

Network topologies for large-scale compute centers: It's the diameter, stupid!

TORSTEN HOEFLER

with support of Maciej Besta @ SPCL
presented at Hot Interconnects 2016, San Jose, CA, USA





50% [1]



33% [2]

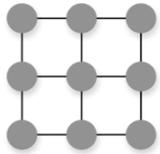
[1] D. Abts et al. (2010), *Energy Proportional Datacenter Networks*, ISCA'10

[2] J. Kim et al. (2007), *Flattened Butterfly: A Cost-Efficient Topology for High-Radix Networks*, ISCA'07

A BRIEF HISTORY OF NETWORK TOPOLOGIES

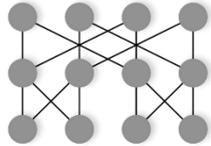
copper cables, small radix switches

fiber, high-radix switches

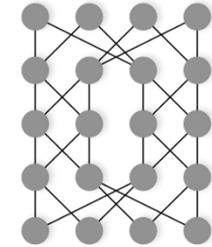


Mesh

1980's

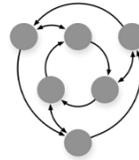


Butterfly



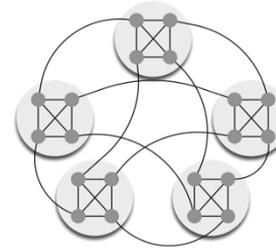
Clos/Benes

2000's



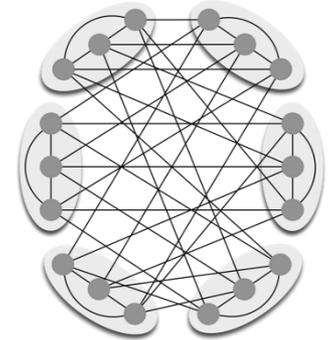
Kautz

~2005



Dragonfly

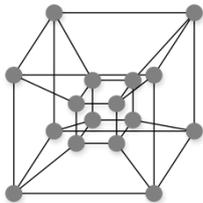
2008



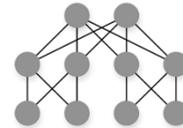
Slim Fly

2014

Hypercube

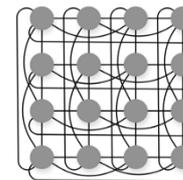


Fat Trees



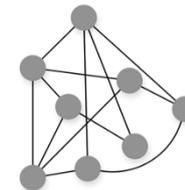
2007

Flat Fly



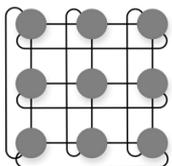
2008

Random

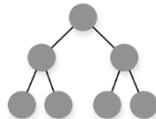


????

Torus



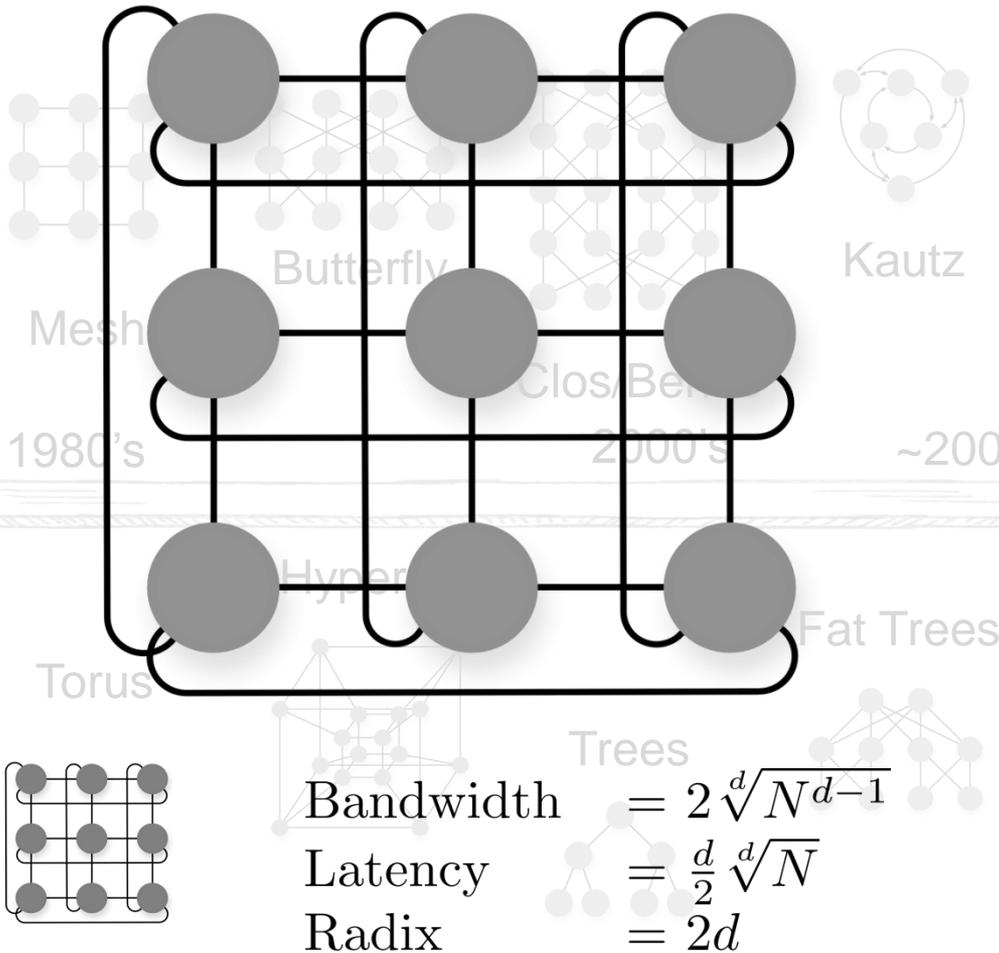
Trees



A BRIEF HISTORY OF NETWORK TOPOLOGIES

copper cables, small radix switches

fiber, high-radix switches



2008

2014



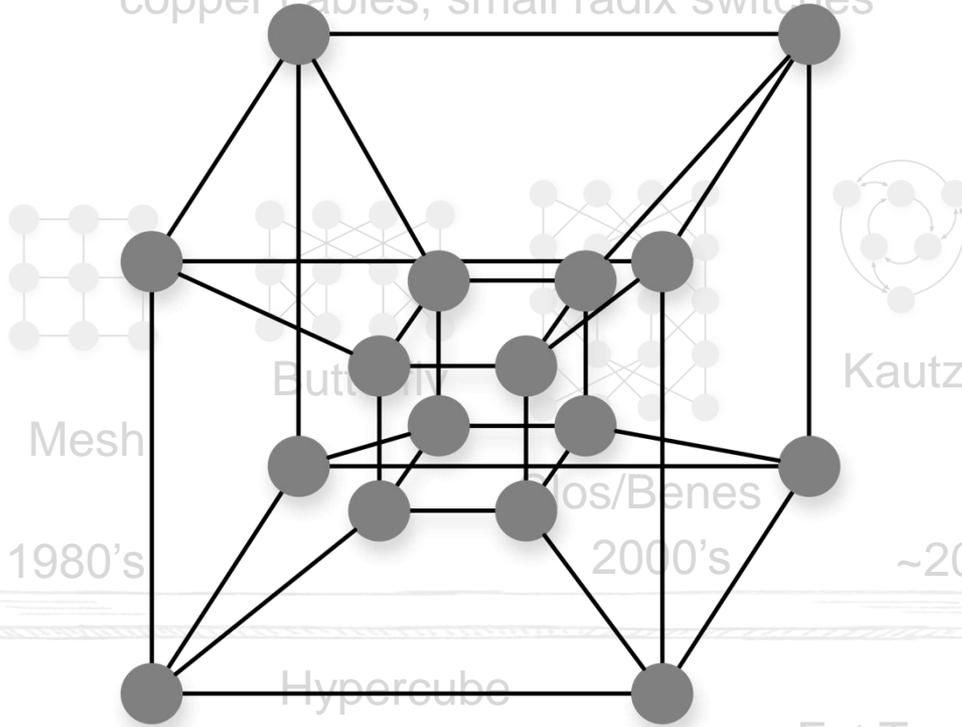
2014

2014

A BRIEF HISTORY OF NETWORK TOPOLOGIES

copper cables, small radix switches

fiber, high-radix switches



1980's

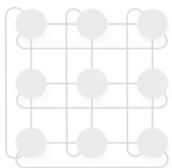
2000's

~2005

2008

2014

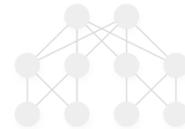
Torus



Bandwidth
Latency
Radix

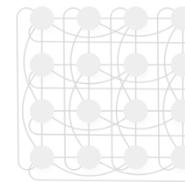
$$\begin{aligned} &= \frac{N}{2} \\ &= \log_2 N \\ &= \log_2 N \end{aligned}$$

Fat Trees



2007

Flat Fly



2008

Random



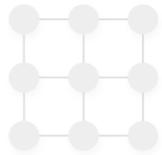
????



A BRIEF HISTORY OF NETWORK TOPOLOGIES

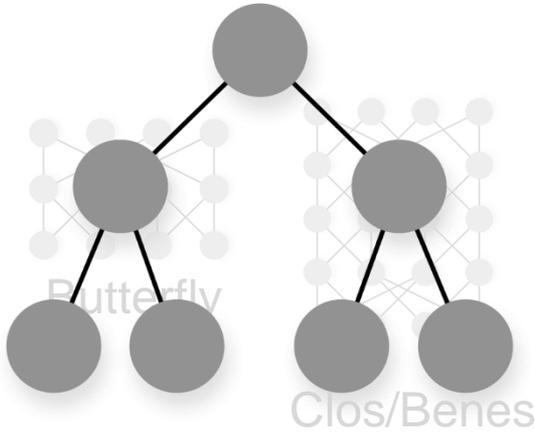
copper cables, small radix switches

fiber, high-radix switches



Mesh

1980's



Butterfly

Clos/Benes

2000's



Kautz

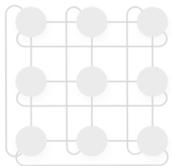
~2005



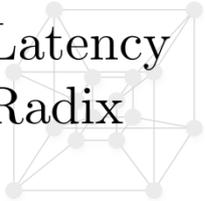
2008

2014

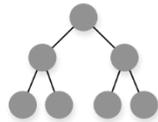
Torus



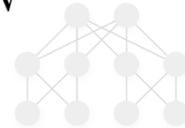
Hypercube
Bandwidth
Latency
Radix



$$\begin{aligned} \text{Bandwidth} &= 1 \\ \text{Latency} &= 2 \log_2 N \\ \text{Radix} &= 2 \end{aligned}$$

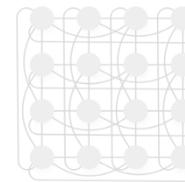


Fat Trees



2007

Flat Fly



2008

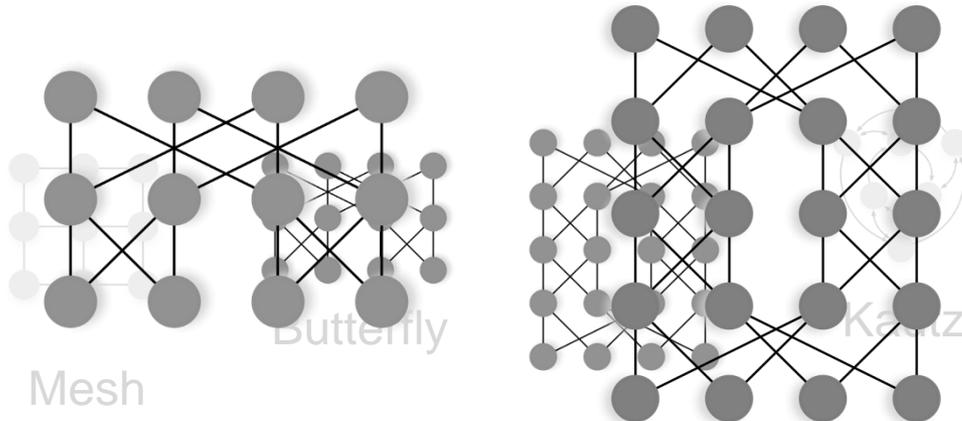
Random



????

