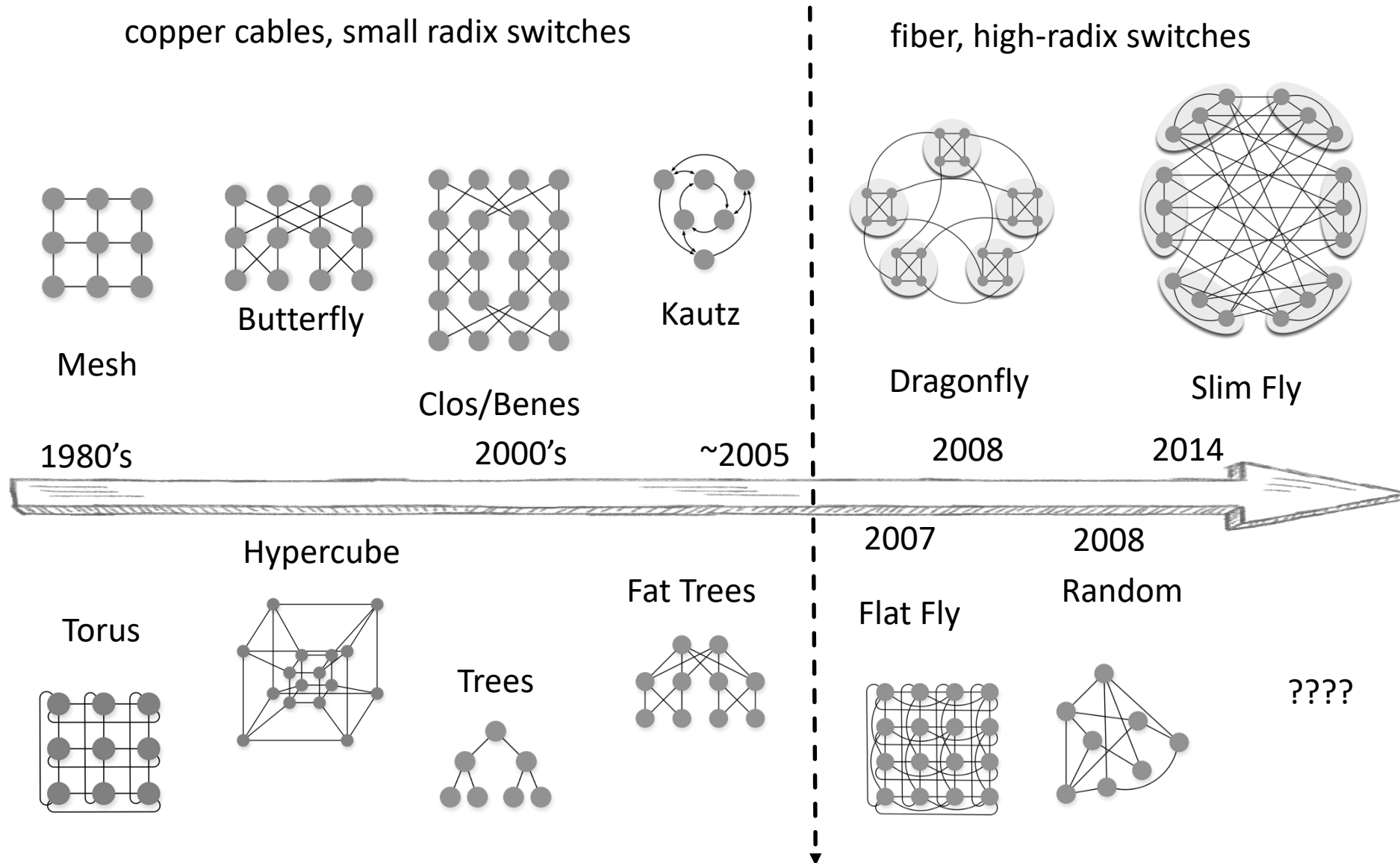


Low Cost



Low Processing Load

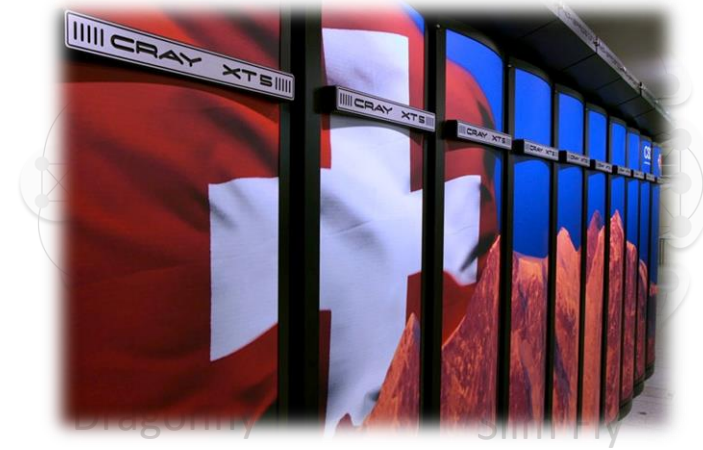
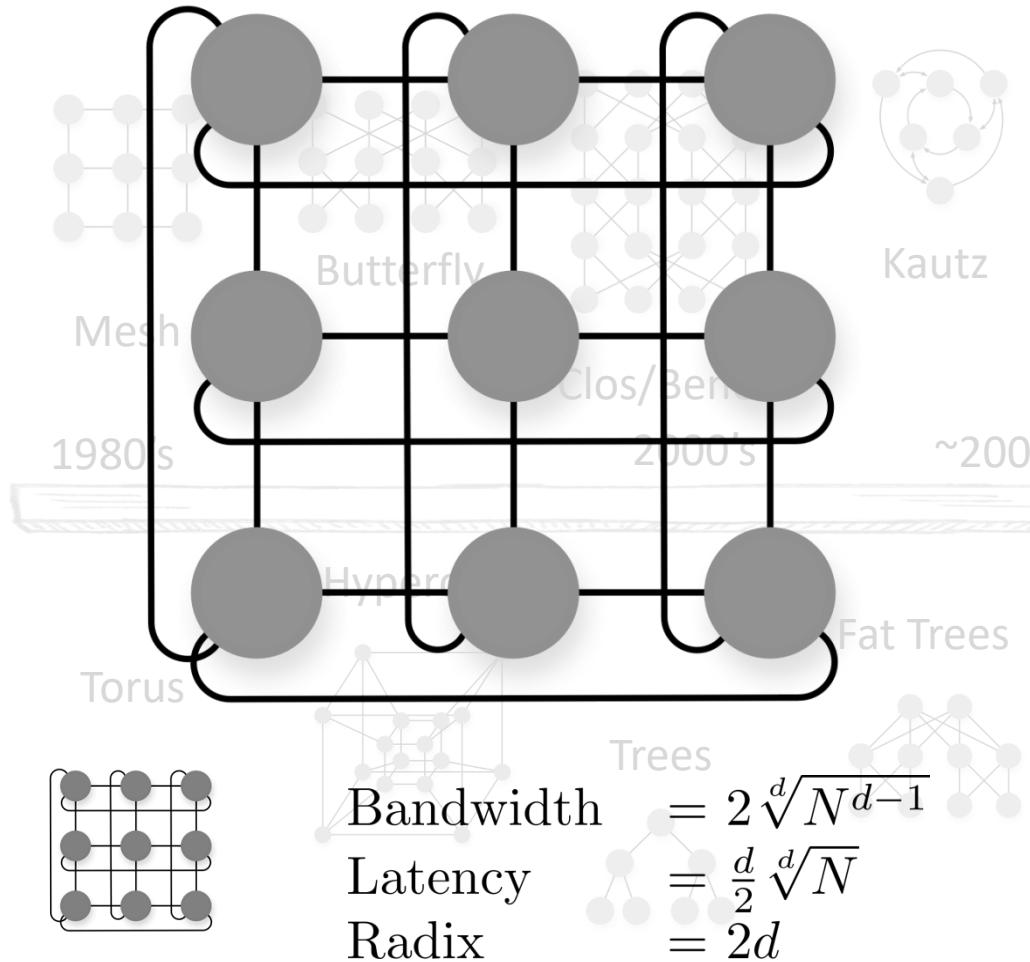
# A BRIEF HISTORY OF NETWORK TOPOLOGIES



# A BRIEF HISTORY OF NETWORK TOPOLOGIES

copper cables, small radix switches

fiber, high-radix switches



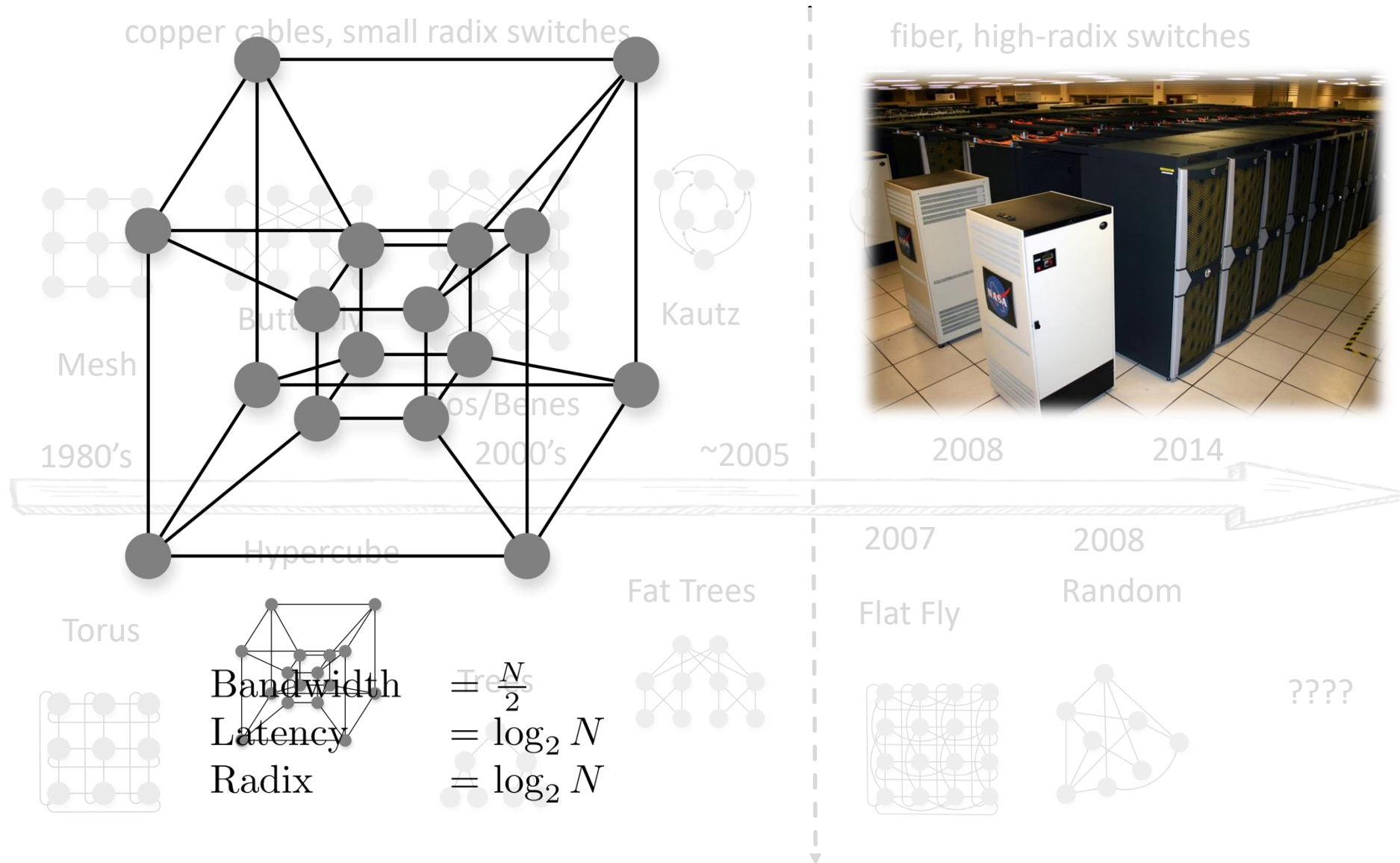
2008

2014



2014

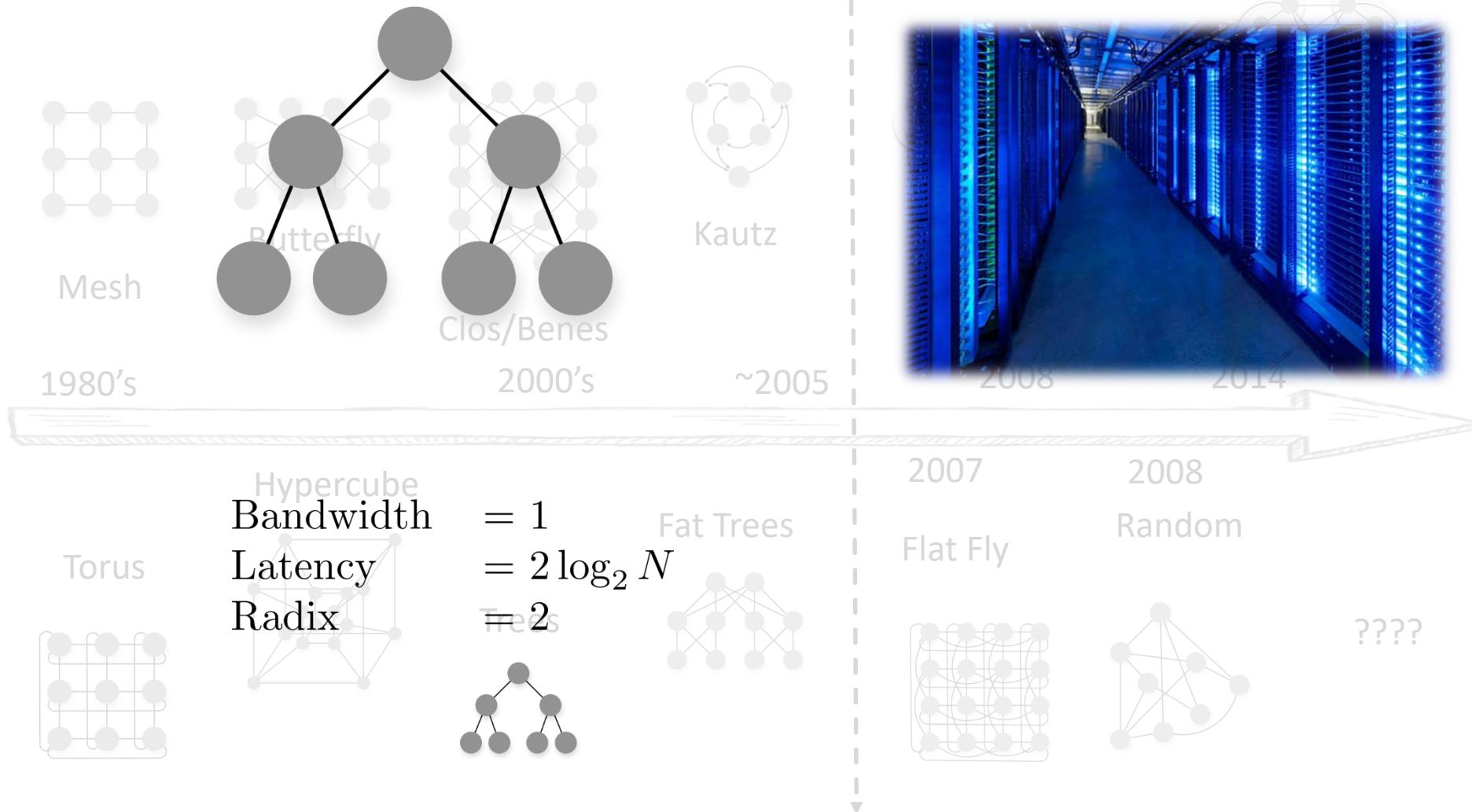
# A BRIEF HISTORY OF NETWORK TOPOLOGIES



# A BRIEF HISTORY OF NETWORK TOPOLOGIES

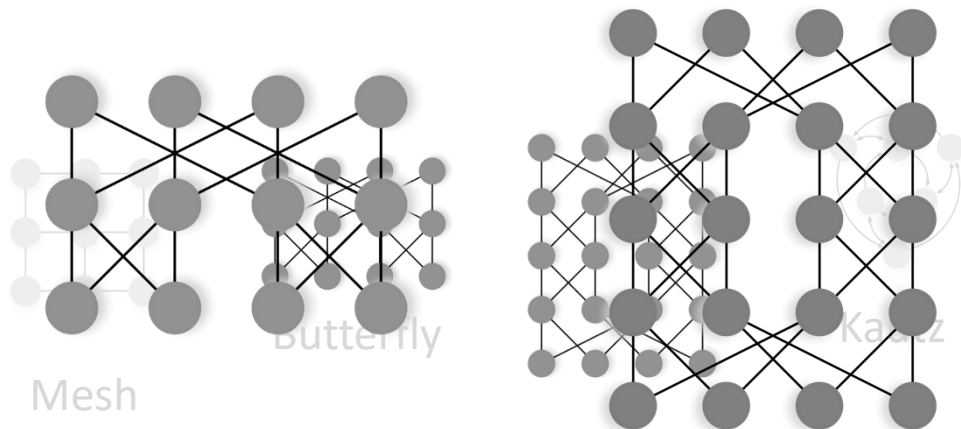
copper cables, small radix switches

fiber, high-radix switches



# A BRIEF HISTORY OF NETWORK TOPOLOGIES

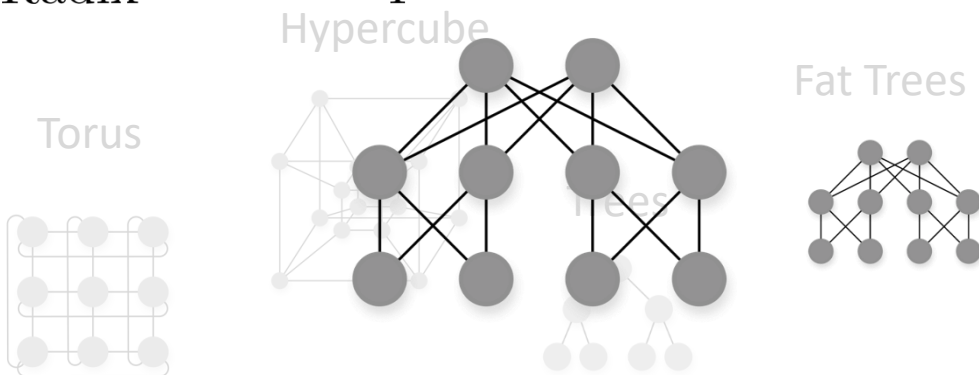
copper cables, small radix switches



Bandwidth =  $\frac{N}{2}$

Latency =  $2 \log_2 N$

Radix = 4



Dragonfly

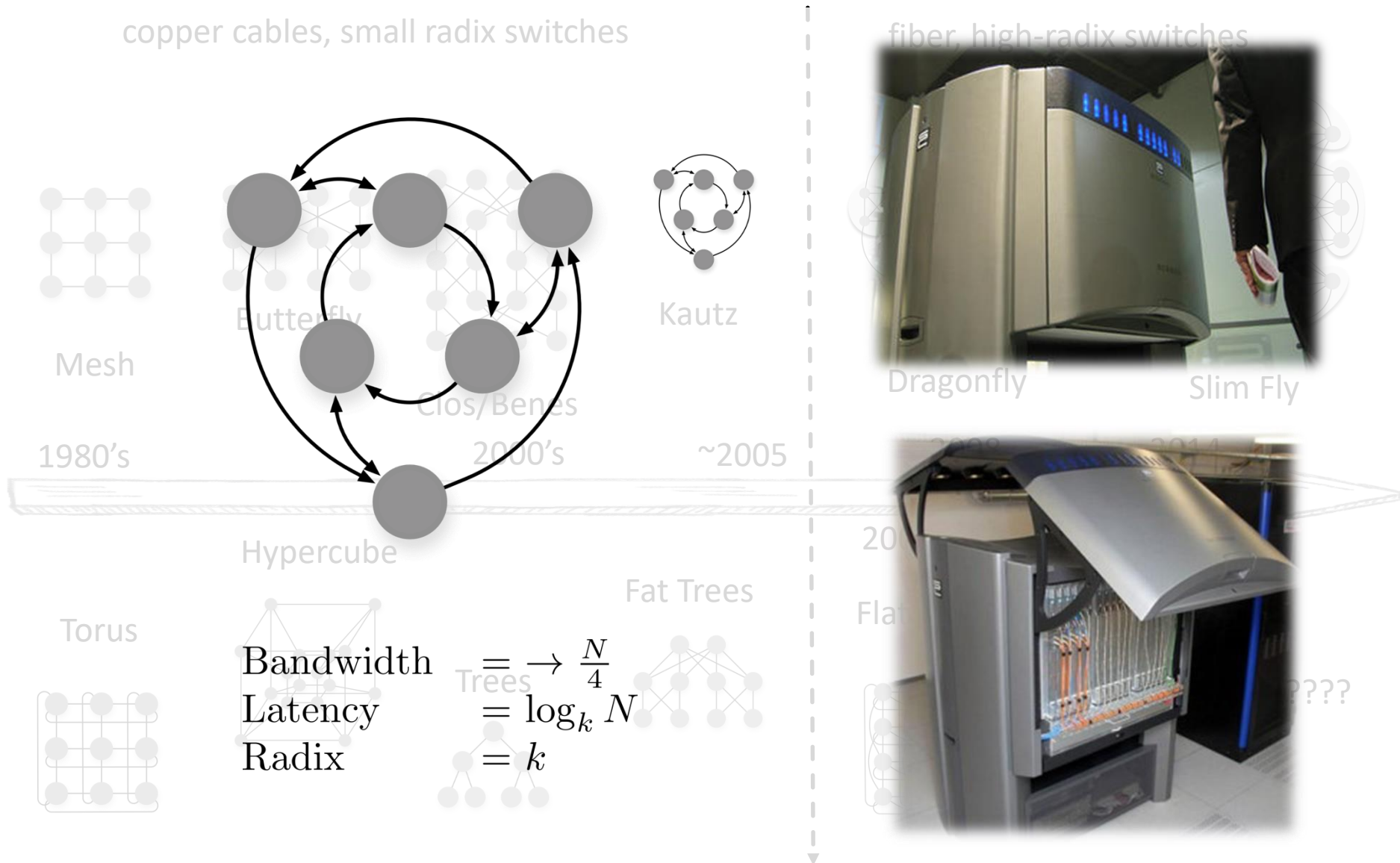
Slim Fly

2008

2014



# A BRIEF HISTORY OF NETWORK TOPOLOGIES

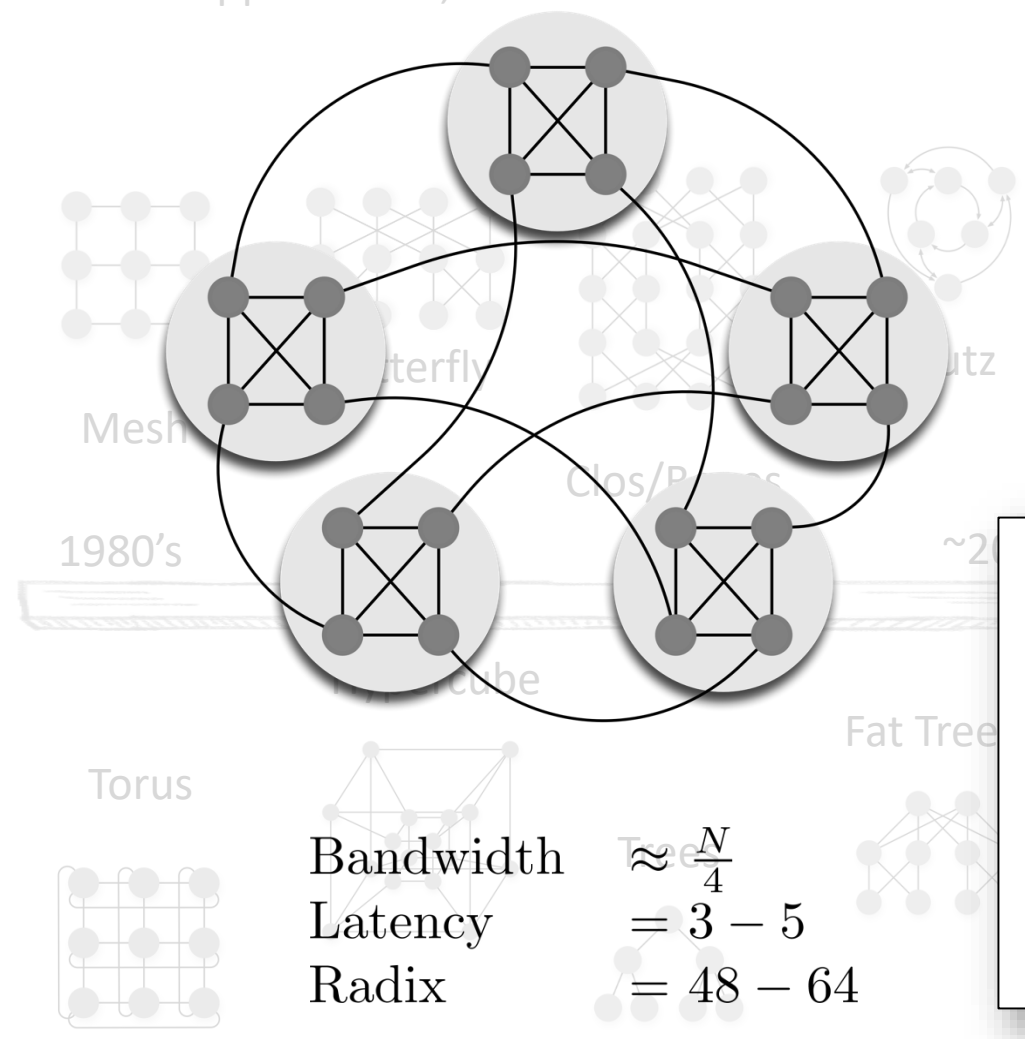




# A BRIEF HISTORY OF NETWORK TOPOLOGIES

copper cables, small radix switches

fiber, high-radix switches



2010 18th IEEE Symposium on High Performance Interconnects

### The PERCS High-Performance Interconnect

Baba Arimilli \*, Ravi Arimilli \*, Vicente Chung \*, Scott Clark \*, Wolfgang Denzel †, Ben Drerup \*, Torsten Hoefler ‡, Jody Joyner \*, Jerry Lewis \*, Jian Li †, Nan Ni \* and Ram Rajamony †

\* IBM Systems and Technology Group, 11501 Burnet Road, Austin, TX 78758  
 † IBM Research (Austin, Zurich), 11501 Burnet Road, Austin, TX 78758  
 ‡ Blue Waters Directorate, NCSA, University of Illinois at Urbana-Champaign, Urbana, IL 61801  
 E-mail: arimilli@us.ibm.com, rajamony@us.ibm.com, htor@illinois.edu

**Abstract**—The PERCS system was designed by IBM in response to a DARPA challenge that called for a high-productivity high-performance computing system. A major innovation in the PERCS design is the network that is built using Hub chips that are integrated into the compute nodes. Each Hub chip is about 580 mm<sup>2</sup> in size, has over 3700 signal I/Os, and is packaged in a module that also contains LGA-attached optical electronic devices.

The Hub module implements five types of high-bandwidth interconnects with multiple links that are fully-connected with a high-performance internal crossbar switch. These links provide over 9 Tbits/second of raw bandwidth and are used to construct a two-level direct-connect topology spanning up to tens of thousands of nodes.

bandwidths do not scale accordingly. For instance, while High Performance Linpack performance [5], [10] shows a steady improvement over time, interconnect-intensive metrics such as G-RandomAccess and G-FFTE [5] show very little improvement.

