

Jellyfish networking data centers randomly



Brighten Godfrey • UIUC

Cisco Systems, September 12, 2013

[Photo: Kevin Raskoff]

Ask me about...

Low latency networked systems

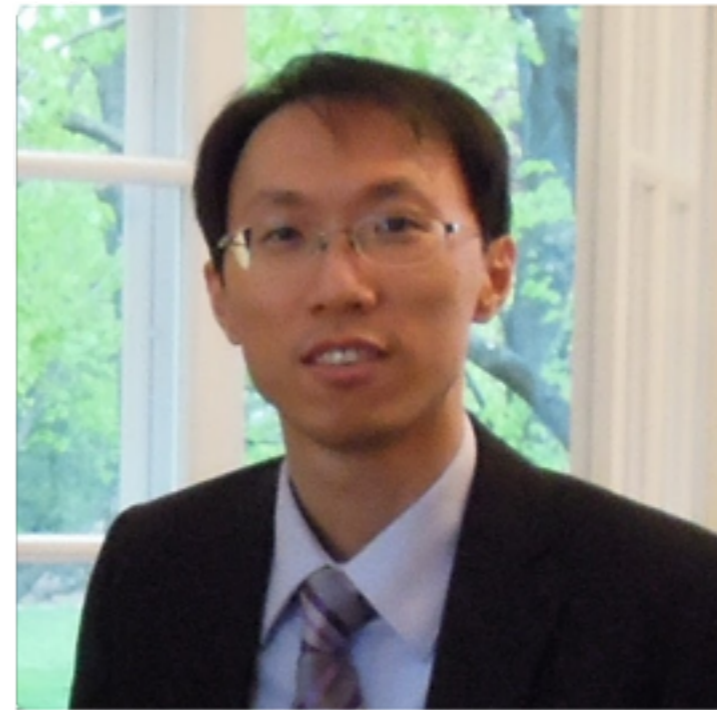
Data plane verification (Veriflow)