A BRIEF HISTORY OF NETWORK TOPOLOGIES

copper cables, small radix switches



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

Latency = $2 \log_2 N$

Radix = 4

1980's

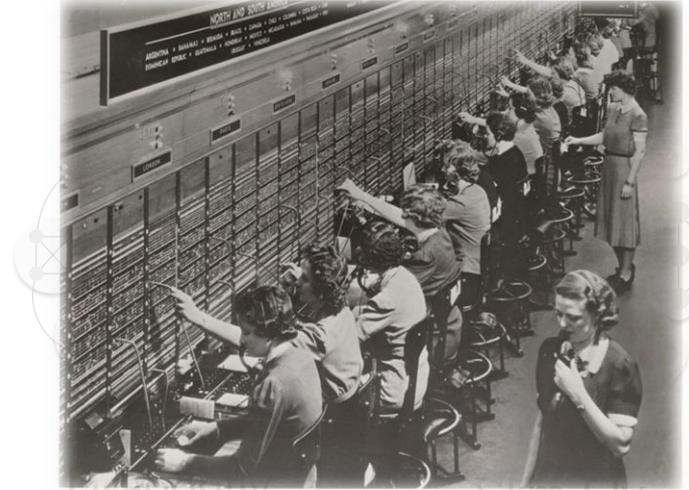
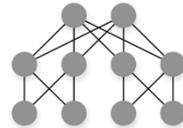
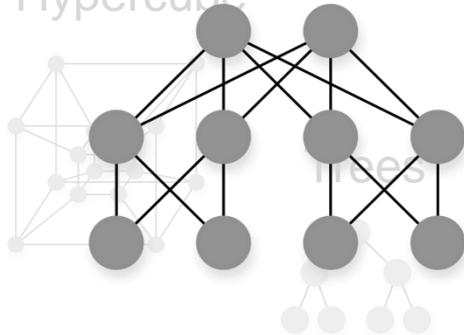
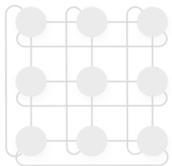
2000's

~2005

Hypercube

Fat Trees

Torus



Dragonfly

Slim Fly

2008

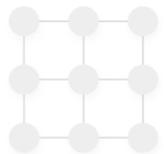
2014



A BRIEF HISTORY OF NETWORK TOPOLOGIES

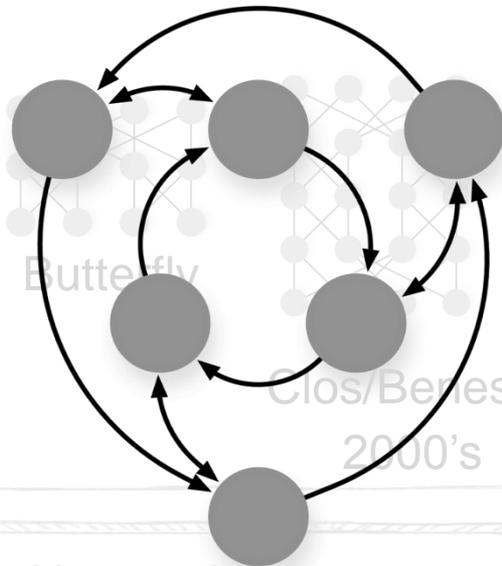
copper cables, small radix switches

fiber, high-radix switches



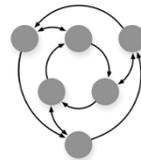
Mesh

1980's



Butterfly

Clos/Benes

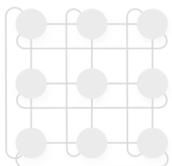


Kautz

~2005

Hypercube

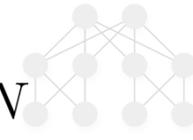
Torus



Bandwidth
Latency
Radix

$$\begin{aligned} &= \rightarrow \frac{N}{4} \\ &= \log_k N \\ &= k \end{aligned}$$

Fat Trees



Dragonfly

Slim Fly



20

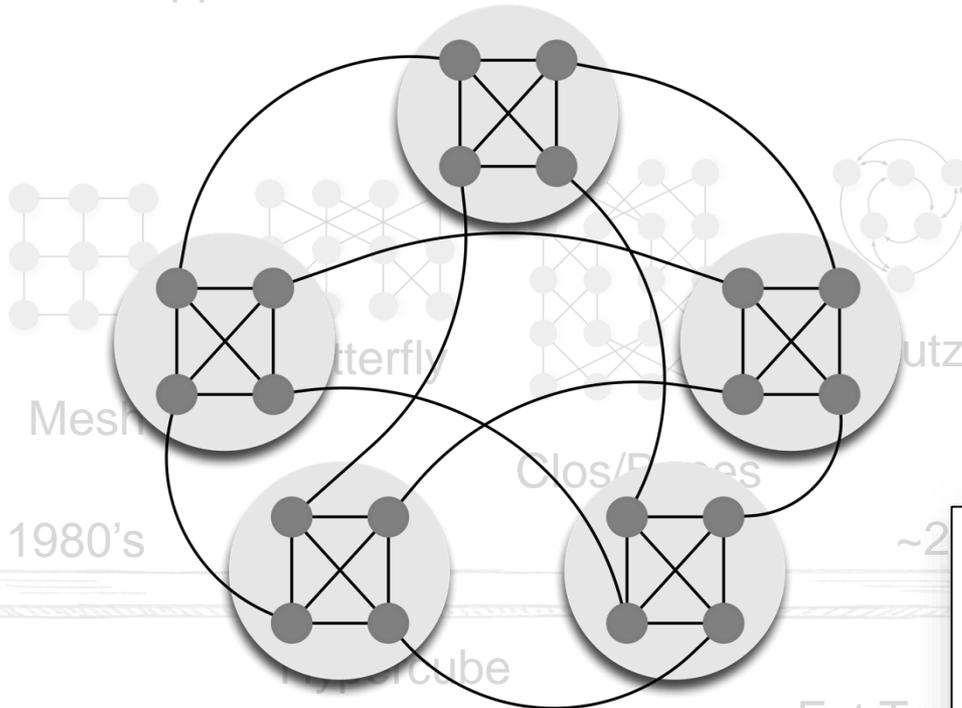
Fla

????

A BRIEF HISTORY OF NETWORK TOPOLOGIES

copper cables, small radix switches

fiber, high-radix switches



Dragonfly

Slim Fly

Torus

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

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

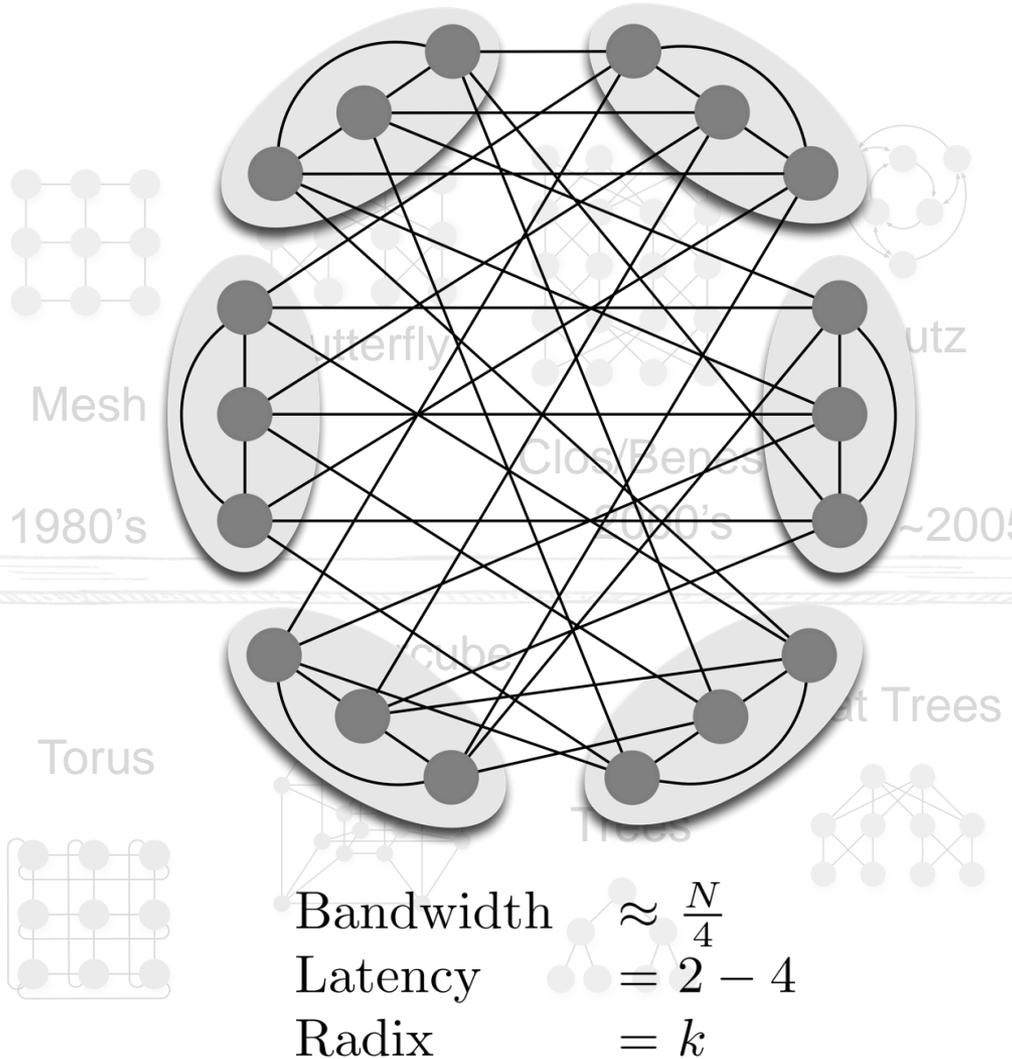
Abstract—The PERCS system was designed by IBM in response to a DARPA challenge that called for a high-productivity high-performance computing system. A major innovation in the PERCS design is the network that is built using Hub chips that are integrated into the compute nodes. Each Hub chip is about 580 mm² 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 thou-

