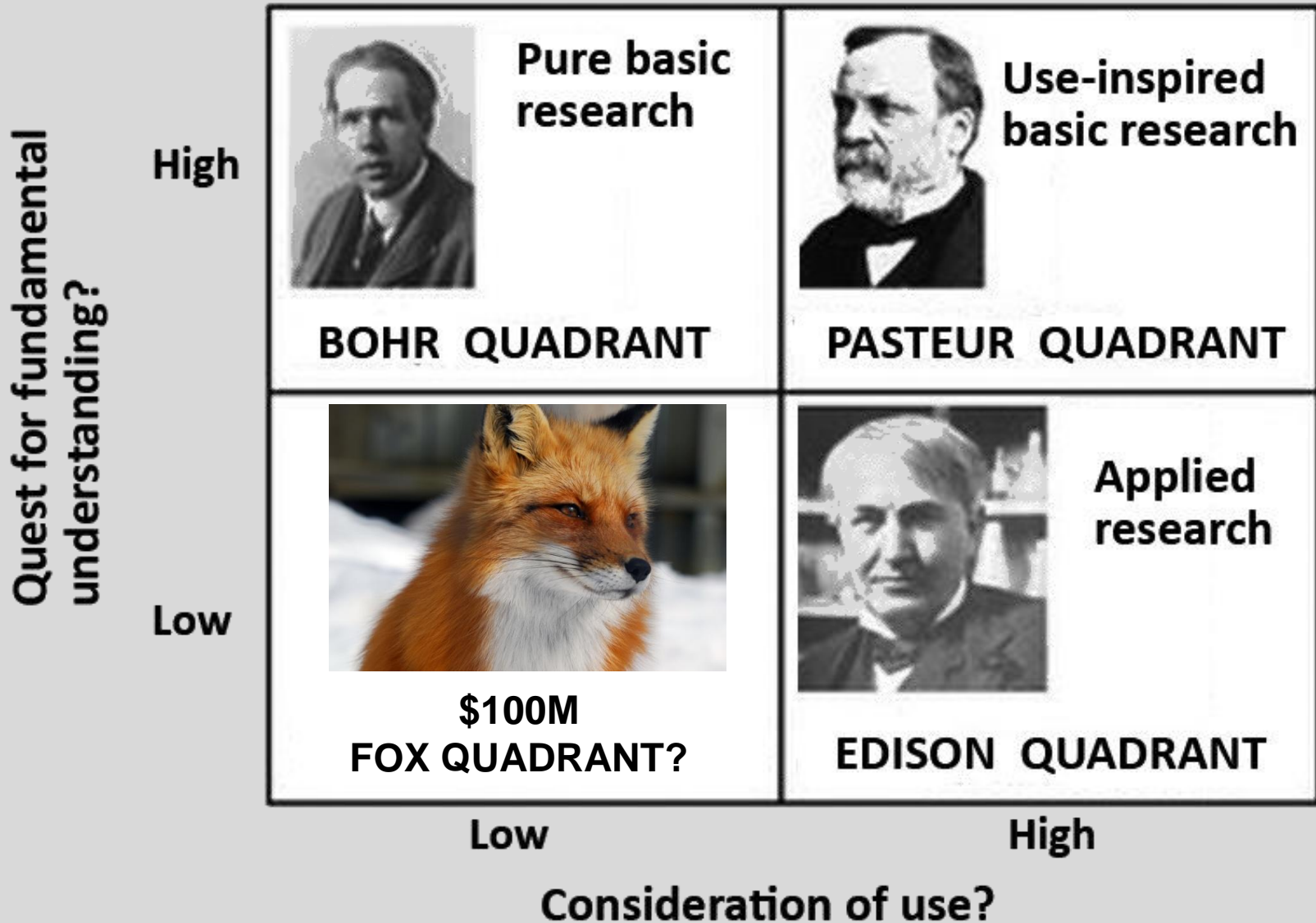


SlimFly: A Cost Effective Low-Diameter Network Topology

TORSTEN HOEFLE
WITH MACIEJ BESTA



Edison's vs. Pasteur's quadrant

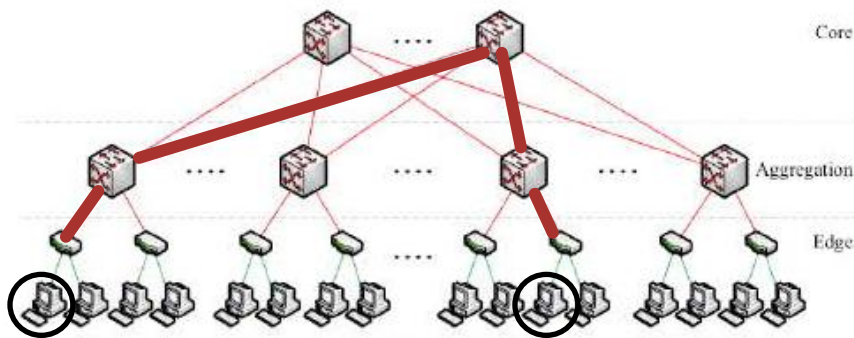


DESIGNING AN EFFICIENT NETWORK TOPOLOGY

- Main intuition/idea: decrease network diameter
 - lower latency
 - smaller cost (fewer routers and cables for same bandwidth)
 - lower power consumption (packets traverse fewer SerDes)

Fat tree:

Diameter == 4

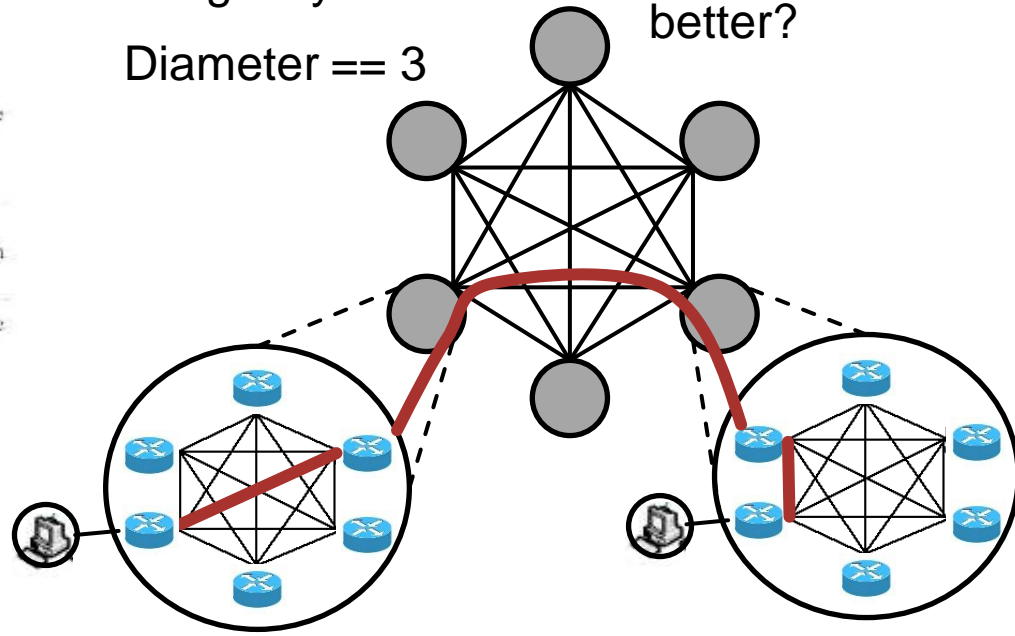


...still high ☹

Dragonfly:

Diameter == 3

...can we do better?