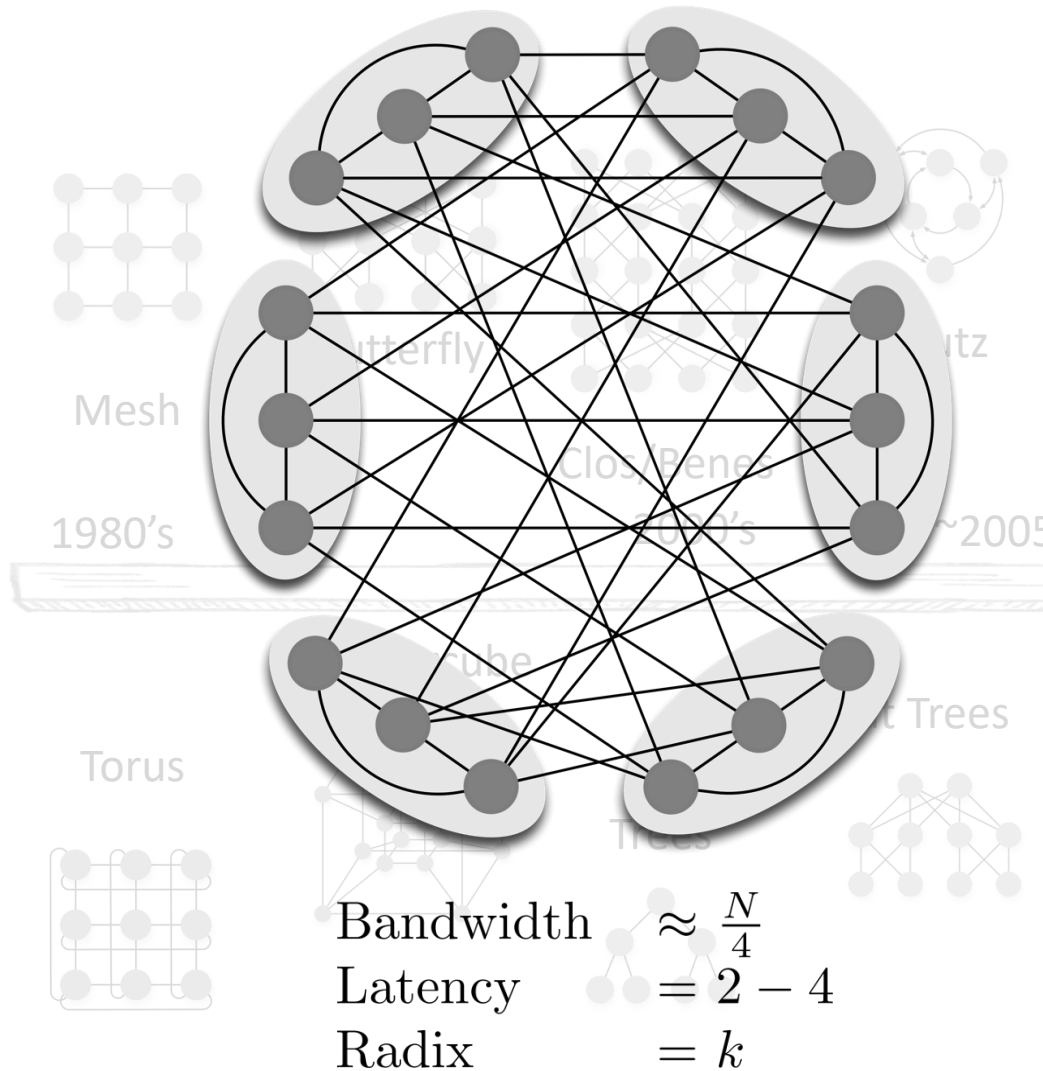
The challenge of building a high-performance, highly productive, multi-Petaflop system forced us to recognize early on that the entire infrastructure had to scale along with the microprocessor's capabilities. A significant component of our scaling solution is a new switchless interconnect with very high fanout organized into a two-level direct connect

Bandwidth  $\approx \frac{N}{4}$   
 Latency  $= 3 - 5$   
 Radix  $= 48 - 64$

# A BRIEF HISTORY OF NETWORK TOPOLOGIES

copper cables, small radix switches

fiber, high-radix switches



Key ideas:

**“It’s the diameter, stupid”**

**Lower diameter:**

- Less cables traversed
- Less cables needed
- Less routers needed

**Cost and energy savings:**

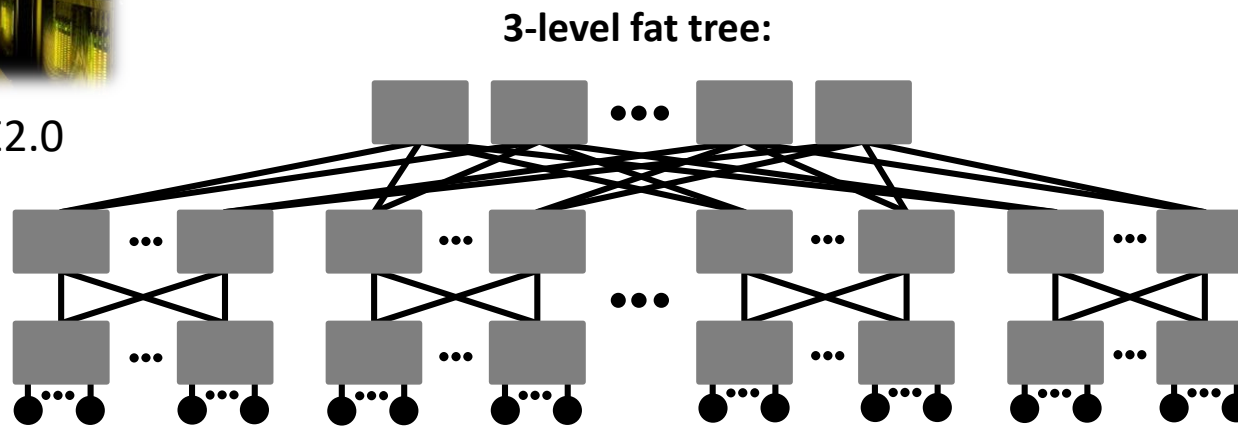
- Up to 50% over Fat Tree
- Up to 33% over Dragonfly

# DESIGNING A LOW-DIAMETER NETWORK TOPOLOGY

EXAMPLE: FULL-BANDWIDTH FAT TREE VS HOFFMAN-SINGLETON GRAPH

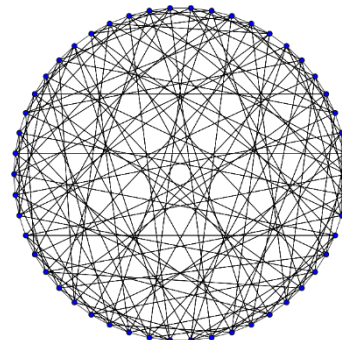


TSUBAME2.0



diameter = 4

Slim Fly based on the Hoffman-Singleton Graph [1]:



diameter = 2  
> ~50% fewer routers  
> ~30% fewer cables

# LIMITS ON LOW DIAMETER NETWORK TOPOLOGIES



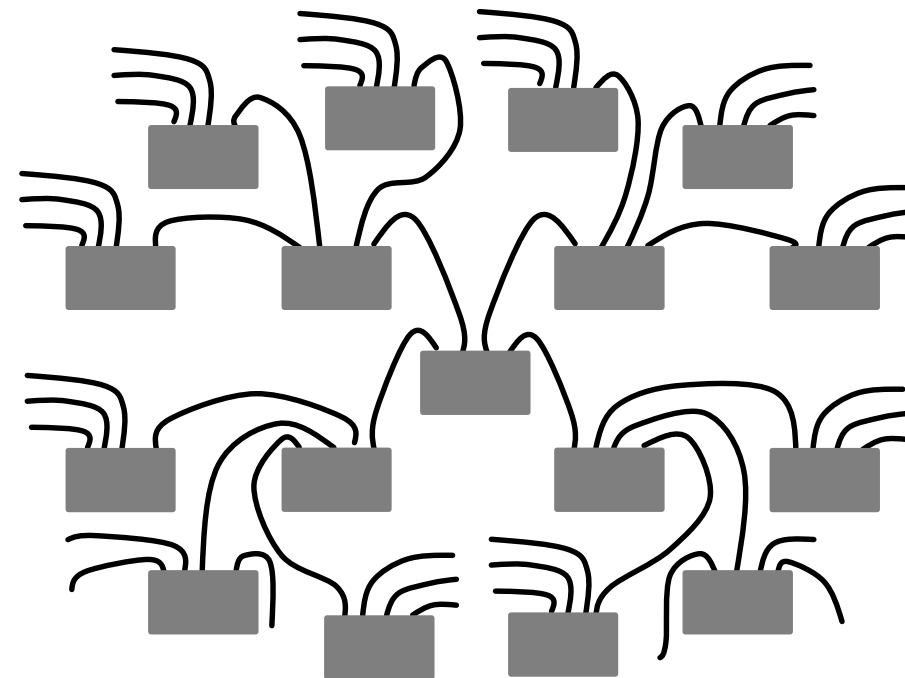
## Key method

### Optimize towards the Moore

**Bound [1]:** the upper bound on the *number of vertices* in a graph with given *diameter*  $D$  and *radix*  $k$ .

$$MB(D, k) = 1 + k + k(k - 1) + k(k - 1)^2 + \dots$$

$$MB(D, k) = 1 + k \sum_{i=0}^{D-1} (k - 1)^i$$

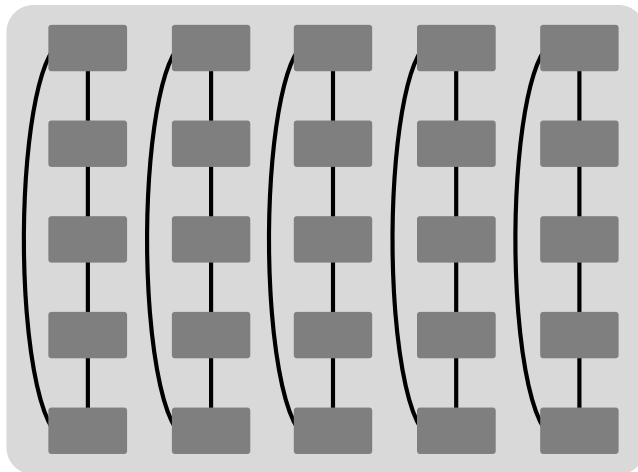


# THE SLIM FLY PRINCIPLE – APPROACHING THE MOORE BOUND!

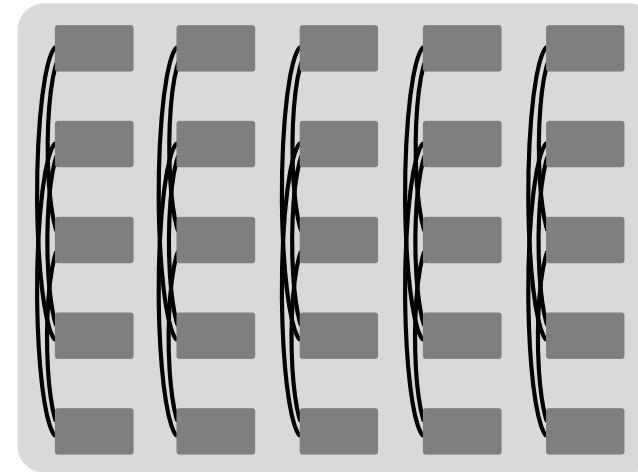
## CONNECTING ROUTERS: DIAMETER 2

Example Slim Fly design for *diameter* = 2: *MMS graphs* [1]

A subgraph with identical groups of routers

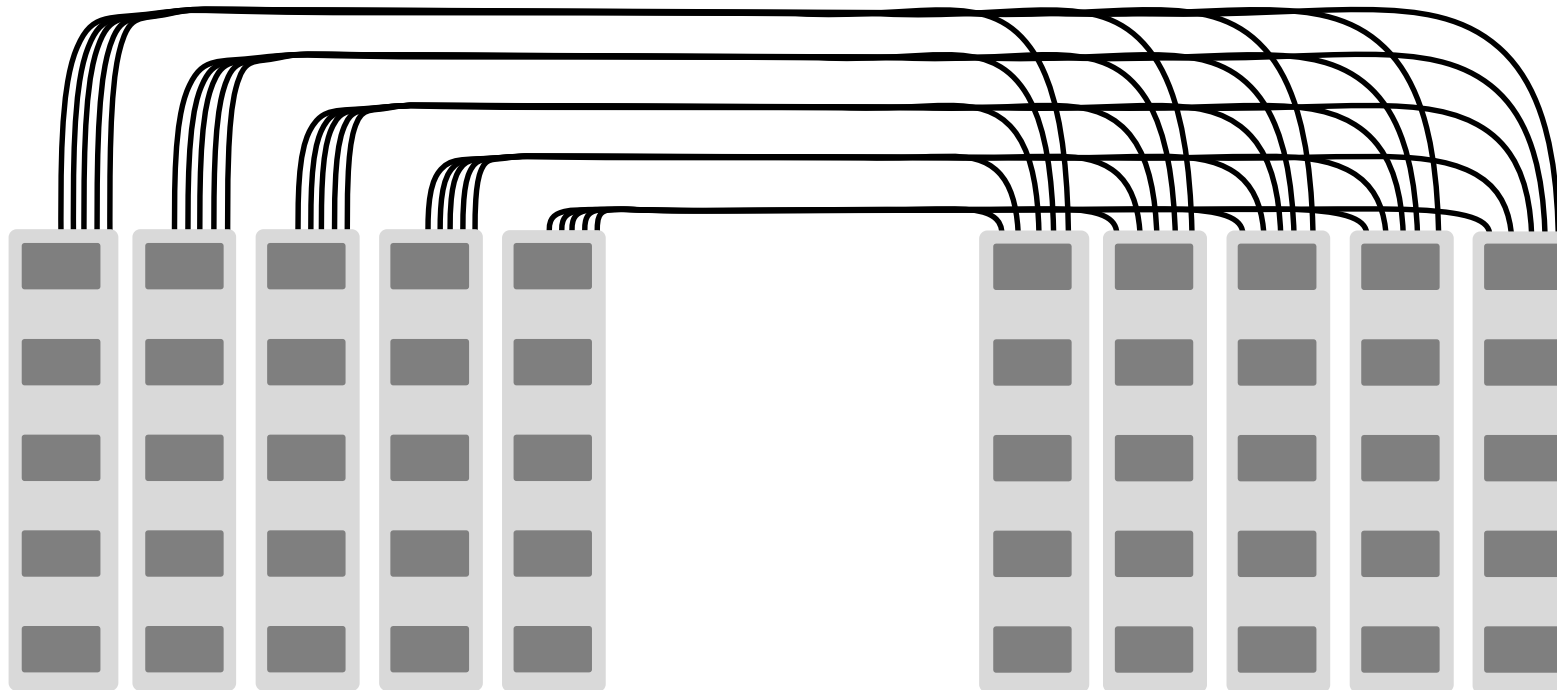


A subgraph with identical groups of routers



# THE SLIM FLY PRINCIPLE – APPROACHING THE MOORE BOUND!

CONNECTING ROUTERS: DIAMETER 2



Groups form a fully-connected bipartite graph



Low Latency



Low Cost

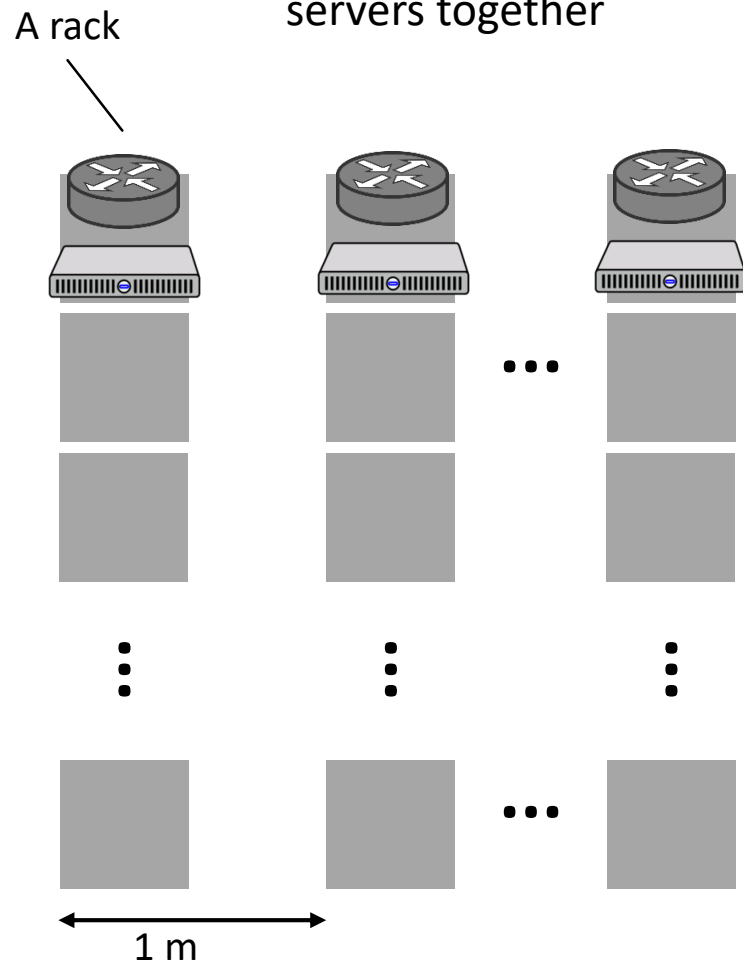


Low Processing Load

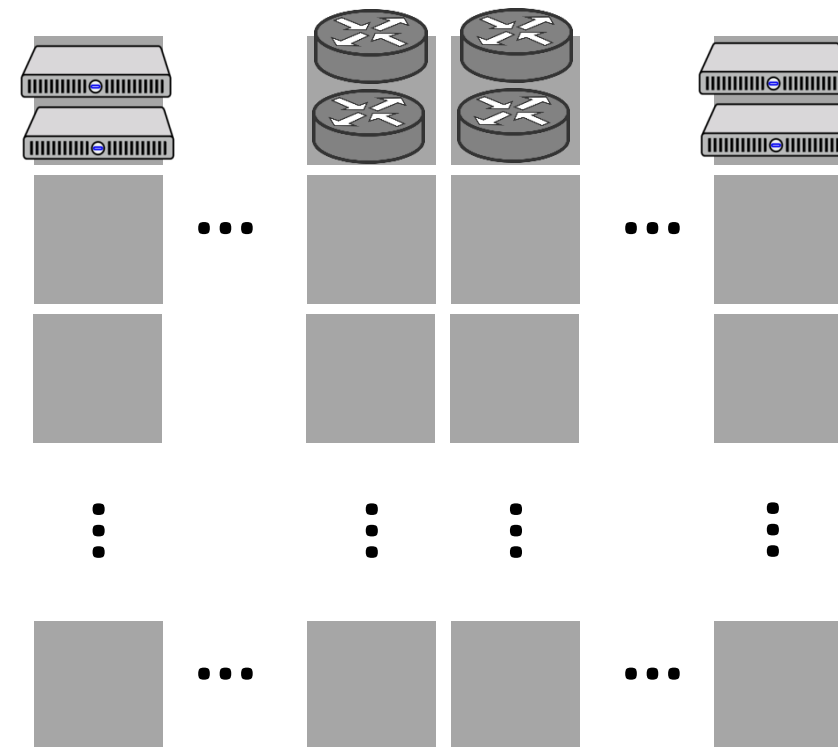
# SLIM FLY MMS - COST COMPARISON

## COST MODELS: VARIANTS

**Variant 1: Routers and servers together**



**Variant 2: Routers and servers separately**





# COST COMPARISON

## CABLE COST MODEL

- Cable cost as a function of distance
  - The functions obtained using linear regression\*
  - Optical transceivers considered
  - Cables used: Mellanox IB FDR10 40Gb/s QSFP
- Other used cables (studies in paper):

Mellanox IB QDR  
56Gb/s QSFP



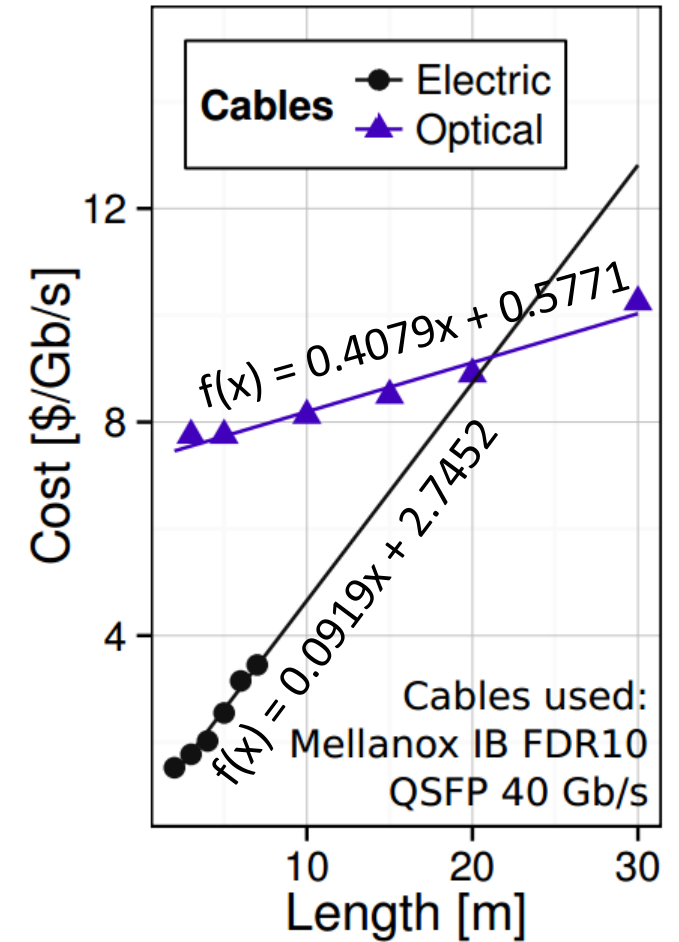
Mellanox Ethernet  
10Gb/s SFP+



Mellanox Ethernet  
40Gb/s QSFP



Elpeus Ethernet  
10Gb/s SFP+



\*Prices based on:



# COST COMPARISON

## ROUTER COST MODEL

- Router cost as a function of radix
  - The function obtained using linear regression\*
  - Routers used:

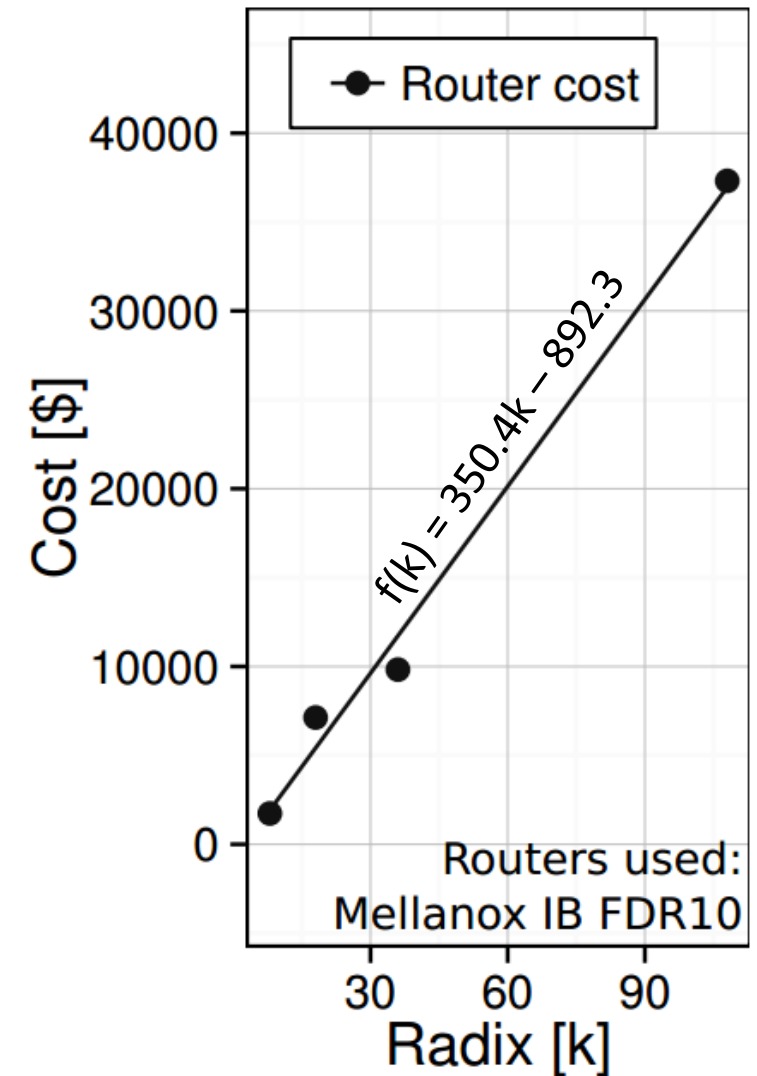
\*Prices based on:



Mellanox IB FDR10

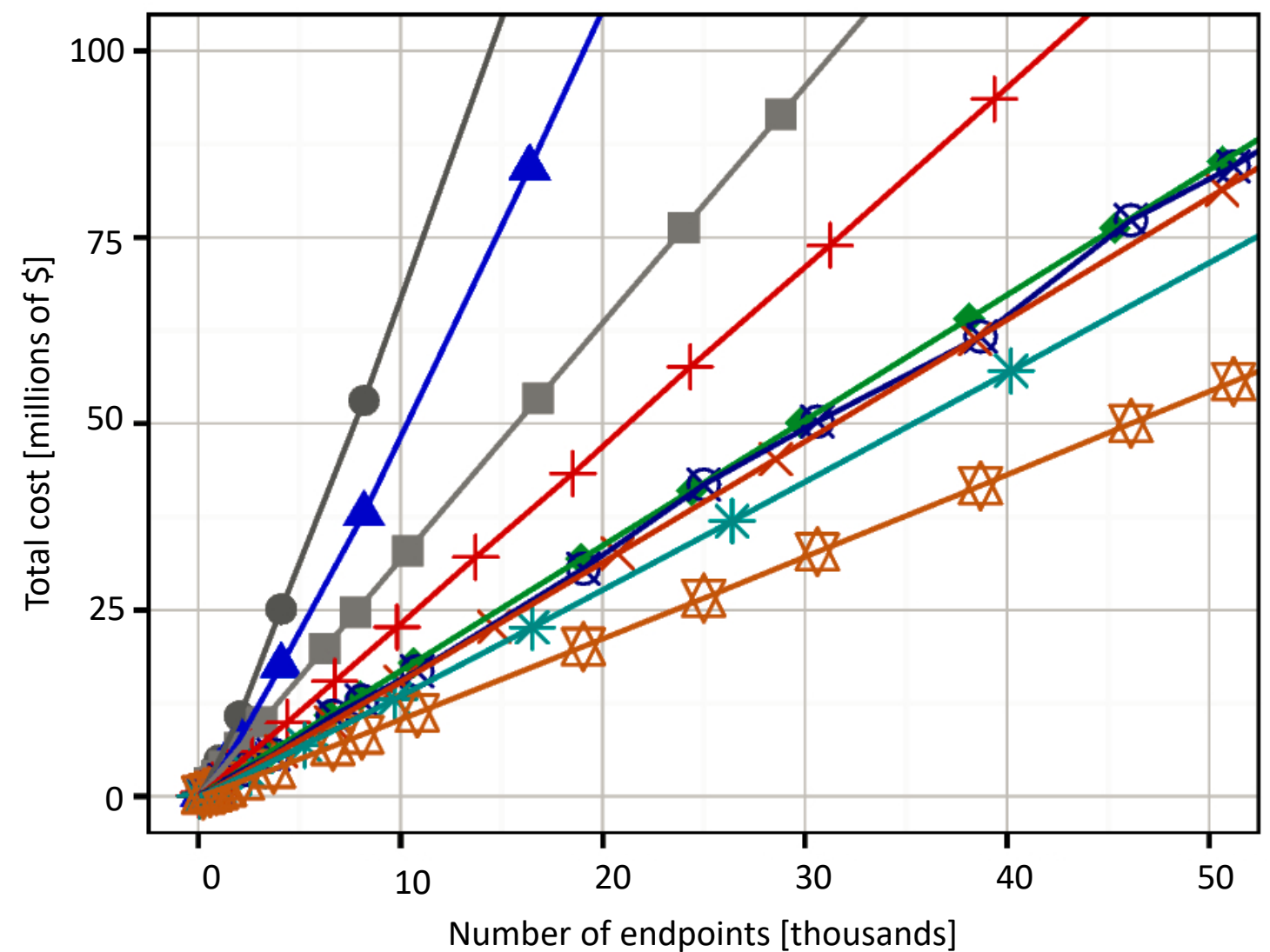


Mellanox Ethernet 10/40 Gb



# COST COMPARISON

## Variant 1



## Variant 2:

SF less expensive than Dragonfly by  
~13% (Mellanox IB routers) up to  
~39% (Mellanox Ethernet routers)