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

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

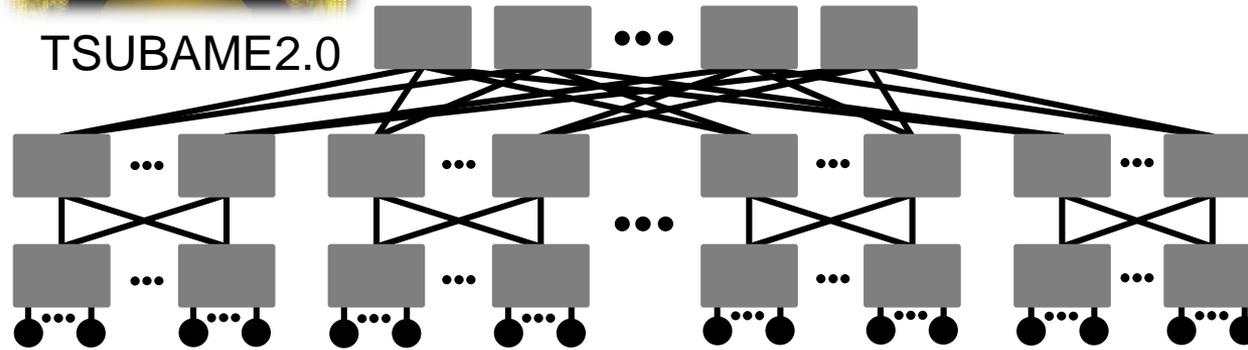
DESIGNING AN EFFICIENT NETWORK TOPOLOGY

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



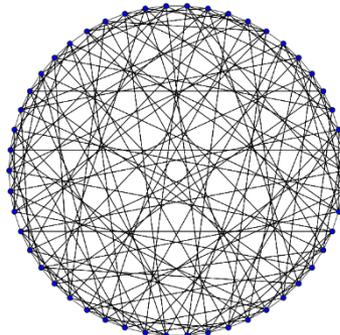
3-level fat tree:

TSUBAME2.0



diameter = 4

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



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

DESIGNING AN EFFICIENT NETWORK TOPOLOGY

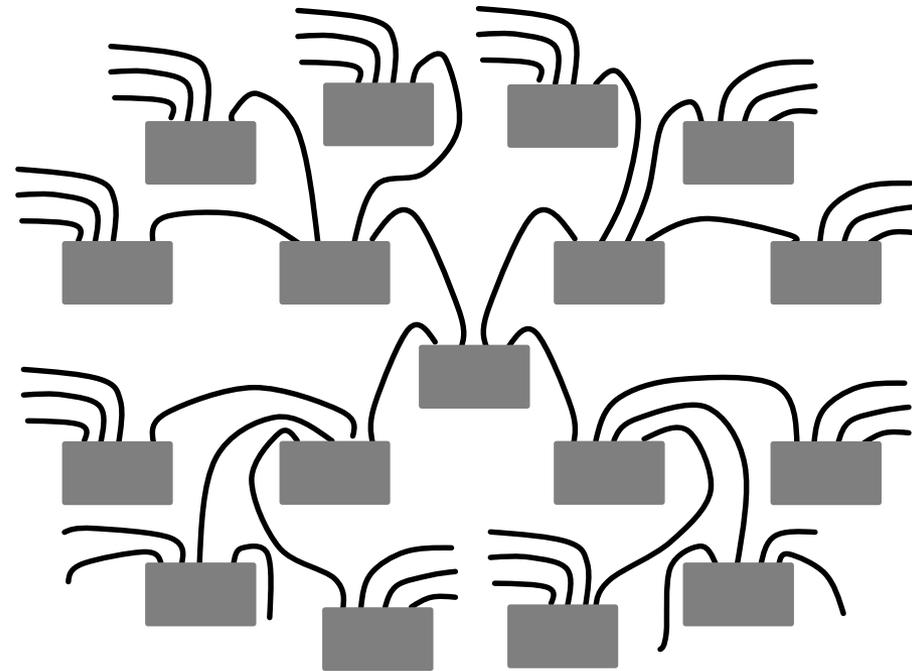


Key method

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

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

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

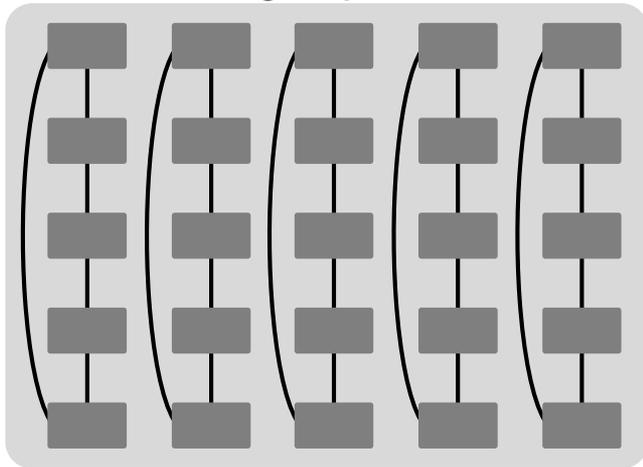


DESIGNING AN EFFICIENT NETWORK TOPOLOGY

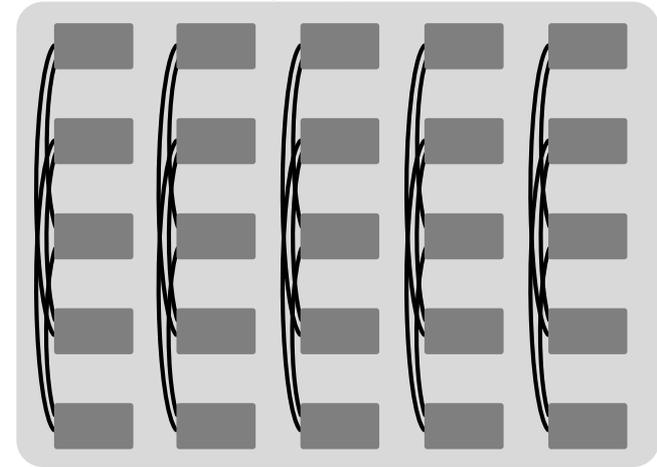
CONNECTING ROUTERS: DIAMETER 2

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

A subgraph with
identical groups of routers

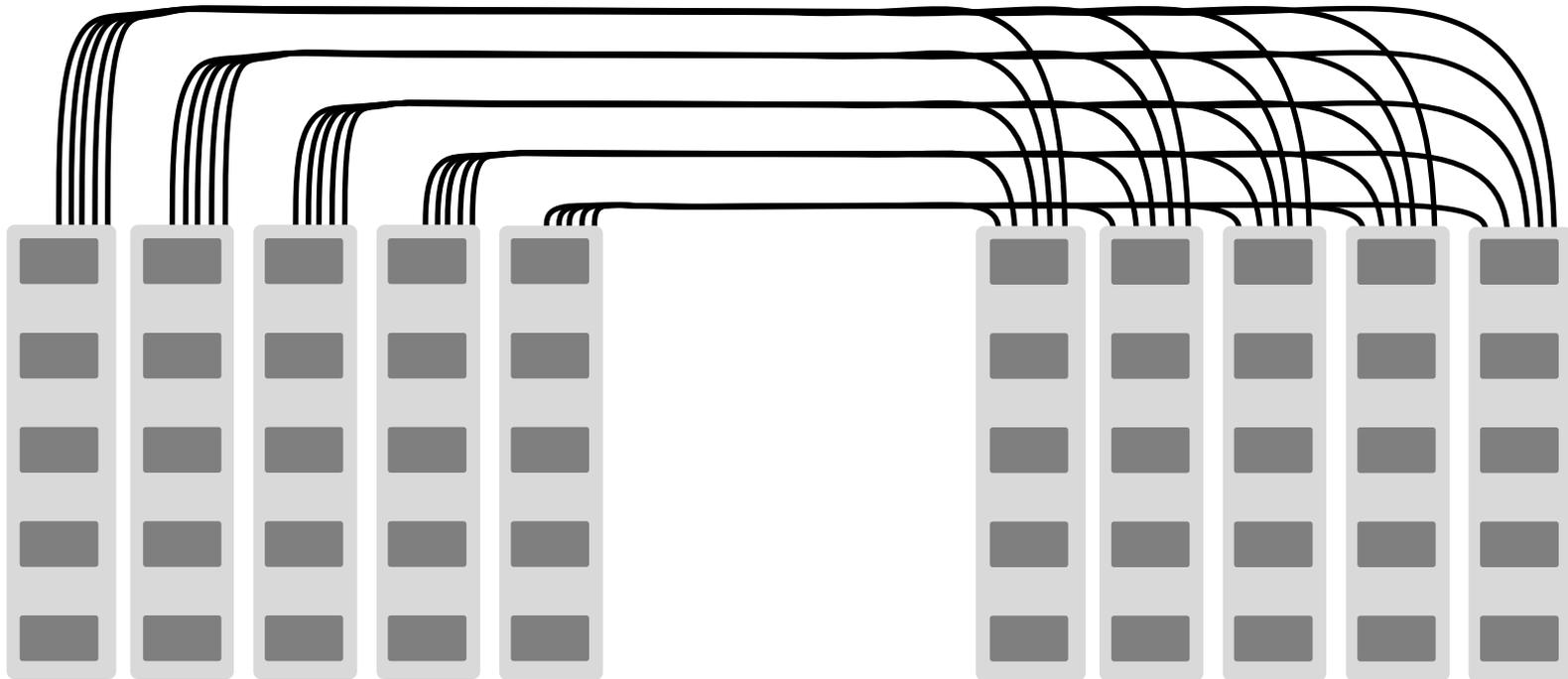


A subgraph with
identical groups of routers



DESIGNING AN EFFICIENT NETWORK TOPOLOGY

CONNECTING ROUTERS: DIAMETER 2



Groups form a fully-connected bipartite graph

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:

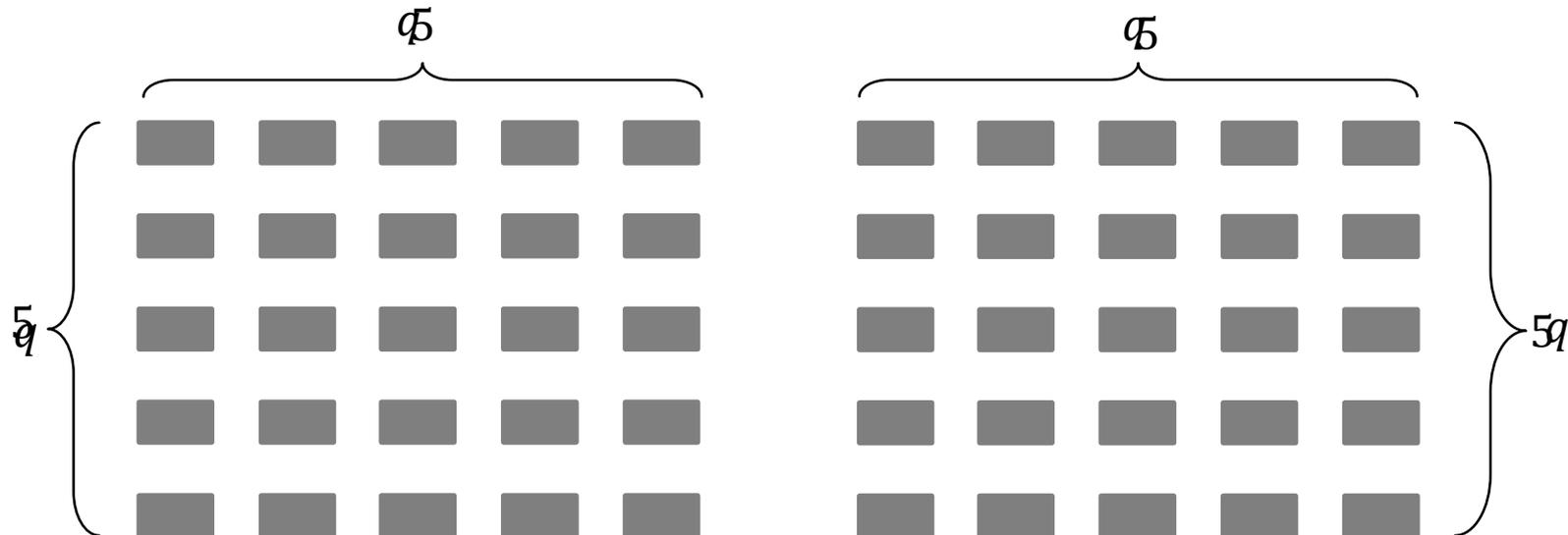
$$\mathcal{F}_q = \mathbb{Z}/q\mathbb{Z} = \{0, 1, \dots, q - 1\}$$

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

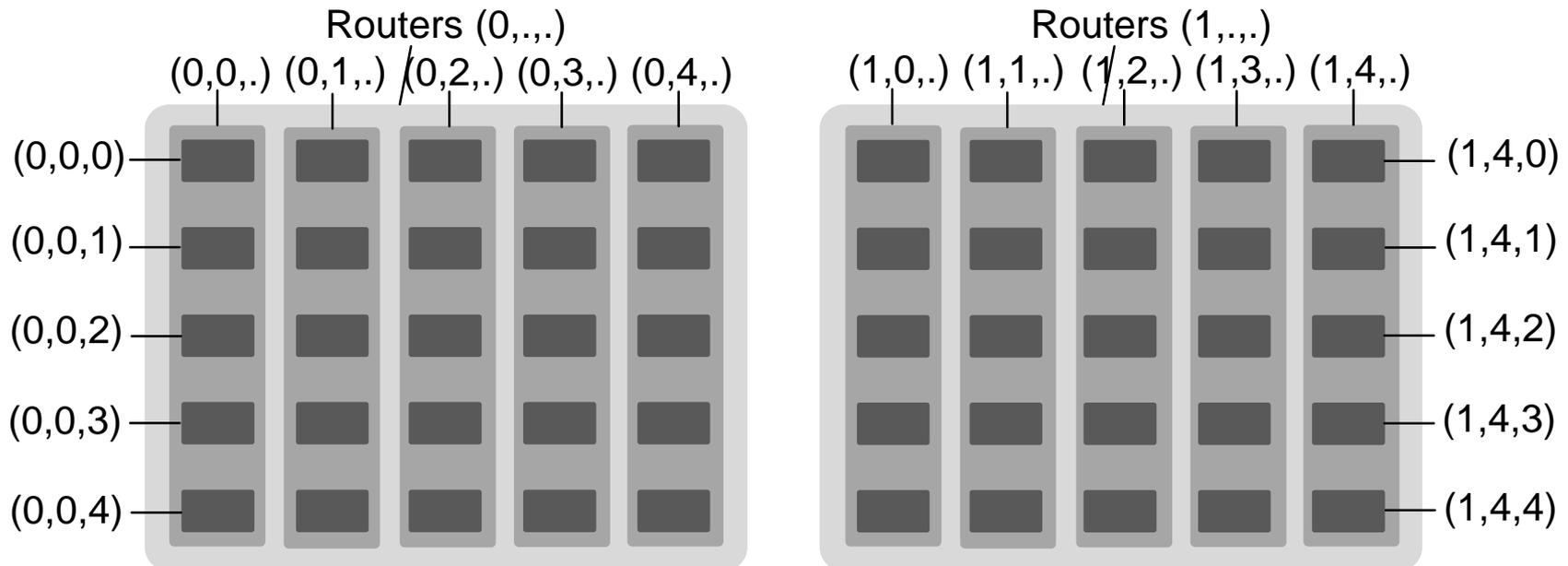
3 Label the routers

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 ξ

$\xi \in \mathcal{F}_q$ generates \mathcal{F}_q :

All non-zero elements of \mathcal{F}_q
 can be written as ξ^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: $q = 5$

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

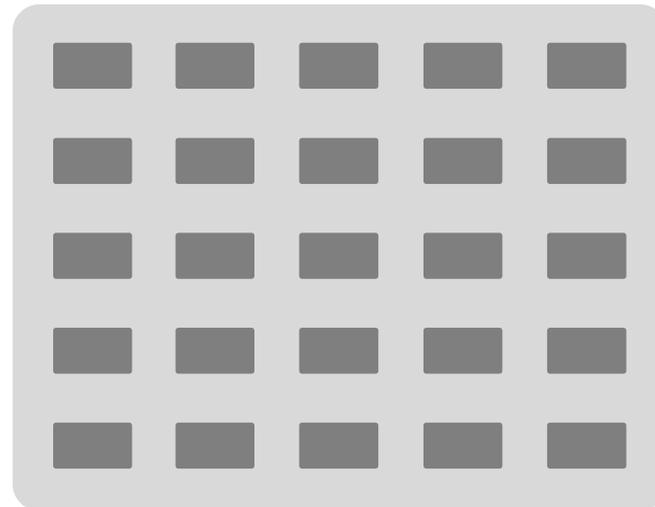
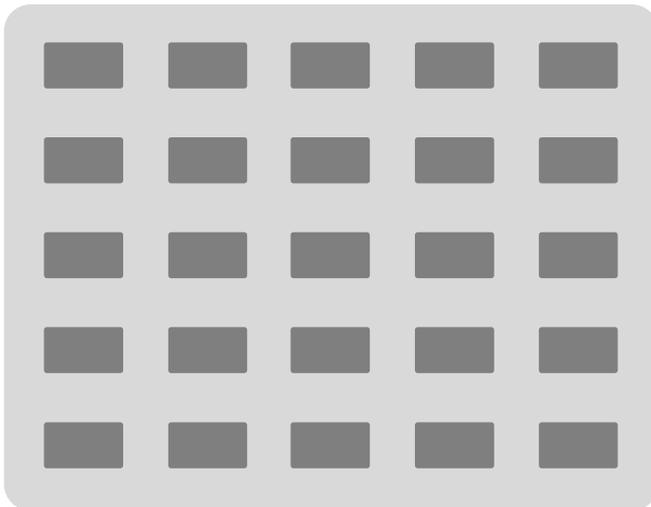
$$\xi = 2$$

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

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

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

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



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:

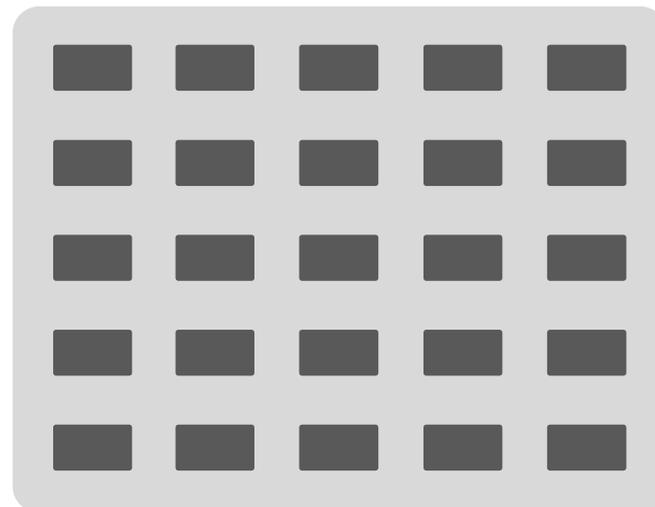
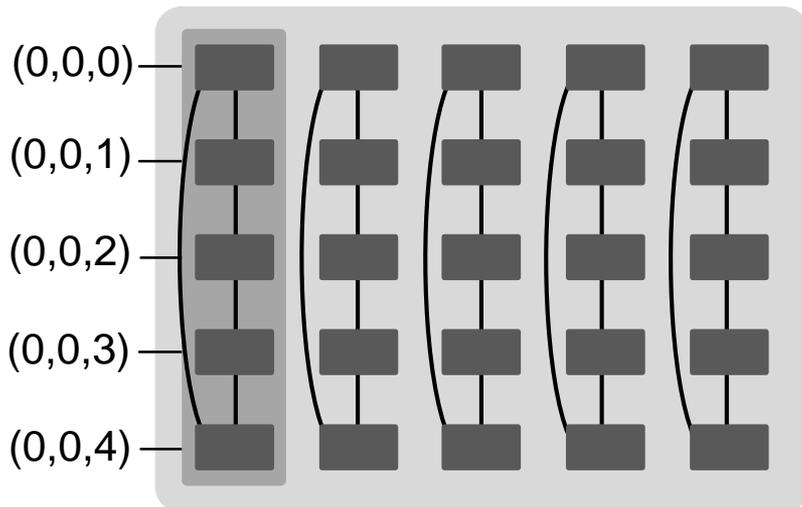
$$X = \{1, \xi^2, \dots, \xi^{q-3}\} \text{ (for subgraph 0)}$$