Ankit Singla
UIUC

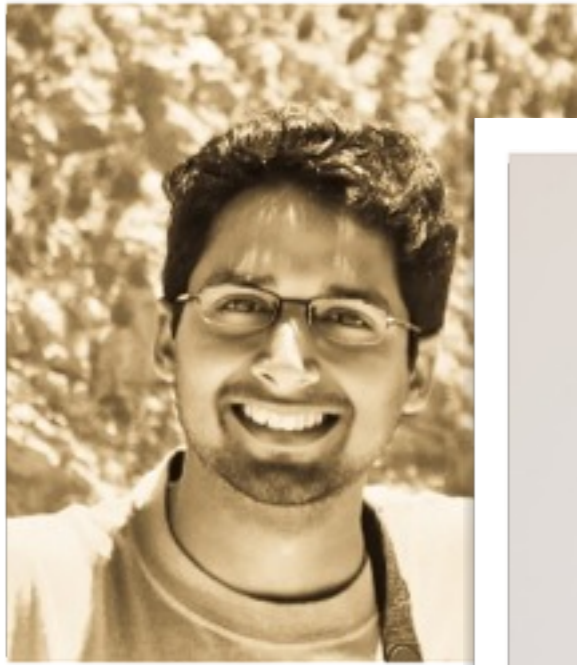


**Chi-Yao
Hong**
UIUC



Kyle Jao
UIUC

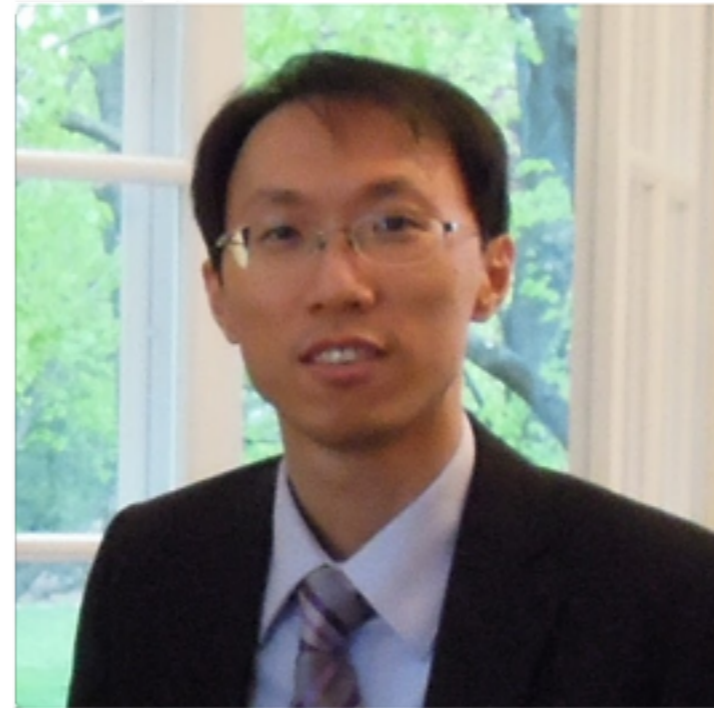
Sangeetha Abdu Jyothi
UIUC



Ankit Singla



Chi-Yao

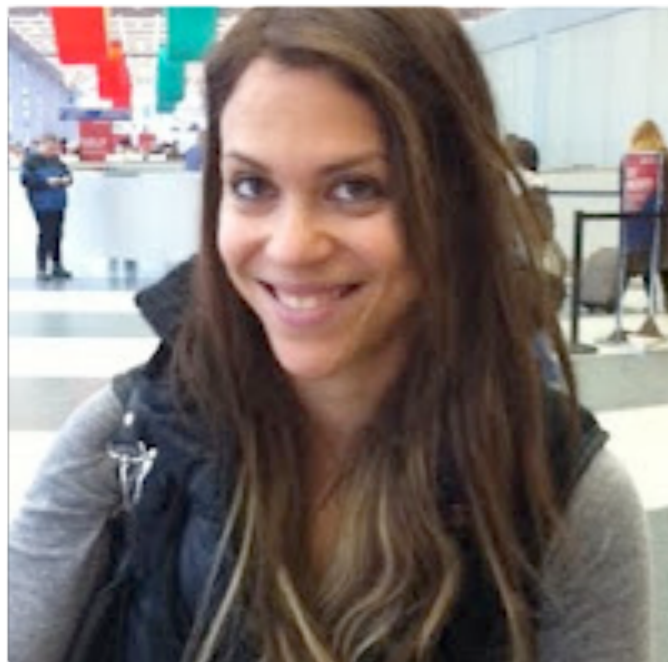


Kyle Jao

UIUC

**Sangeetha
Abdu Jyothi**

UIUC