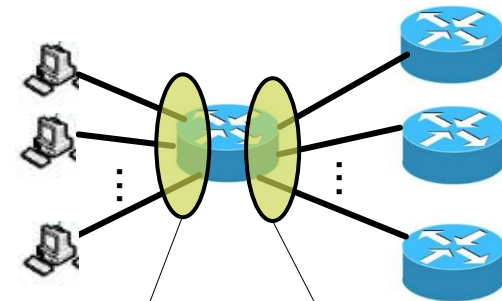
DESIGNING AN EFFICIENT NETWORK TOPOLOGY

- Goal: find a close-to-optimal topology that maximizes the number of endpoints (N) for a given diameter (D) and degree (k)
- Moore Bound: upper bound on the number of routers (N_r) in a graph with given D and k' .

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

$D = 2$: $N_r \approx k'^2$
(~200,000 endpoints
with 108-port switches)

$D = 3$: $N_r \approx k'^3$
(>10,000,000 endpoints
with 108-port switches)



Number of ports
to endpoints (p)

Number of ports
to routers (k')

DESIGNING AN EFFICIENT NETWORK TOPOLOGY

- Degree-Diameter problem

Graph with the maximum N_r for a given D and k

Diameter (D)

→

Degree (k)

↓

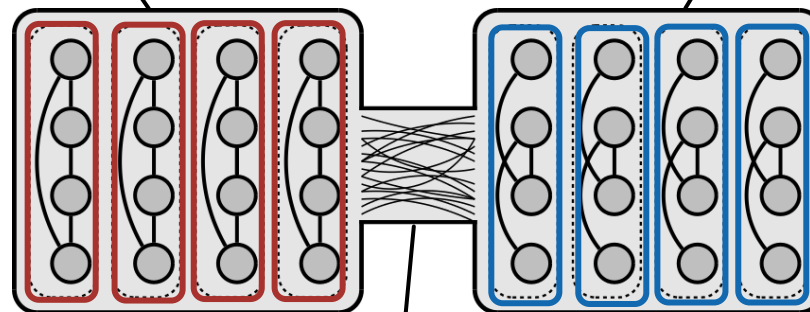
	2	3	4	5	6	7	8	9	10
3	<div>10</div> <div>100%</div>	<div>20</div> <div>100%</div>	<div>38</div> <div>100%</div>	<div>92</div> <div>76.08%</div>	<div>188</div> <div>70.21%</div>	<div>380</div> <div>51.57%</div>	<div>764</div> <div>43.97%</div>	<div>1532</div> <div>39.16%</div>	<div>3068</div> <div>40.74%</div>
4	<div>15</div> <div>100%</div>	<div>52</div> <div>78.84%</div>	<div>160</div> <div>61.25%</div>	<div>484</div> <div>75.20%</div>	<div>1456</div> <div>50.82%</div>	<div>4372</div> <div>30.19%</div>	<div>13120</div> <div>24.71%</div>	<div>39364</div> <div>19.24%</div>	<div>118096</div> <div>14.99%</div>
5	<div>24</div> <div>100%</div>	<div>104</div> <div>69.23%</div>	<div>424</div> <div>50%</div>	<div>1704</div> <div>36.61%</div>	<div>6824</div> <div>40.62%</div>	<div>27304</div> <div>20.20%</div>	<div>109224</div> <div>15.59%</div>	<div>436904</div> <div>13.23%</div>	<div>1747624</div> <div>10.70%</div>
6	<div>32</div> <div>100%</div>	<div>186</div> <div>59.67%</div>	<div>936</div> <div>41.66%</div>	<div>4686</div> <div>29.96%</div>	<div>23436</div> <div>33.78%</div>	<div>117186</div> <div>16.54%</div>	<div>585936</div> <div>13.04%</div>	<div>2929686</div> <div>10.50%</div>	<div>14648436</div> <div>8.55%</div>
7	<div>50</div> <div>100%</div>	<div>301</div> <div>55.81%</div>	<div>1813</div> <div>37.06%</div>	<div>10885</div> <div>25.31%</div>	<div>65317</div> <div>18.35%</div>	<div>391909</div> <div>13.46%</div>	<div>2351461</div> <div>10.61%</div>	<div>14108773</div> <div>8.66%</div>	<div>84652645</div> <div>7.09%</div>
8	<div>63</div> <div>90.47%</div>	<div>456</div> <div>55.48%</div>	<div>3200</div> <div>34.37%</div>	<div>22408</div> <div>22.58%</div>	<div>156864</div> <div>25.29%</div>	<div>1098056</div> <div>11.94%</div>	<div>7686400</div> <div>9.56%</div>	<div>53804808</div> <div>7.88%</div>	<div>376633664</div> <div>6.61%</div>
9	<div>80</div> <div>92.50%</div>	<div>657</div> <div>89.04%</div>	<div>5265</div> <div>29.43%</div>	<div>42129</div> <div>19.46%</div>	<div>337041</div> <div>22.51%</div>	<div>2696337</div> <div>10.37%</div>	<div>21570705</div> <div>7.81%</div>	<div>172565649</div> <div>7.02%</div>	<div>1380525201</div> <div>4.77%</div>
10	<div>99</div> <div>91.91%</div>	<div>910</div> <div>71.42%</div>	<div>8200</div> <div>27.87%</div>	<div>73810</div> <div>17.80%</div>	<div>664300</div> <div>20.27%</div>	<div>5978710</div> <div>9.75%</div>	<div>53808400</div> <div>7.97%</div>	<div>484275610</div> <div>5.78%</div>	<div>4358480500</div> <div>4.61%</div>

DESIGNING AN EFFICIENT NETWORK TOPOLOGY

- Degree-Diameter problem
 - Construct a graph with the maximum N_r for a given D and k'
- We use a result from McKay, Miller, Siran (MMS graphs) [1]; $D = 2$

A subgraph with identical groups of routers

A subgraph with identical groups of routers

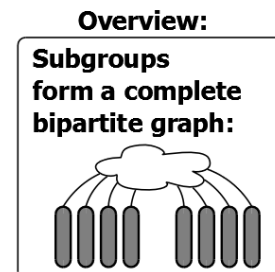
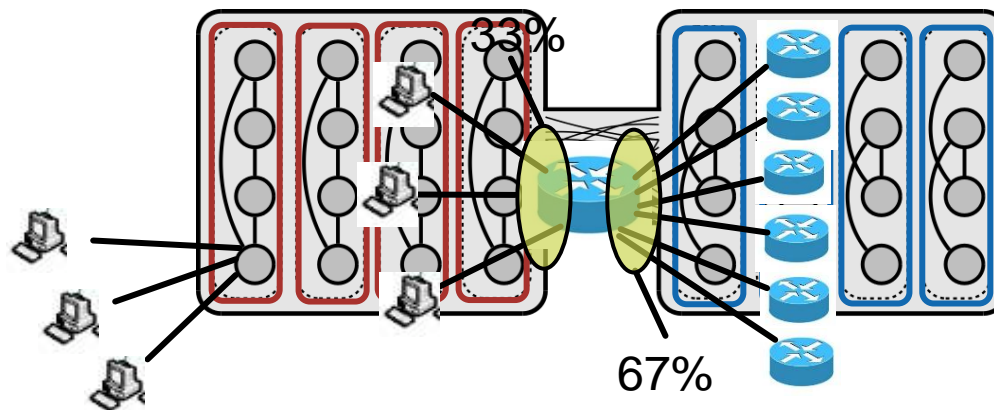