### Topology

- Long Hop
- ▲ Hypercube
- Torus 5D
- + Fat Tree
- ◆ Torus 3D
- ⊗ Random Top..
- × Flat. Butterfly
- \* Dragonfly
- ⊠ Slim Fly

# COST COMPARISON

## DETAILED CASE-STUDY

- A Slim Fly with;
  - $N = 10,830$
  - $k = 43$
  - $N_r = 722$

# COST & POWER COMPARISON

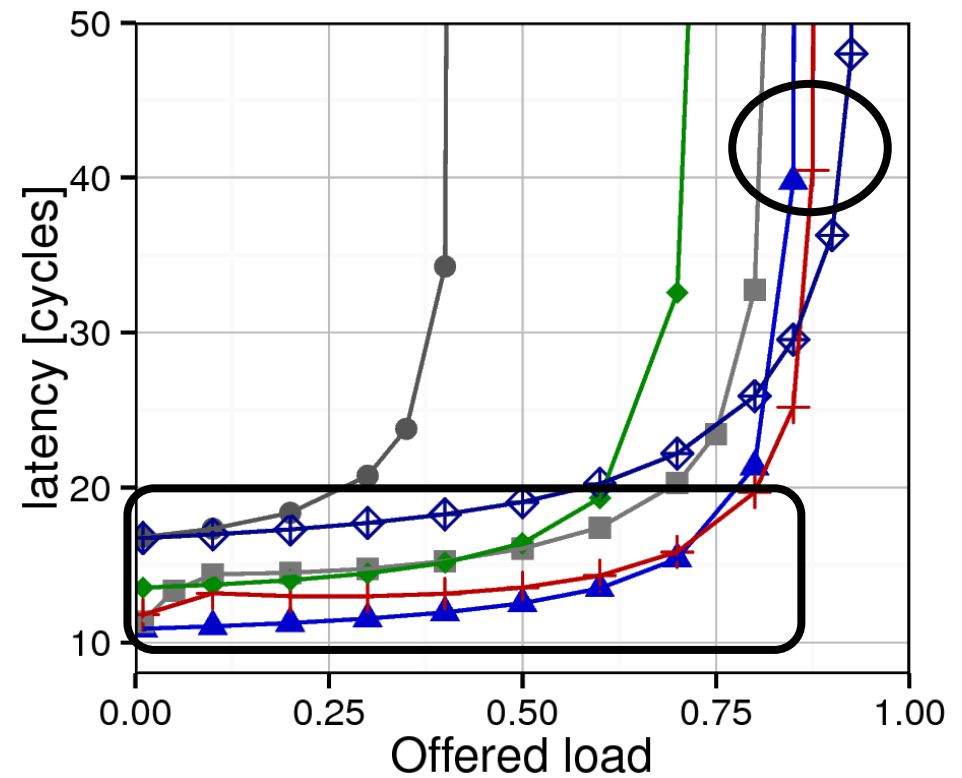
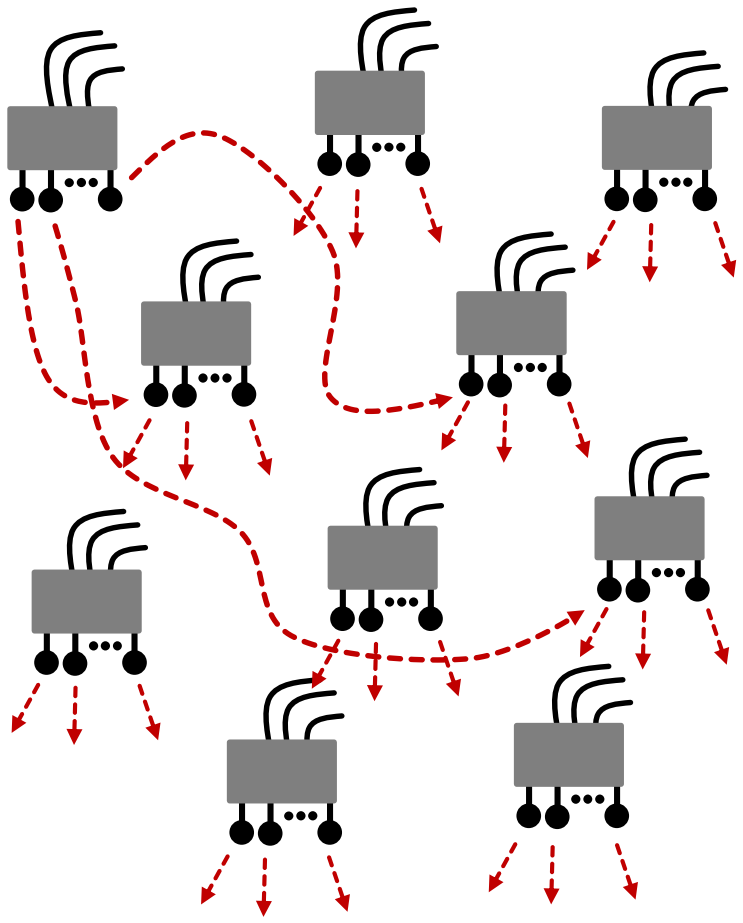
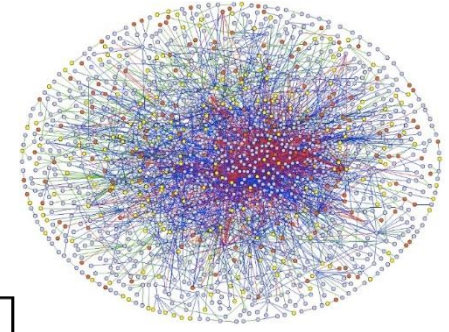
## DETAILED CASE-STUDY: HIGH-RADIX TOPOLOGIES

Topology	Fat tree	Random	Flat. Butterfly	Dragonfly	Slim Fly
Endpoints ( $N$ )	19,876	40,200	20,736	58,806	<b>10,830</b>
Routers ( $N_r$ )	2,311	4,020	1,728	5,346	<b>722</b>
Radix ( $k$ )	<b>43</b>	<b>43</b>	<b>43</b>	<b>43</b>	<b>43</b>
Electric cables	19,414	32,488	9,504	56,133	<b>6,669</b>
Fiber cables	40,215	33,842	20,736	29,524	<b>6,869</b>
Cost per node [\$]	2,346	1,743	1,570	1,438	<b>1,033</b>
Power per node [W]	14.0	12.04	10.8	10.9	<b>8.02</b>

Topology	Fat tree	Random	Flat. Butterfly	Dragonfly	Slim Fly
Endpoints ( $N$ )	<b>10,718</b>	<b>9,702</b>	<b>10,000</b>	<b>9,702</b>	<b>10,830</b>
Routers ( $N_r$ )	1,531	1,386	1,000	1,386	<b>722</b>
Radix ( $k$ )	35	28	33	27	<b>43</b>
Electric cables	7,350	6,837	4,500	9,009	<b>6,669</b>
Fiber cables	24,806	7,716	10,000	4,900	<b>6,869</b>
Cost per node [\$]	2,315	1,566	1,535	1,342	<b>1,033</b>
Power per node [W]	14.0	11.2	10.8	10.8	<b>8.02</b>

# PERFORMANCE & ROUTING

## RANDOM UNIFORM TRAFFIC



- Routing protocol**
- Slim Fly (Valiant)
  - ▲ Slim Fly (Minimum)
  - Slim Fly (UGAL-L)
  - ✚ Slim Fly (UGAL-G)
  - ◆ Dragonfly (UGAL-L)
  - ◇ Fat Tree (ANCA)

# Tuning VNs and VCs to avoid HoL blocking – uniform traffic (no VOQ)

Simple 2VC deadlock avoidance

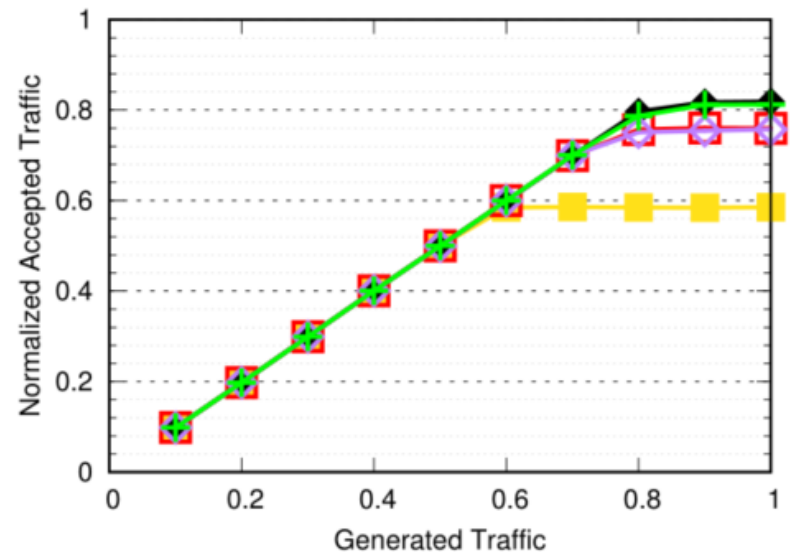
Destination-based buffer mgmt.

SF-optimized 4 VN and 8 VC

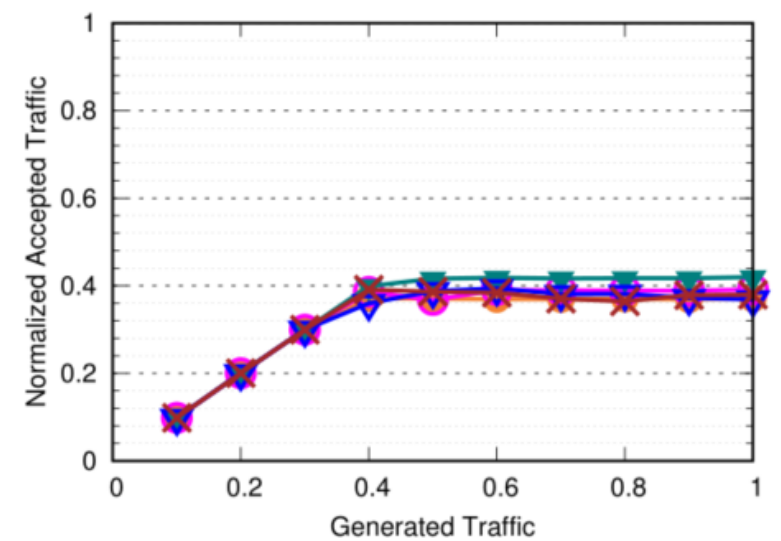
- MIN-DLA-2
- MIN-DBBM-8
- MIN-DBBM-16
- MIN-SF2LQ-8
- MIN-SF2LQ-16

- VAL-DLA-4
- VAL-DBBM-8
- VAL-DBBM-16
- VAL-SF4LQ-8
- VAL-SF4LQ-16

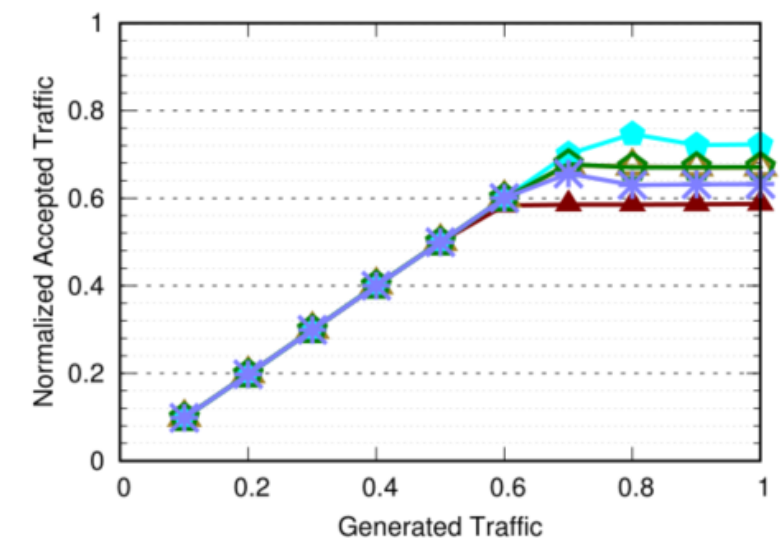
- UGAL-DLA-4
- UGAL-DBBM-8
- UGAL-DBBM-16
- UGAL-SF4LQ-8
- UGAL-SF4LQ-16



MIN Routing. Slim Fly 19\_10



VAL Routing. Slim Fly 19\_10



UGAL Routing. Slim Fly 19\_10

# Tuning VNs and VCs to avoid HoL blocking – hotspot traffic (no VOQ)

Simple 2VC deadlock avoidance

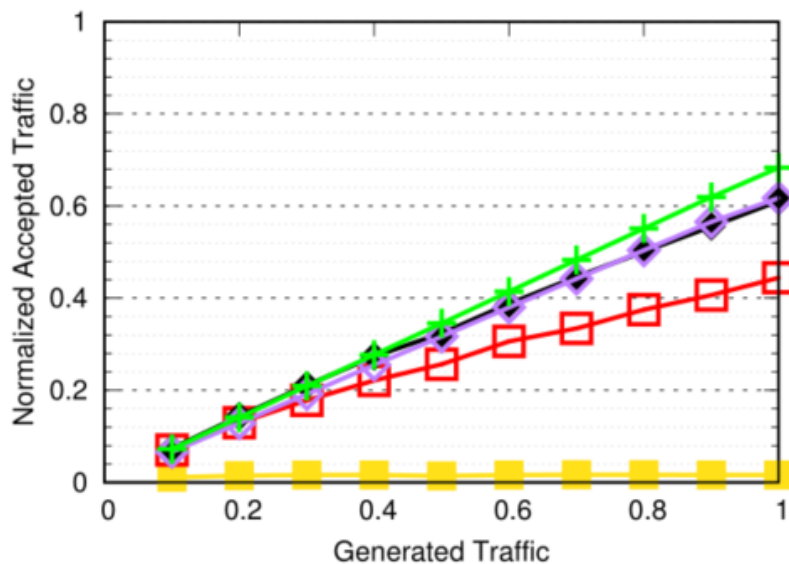
Destination-based buffer mgmt.

SF-optimized 4 VN and 8 VC

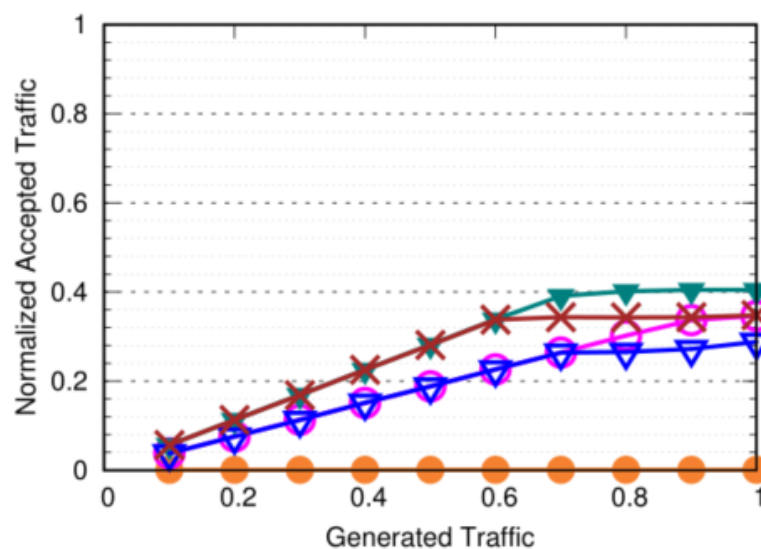
- MIN-DLA-2
- MIN-DBBM-8
- MIN-DBBM-16
- MIN-SF2LQ-8
- MIN-SF2LQ-16

- VAL-DLA-4
- VAL-DBBM-8
- VAL-DBBM-16
- VAL-SF4LQ-8
- VAL-SF4LQ-16

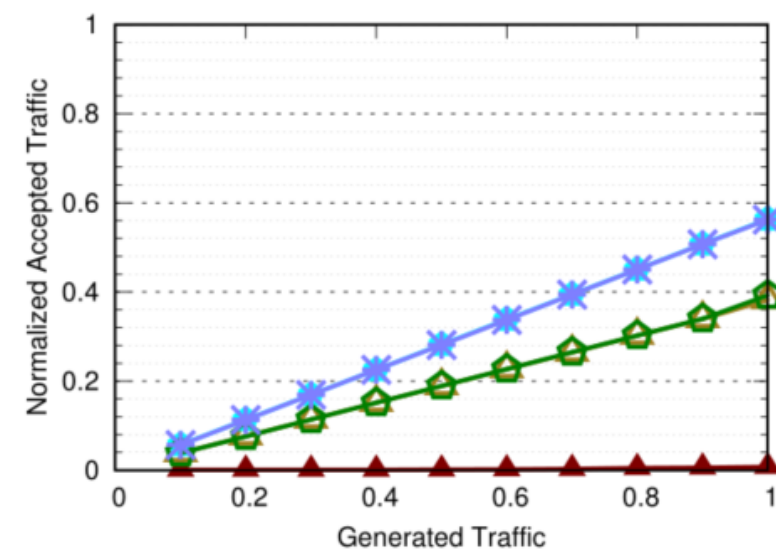
- UGAL-DLA-4
- UGAL-DBBM-8
- UGAL-DBBM-16
- UGAL-SF4LQ-8
- UGAL-SF4LQ-16



MIN Routing. Slim Fly 19\_10



VAL Routing. Slim Fly 19\_10

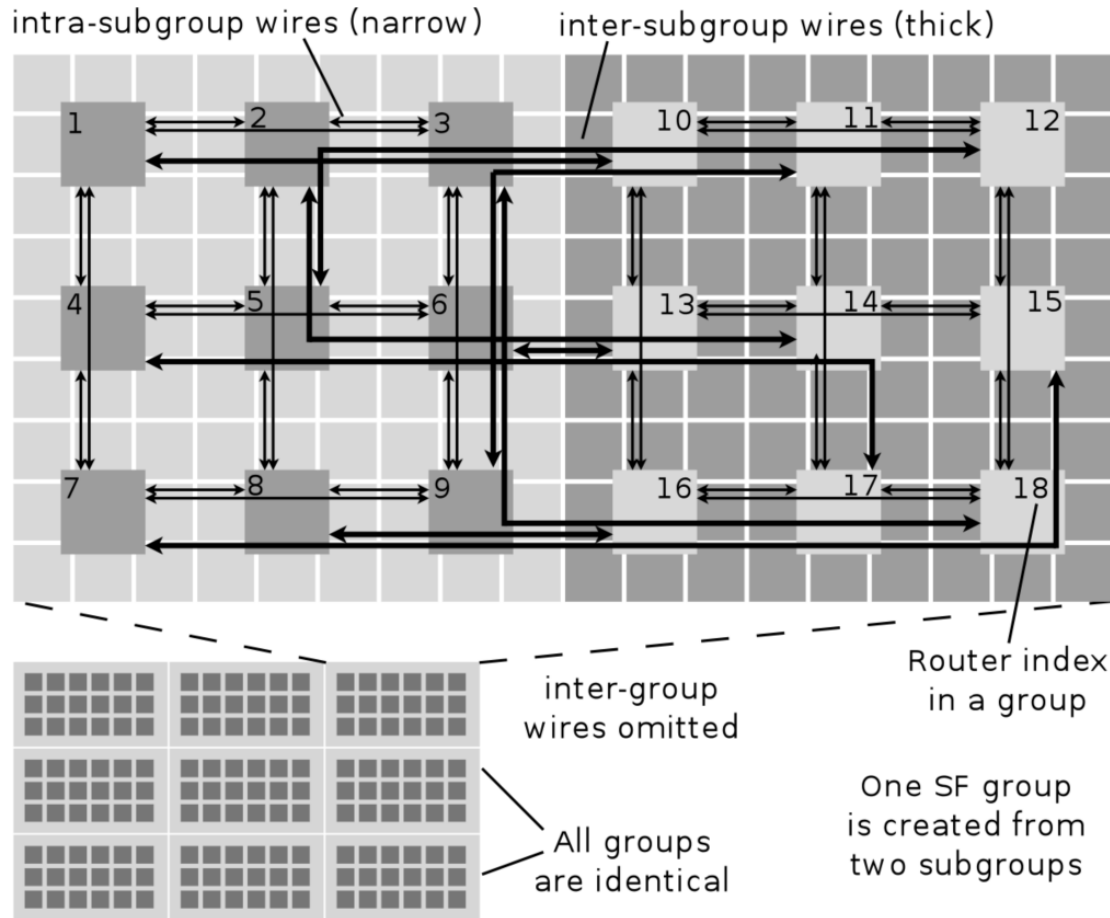


UGAL Routing. Slim Fly 19\_10

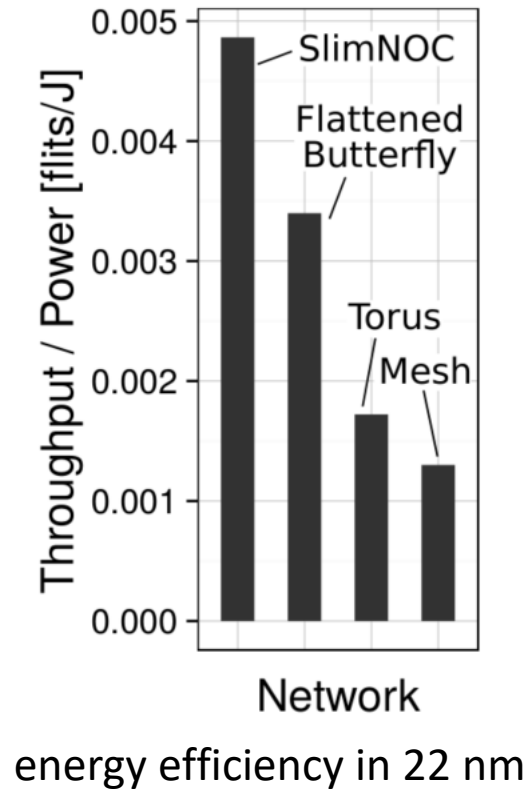


# Slim NoC – Slim Fly topologies for on chip networks

- New challenges – layout in the chip’s metal layers



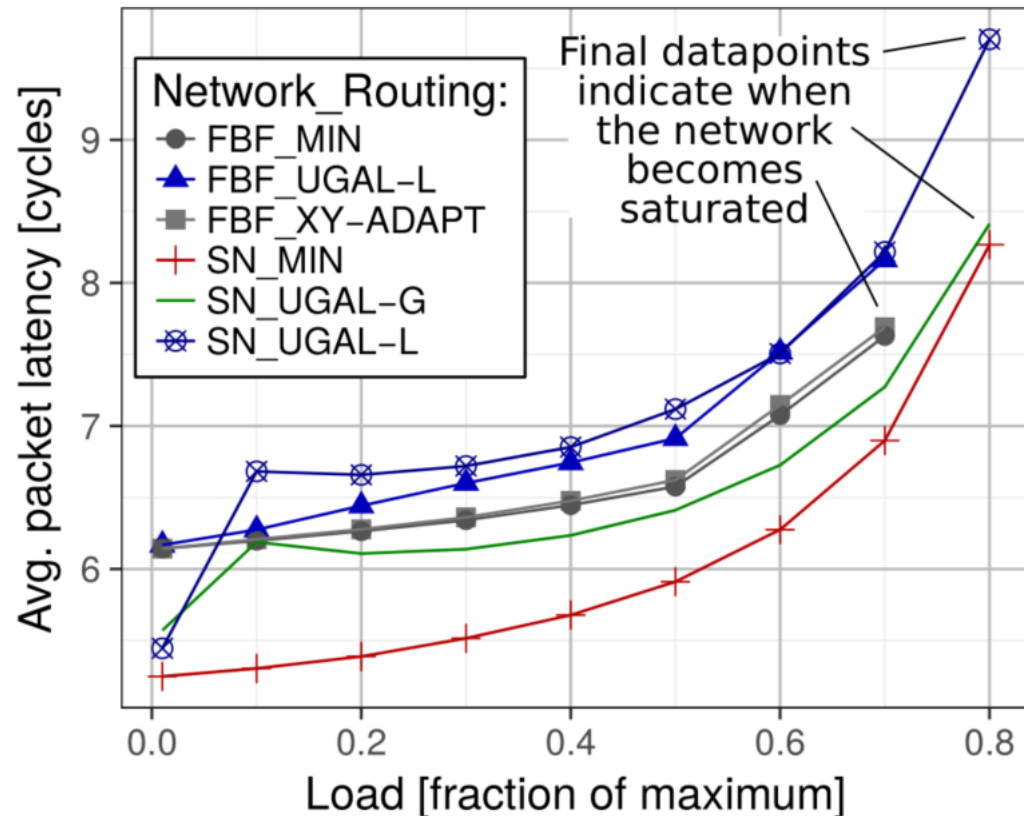
1296 cores, 162 routers



# Slim NoC performance

uniform random load

Flattened Butterfly vs. Slim Fly



Energy-Delay Product  
PARSEC/SPLASH

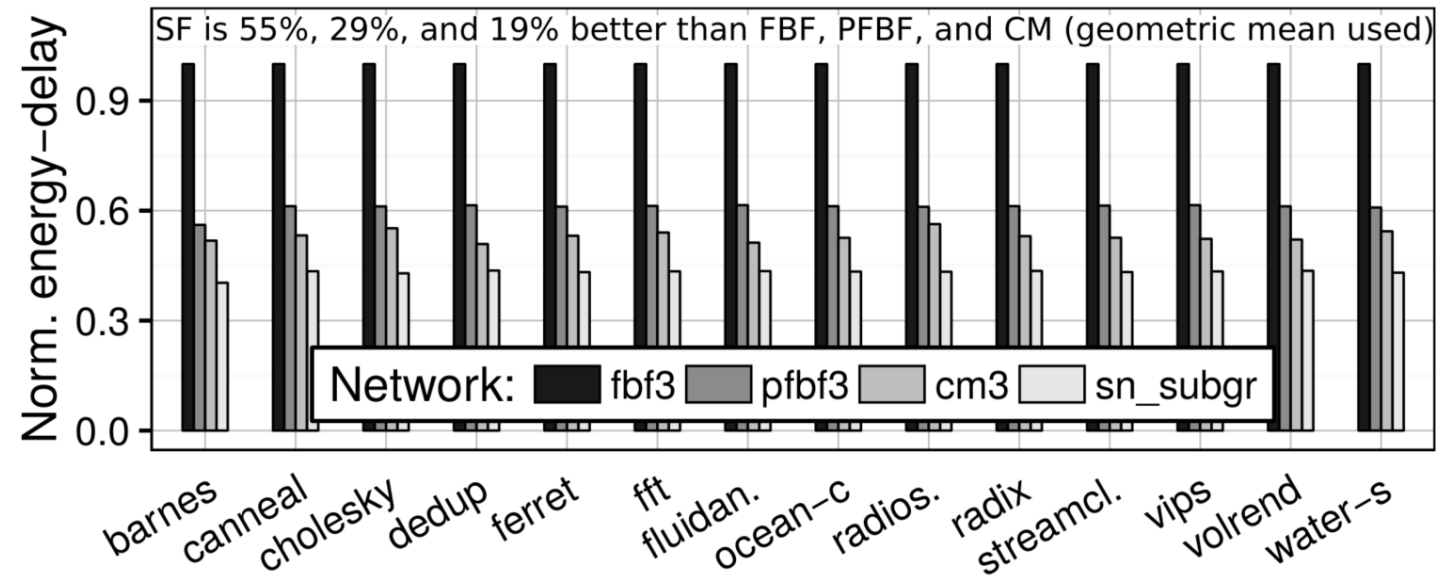


Figure Legend:

FBF/PFBF – Flattened Butterfly

CM – Concentrated Mesh

SN\_SUBGR – Slim NOC



Low Latency

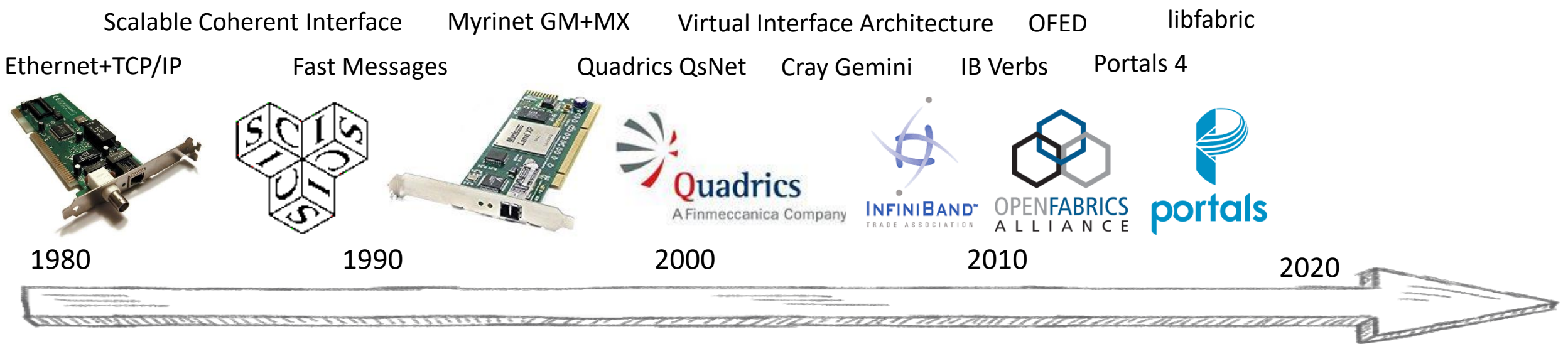


Low Cost



Low Processing Load

# The Development of High-Performance Networking Interfaces



sockets

(active) message based

protocol offload

remote direct memory access (RDMA)

coherent memory access

OS bypass

zero copy

triggered operations

**InfiniBand Trade Association Launches the RoCE Initiative to Advance RDMA over Converged Ethernet Solutions**

*RoCE delivers significant performance and efficiency gains to cloud, storage, virtualization and hyper-converged infrastructures*

businessinsider.com

**Microsoft to Drive RDMA Into Datacenters and Clouds**

November 18, 2013 by Timothy Prickett Morgan

**RDMA over Ethernet - the Rocky road to convergence**

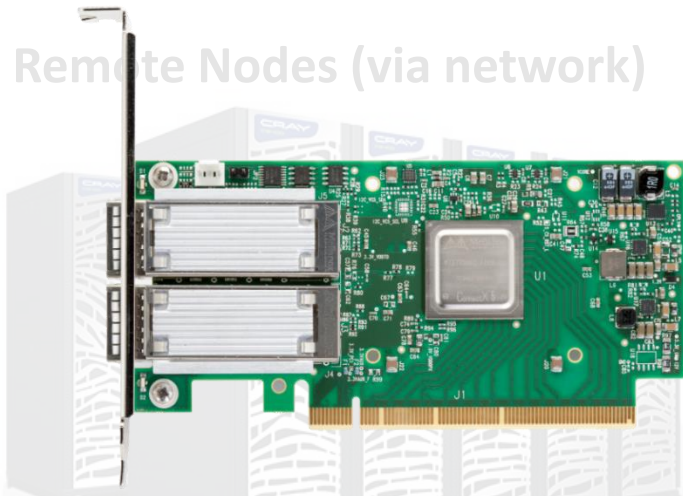
17 November 2015 | By Brandon Hoff

June 2017

95 / top-100 systems use RDMA

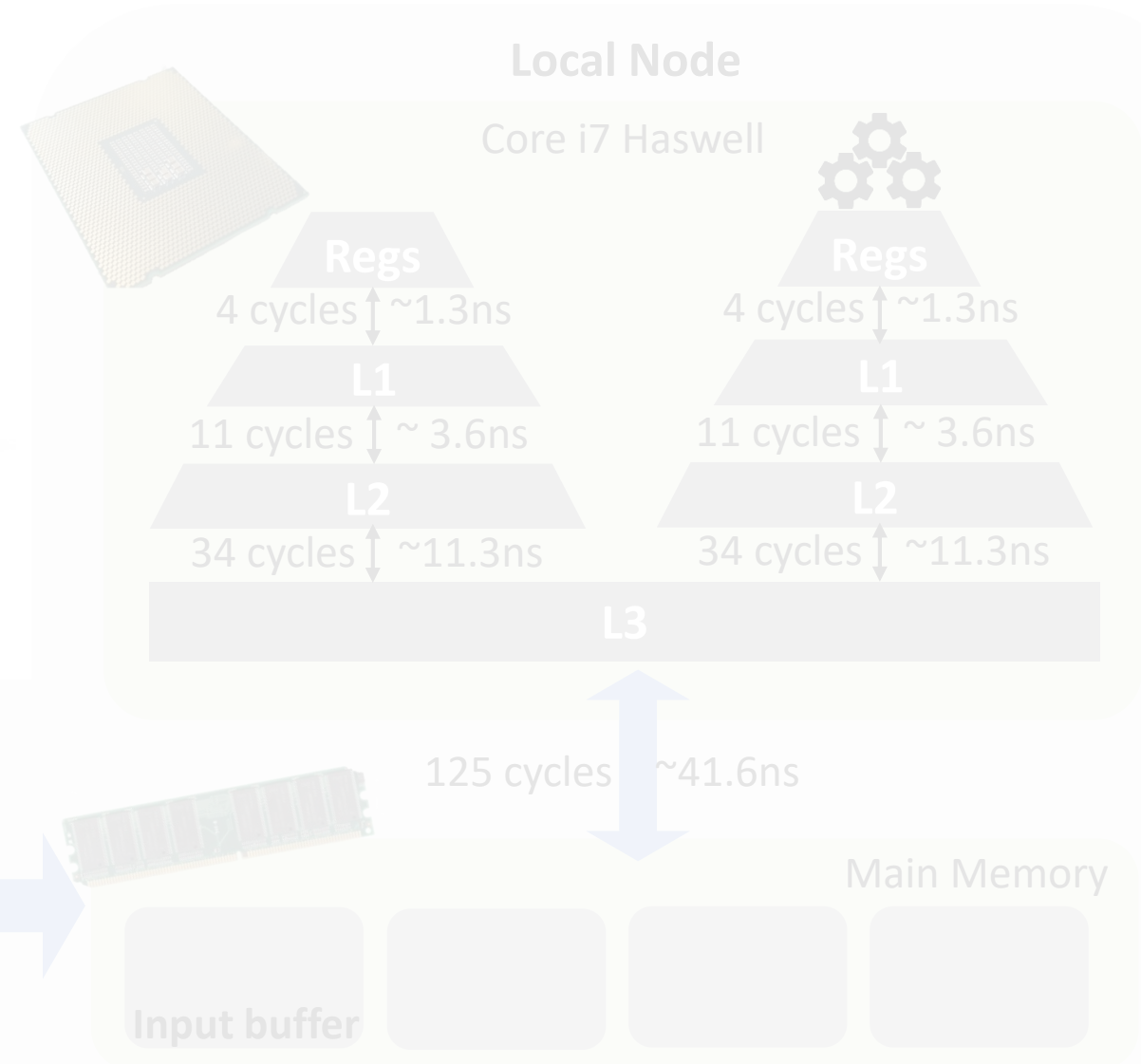
>285 / top-500 systems use RDMA

# Data Processing in modern RDMA networks



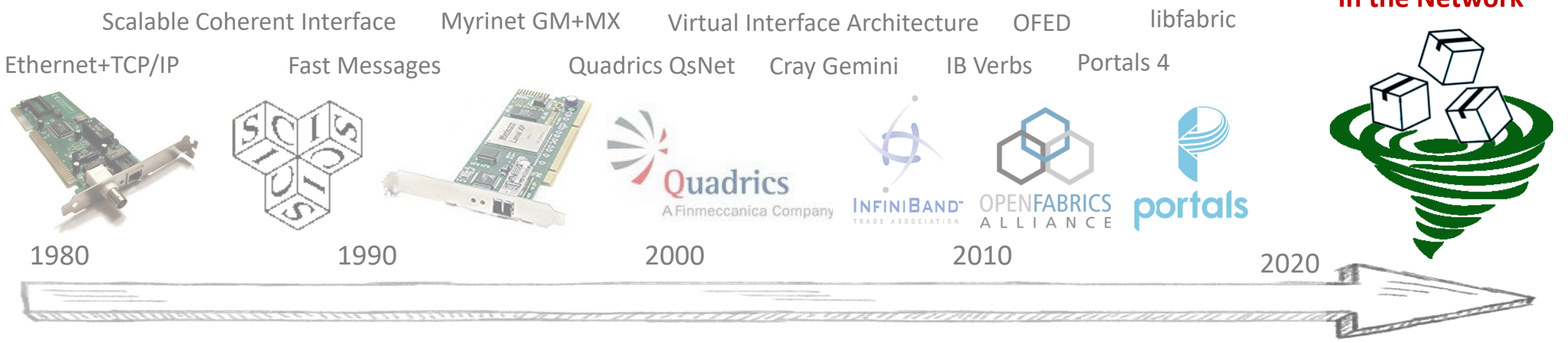
Remote Nodes (via network)

**Mellanox Connect-X5: 1 msg/5ns**  
**Tomorrow (400G): 1 msg/1.2ns**



# The future of High-Performance Networking Interfaces

**SPIN**  
Streaming Processing  
In the Network



sockets    (active) message based    protocol offload    remote direct memory access (RDMA)    fully programmable NIC acceleration

coherent memory access    OS bypass    zero copy    triggered operations

## Established Principles for Compute Acceleration

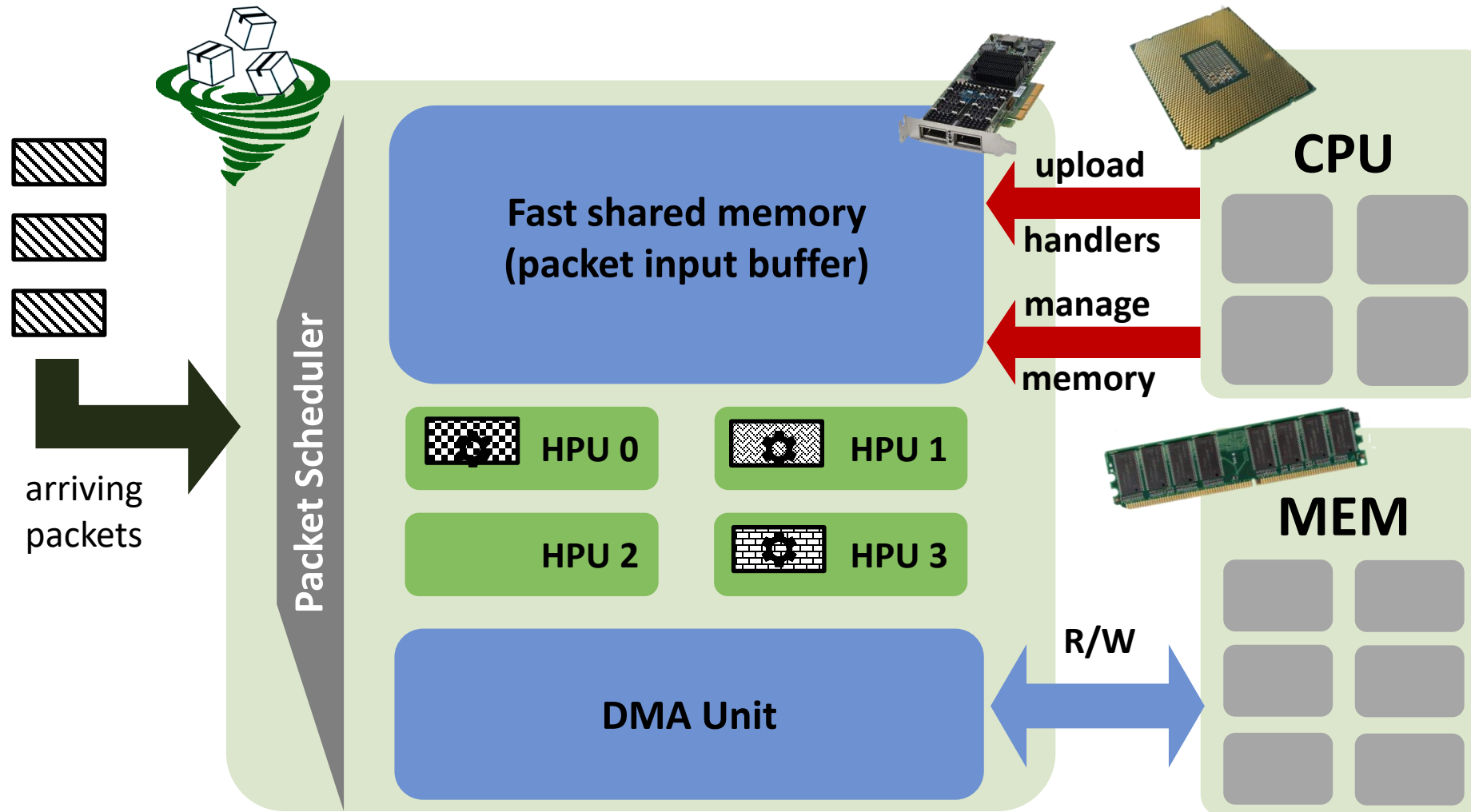
Specialization    Programmability    Libraries  
Ease-of-use    Portability    Efficiency



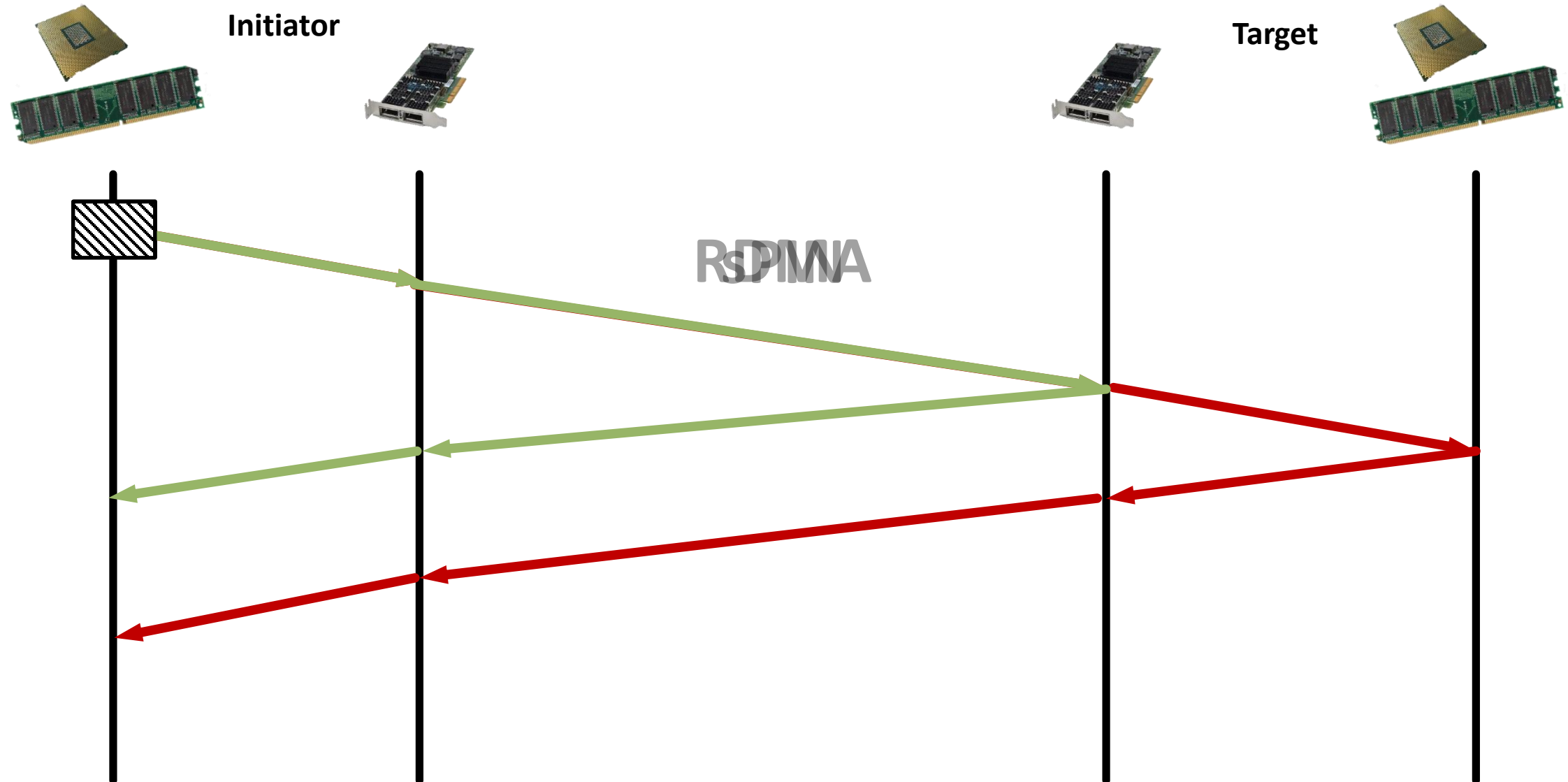
TOP 500<sup>®</sup> SUPERCOMPUTER SITES    June 2017

95 / top-100 systems use RDMA  
>285 / top-500 systems use RDMA

# sPIN NIC - Abstract Machine Model

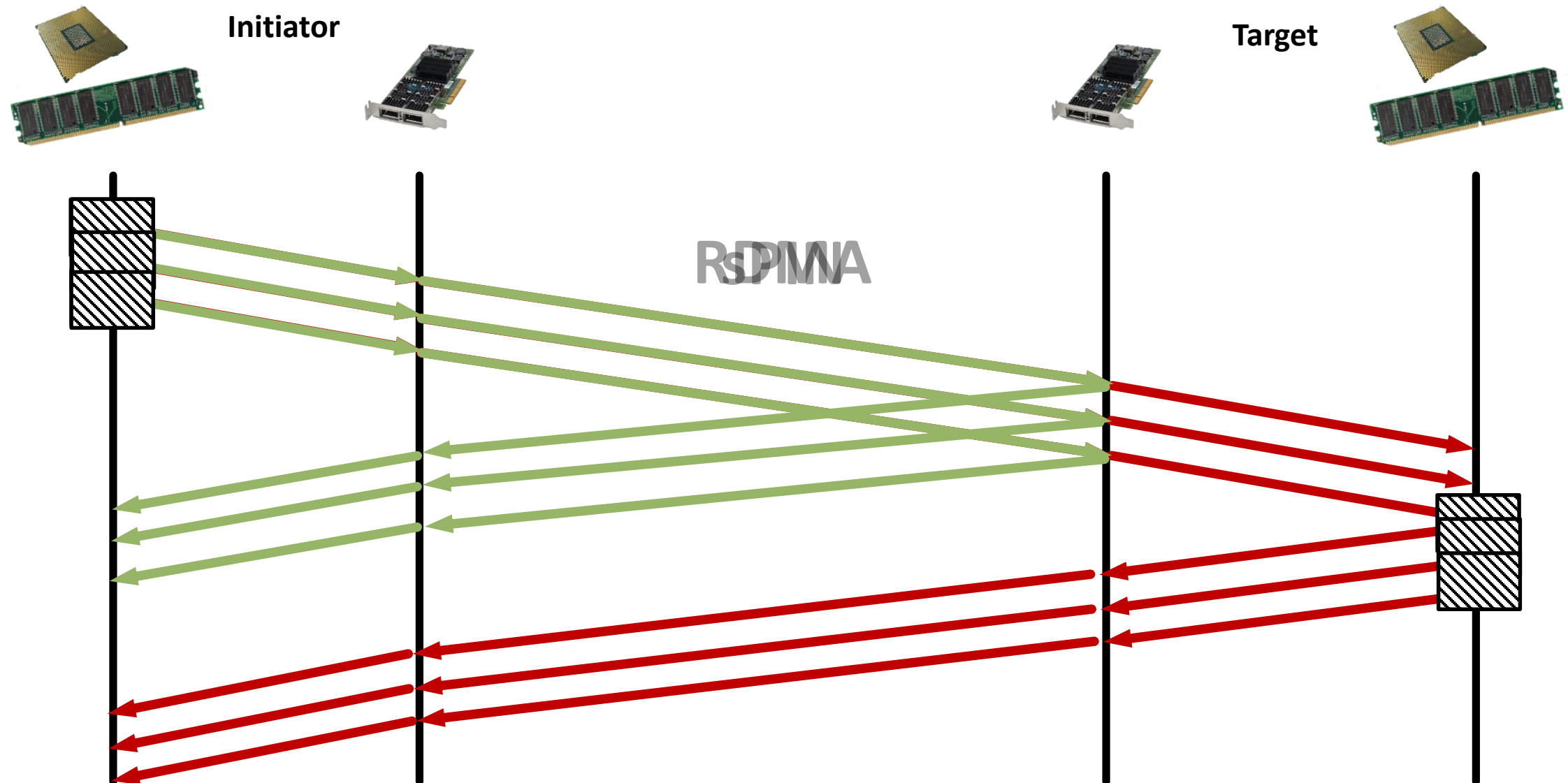


# RDMA vs. sPIN in action: Simple Ping Pong

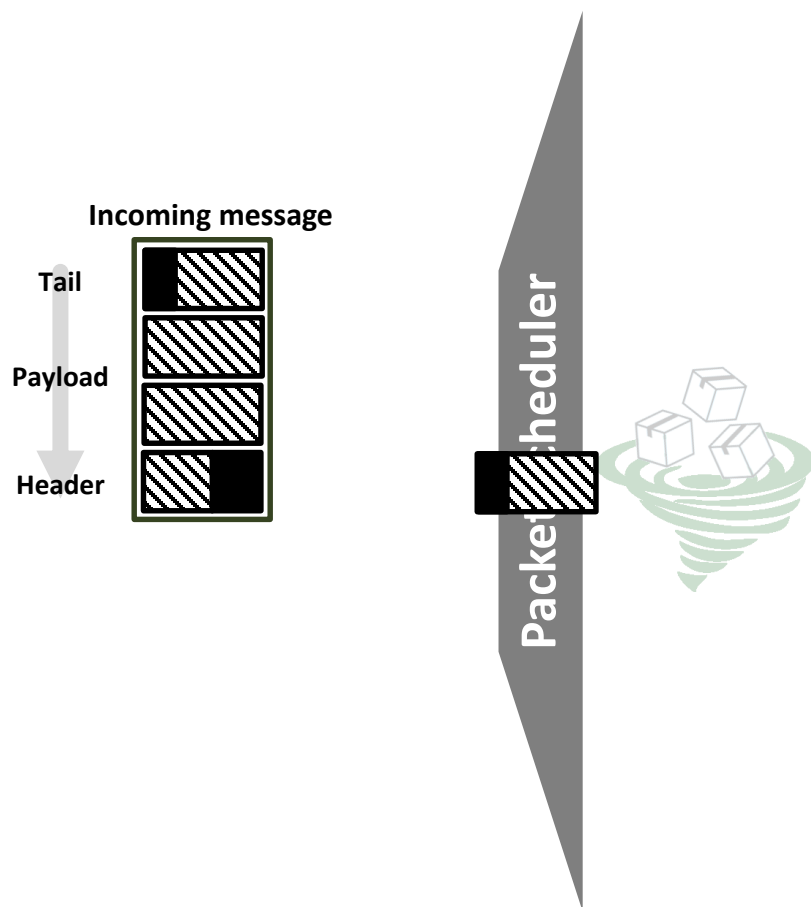




# RDMA vs. sPIN in action: Streaming Ping Pong



# sPIN – Programming Interface



## Header handler

```
__handler int pp_header_handler(const ptl_header_t h, void *state) {
    pingpong_info_t *i = state;
    i->source = h.source_id;
    return PROCESS_DATA; // execute payload handler to put from device
}
```

## Payload handler

```
__handler int pp_payload_handler(const ptl_payload_t p, void *state) {
    pingpong_info_t *i = state;
    PtlHandlerPutFromDevice(p.base, p.length, 1, 0, i->source, 10, 0, NULL, 0);
    return SUCCESS;
}
```

## Completion handler

```
__handler int pp_completion_handler(int dropped_bytes,
                                     bool flow_control_triggered, void *state) {
    return SUCCESS;
}
```

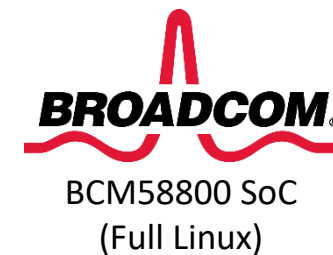
```
connect(peer, /* ... */, &pp_header_handler, &pp_payload_handler, &pp_completion_handler);
```

# Possible sPIN implementations



- **sPIN is a programming abstraction, similar to CUDA or OpenCL combined with OFED or Portals 4**
  - It enables a large variety of NIC implementations!
  - For example, massively multithreaded HPUs  
*Including warp-like scheduling strategies*
- **Main goal: sPIN must not obstruct line-rate**
  - Programmer must limit processing time per packet  
*Little's Law: 500 instructions per handler, 2.5 GHz, IPC=1, 1 Tb/s → 25 kiB memory*
  - Relies on fast shared memory (processing in packet buffers)  
*Scratchpad or registers*
  - Quick (single-cycle) handler invocation on packet arrival  
*Pre-initialized memory & context*
- **Can be implemented in most RDMA NICs with a firmware update**
  - Or in software in programmable (Smart) NICs

at 400G, process more than  
833 million messages/s

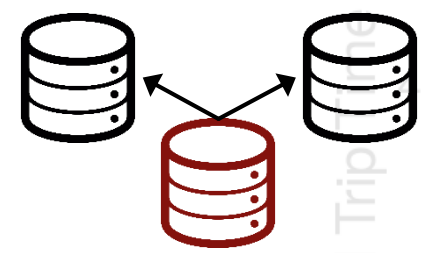


# Simulating a sPIN NIC – Ping Pong

- LogGOPSim v2 [1]: combine LogGOPSim (packet-level network) with gem5 (cycle accurate CPU simulation)

## Network (LogGOPSim):

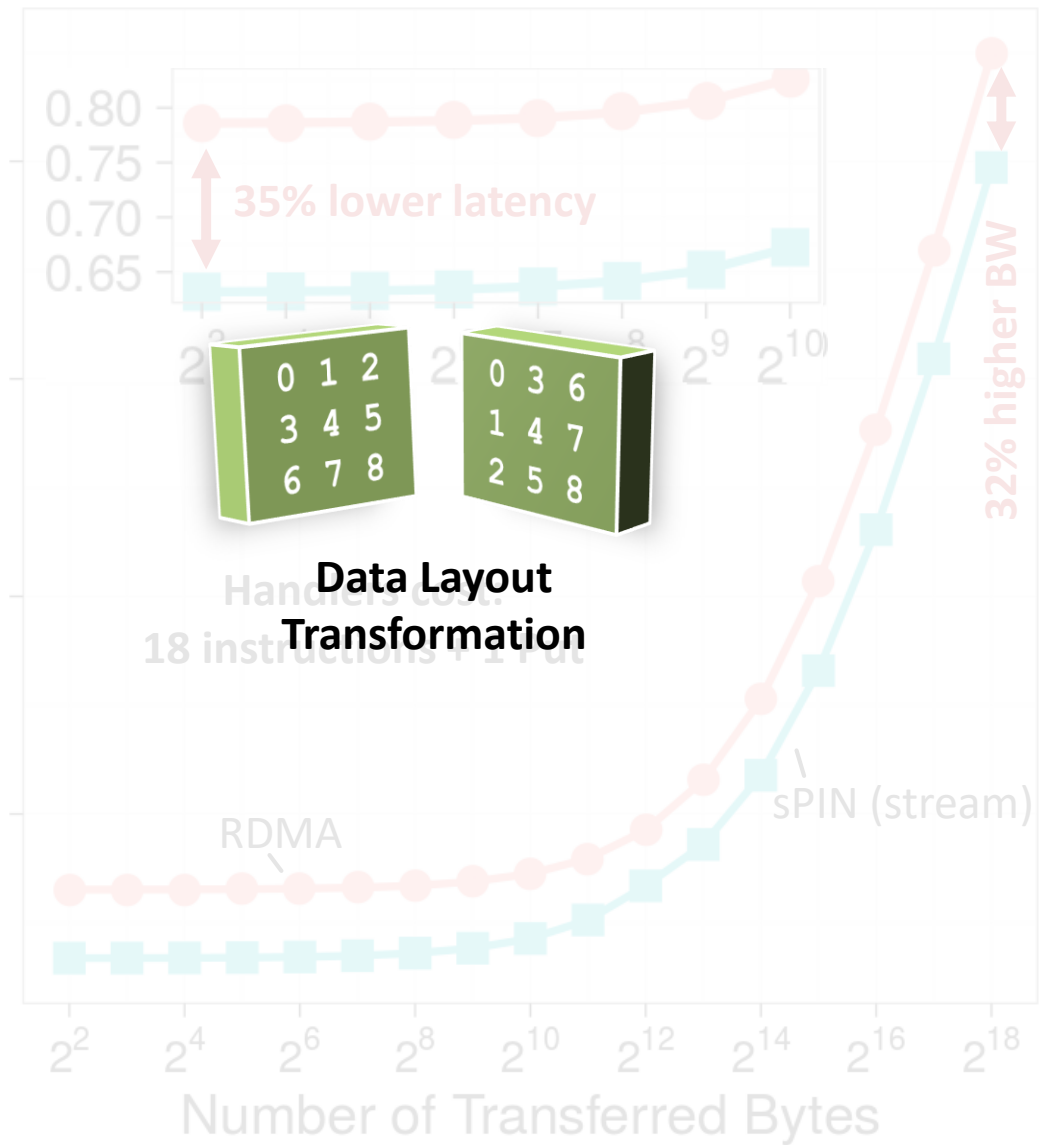
- Supports Portals 4 and MPI
- Parametrized for future InfiniBand
- $\sigma=65ns$  (measured)
- $g=6.7ns$  (150 MM/s)
- $G=2.5ps$  (400 Gib/s)
- Switch  $L=50ns$  (measured)
- Wire  $L=33.4ns$  (100 Gb/s)



Distributed Data Management

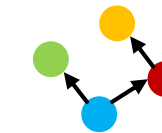
## NIC HPU

- 2.5 GHz ARM Cortex A15 OOO
- ARMv8-A 32 bit ISA
- Single-cycle access SRAM (no DRAM)
- Header matching  $m=30ns$ , per packet 2ns
- In parallel with  $g$ !