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

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:

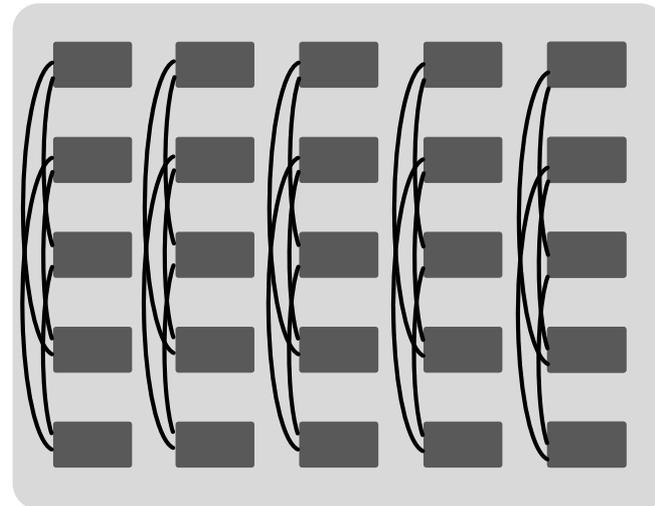
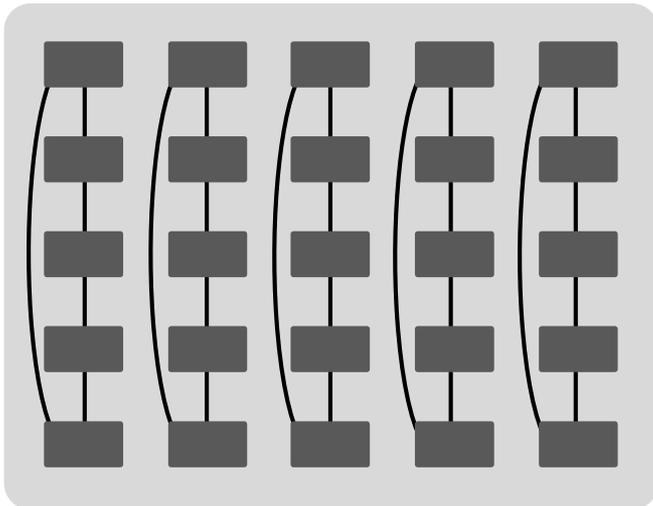
$$X = \{1, \xi^2, \dots, \xi^{q-3}\} \text{ (for subgraph 0)}$$

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

E Example: $q = 5$

Take Routers (1,4,.)

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



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 $(1, 0, 0)$

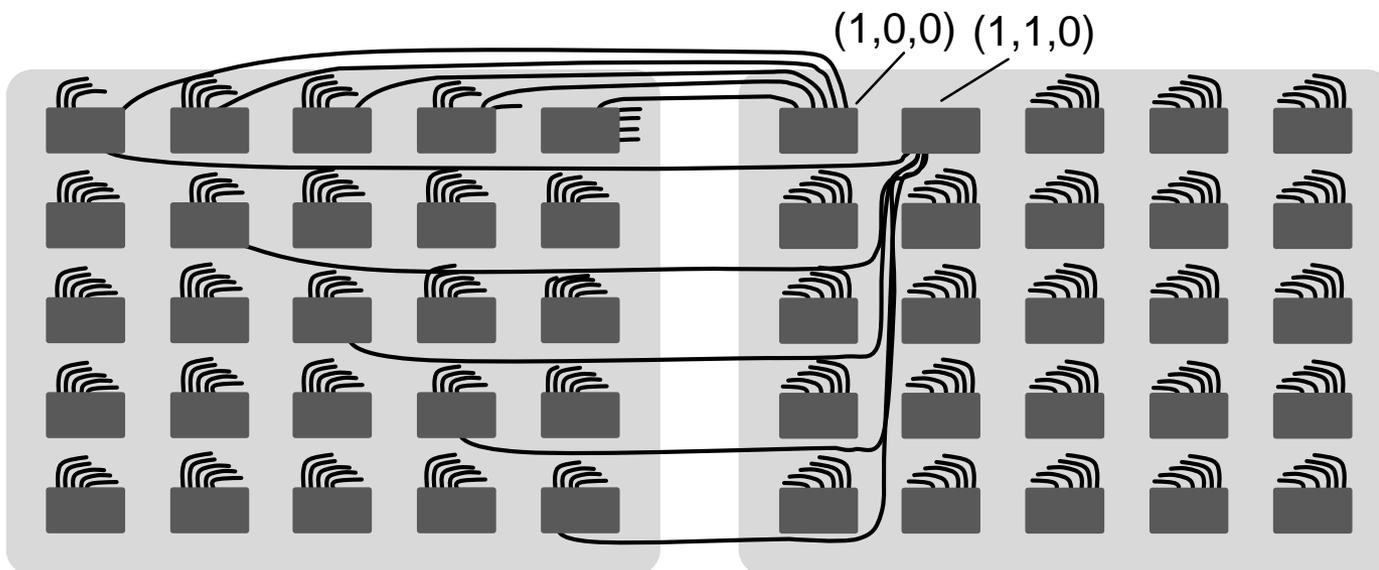
$(1, 0, 0) \leftrightarrow (0, x, 0)$

$m = 0, c = 0$

Take Router $(1, 1, 0)$

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

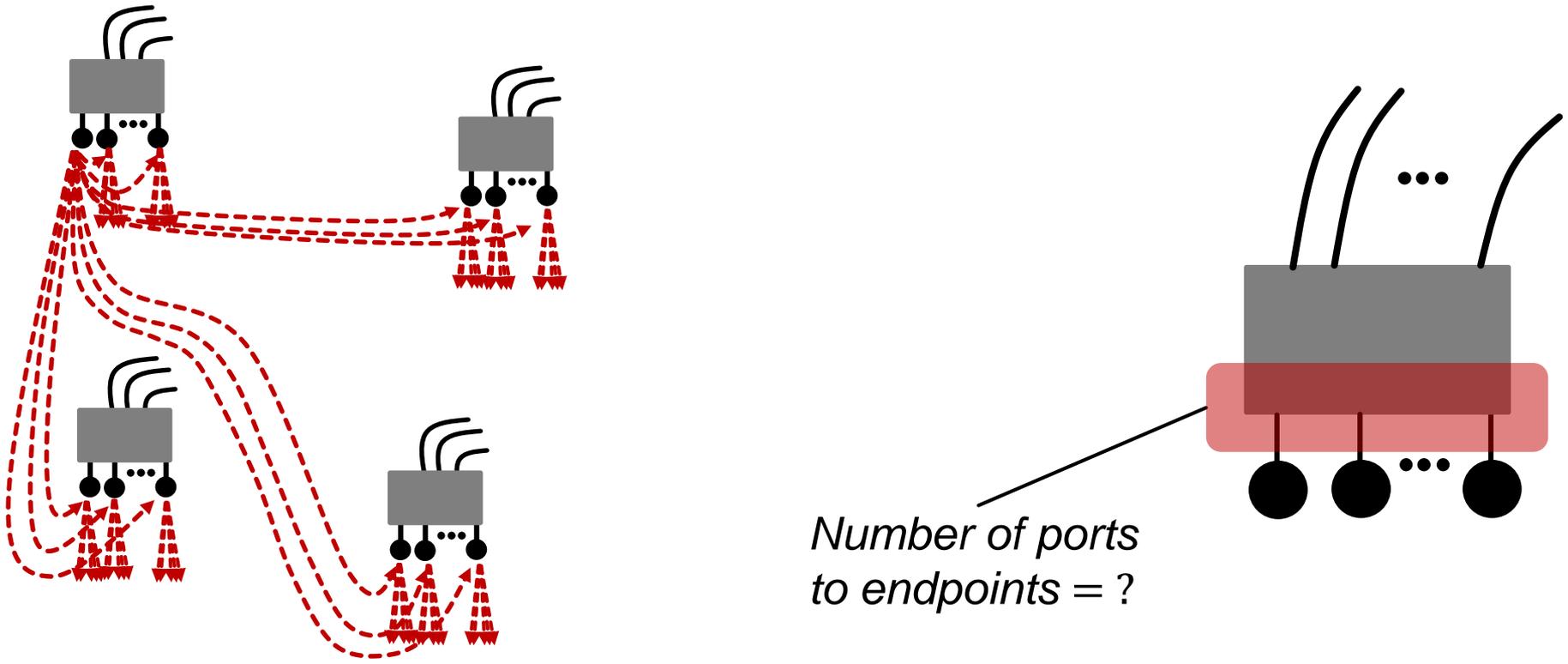
$m = 1, c = 0$



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



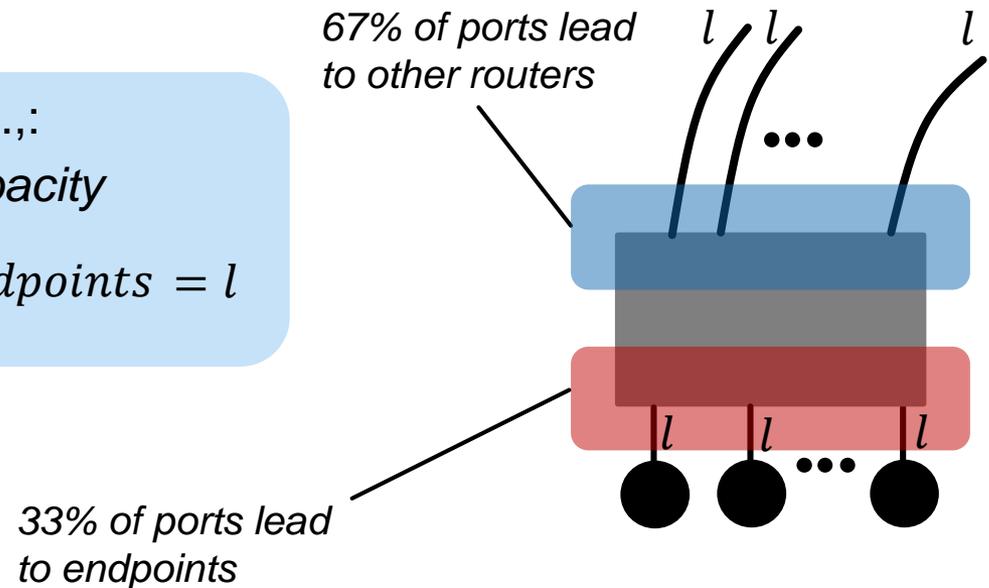
DESIGNING AN EFFICIENT NETWORK TOPOLOGY

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
local uplink load = number of endpoints = l