Connections between subgraphs
(details skipped for clarity)

ATTACHING ENDPOINTS

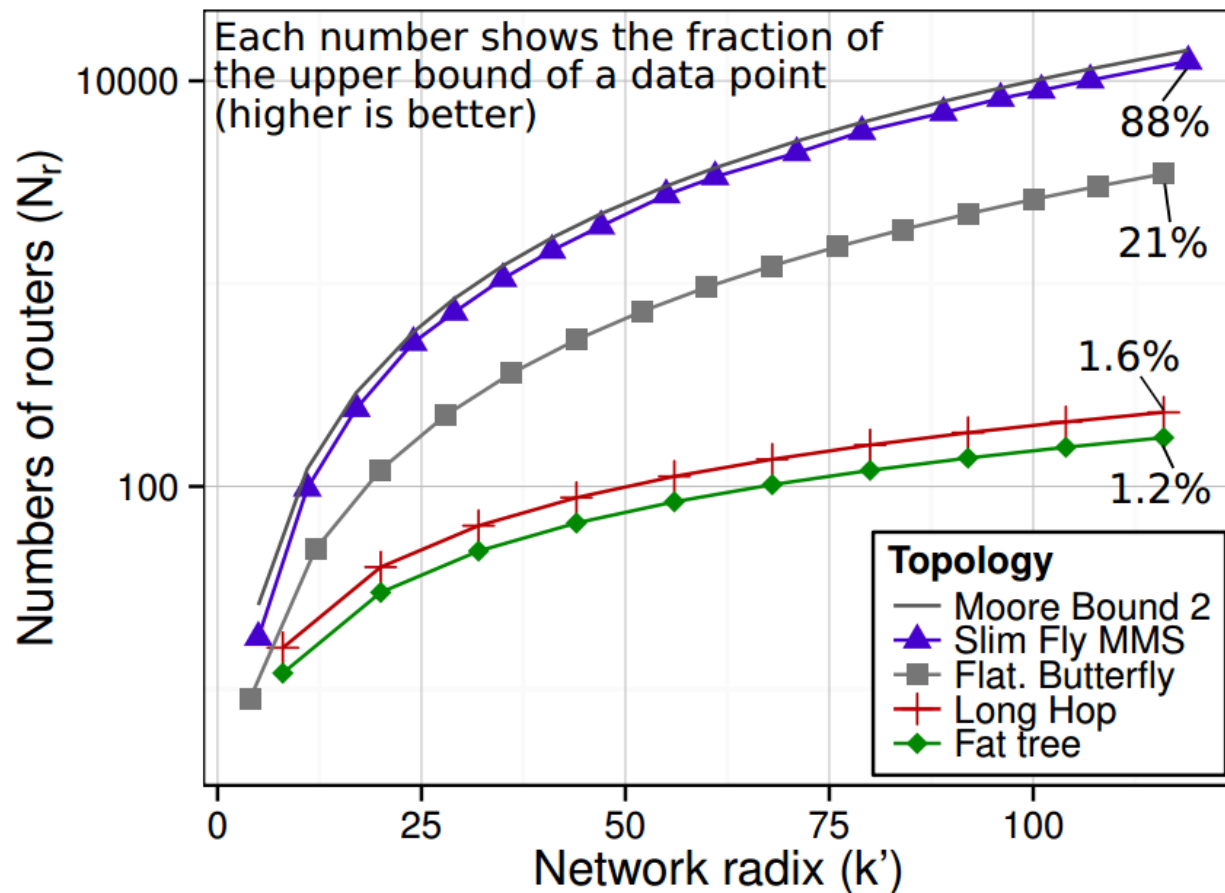
- How many endpoints do we attach to each router?
- Maximize for p while maintaining full global bandwidth
 - Global bandwidth: the theoretical cumulative throughput if all processes simultaneously communicate with all other processes in a steady state
- Result:

$$p = \left\lceil \frac{k'}{2} \right\rceil$$



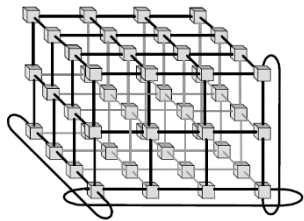
COMPARISON TO OPTIMALITY

- How close is SlimFly MMS to the Moore Bound ($D=2$)?

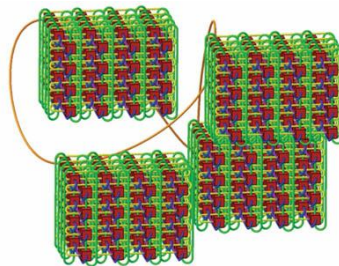


STRUCTURE ANALYSIS

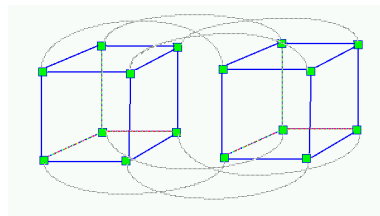
COMPARISON TARGETS



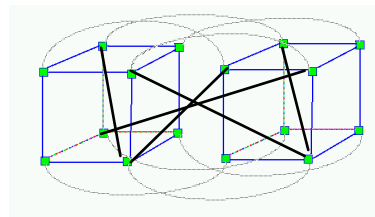
Torus 3D



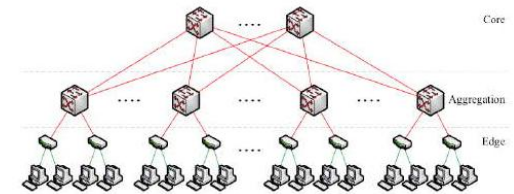
Torus 5D



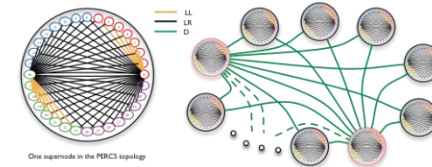
Hypercube



Long Hop [1]

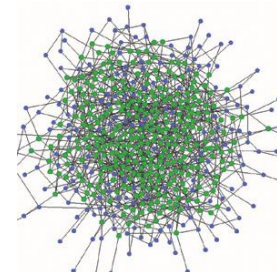
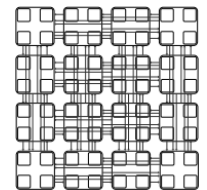