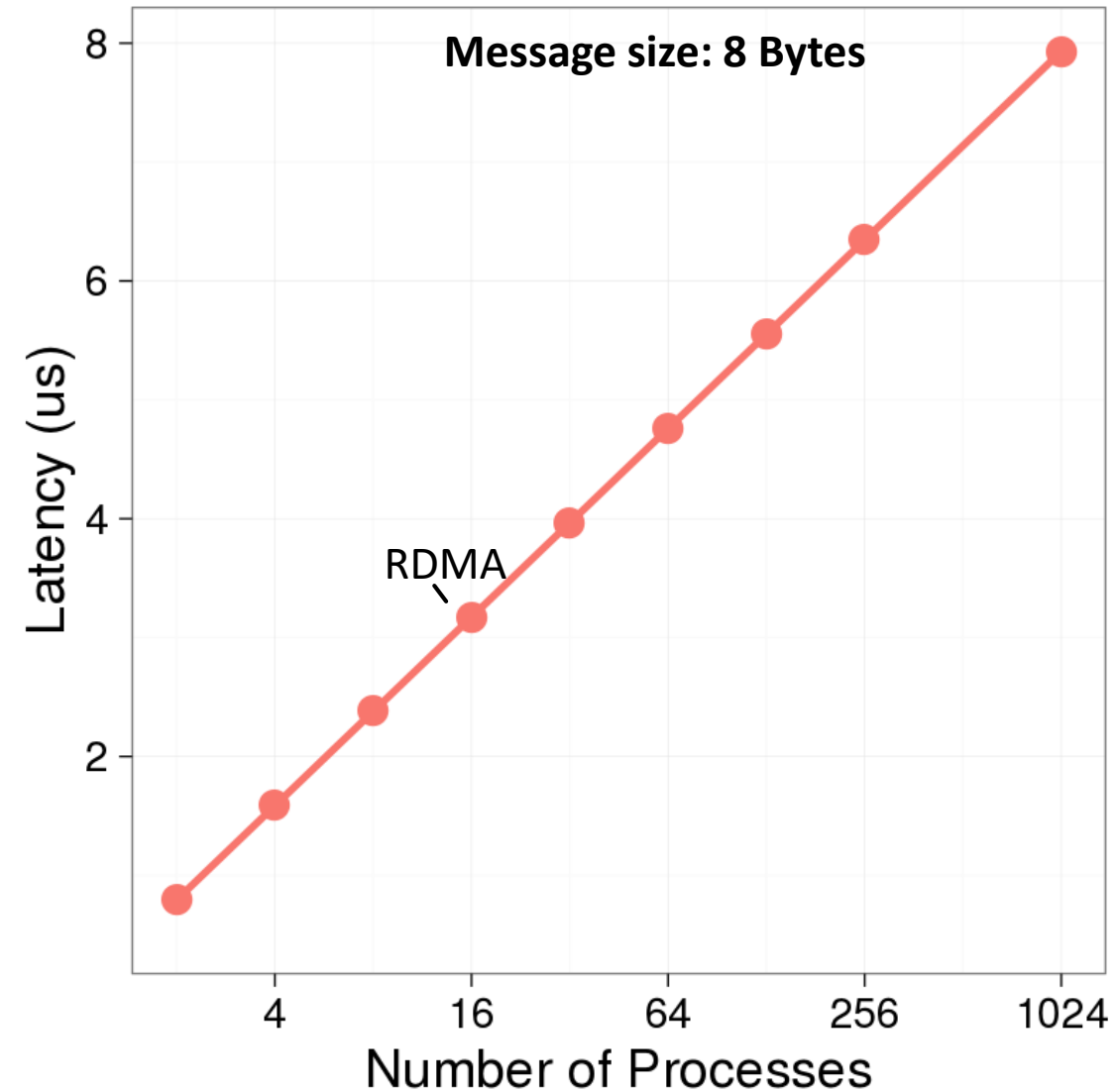
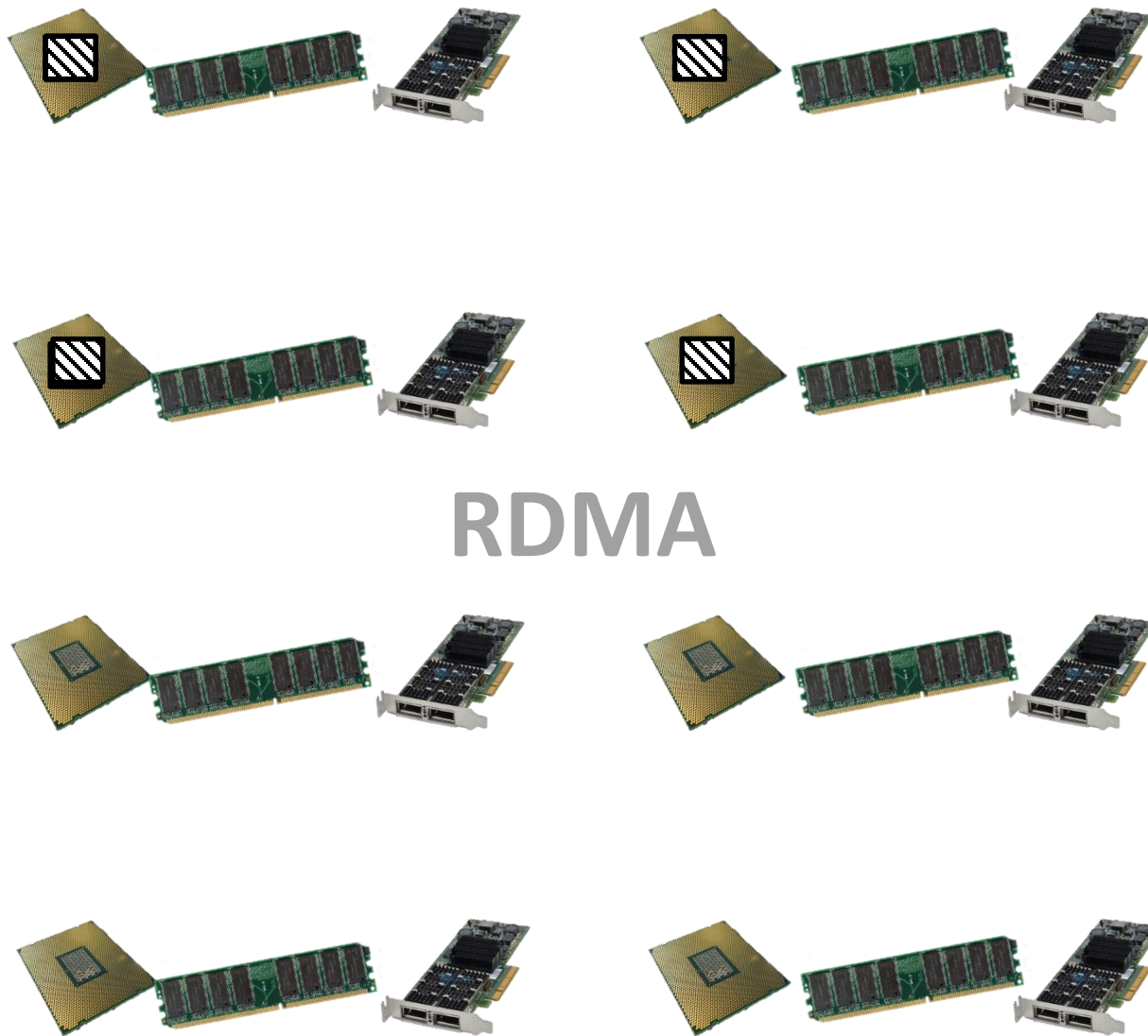


[1] S. Di Girolamo, K. Taranov, T. Schneider, E. Stalder, T. Hoefler, LogGOPSim+gem5: Simulating Network Offload Engines Over Packet-Switched Networks. Presented at ExaMPI'17

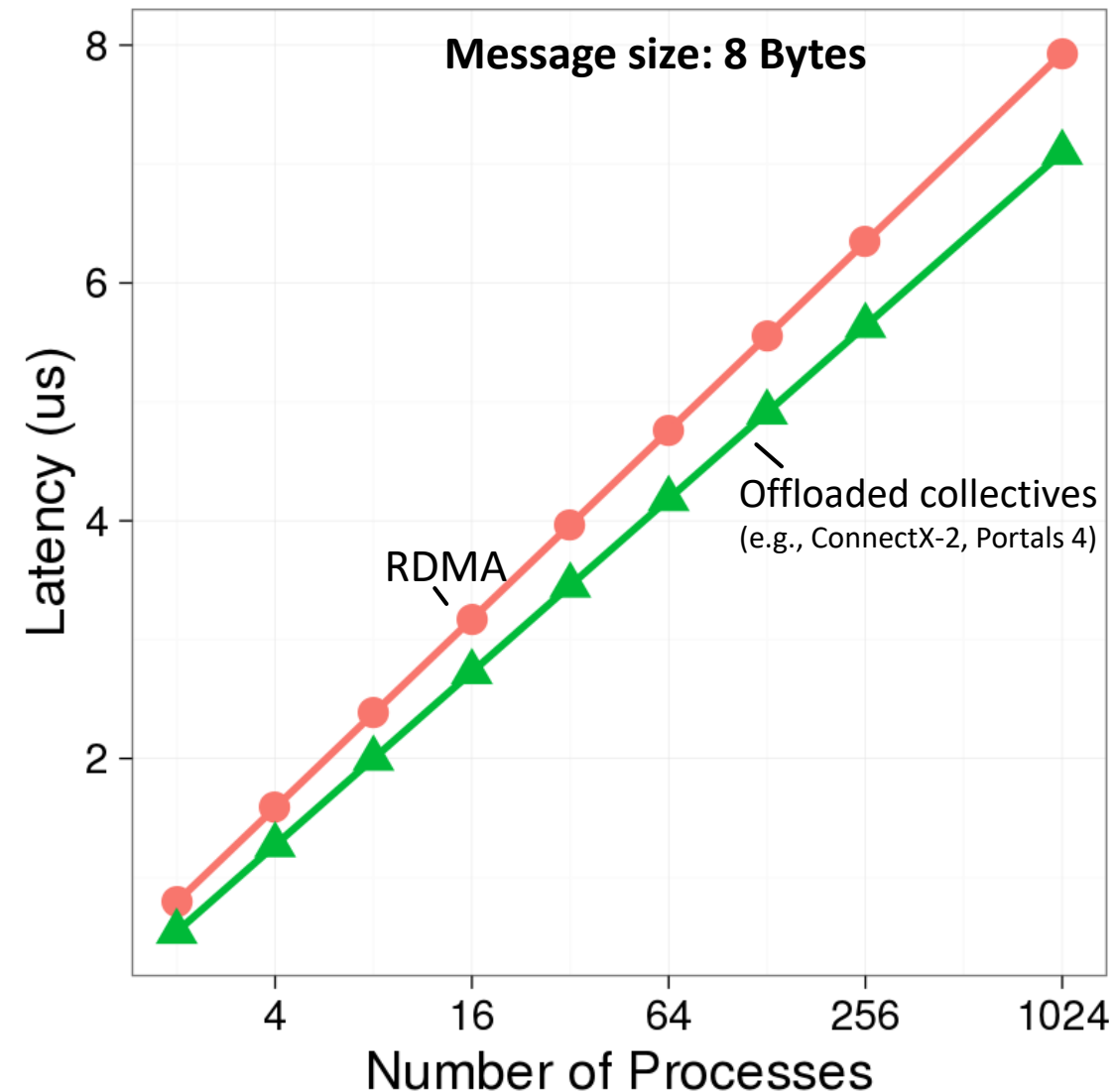
# Use Case 1: Broadcast acceleration



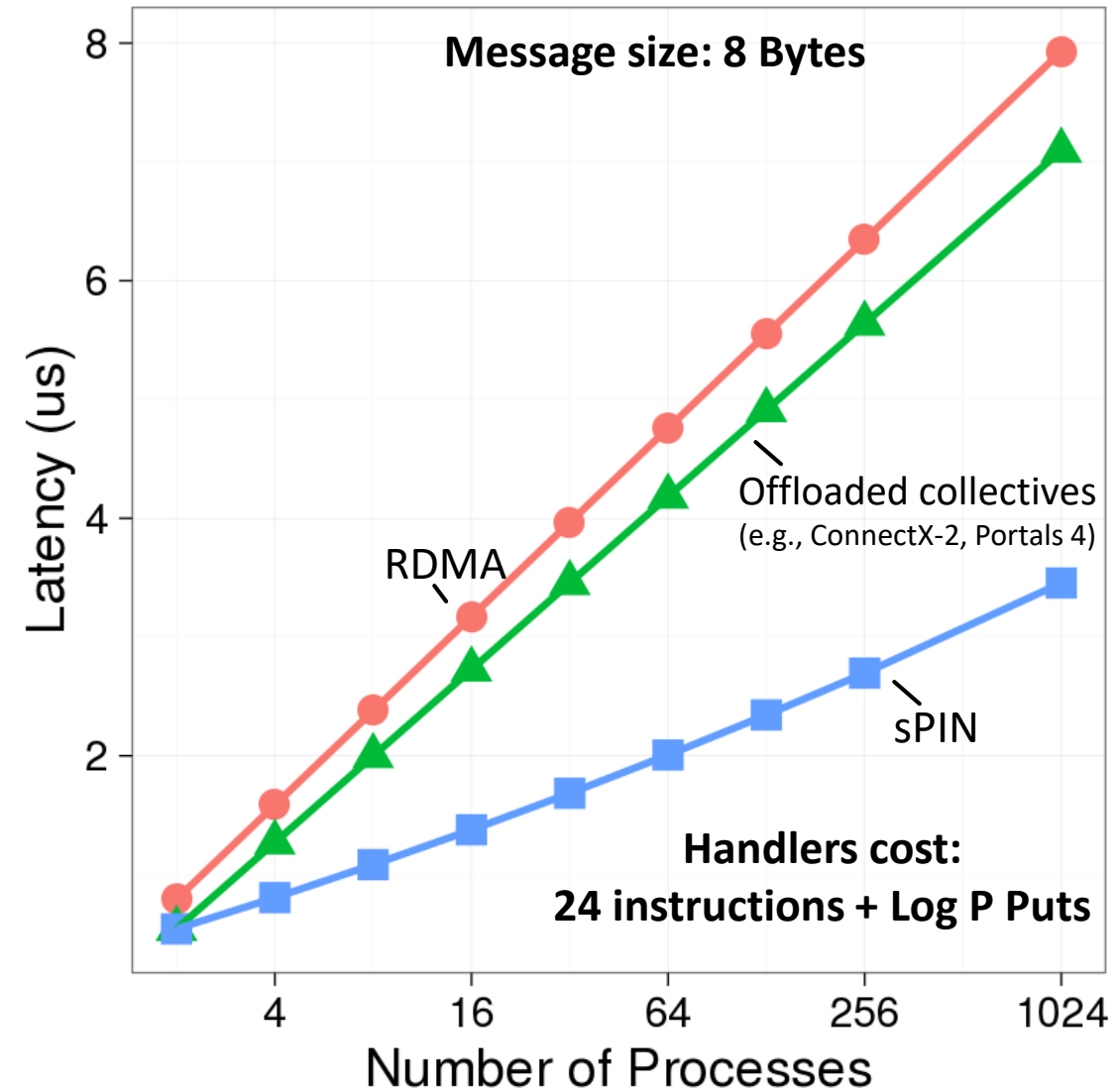
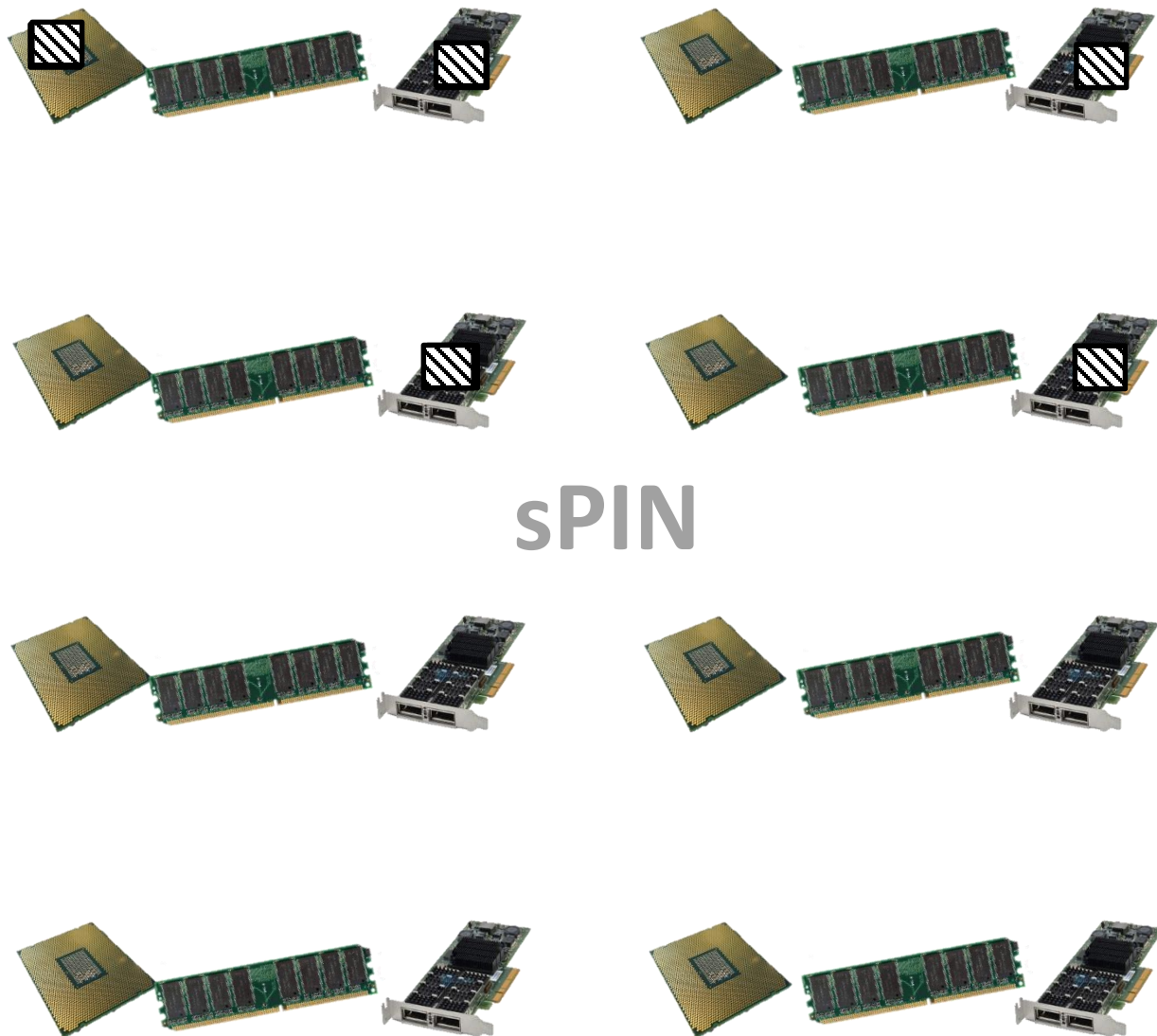
Network Group Communication



# Use Case 1: Broadcast acceleration



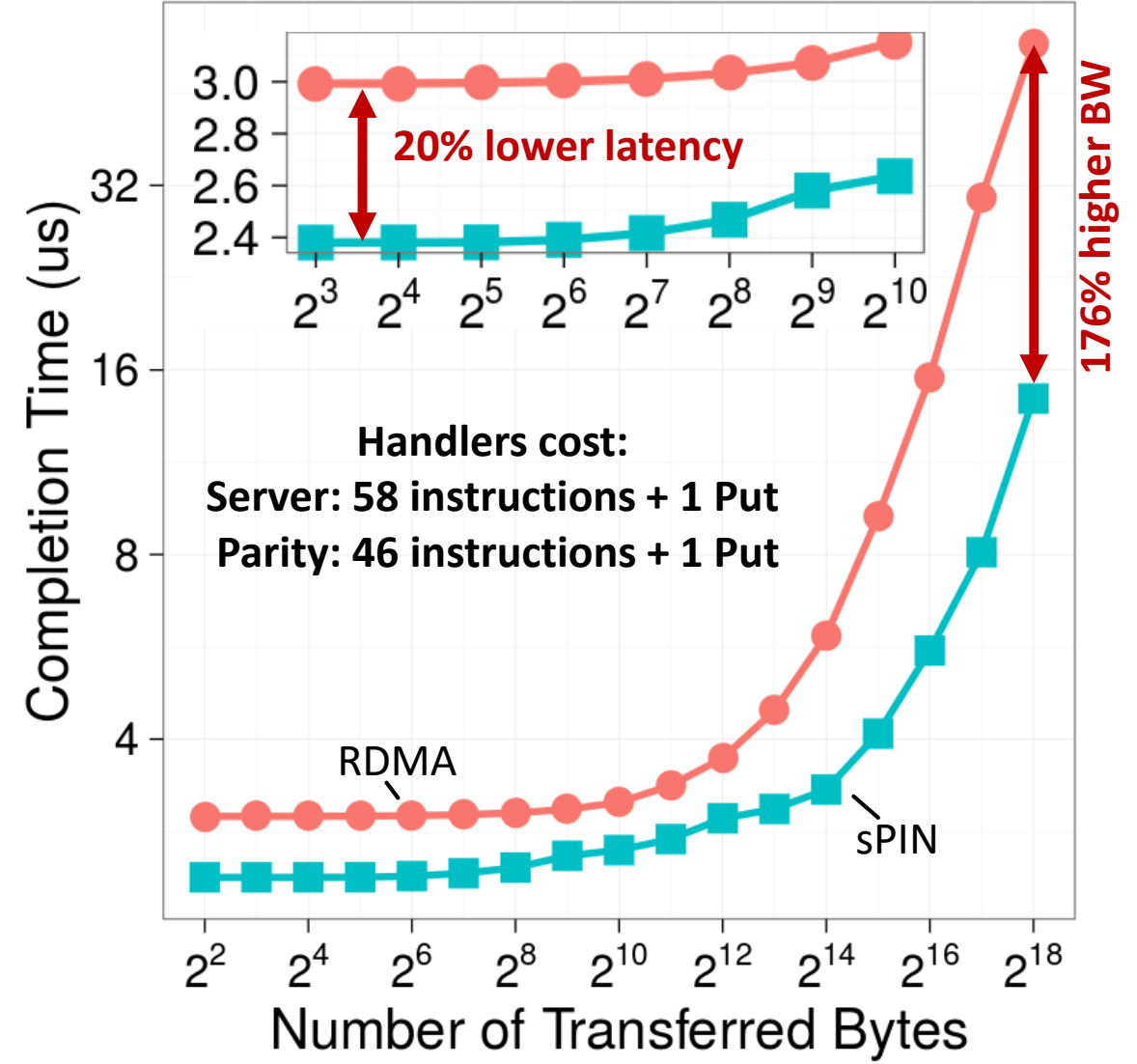
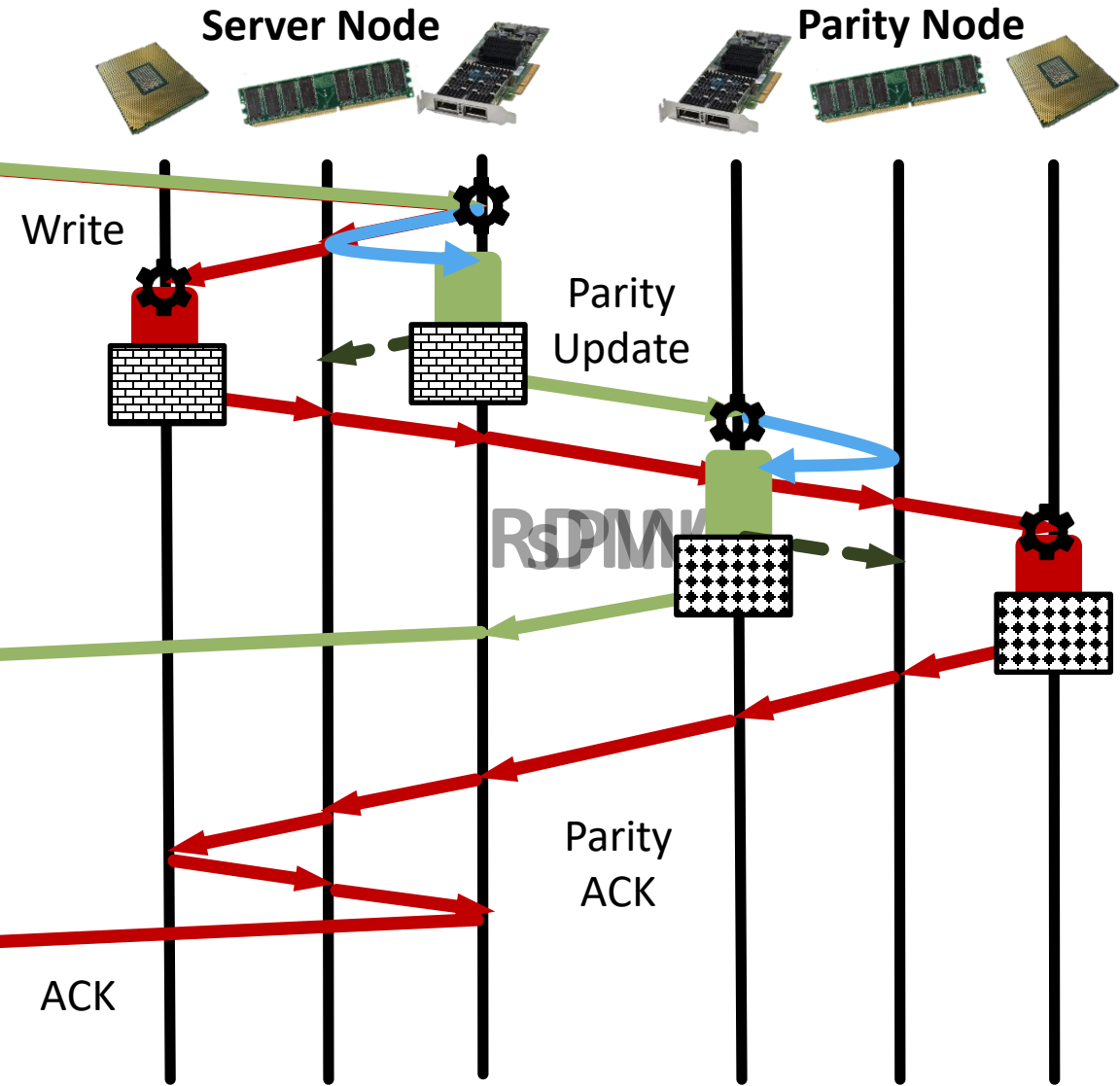
# Use Case 1: Broadcast acceleration