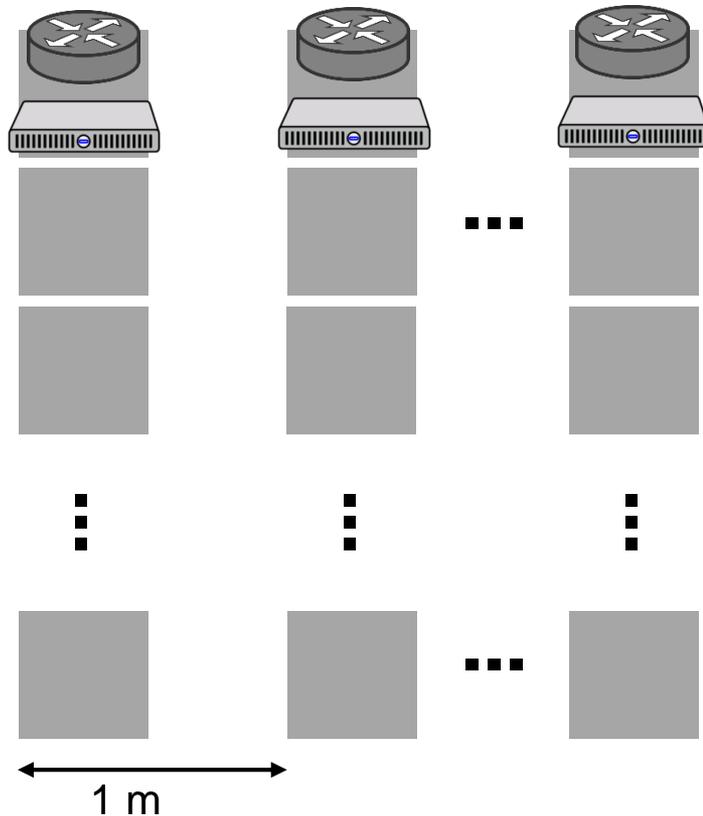


COST COMPARISON

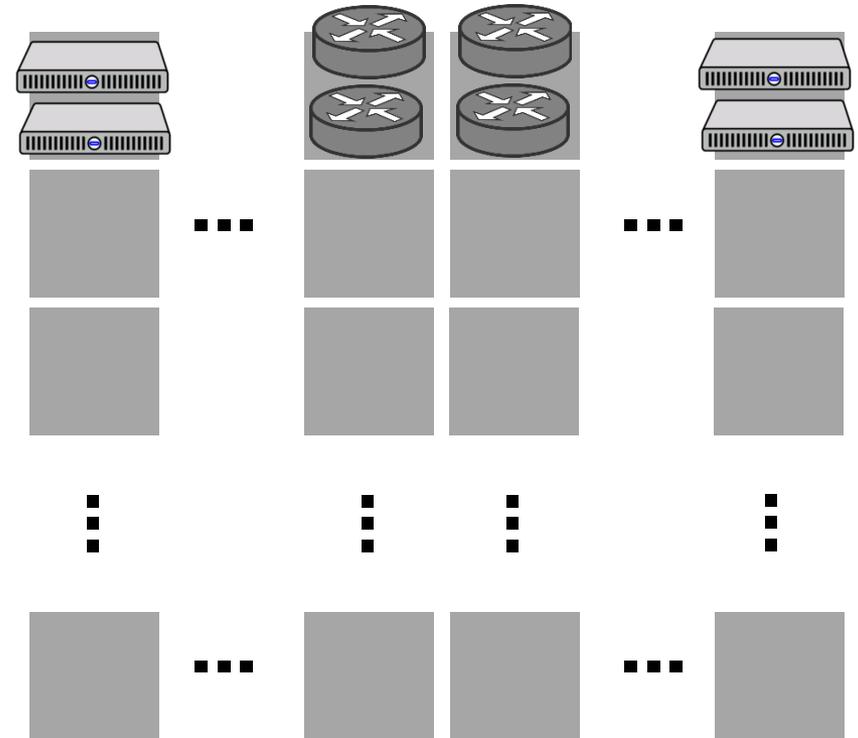
COST MODELS: VARIANTS

Variant 1: Routers and servers together

A rack



Variant 2: Routers and servers separately



COST COMPARISON

CABLE COST MODEL

*Prices based on:  COLFAX DIRECT
HPC and Data Center Gear

- 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:

Mellanox IB QDR
56Gb/s QSFP



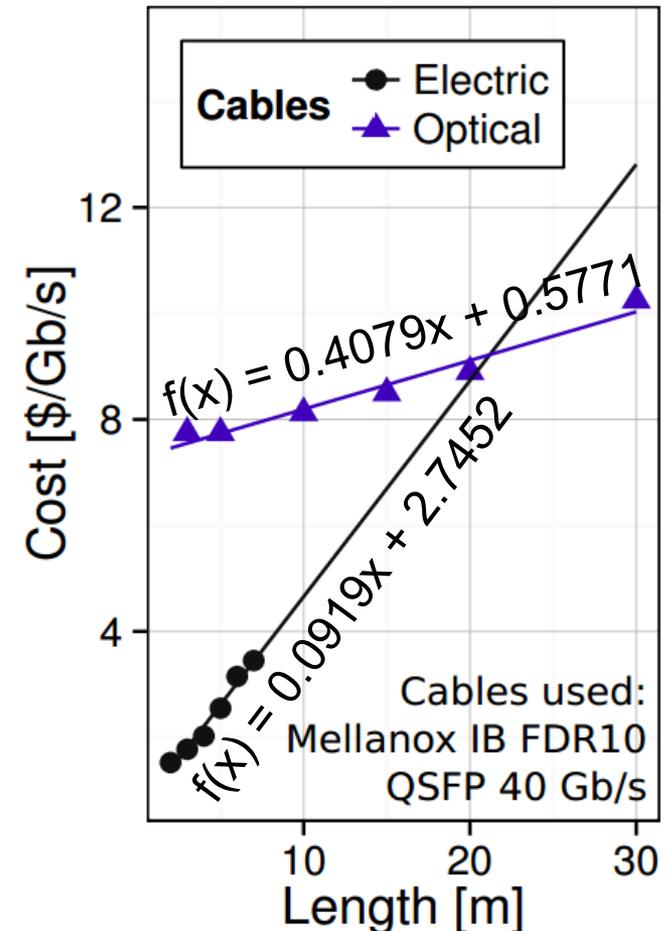
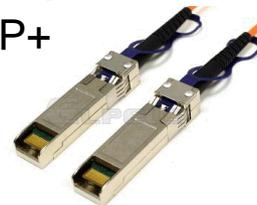
Mellanox Ethernet
40Gb/s QSFP