**Alexandra
Kolla**

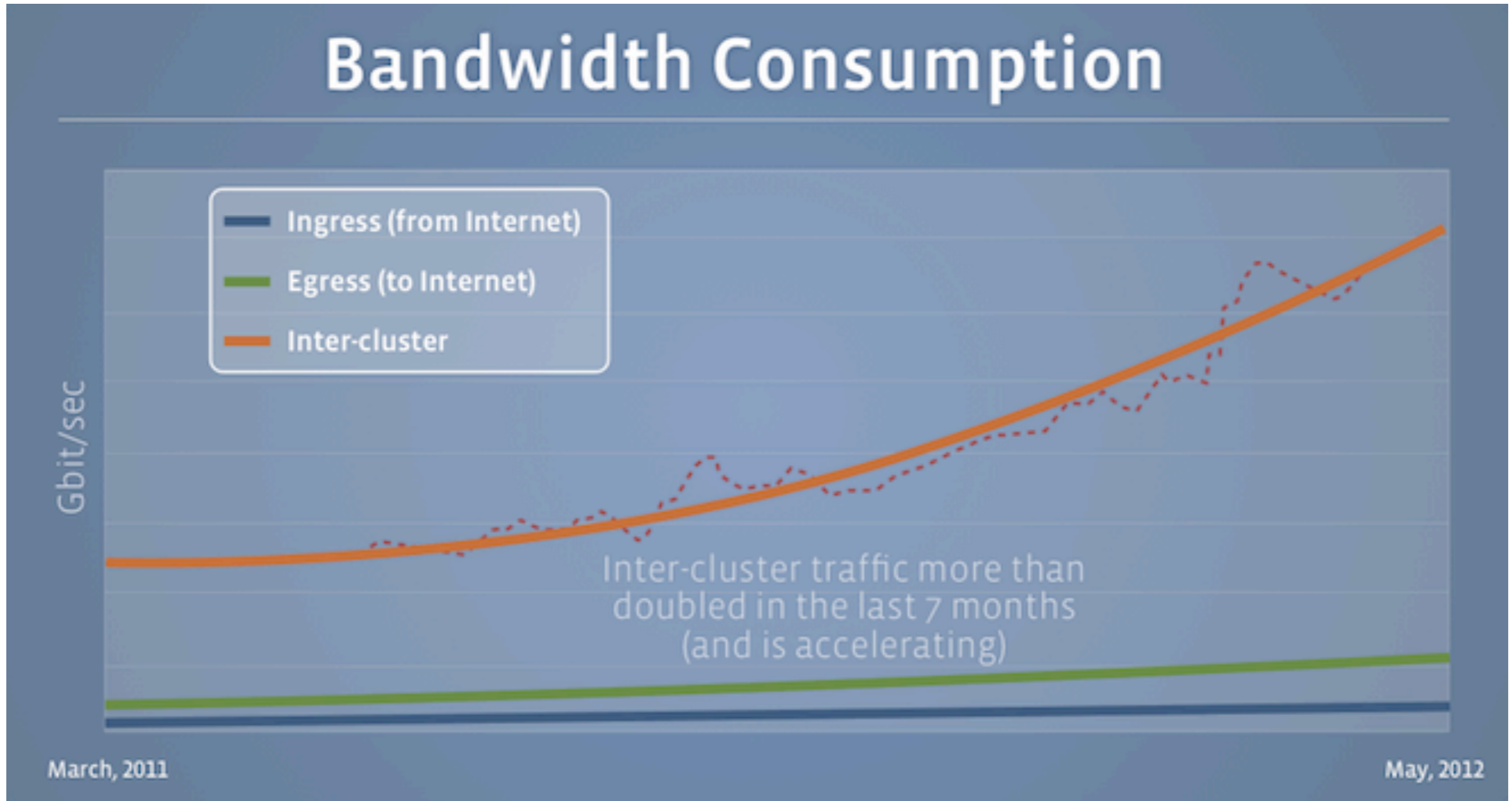
UIUC



Lucian Popa

HP Labs

The need for throughput



March
2011

May
2012

[Facebook, via Wired]

Difficult goals

High throughput
with minimal cost

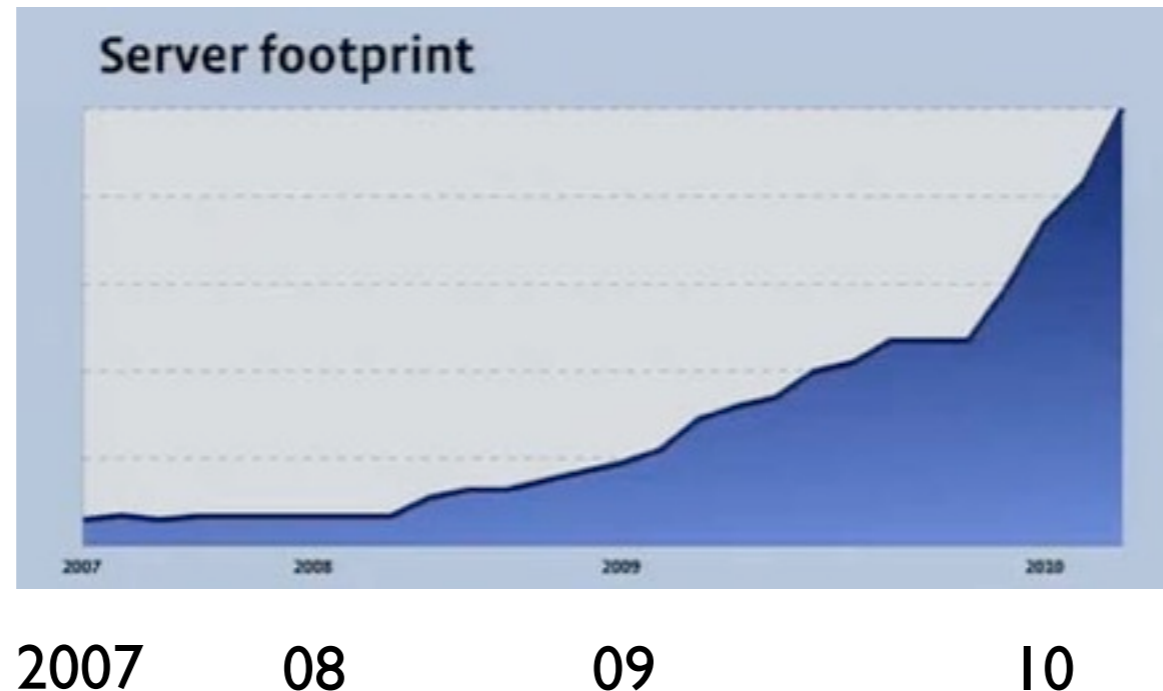
Support big data analytics
Agile placement of VMs