Underwood, K.D., et al., Enabling flexible collective communication offload with triggered operations. *HOTI'11*

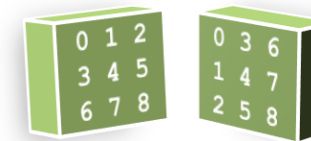
Liu, J., et al., High performance RDMA-based MPI implementation over InfiniBand. *International Journal of Parallel Programming* 2004

# Use Case 2: RAID acceleration

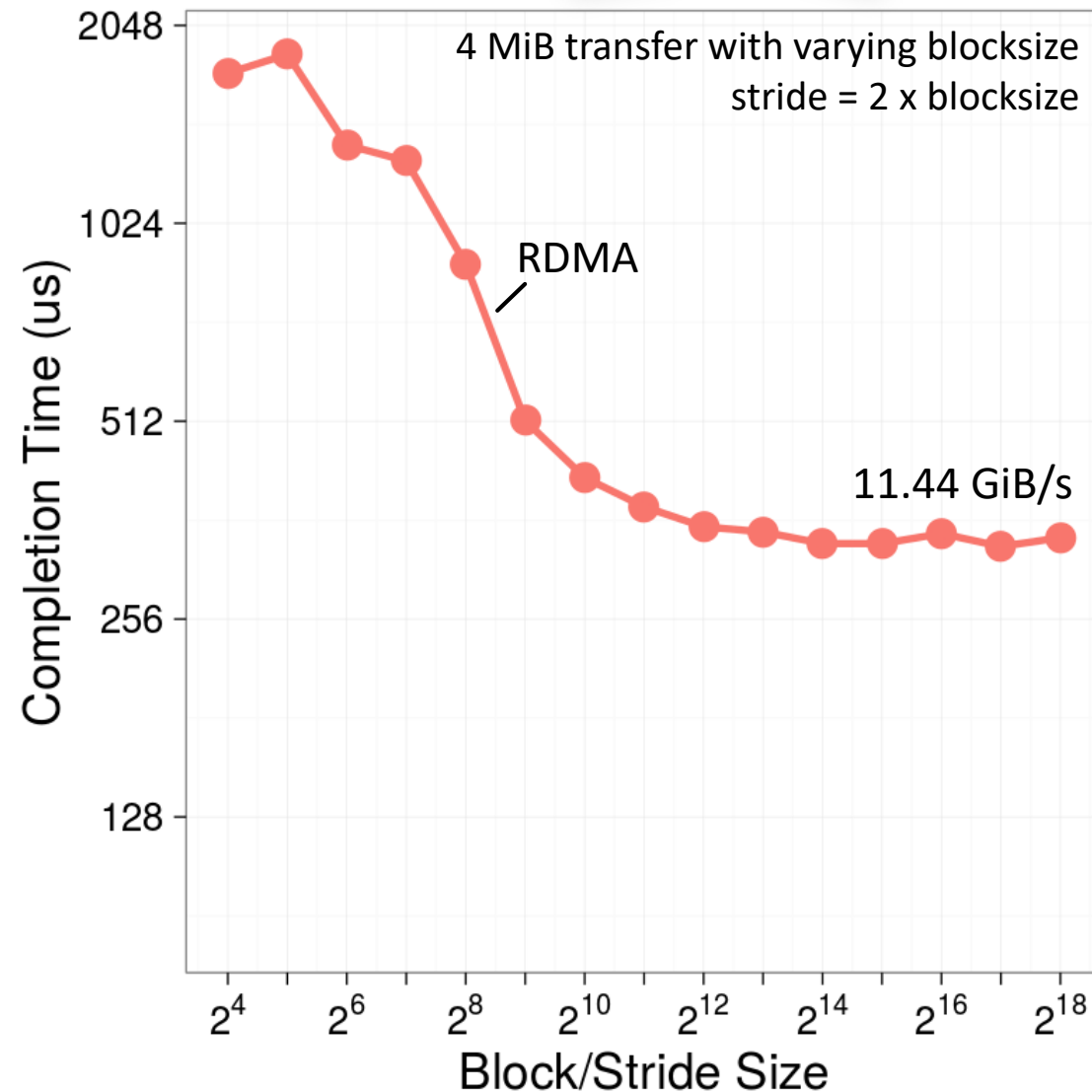
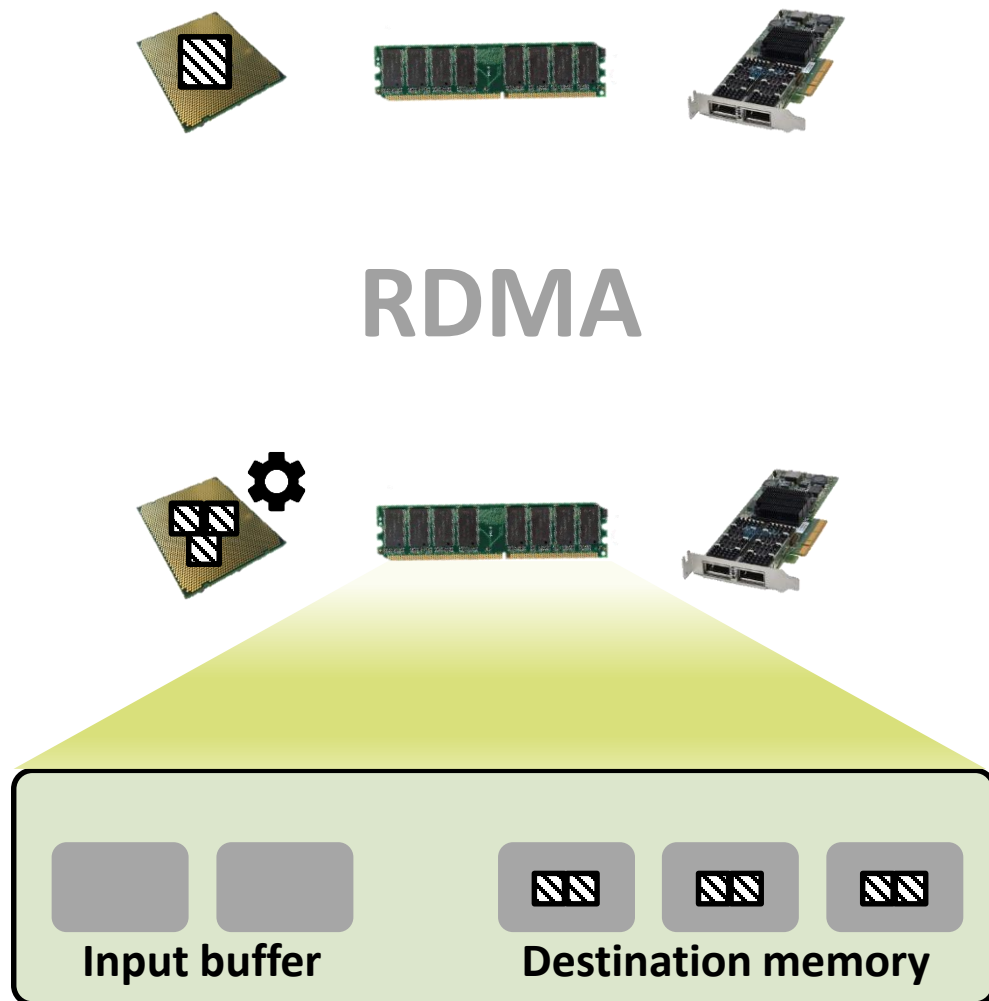




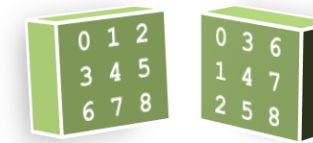
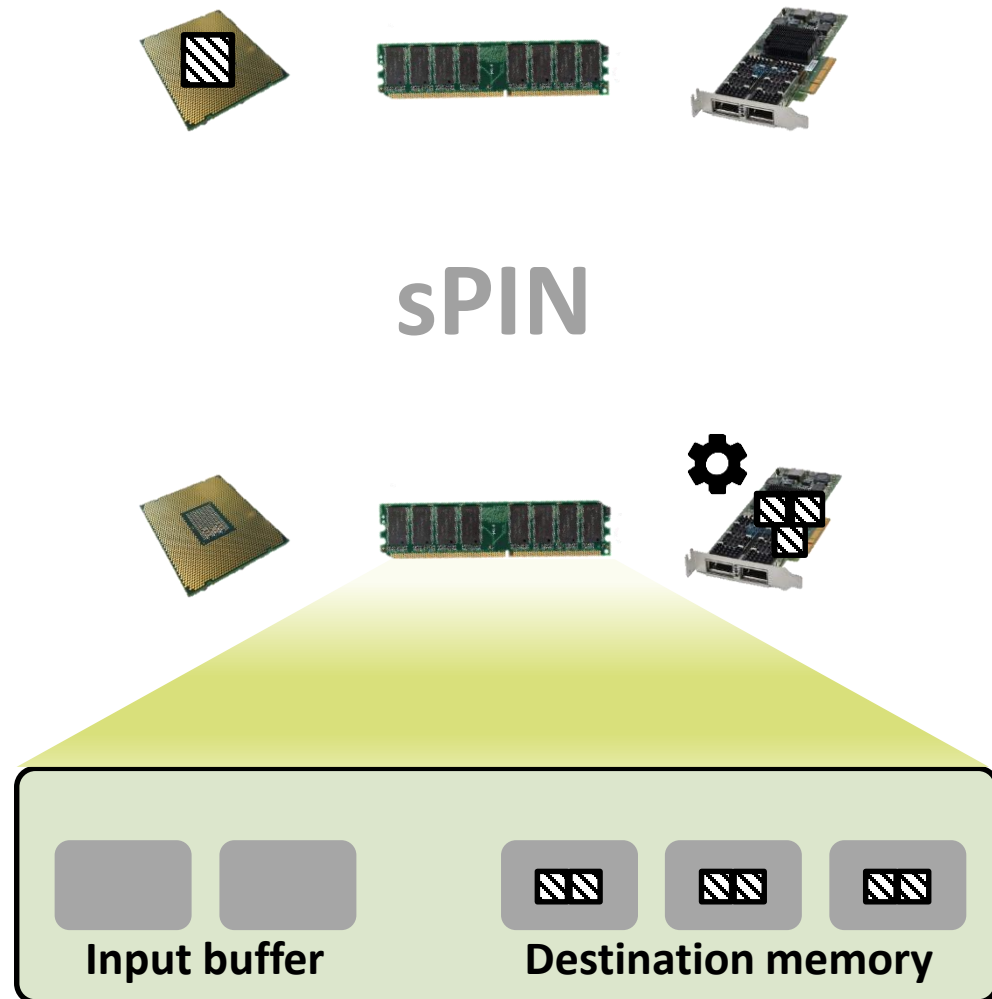
# Use Case 3: MPI Datatypes acceleration



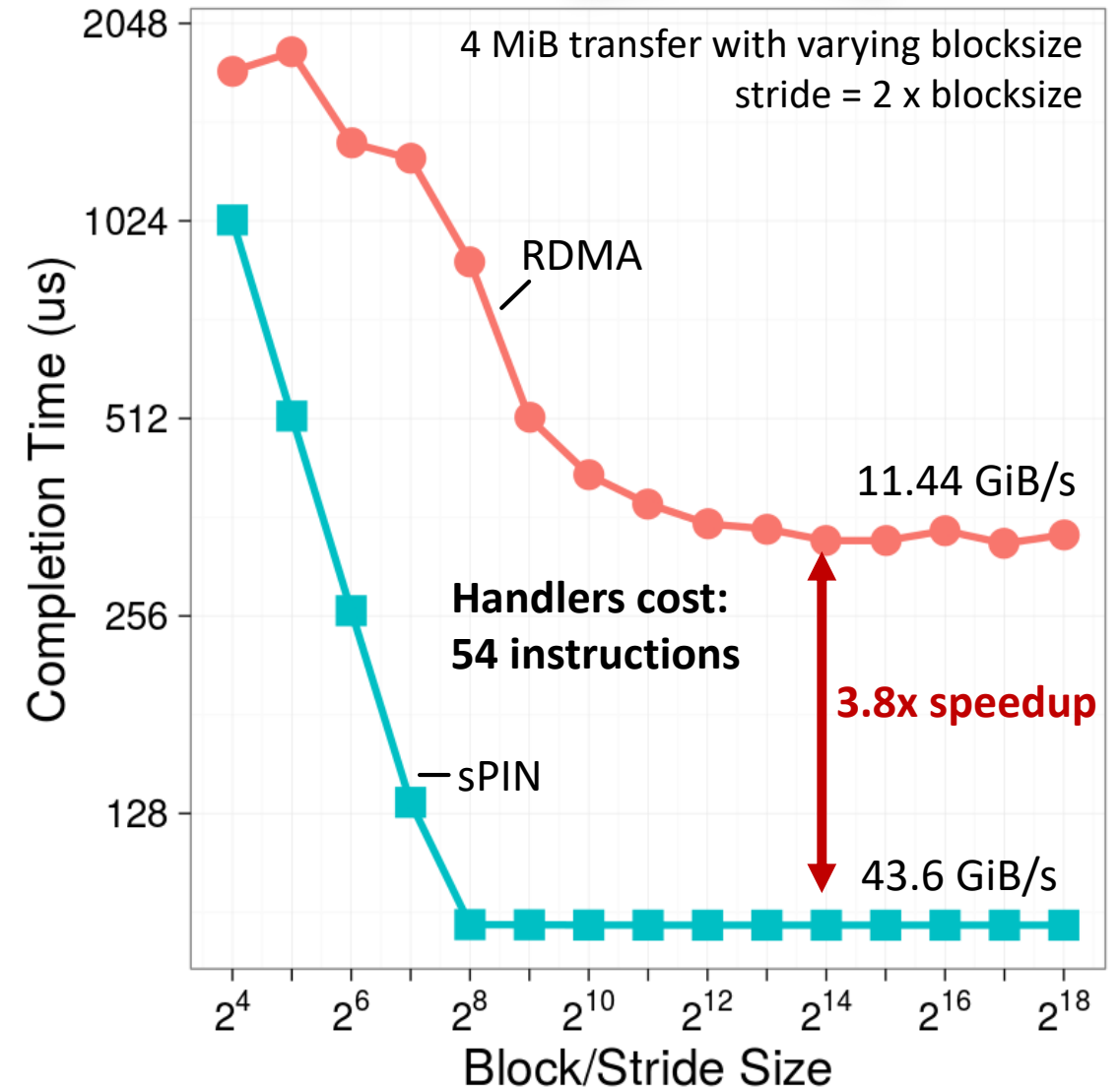
Data Layout Transformation



# Use Case 3: MPI Datatypes acceleration



Data Layout Transformation




## Further results and use-cases

# Further results and use-cases

SPCL ETH zürich


### Use Case 4: MPI Rendezvous Protocol



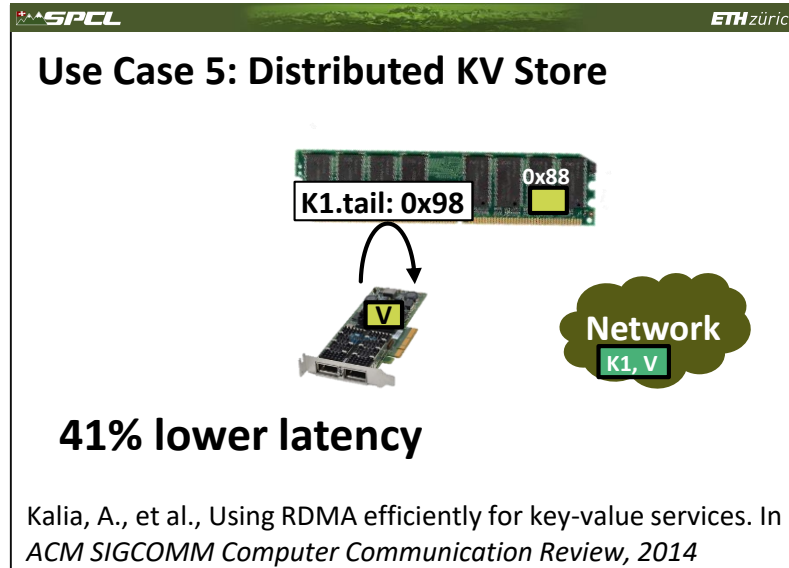
program	p	msgs	ovhd	ovhd	red
MILC	64	5.7M	5.5%	1.9%	65%
POP	64	772M	3.1%	2.4%	22%
coMD	72	5.3M	6.1%	2.4%	60%
coMD	360	28.1M	6.5%	2.8%	58%
Cloverleaf	72	2.7M	5.2%	2.4%	53%
Cloverleaf	360	15.3M	5.6%	3.2%	42%

# Further results and use-cases

**Use Case 4: MPI Rendezvous Protocol**



program	p	msgs	ovhd	ovhd	red
MILC	64	5.7M	5.5%	1.9%	65%
POP	64	772M	3.1%	2.4%	22%
coMD	72	5.3M	6.1%	2.4%	60%
coMD	360	28.1M	6.5%	2.8%	58%
Cloverleaf	72	2.7M	5.2%	2.4%	53%
Cloverleaf	360	15.3M	5.6%	3.2%	42%



# Further results and use-cases

**Use Case 4: MPI Rendezvous Protocol**

program	p	msgs	ovhd	ovhd	red
MILC	64	5.7M	5.5%	1.9%	65%
POP	64	772M	3.1%	2.4%	22%
coMD	72	5.3M	6.1%	2.4%	60%
coMD	360	28.1M	6.5%	2.8%	58%
Cloverleaf	72	2.7M	5.2%	2.4%	53%
Cloverleaf	360	15.3M	5.6%	3.2%	42%

**Use Case 5: Distributed KV Store**

**41% lower latency**

Kalia, A., et al., Using RDMA efficiently for key-value services. In *ACM SIGCOMM Computer Communication Review*, 2014

**Use Case 6: Conditional Read**

Discarded data: 80%

Speedup

Data Size

Barthels, C., et al., Designing Databases for Future High-Performance Networks. *IEEE Data Eng. Bulletin*, 2017

# Further results and use-cases

**Use Case 4: MPI Rendezvous Protocol**

program	p	msgs	ovhd	ovhd	red
MILC	64	5.7M	5.5%	1.9%	65%
POP	64	772M	3.1%	2.4%	22%
coMD	72	5.3M	6.1%	2.4%	60%
coMD	360	28.1M	6.5%	2.8%	58%
Cloverleaf	72	2.7M	5.2%	2.4%	53%
Cloverleaf	360	15.3M	5.6%	3.2%	42%

**Use Case 5: Distributed KV Store**

**41% lower latency**

Kalia, A., et al., Using RDMA efficiently for key-value services. In *ACM SIGCOMM Computer Communication Review*, 2014

**Use Case 6: Conditional Read**

Discarded data: 80%

Speedup

Data Size

Barthels, C., et al., Designing Databases for Future High-Performance Networks. *IEEE Data Eng. Bulletin*, 2017

**Use Case 7: Distributed Transactions**

data pkts

log pkts

Dragojević, A, et al., No compromises: distributed transactions with consistency, availability, and performance. SOSP'15

# Further results and use-cases

**Use Case 4: MPI Rendezvous Protocol**

program	p	msgs	ovhd	ovhd	red
MILC	64	5.7M	5.5%	1.9%	65%
POP	64	772M	3.1%	2.4%	22%
coMD	72	5.3M	6.1%	2.4%	60%
coMD	360	28.1M	6.5%	2.8%	58%
Cloverleaf	72	2.7M	5.2%	2.4%	53%
Cloverleaf	360	15.3M	5.6%	3.2%	42%

**Use Case 5: Distributed KV Store**

**41% lower latency**

Kalia, A., et al., Using RDMA efficiently for key-value services. In *ACM SIGCOMM Computer Communication Review*, 2014

**Use Case 6: Conditional Read**

Barthels, C., et al., Designing Databases for Future High-Performance Networks. *IEEE Data Eng. Bulletin*, 2017

**Use Case 7: Distributed Transactions**

Dragojević, A, et al., No compromises: distributed transactions with consistency, availability, and performance. *SOSP'15*

**Use Case 8: FT Broadcast**

Bosilca, G., et al., Failure Detection and Propagation in HPC systems. *SC'16*



# Further results and use-cases

### Use Case 4: MPI Rendezvous Protocol

program	p	msgs	ovhd	ovhd	red
MILC	64	5.7M	5.5%	1.9%	65%
POP	64	772M	3.1%	2.4%	22%
coMD	72	5.3M	6.1%	2.4%	60%
coMD	360	28.1M	6.5%	2.8%	58%
Cloverleaf	72	2.7M	5.2%	2.4%	53%
Cloverleaf	360	15.3M	5.6%	3.2%	42%

### Use Case 5: Distributed KV Store

## The Next 700 sPIN use-cases

... just think about sPIN graph kernels ...

**41% lower latency**

Kalia, A., et al., Using sPIN for distributed KV services. In ACM SIGCOMM Conference on Data Communication Systems, 2017.