Mellanox Ethernet
10Gb/s SFP+



Elpeus Ethernet
10Gb/s SFP+



COST COMPARISON

ROUTER COST MODEL

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

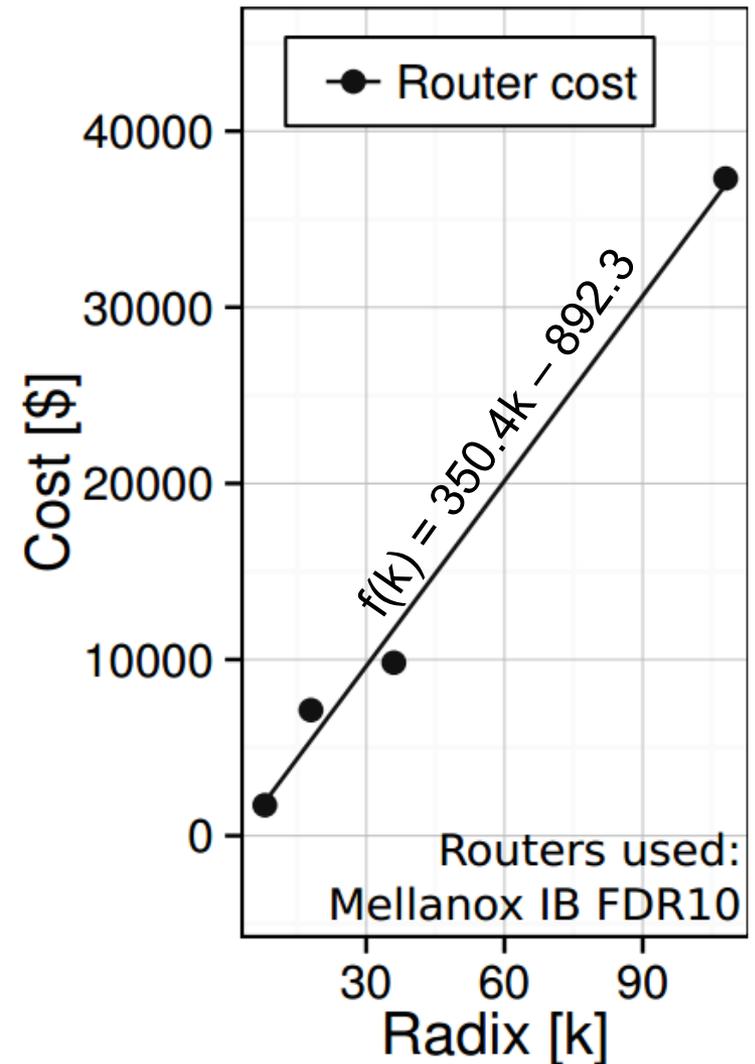
Mellanox IB FDR10



Mellanox Ethernet 10/40 Gb



*Prices based on:  COLFAX DIRECT
HPC and Data Center Gear



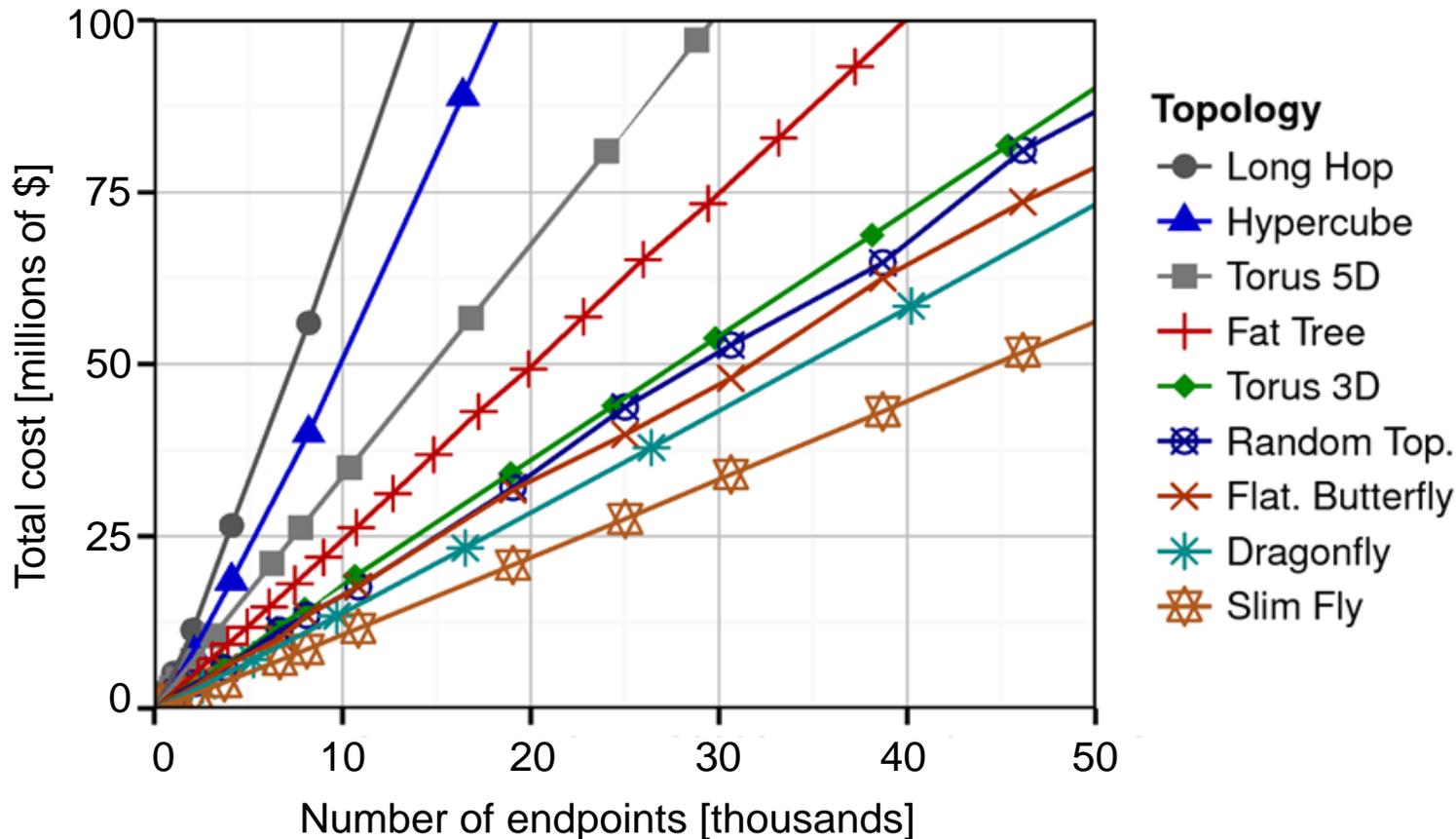
COST COMPARISON

RESULTS

Variant 2:

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

Variant 1:



COST & POWER 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

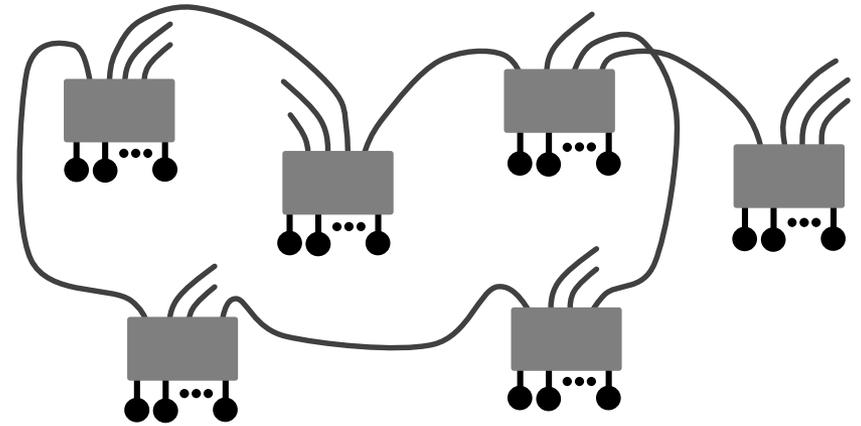
Topology	Fat tree	Random	Flat. Butterfly	Dragonfly	Slim Fly
Endpoints (N)	19,876	40,200	20,736	58,806	10,830
Routers (N_r)	2,311	4,020	1,728	5,346	722
Radix (k)	43	43	43	43	43
Electric cables	19,414	32,488	9,504	56,133	6,669
Fiber cables	40,215	33,842	20,736	29,524	6,869
Cost per node [\$]	2,346	1,743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

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

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%	60%
1024	25%	40%	40%	55%	40%	50%	60%	-	-
2048	20%	-	40%	55%	40%	55%	65%	65%	65%
4096	15%	-	45%	55%	55%	60%	70%	70%	70%
8192	10%	35%	45%	55%	60%	65%	-	75%	75%

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

PERFORMANCE & ROUTING

MINIMUM ROUTING

1 Intra-group connections

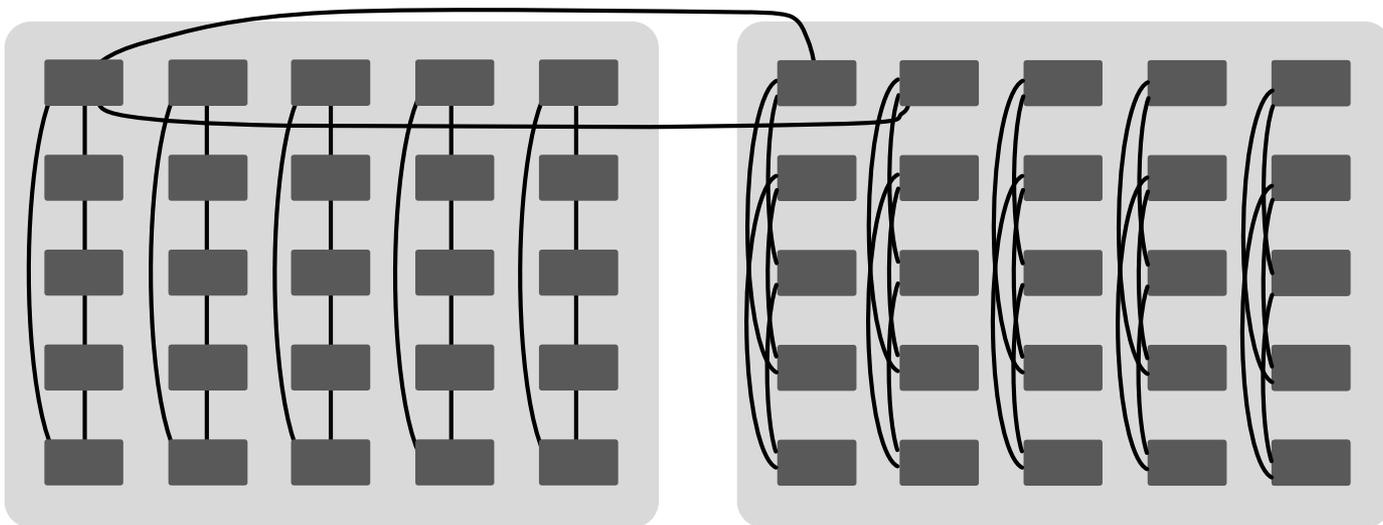
- ⊃ Path of length 1 or 2 between two routers

2 Inter-group connections (different types of groups)

- ⊃ Path of length 1 or 2 between two routers

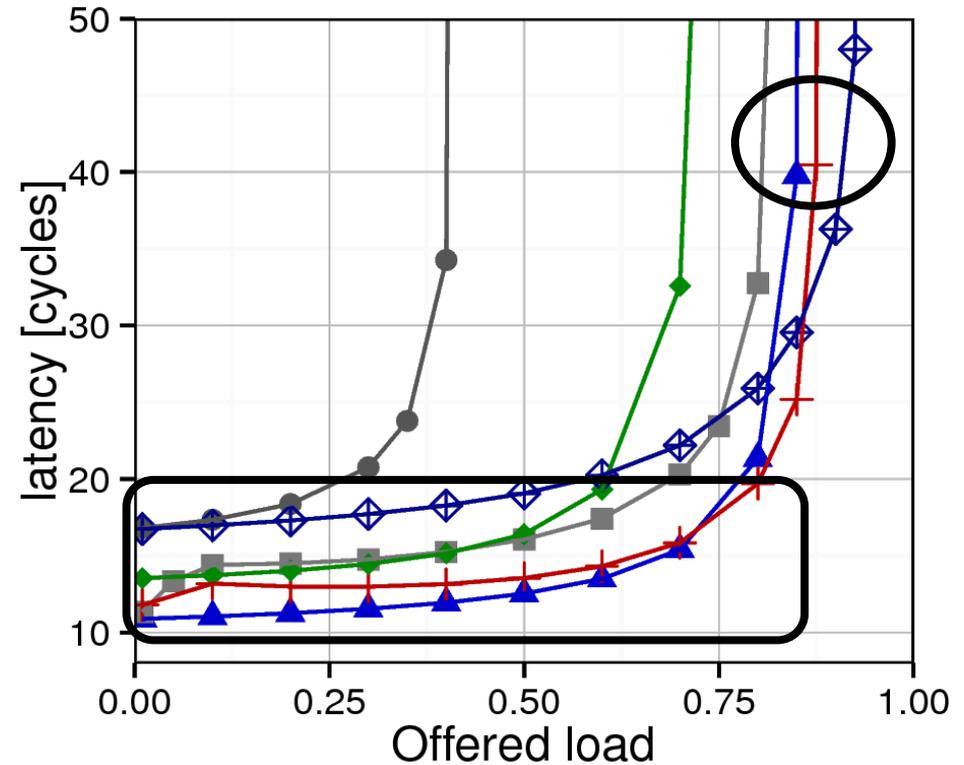
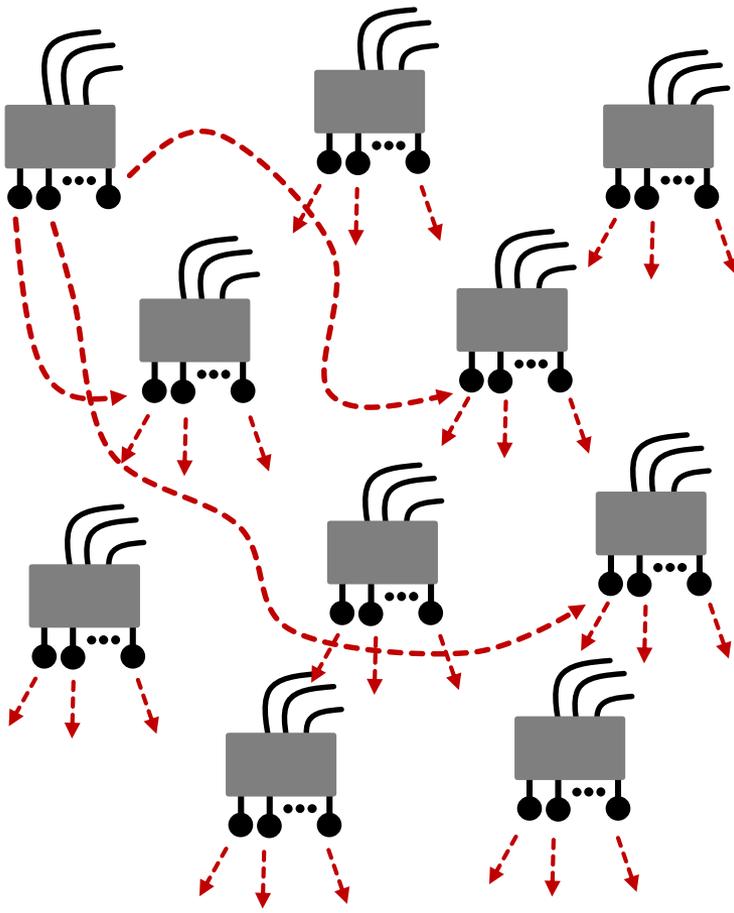
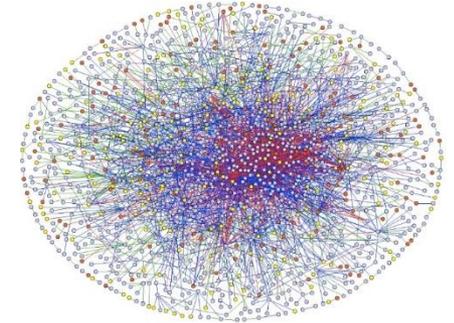
3 Inter-group connections (identical types of groups)