Flexible incremental
expandability

Easily add/replace
servers & switches

Incremental expansion

Facebook “adding capacity on a daily basis”

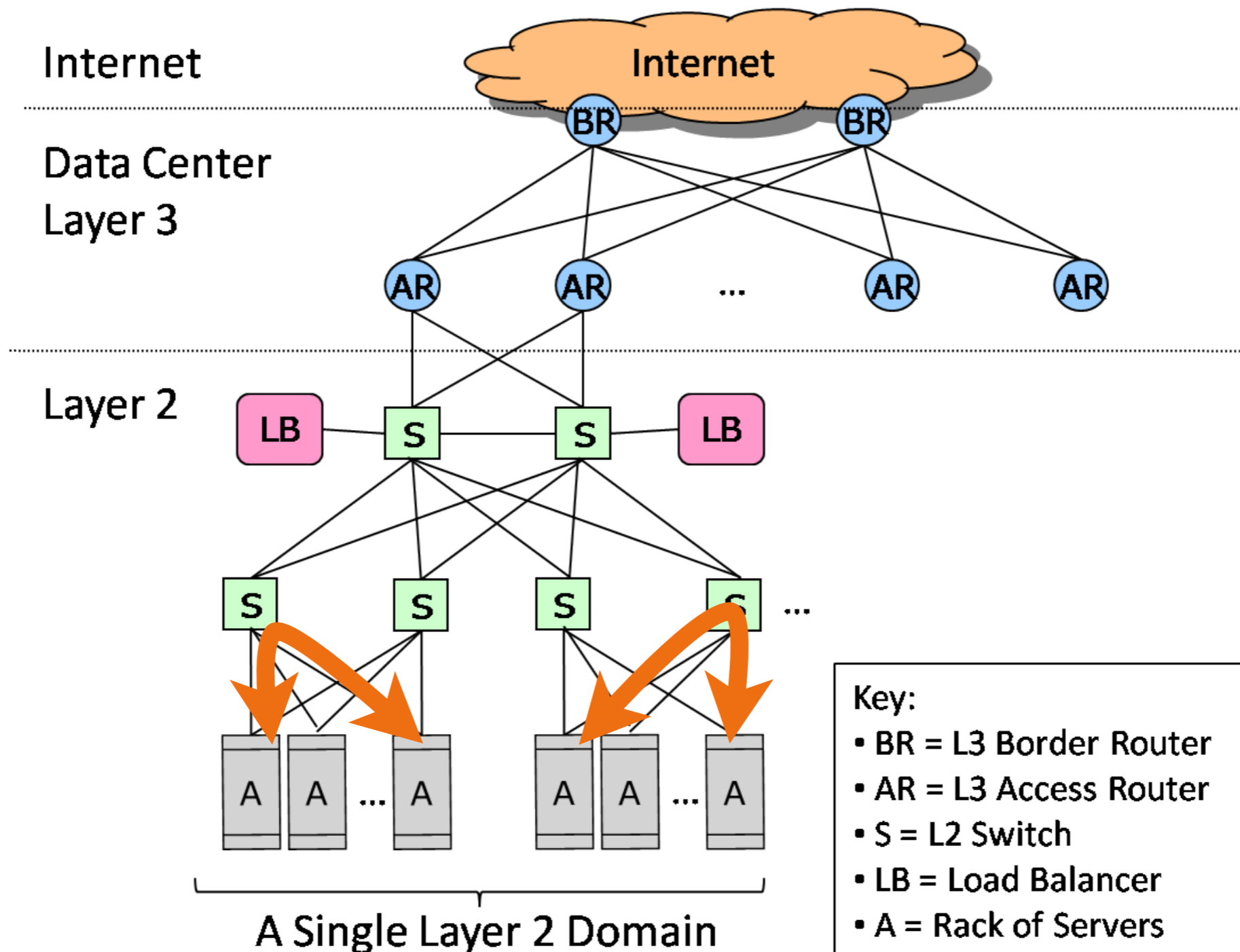


Reduces up-front capital expenditure