Fat tree



Dragonfly

Flattened
Butterfly



Random
networks

STRUCTURE ANALYSIS

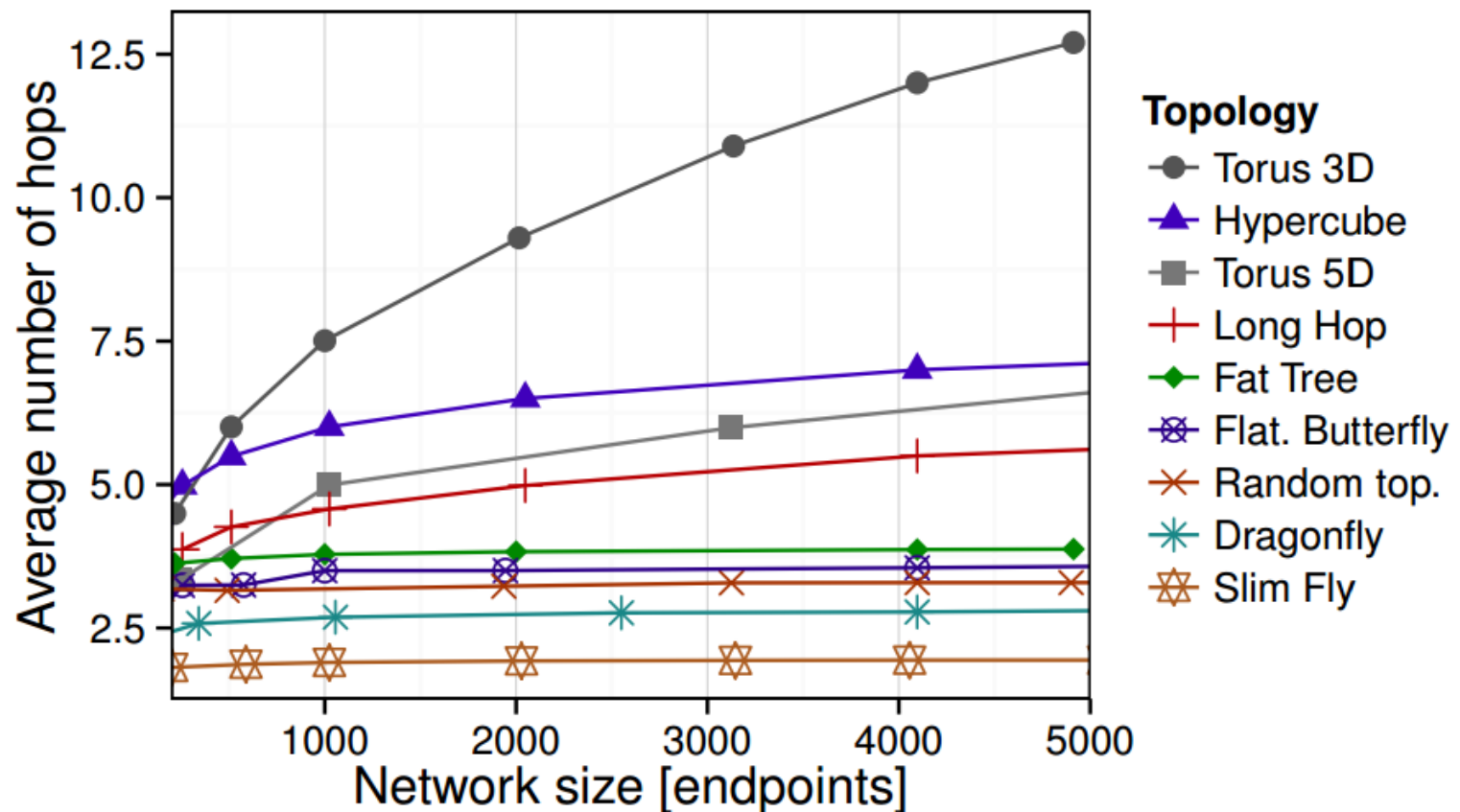
DIAMETER

Topology	Symbol	Example System	Diameter
3-dimensional torus	T3D	Cray Gemini	$\lceil 3/2 \sqrt[3]{N_r} \rceil$
5-dimensional torus	T5D	IBM BlueGene/Q	$\lceil 5/2 \sqrt[5]{N_r} \rceil$
Hypercube	HC	NASA Pleiades	$\lceil \log_2 N_r \rceil$
3-level fat tree	FT-3	Tianhe-2	4
3-level Flat. Butterfly	FBF-3	-	3
Dragonfly topologies	DF	Cray Cascade	3
Random topologies	DLN	-	3–10
Long Hop topologies	LH-HC	Infinetics Systems	4–6
Slim Fly MMS	SF	-	2

STRUCTURE ANALYSIS

AVERAGE DISTANCE

Random uniform traffic
using minimum path routing



STRUCTURE ANALYSIS

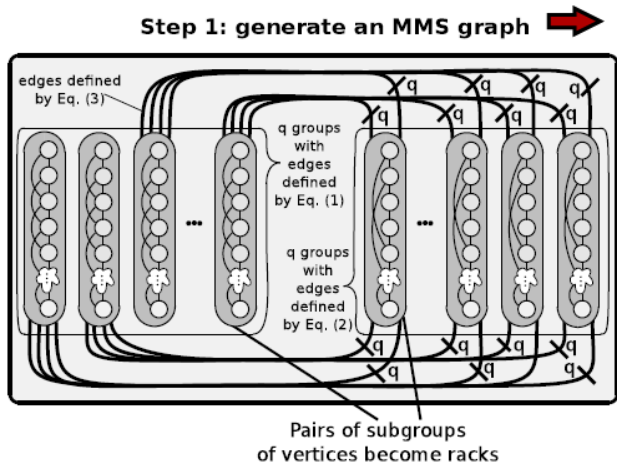
RESILIENCY

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

- Disconnection metrics*
- Other studied metrics (N≈8192):
 - Diameter (increase by 2) [1]; SF: 40%, DF: 25%, DLN: 60%
 - Average path length (increase by 2); SF: 55%, DF: 45%, DLN: 60%

$\approx N$	T3D	T5D	HC	LH-HC	FT-3	DF	FBF-3	DLN	SF
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%

PHYSICAL LAYOUT



COMPARISON TO DRAGONFLY

SlimFly:

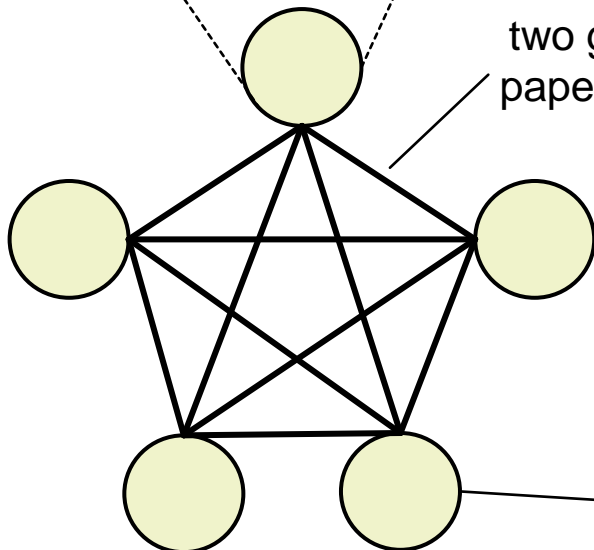
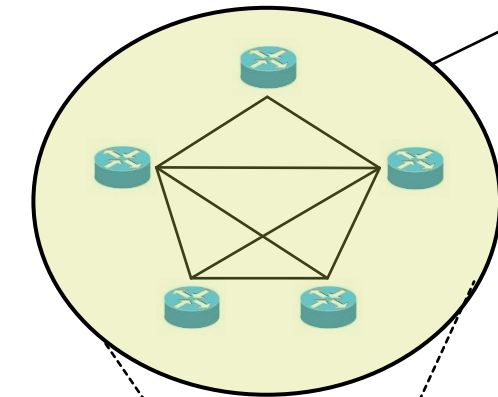
Groups not necessarily fully connected