### Use Case 6: Conditional Read

Discarded data: 80%

Speedup

Data Size

Barthels, C., et al., Designing Databases for Future High-Performance Networks. *IEEE Data Eng. Bulletin*, 2017

### Use Case 7: Distributed Transactions

Dragojević, A, et al., No compromises: distributed transactions with consistency, availability, and performance. SOSP'15

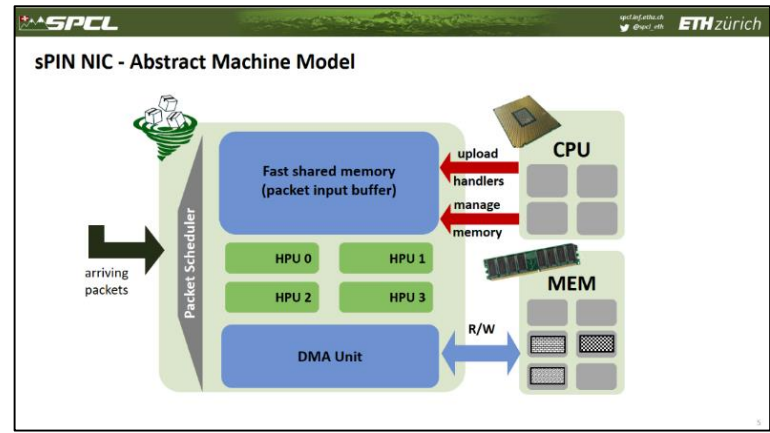
### Use Case 8: Distributed Graph Kernels

Bosilca, G., et al., Failure Detection and Propagation in HPC systems. SC'16

### Use Case 9: Distributed Consensus

István, Z., et al., Consensus in a Box: Inexpensive Coordination in Hardware. NSDI'16

# sPIN Streaming Processing in the Network for Network Acceleration



### sPIN – Programming Interface

```

Header handler
_handler int pp_header_handler(const pti_header_t h, void *state) {
    pingpong_info_t *i = state;
    i->source = h.source_id;
    return PROCESS_DATA; // execute payload handler to put from device
}

Payload handler
_handler int pp_payload_handler(const pti_payload_t b, void *state) {
    pingpong_info_t *i = state;
    PtiHandlerPutFromDevice(p.base, p.length, 1, 0, i->source, 10, 0, NULL, 0);
    return SUCCESS;
}

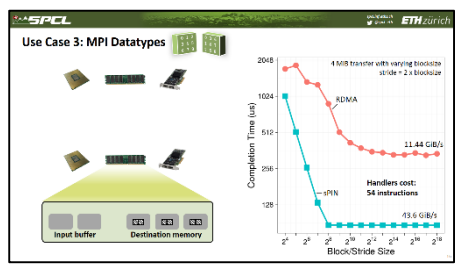
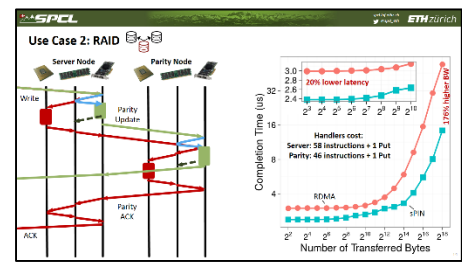
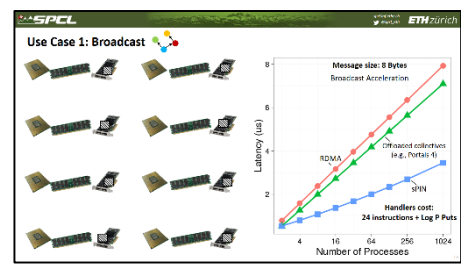
Completion handler
_handler int pp_completion_handler(int dropped_bytes,
    bool flow_control_triggered, void *state) {
    return SUCCESS;
}
    
```

connect(peer, /\* ... \*/, &pp\_header\_handler, &pp\_payload\_handler, &pp\_completion\_handler);

sPIN



beyond RDMA



### Possible sPIN implementations

- sPIN is a programming abstraction, similar to CUDA or OpenCL combined with OFED or Portals 4
- It enables a large variety of NIC implementations!
- For example, massively multithreaded HPUs
  - Including warp-like scheduling strategies
  - at 400G, process more than 833 million messages/s
- Main goal: sPIN must not obstruct line-rate
  - Programmer must limit processing time per packet
  - Little's Law: 500 instructions per handler, 2.5 GHz, IPC=1, 1 Tb/s → 25 kiB memory
  - Relies on fast shared memory (processing in packet buffers)
  - Scratchpad or registers
  - Quick (single-cycle) handler invocation on packet arrival
  - Pre-initialized memory & context
- Can be implemented in most RDMA NICs with a firmware update
  - Or in software in programmable (Smart) NICs
- Two implementation modes: integrated and discrete
  - Integrated on the same SoC (e.g., through CC mechanism)
  - Discrete is connected via bus interface (e.g., PCIe)

### A BRIEF HISTORY OF NETWORK TOPOLOGIES

copper cables, small radix switches | fiber, high-radix switches

**Key ideas:**

**"It's the diameter, stupid"**

**Lower diameter:**

- Less cables traversed
- Less cables needed
- Less routers needed

**Cost and energy savings:**

- Up to 50% over Fat Tree
- Up to 33% over Dragonfly

Bandwidth  $\approx \frac{N}{4}$   
 Latency  $= 2 - 4$   
 Radix  $= k$

### Slim NoC – Slim Fly topologies for on chip networks

**New challenges – layout in the chip's metal layers**

1296 cores, 162 routers

Network	Throughput / Power [flits/J]
SlimNOC	~0.0048
Flattened Butterfly	~0.0035
Torus	~0.0018
Mesh	~0.0012



### sPIN NIC - Abstract Machine Model

upload handlers manage memory

R/W

arriving packets

### sPIN Streaming Processing in the Network for Network Acceleration

**beyond RDMA**

Full paper: <https://arxiv.org/abs/1709.05483>

Try it out: [https://spcl.inf.ethz.ch/Research/Parallel\\_Programming/sPIN/](https://spcl.inf.ethz.ch/Research/Parallel_Programming/sPIN/)

## Backup Slides

# Slim Fly Backup

# DESIGNING AN EFFICIENT NETWORK TOPOLOGY

## CONNECTING ROUTERS: DIAMETER 2

1 Select a prime power  $q$

$$q = 4w + \delta;$$

$$w \in \mathbb{N} \quad \delta \in \{-1, 0, 1\},$$

A Slim Fly based on  $q$   
 Number of routers:  $2q^2$   
 Network radix:  $(3q - \delta)/2$

2 Construct a finite field  $\mathcal{F}_q$

Assuming  $q$  is prime:

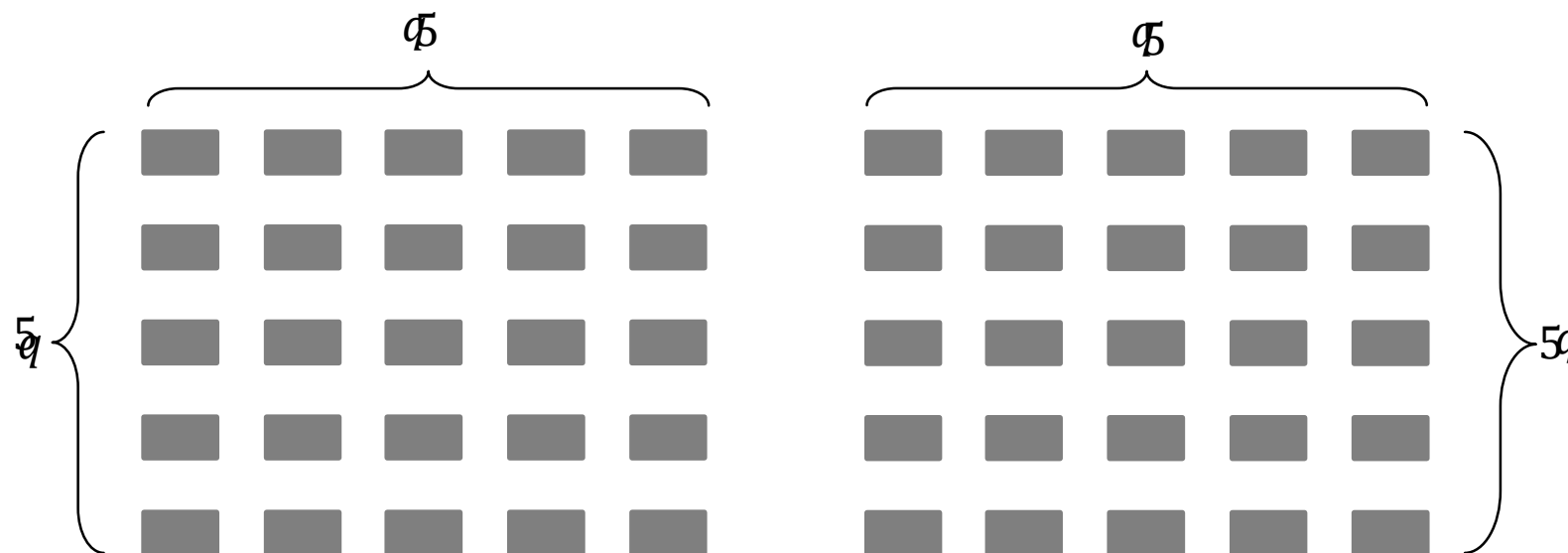
$$\mathcal{F}_q = \mathbb{Z}/q\mathbb{Z} = \{0, 1, \dots, q - 1\}$$

with modular arithmetic.

E Example:  $q = 5$

50 routers  
 network radix: 7

$$\mathcal{F}_5 = \{0, 1, 2, 3, 4\}$$



# DESIGNING AN EFFICIENT NETWORK TOPOLOGY

## CONNECTING ROUTERS: DIAMETER 2

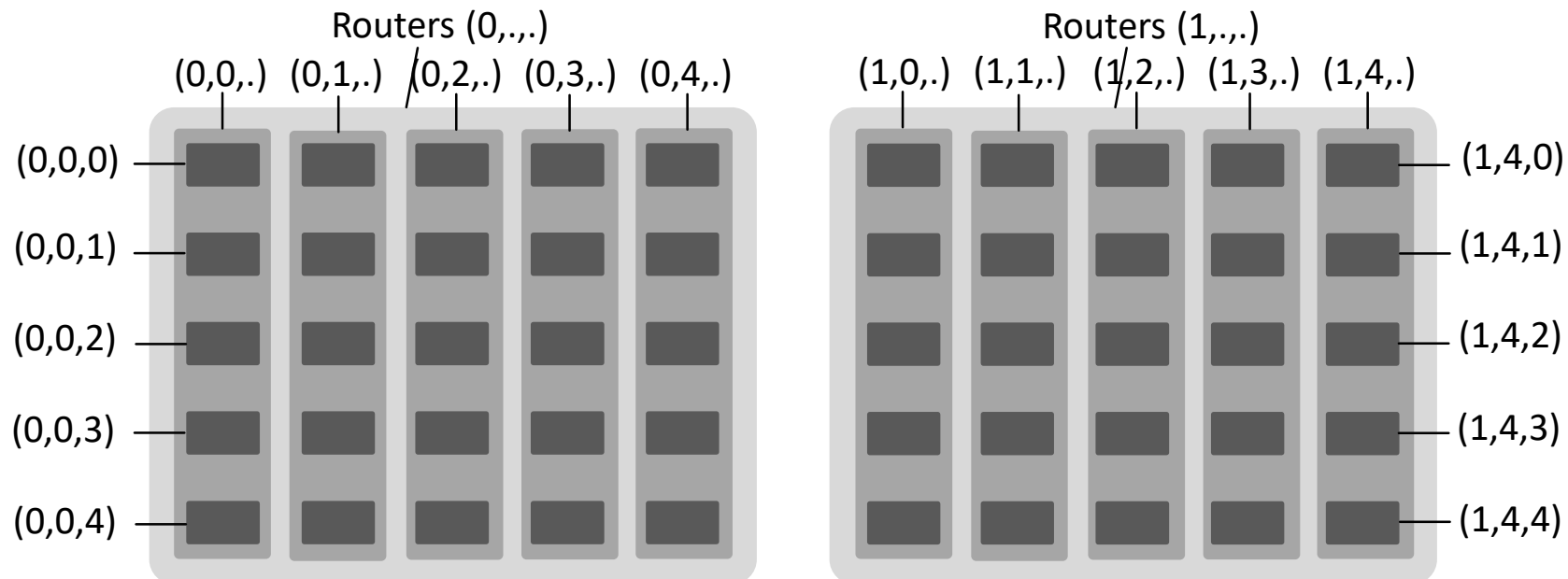
3 Label the routers

Set of routers:

$$\{0,1\} \times \mathcal{F}_q \times \mathcal{F}_q$$

E Example:  $q = 5$

...



# DESIGNING AN EFFICIENT NETWORK TOPOLOGY

## CONNECTING ROUTERS: DIAMETER 2

**4** Find primitive element  $\xi$   
 $\xi \in \mathcal{F}_q$  generates  $\mathcal{F}_q$   
 All non-zero elements of  $\mathcal{F}_q$   
 can be written as  $\xi^i; i \in \mathbb{N}$

**5** Build Generator Sets  
 $X = \{1, \xi^2, \dots, \xi^{q-3}\}$   
 $X' = \{\xi, \xi^3, \dots, \xi^{q-2}\}$

**E** Example:  $q = 5$

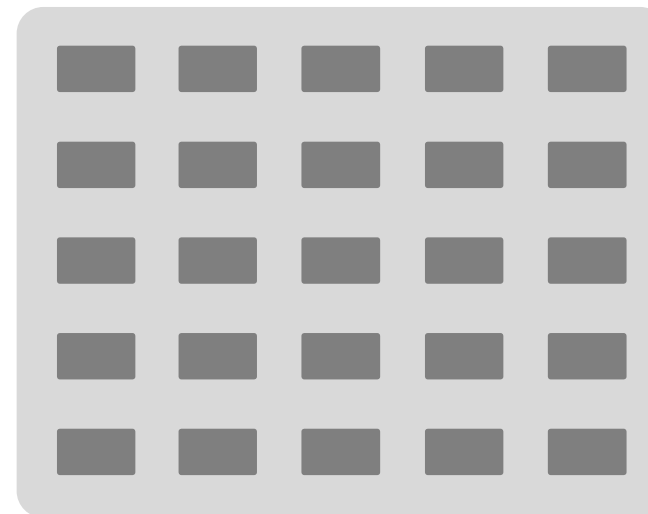
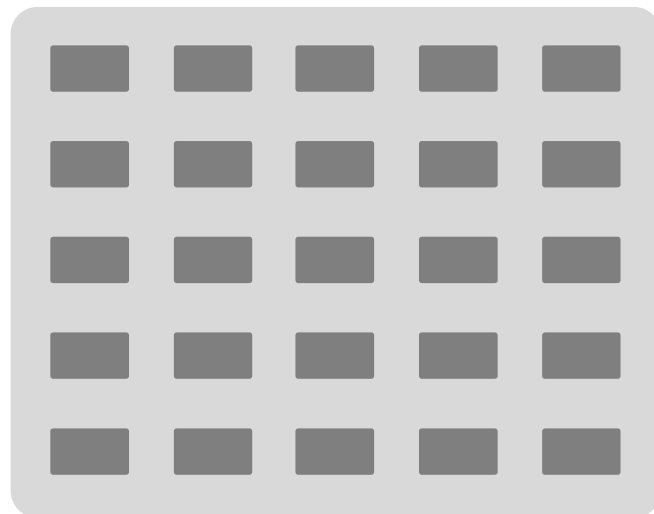
$$\mathcal{F}_5 = \{0, 1, 2, 3, 4\}$$

$$\xi = 2$$

$$1 = \xi^4 \bmod 5 = 2^4 \bmod 5 = 16 \bmod 5$$

$$X = \{1, 4\}$$

$$X' = \{2, 3\}$$





# DESIGNING AN EFFICIENT NETWORK TOPOLOGY

## CONNECTING ROUTERS: DIAMETER 2

### 6 Intra-group connections

Two routers in one group are connected iff their “vertical Manhattan distance” is an element from:

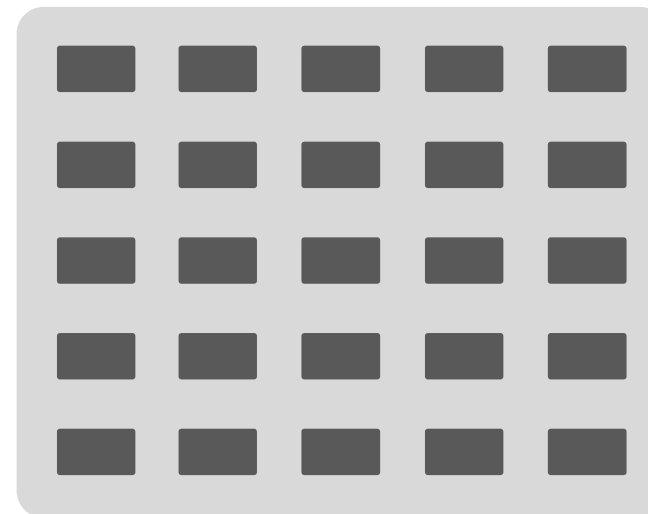
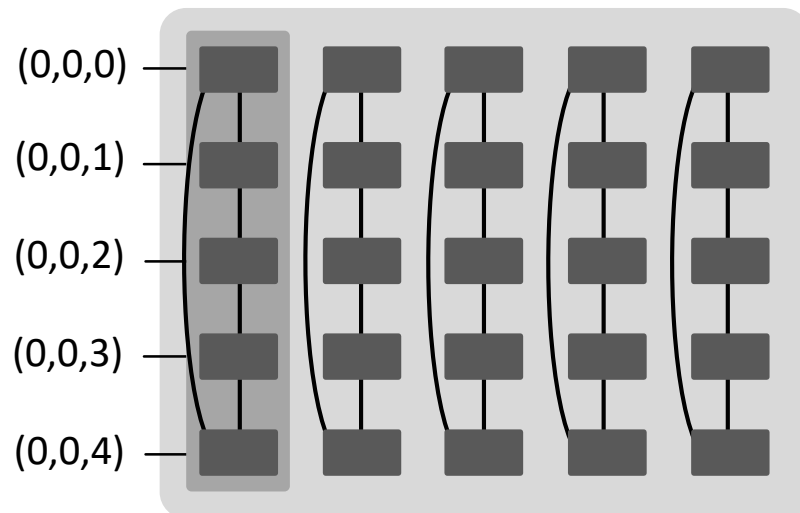
$$X = \{1, \xi^2, \dots, \xi^{q-3}\} \text{ (for subgraph 0)}$$

$$X' = \{\xi, \xi^3, \dots, \xi^{q-2}\} \text{ (for subgraph 1)}$$

### E Example: $q = 5$

Take Routers  $(0,0,.)$

$$X = \{1, 4\}$$



# DESIGNING AN EFFICIENT NETWORK TOPOLOGY

## CONNECTING ROUTERS: DIAMETER 2

### 6 Intra-group connections

Two routers in one group are connected iff their “vertical Manhattan distance” is an element from:

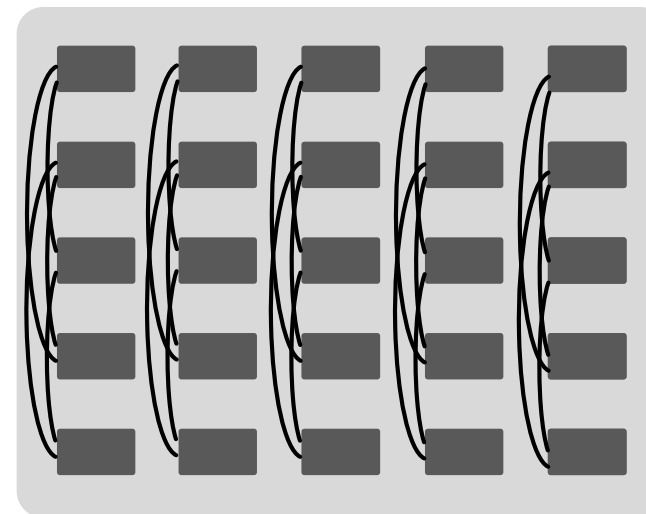
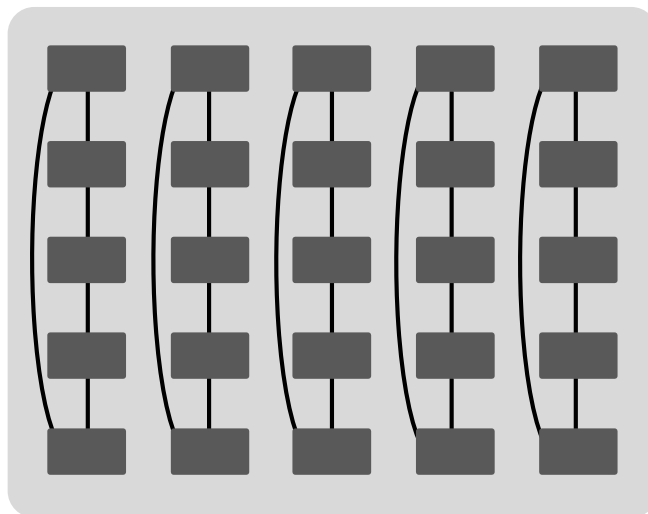
$$X = \{1, \xi^2, \dots, \xi^{q-3}\} \text{ (for subgraph 0)}$$

$$X' = \{\xi, \xi^3, \dots, \xi^{q-2}\} \text{ (for subgraph 1)}$$

E Example:  $q = 5$

Take Routers  $(1, 4, \dots)$

$$X' = \{2, 3\}$$



# DESIGNING AN EFFICIENT NETWORK TOPOLOGY

## CONNECTING ROUTERS: DIAMETER 2

### 7 Inter-group connections

Router  $(0, x, y) \leftrightarrow (1, m, c)$

iff  $y = mx + c$

### E Example: $q = 5$

Take Router  $(1,0,0)$

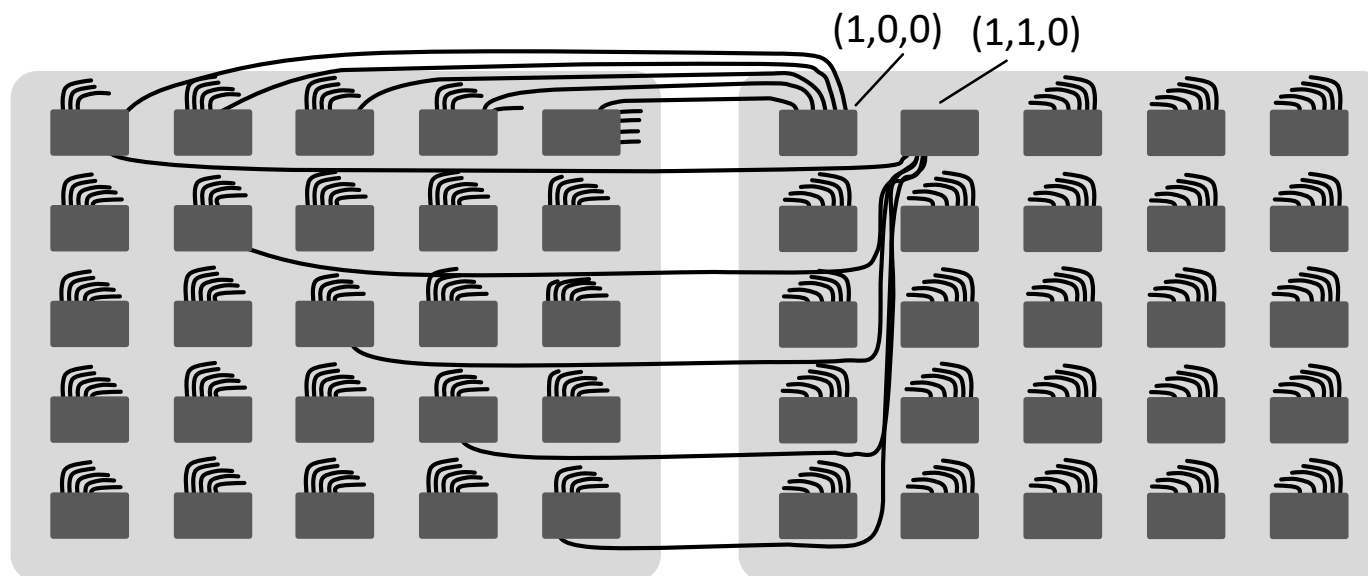
$m = 0, c = 0$

$(1,0,0) \leftrightarrow (0, x, 0)$

Take Router  $(1,1,0)$

$m = 1, c = 0$

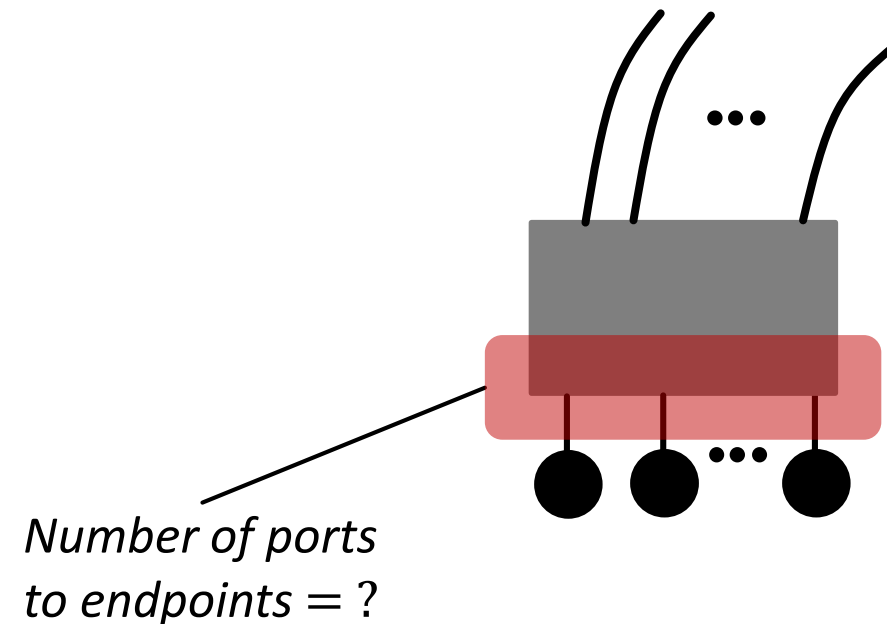
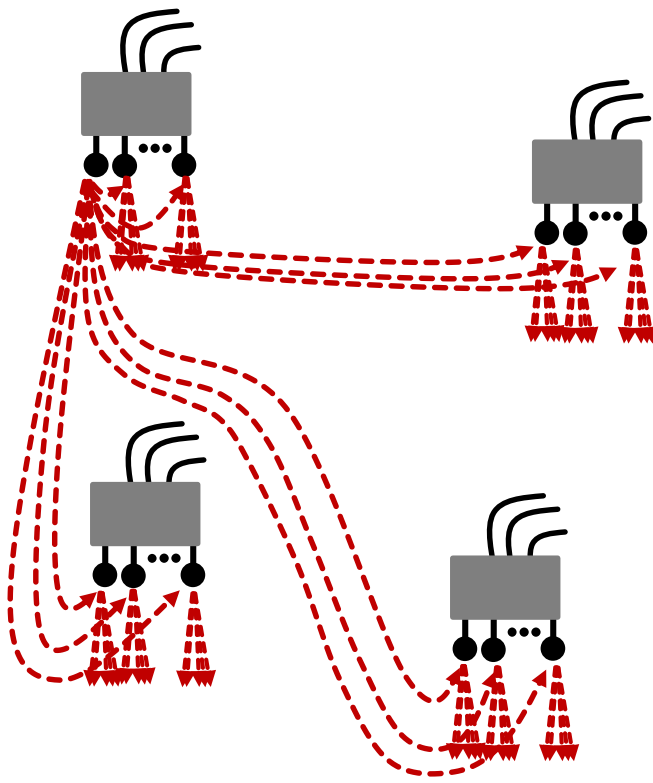
$(1,0,0) \leftrightarrow (0, x, x)$



# DESIGNING AN EFFICIENT NETWORK TOPOLOGY

## ATTACHING ENDPOINTS: DIAMETER 2

- How many endpoints do we attach to each router?
- As many to ensure *full global bandwidth*:
  - Global bandwidth: the theoretical cumulative throughput in all-to-all in a steady state



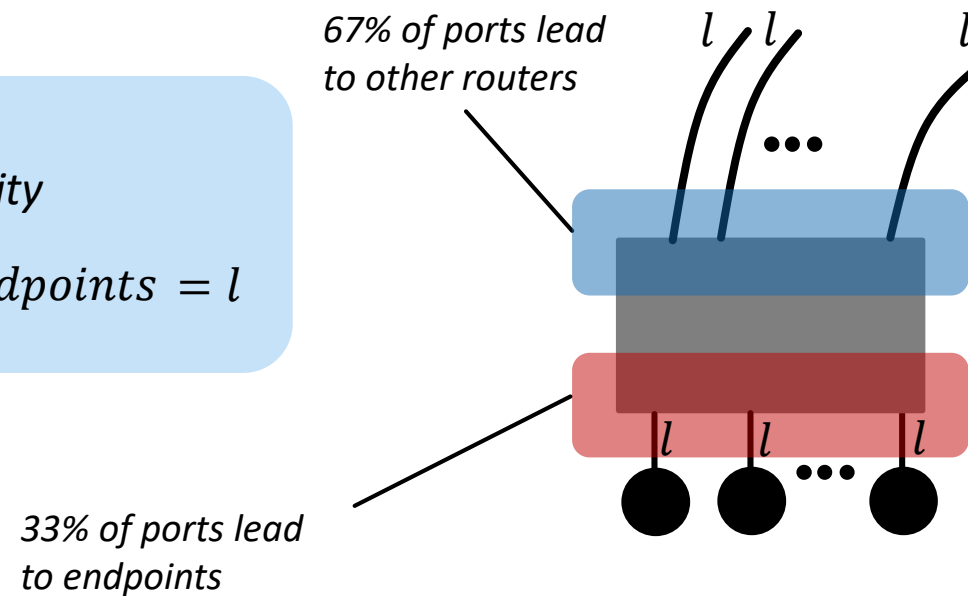
# DESIGNING AN EFFICIENT NETWORK TOPOLOGY

## ATTACHING ENDPOINTS: DIAMETER 2

- 1 Get load  $l$  per router-router channel (average number of routes per channel)

$$l = \frac{\text{total number of routes}}{\text{total number of channels}}$$

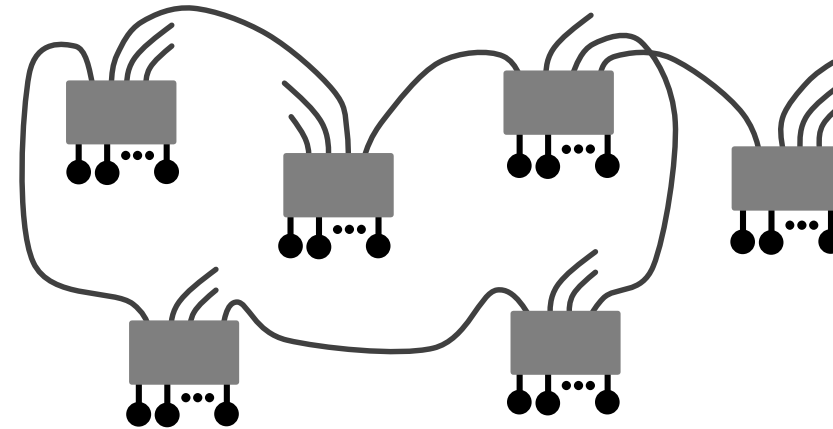
- 2 Make the network balanced, i.e.,:  
each endpoint can inject at full capacity  
 $\text{local uplink load} = \text{number of endpoints} = l$