- ⊃ Path of length 2 between two routers



PERFORMANCE & ROUTING

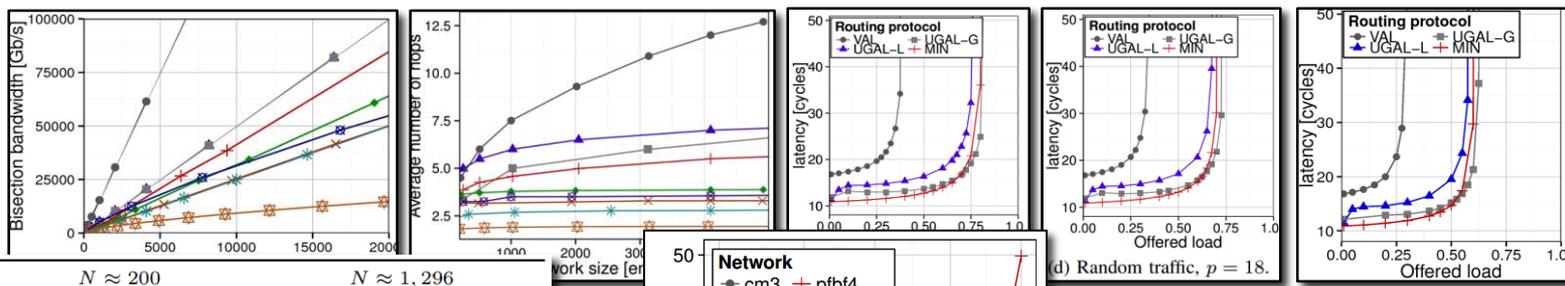
RANDOM UNIFORM TRAFFIC

Routing protocol

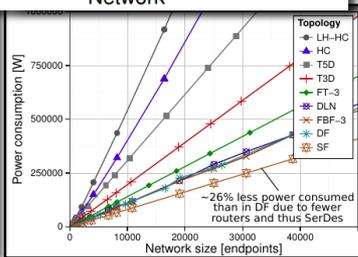
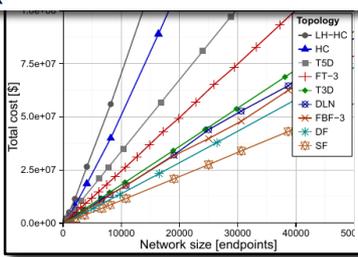
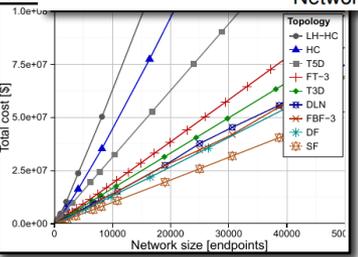
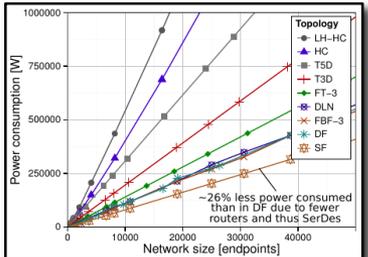
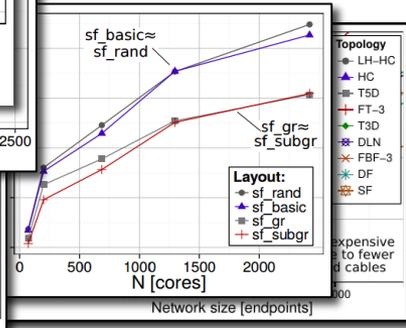
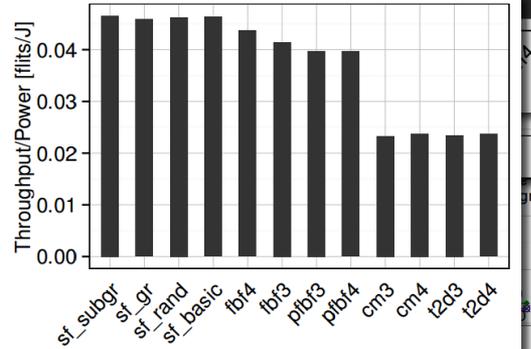
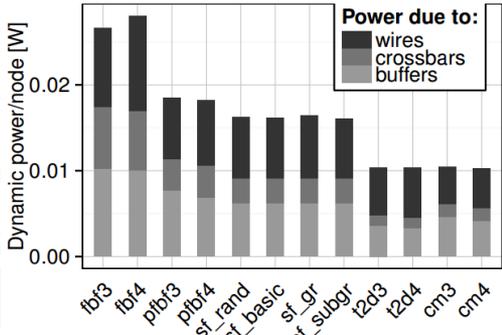
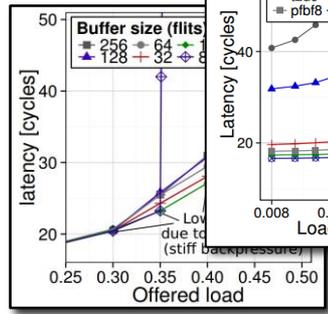
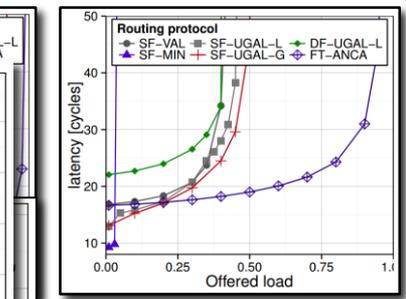
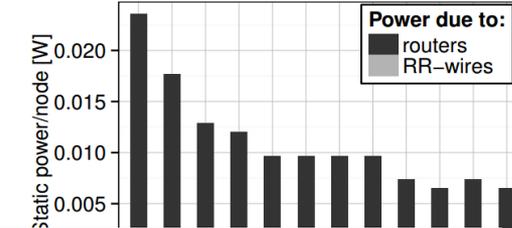
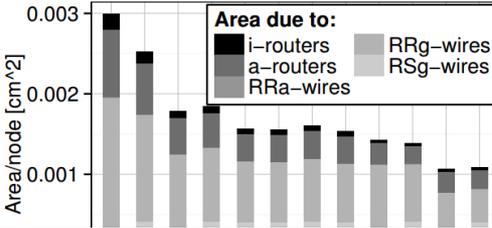
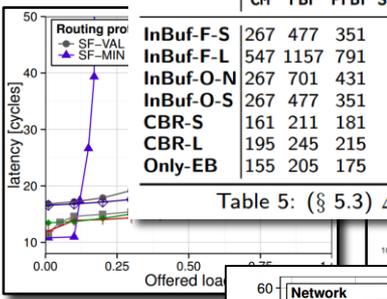
- Slim Fly (Valiant)
- ▲ Slim Fly (Minimum)
- Slim Fly (UGAL-L)
- ✦ Slim Fly (UGAL-G)
- ◆ Dragonfly (UGAL-L)
- ◇ Fat Tree (ANCA)

OTHER RESULTS



Strategy	$N \approx 200$					$N \approx 1,296$				
	CM	FBF	PFBF	SF_subgr	SF_gr	CM	FBF	PFBF	SF_subgr	SF_gr
InBuf-F-S	267	477	351	305	305	1424	2275	1802	1926	1926
InBuf-F-L	547	1157	791	649	649	2594	5065	3692	4052	4052
InBuf-O-N	267	701	431							
InBuf-O-S	267	477	351							
CBR-S	161	211	181							
CBR-L	195	245	215							
Only-EB	155	205	175							

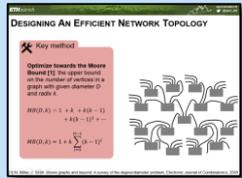
Table 5: (§ 5.3)



Topology	Dragonfly	Slim Fly
Endpoints (N)	10,890	10,830
Routers (N_r)	990	722
Radix (k)	43	43
Electric cables	6,885	6,669
Fiber cables	1,012	6,869
Cost per node [€]	1,365	1,033
Power per node [W]	10.9	8.02

A LOWEST-DIAMETER TOPOLOGY

- Approaching the Moore Bound
- Resilient



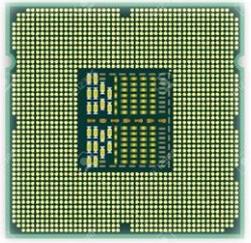
STRUCTURE ANALYSIS

Resiliency*

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

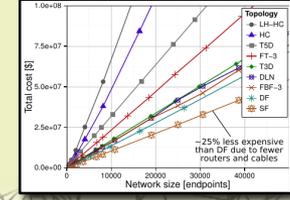
No. N	Health	Health	Health	Long Hop	Fatness	Outage	Ext. Industry	Resiliency	Min. Pk
512	30%	-	40%	55%	35%	-	55%	60%	60%
1024	23%	40%	40%	55%	40%	50%	60%	-	-
2048	20%	-	40%	55%	40%	55%	65%	65%	65%
4096	15%	-	45%	55%	55%	60%	70%	70%	70%
8192	10%	35%	45%	55%	60%	65%	75%	75%	75%

*Missing values indicate the inapplicability of a balanced topology variant for a given N.



A COST & POWER EFFECTIVE TOPOLOGY

- 25% less expensive than Dragonfly,
- 26% less power-hungry than Dragonfly



Thank you
for your attention
<http://spcl.inf.ethz.ch/SlimFly>



A HIGH-PERFORMANCE TOPOLOGY

- Lowest latency
- Full global bandwidth