~50% fewer intra-group cables

$2q$ inter-group cable between two groups (see paper for details)

~25% fewer routers

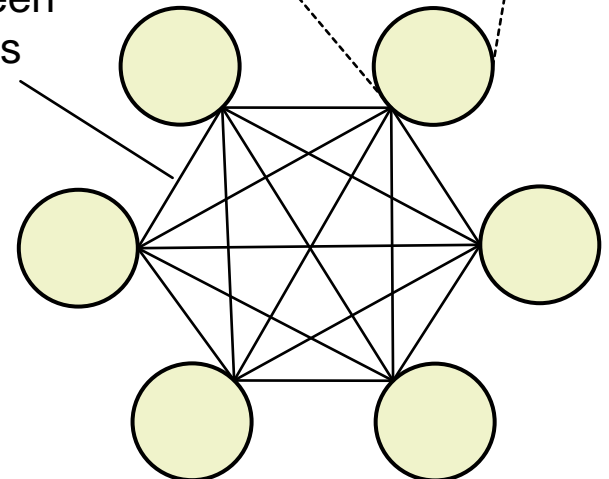
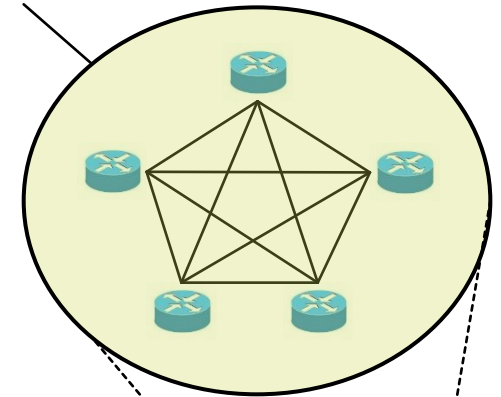
~33% higher endpoint density



Groups fully connected

Dragonfly:

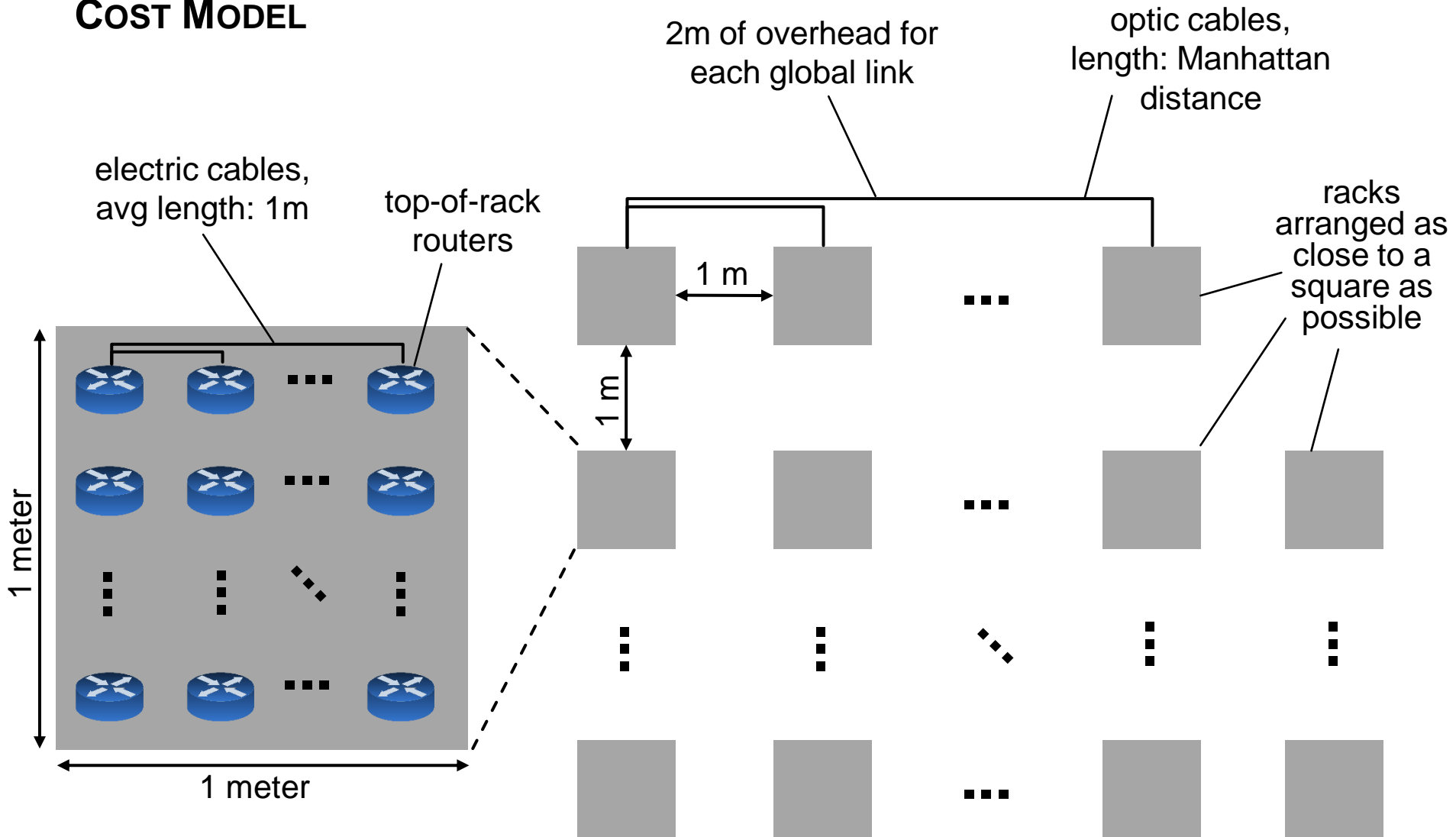
One inter-group cable between two groups



COST COMPARISON

COST MODEL

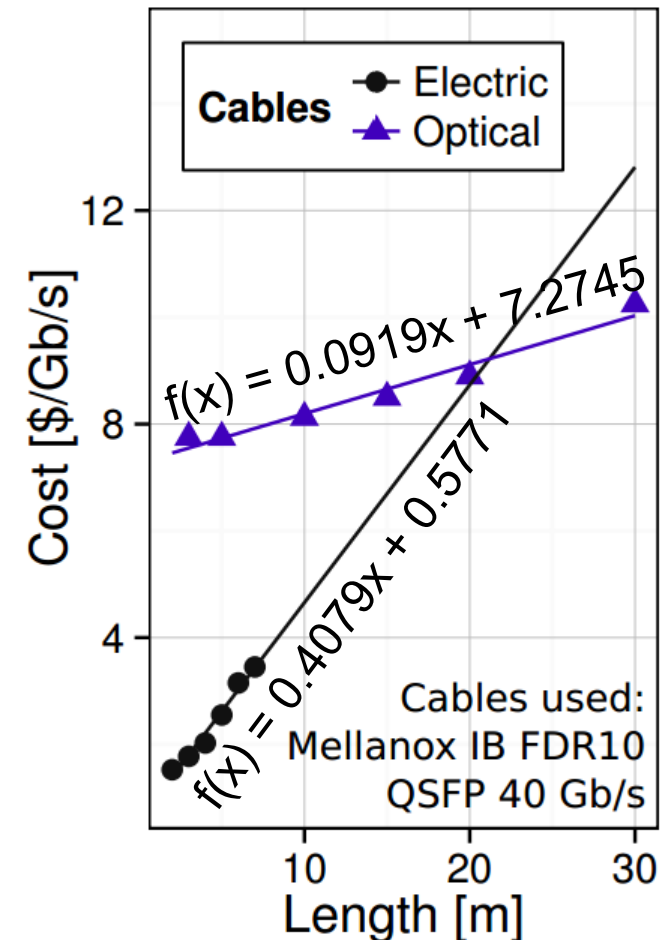
*Most cables skipped for clarity



COST COMPARISON

CABLE COST MODEL

- Bandwidth cost as a function of distance
 - The functions obtained using linear regression*
- Used cables:
 - Mellanox IB QDR 56Gb/s QSFP
 - Mellanox Ethernet 40Gb/s QSFP
 - Mellanox Ethernet 10Gb/s SFP+
 - Elpeus Ethernet 10Gb/s SFP+