# STRUCTURE ANALYSIS

## RESILIENCY

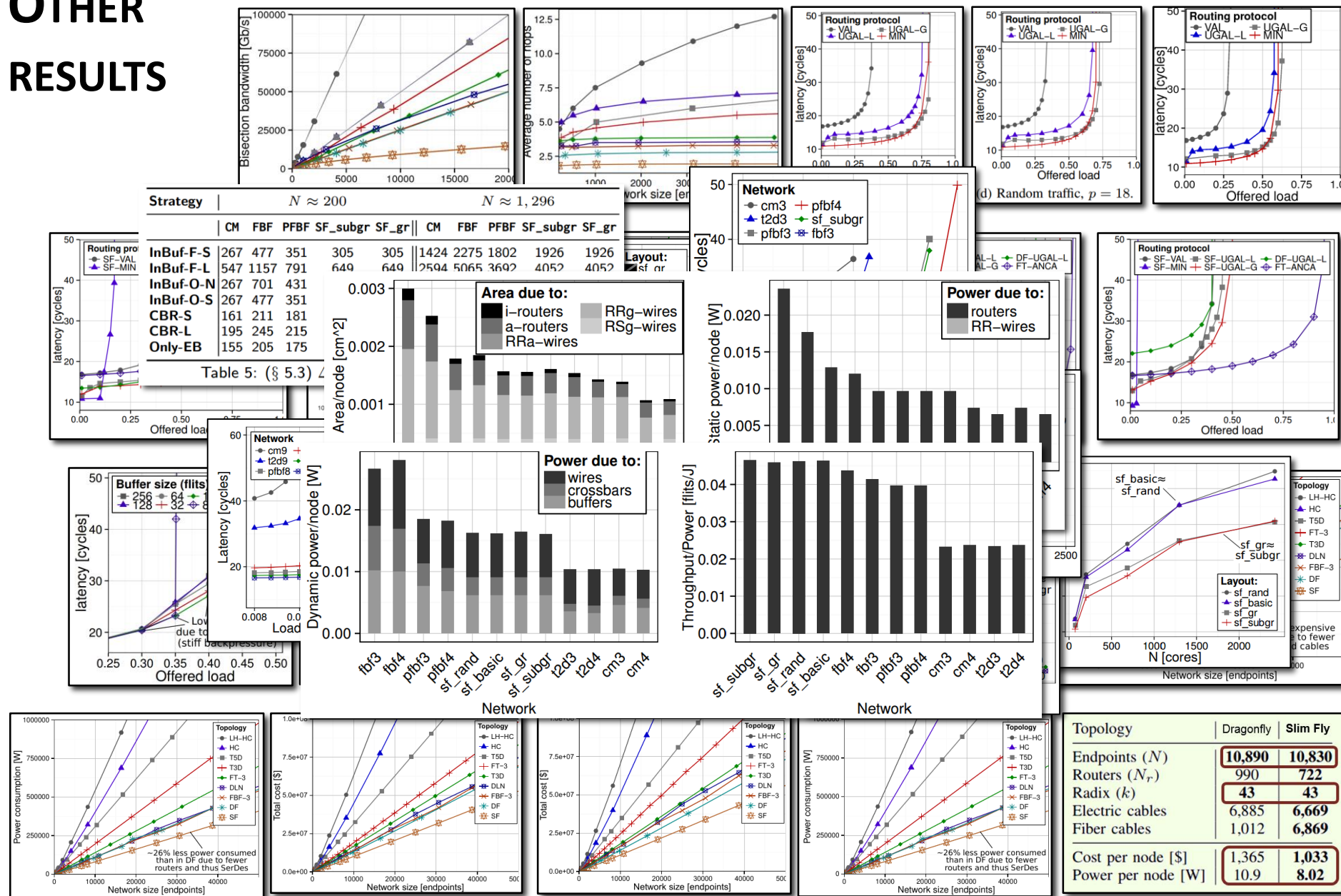
- Disconnection metrics
- Other studied metrics:
  - Average path length (increase by 2); SF is 10% more resilient than DF



Number of endpoints	Torus3D	Torus5D	Hypercube	Long Hop	Fat tree	Dragonfly	Flat. Butterfly	Random	Slim Fly
512	30%	-	40%	55%	35%	-	55%	60%	<b>60%</b>
1024	25%	40%	40%	55%	40%	50%	60%	-	-
2048	20%	-	40%	55%	40%	55%	65%	65%	<b>65%</b>
4096	15%	-	45%	55%	55%	60%	70%	70%	<b>70%</b>
8192	10%	35%	45%	55%	60%	65%	-	75%	<b>75%</b>

“-” means that a given topology does not have a variant of a given size

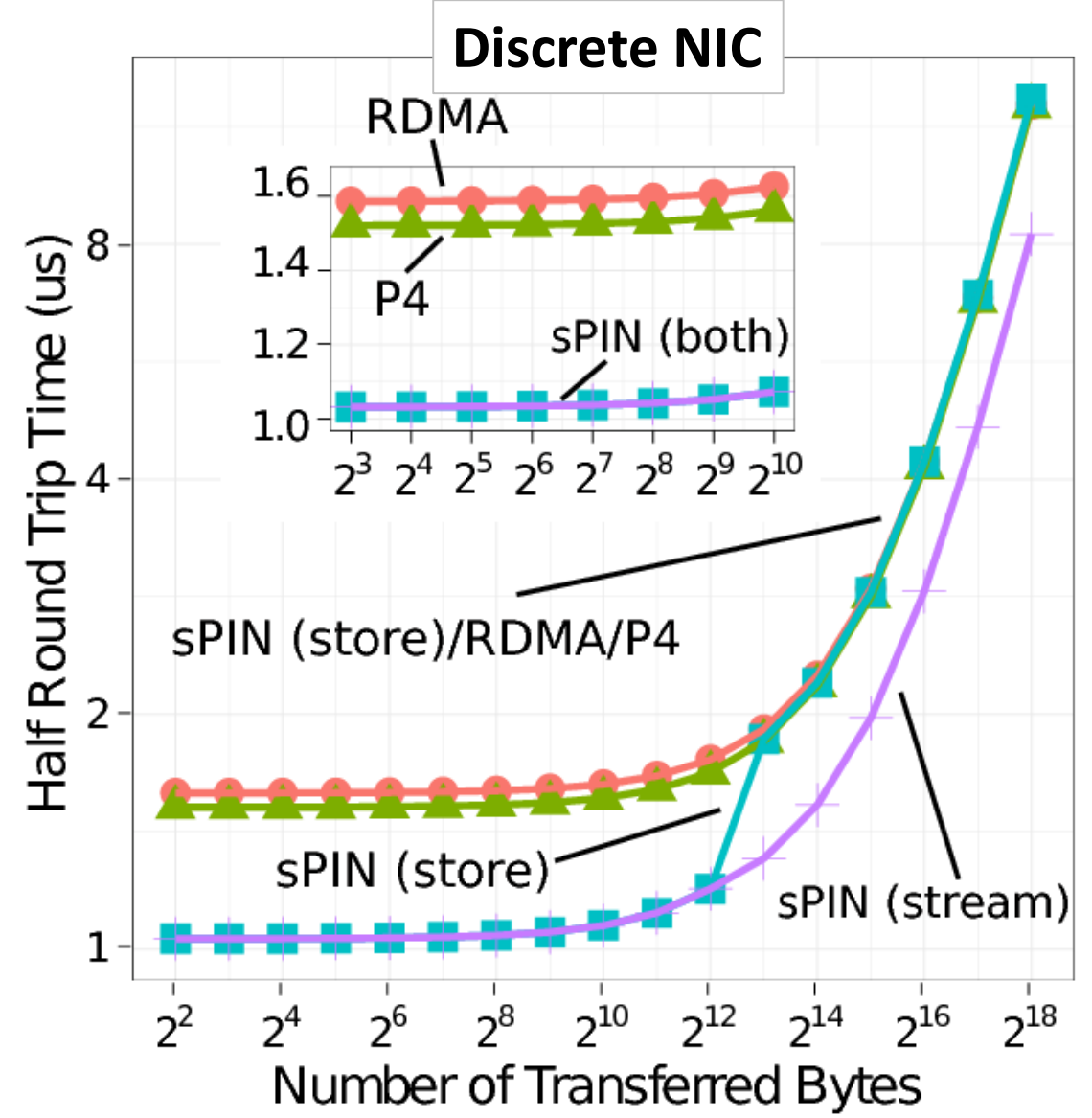
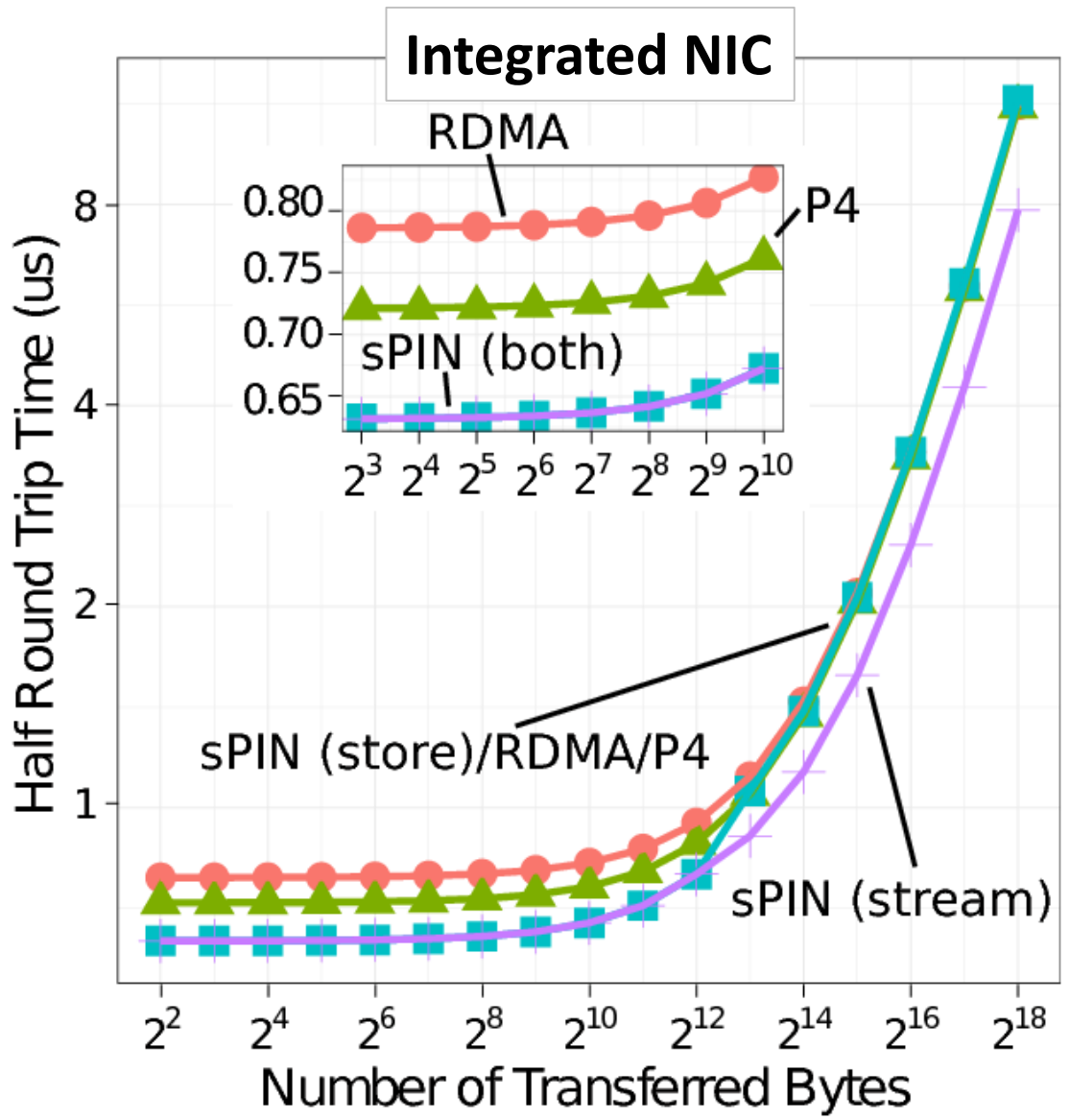
# OTHER RESULTS



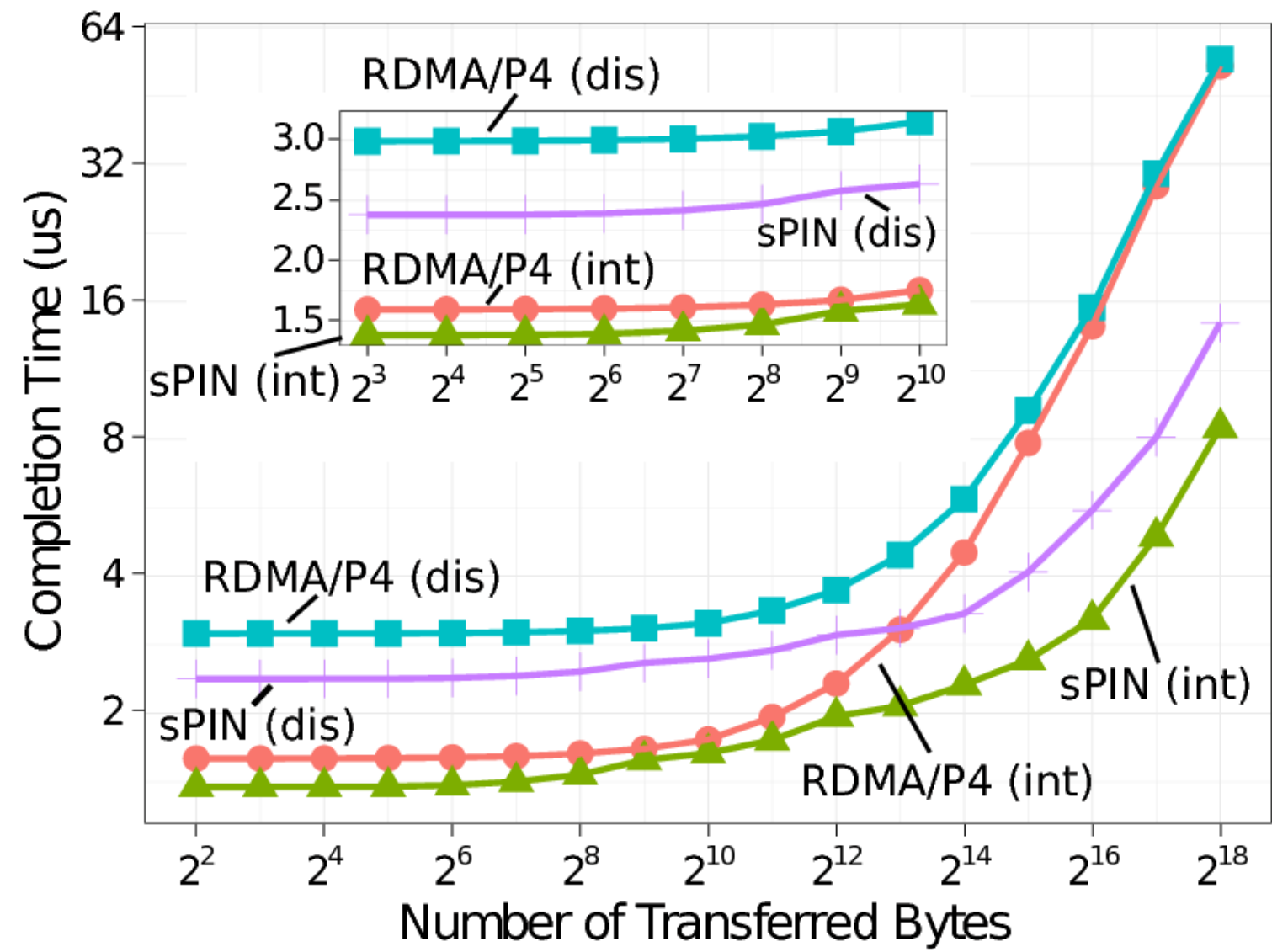
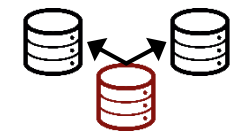
# sPIN backup



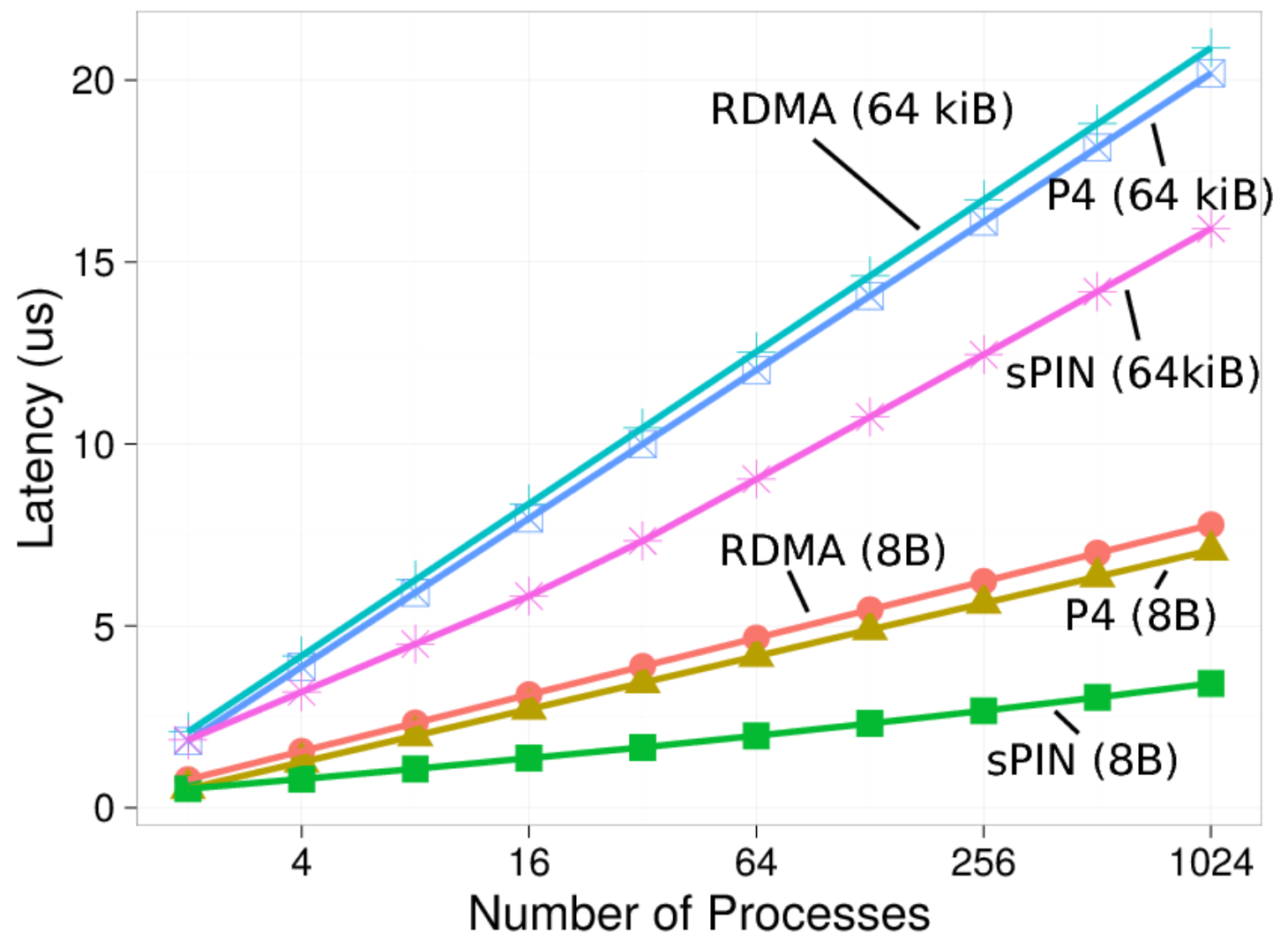
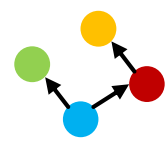
# Ping-Pong results (integrated/discrete)



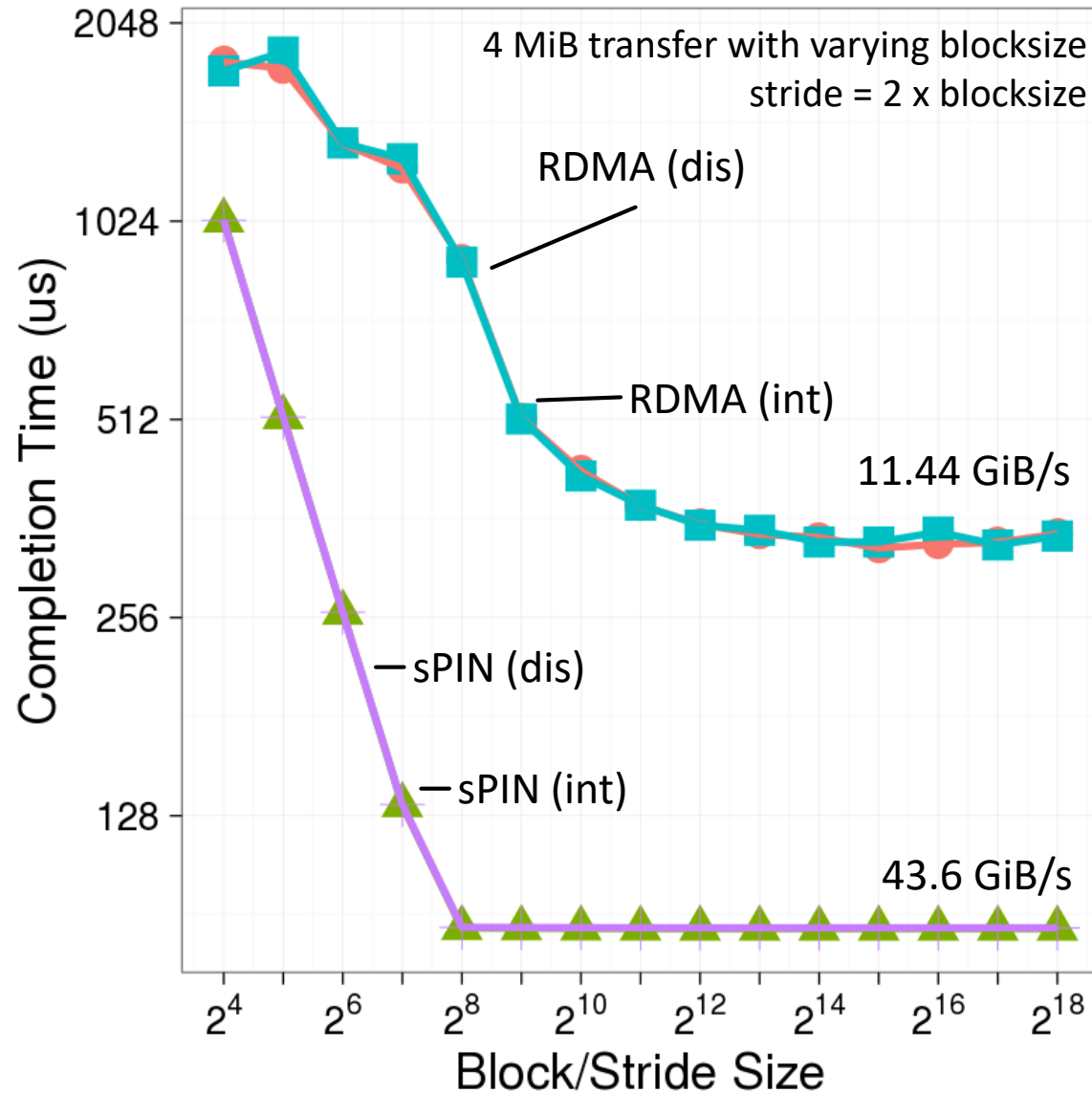
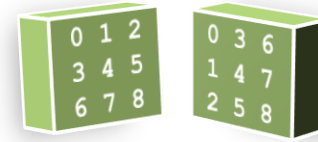
# RAID acceleration (integrated/discrete)



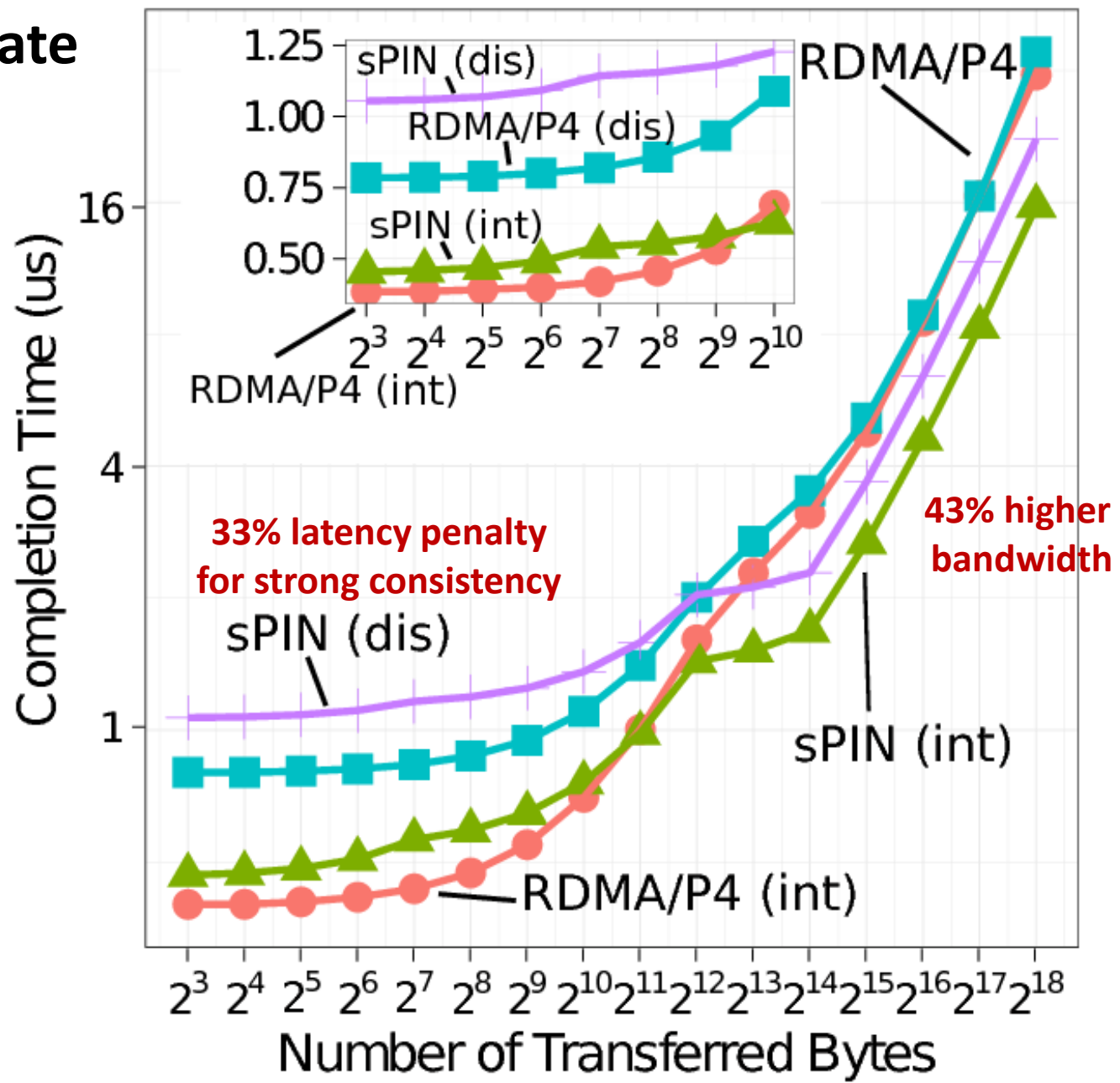
# Broadcast acceleration for large messages



# MPI Datatypes acceleration (integrated/discrete)



# Remote Accumulate



# HPUs needed depending on packet size and execution time per packet