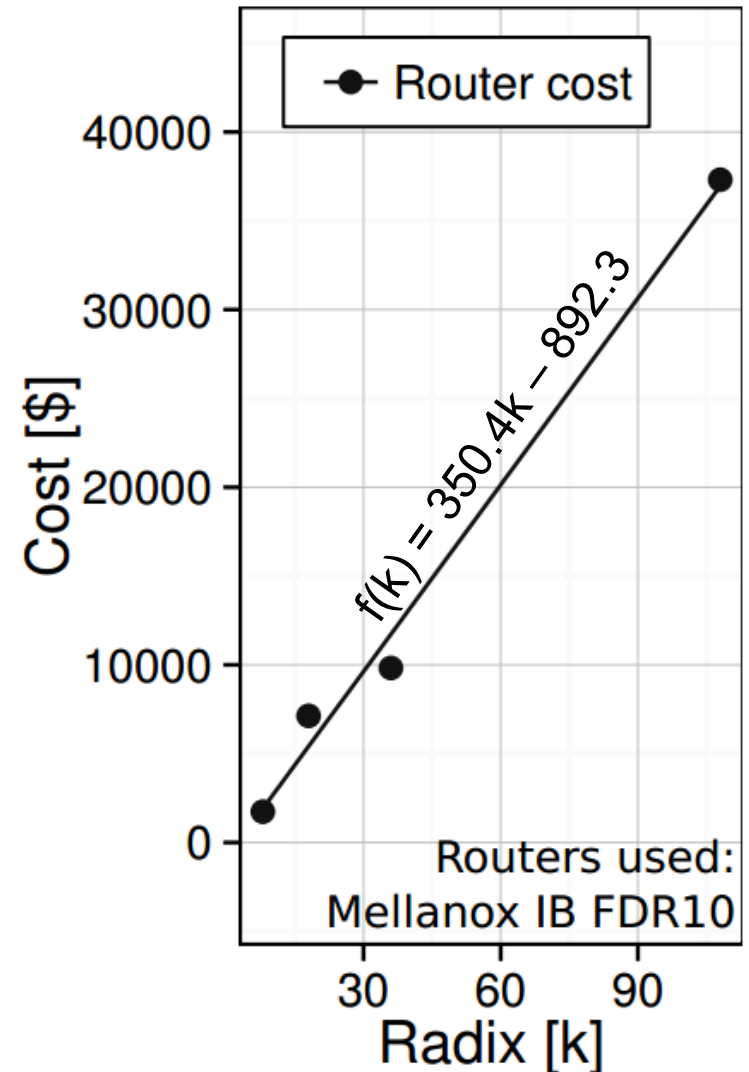


*Prices based on ColfaxDirect, June 2014

COST COMPARISON

ROUTER COST MODEL

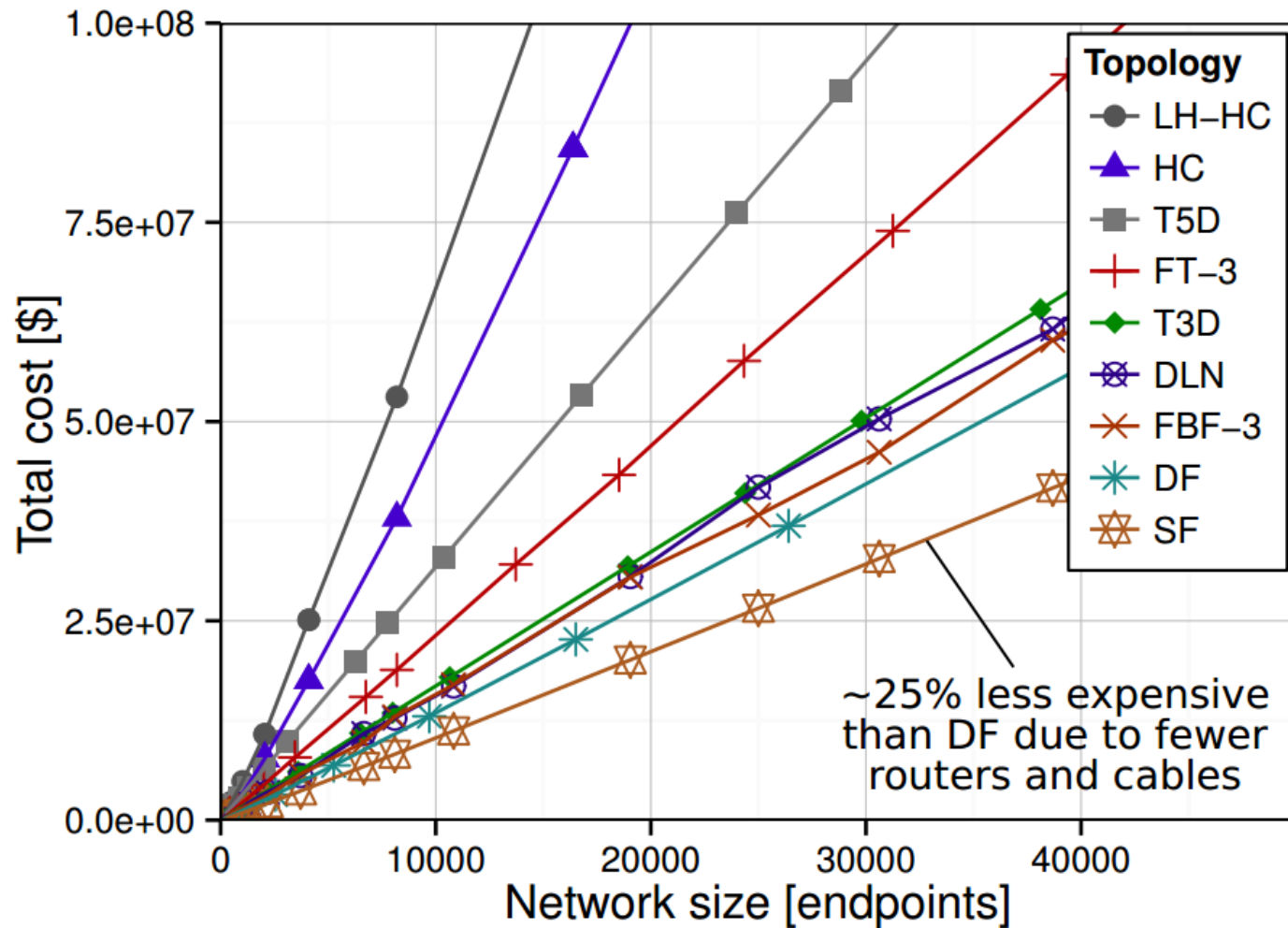
- Router cost as a function of radix
 - The function obtained using linear regression*
- Used routers:
 - Mellanox Ethernet 10/40Gb



*Prices based on ColfaxDirect, June 2014

COST COMPARISON

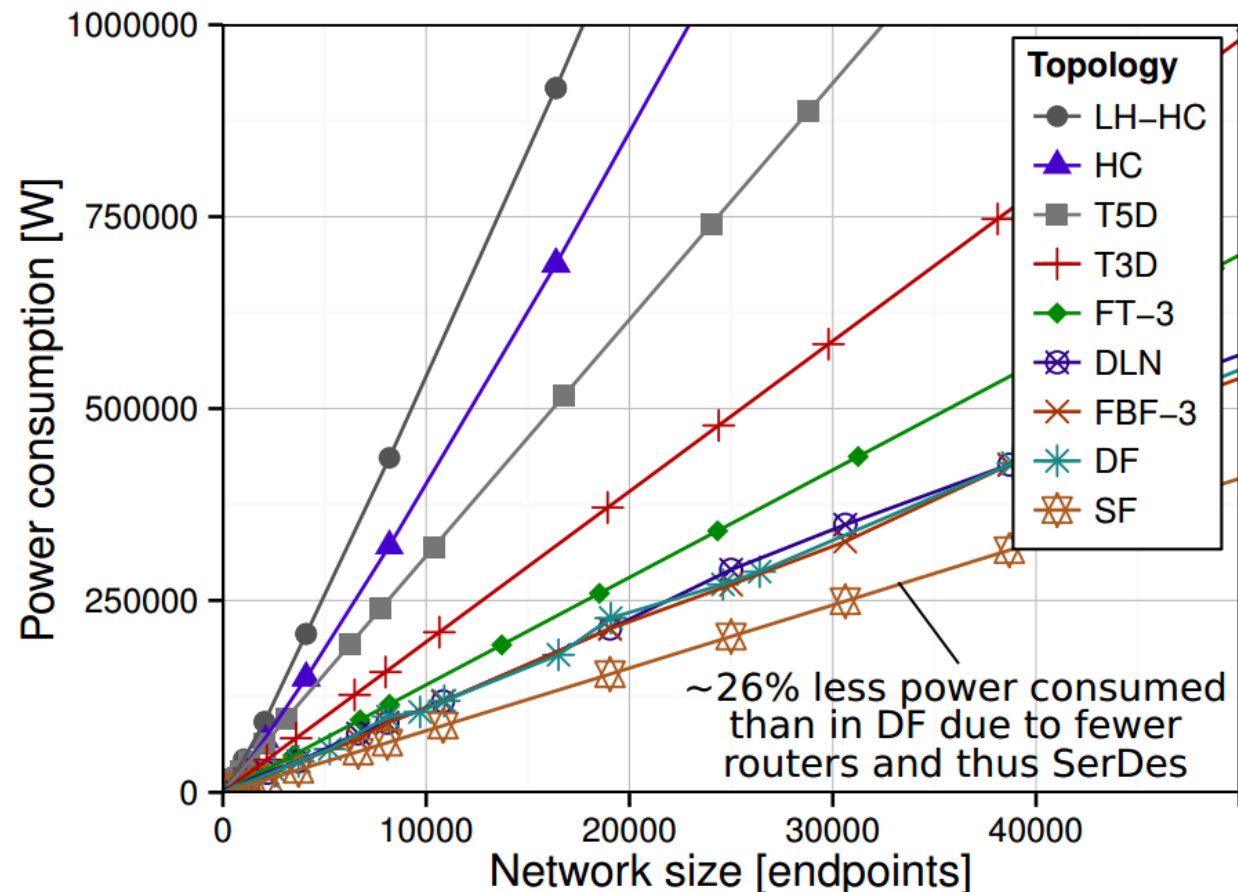
RESULTS



POWER COMPARISON

POWER MODEL

- Model similar to [1],
 - Each router port has four lanes,
 - Each lane has one SerDes,
 - Each SerDes consumes 0.7 W
 - Other parameters as in the cost model



COST & POWER COMPARISON

DETAILED CASE-STUDY

	Low-radix topologies				
Topology	T3D	T5D	HC	LH-HC	SF
Endpoints (N)	10,648	10,368	8,192	8,192	10,830
Routers (N_r)	10,648	10,368	8,192	8,192	722
Radix (k)	7	11	14	19	43
Electric cables	31,900	50,688	32,768	53,248	6,669
Fiber cables	0	0	12,288	12,288	6,869
Cost per node [\$]	1,682	3,176	4,631	6,481	1,033
Power per node [W]	19.6	30.8	39.2	53.2	8.02

	High-radix topologies									
Topology	FT-3	DLN	FBF-3	DF	FT-3	DLN	FBF-3	DF	DF	SF
Endpoints (N)	19,876	40,200	20,736	58,806	10,718	9,702	10,000	9,702	10,890	10,830
Routers (N_r)	2,311	4,020	1,728	5,346	1,531	1,386	1,000	1,386	990	722
Radix (k)	43	43	43	43	35	28	33	27	43	43
Electric cables	19,414	32,488	9,504	56,133	7,350	6,837	4,500	9,009	6,885	6,669
Fiber cables	40,215	33,842	20,736	29,524	24,806	7,716	10,000	4,900	1,012	6,869
Cost per node [\$]	2,346	1,743	1,570	1,438	2,315	1,566	1,535	1,342	1,365	1,033
Power per node [W]	14.0	12.04	10.8	10.9	14.0	11.2	10.8	10.8	10.9	8.02

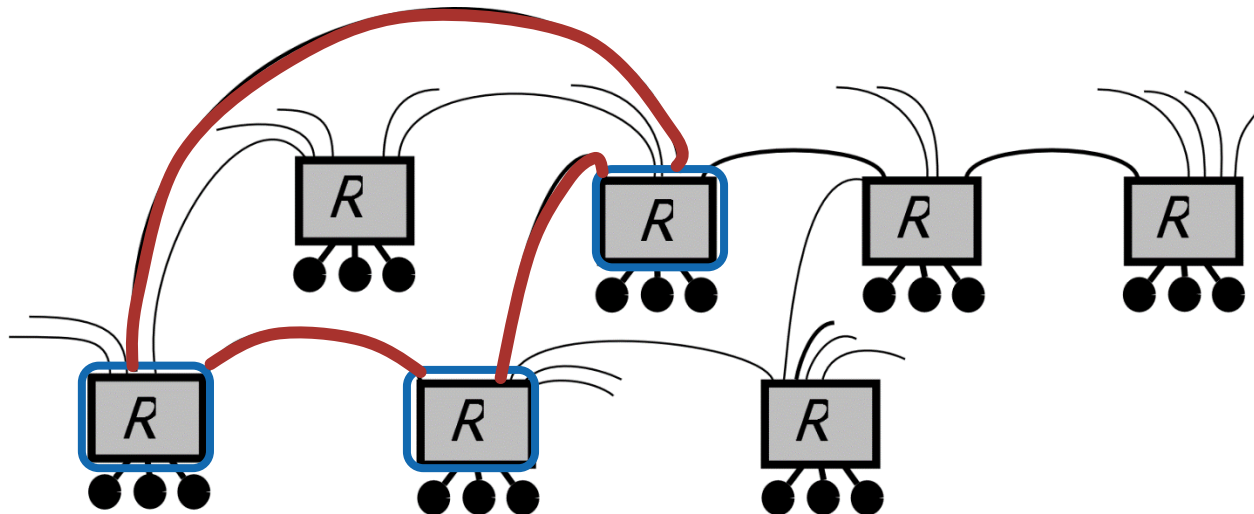
DEADLOCK FREEDOM

MINIMUM STATIC ROUTING

- Assign two virtual channels (VC0 and VC1) to each link
- For a 1-hop path use VC0
- For a 2-hop path use VC0 (hop 1) and VC1 (hop 2)
- One can also use the DFSSSP scheme [1]

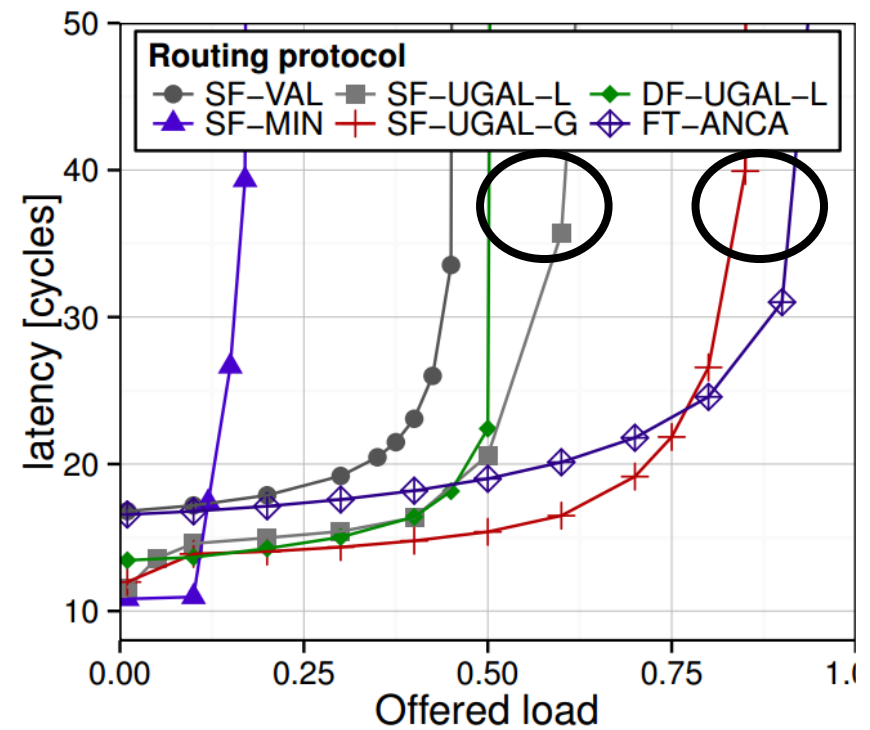
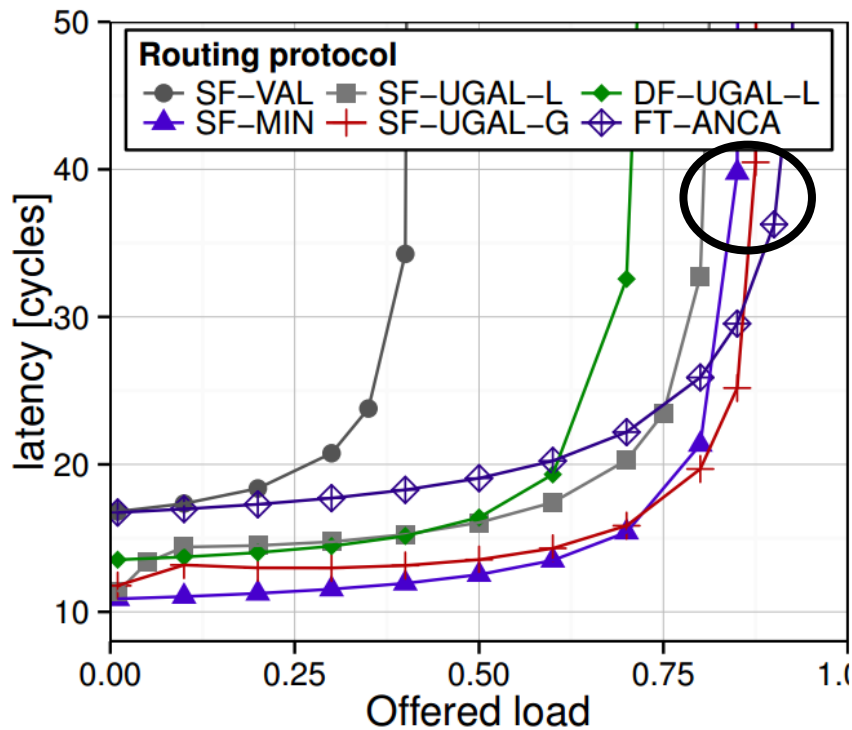
PERFORMANCE

- Cycle-based flit-level simulations (Booksim)
- Routing protocols:
 - Minimum static routing
 - Valiant's random routing
 - Universal Globally-Adaptive Load-Balancing routing
 - UGAL-L*: each router has access to its local output queues
 - UGAL-G*: each router has access to the sizes of all router queues in the network



PERFORMANCE

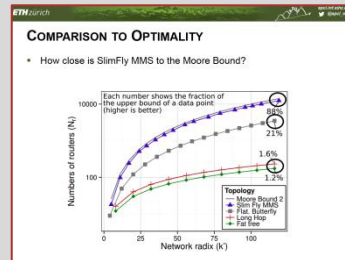
- Random uniform traffic
- Bit permutation (reverse) traffic



CONCLUSIONS

Topology design

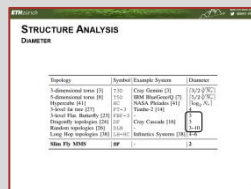
Optimizing towards the Moore Bound reduces expensive network resources



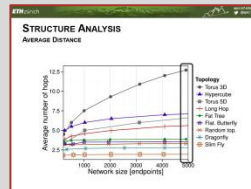
PhD fellowship for parallel computing

Advantages of SlimFly

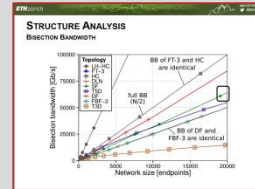
Diameter



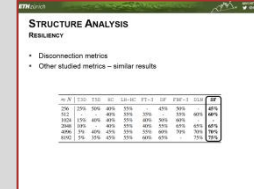
Avg. distance



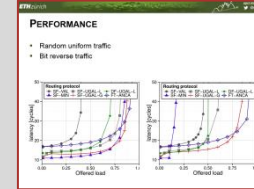
Bisection bandwidth



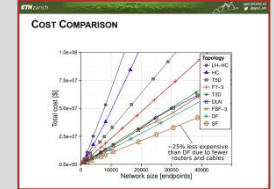
Resilience



Performance



Cost & power



Optimization approach

Combining mathematical optimization and current technology trends effectively tackles challenges in networking