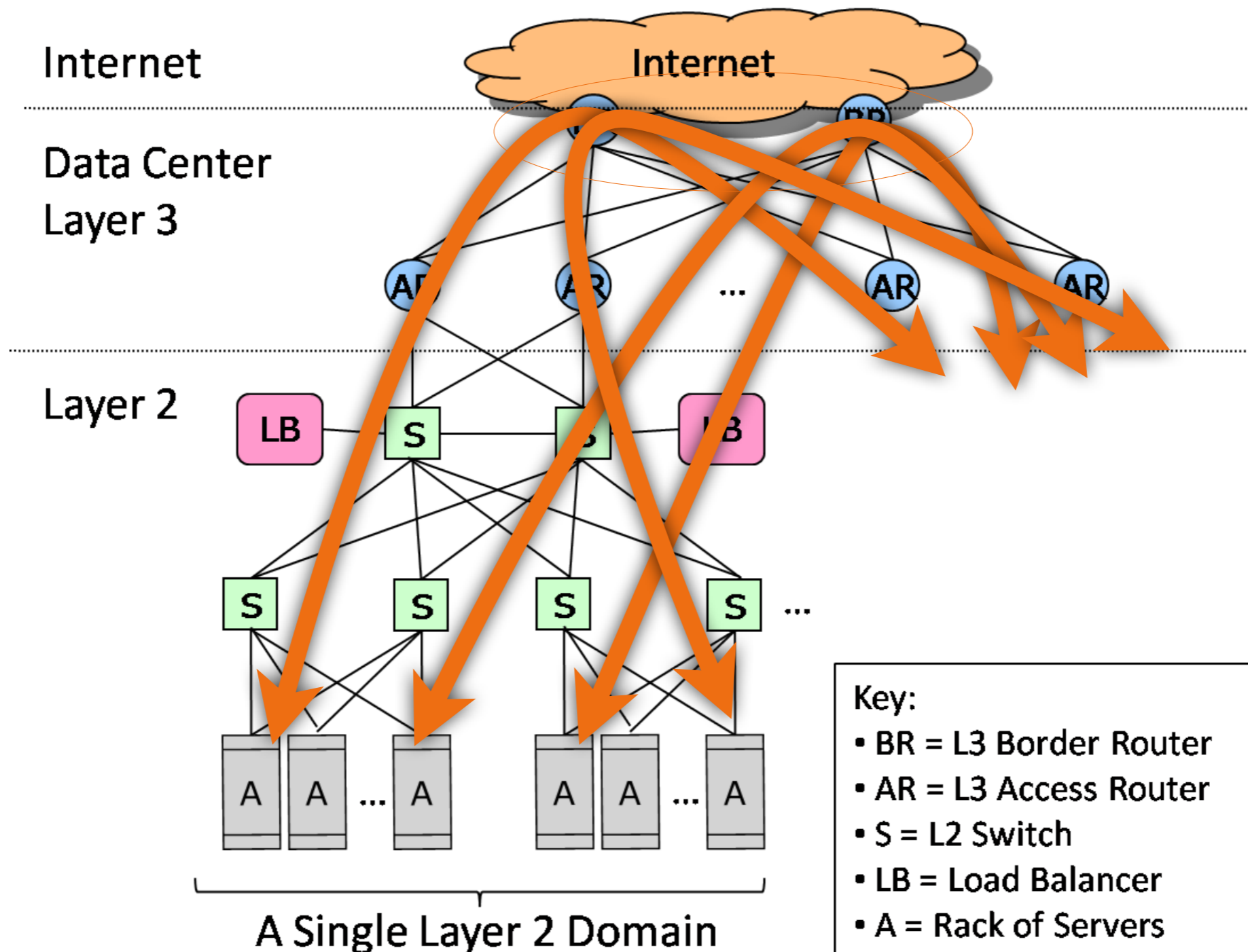
Commercial products expand servers but not the net

- SGI Ice Cube (“Expandable Modular Data Center”)
- HP EcoPod (“Pay-as-you-grow”)

Today's structured networks



Today's structured networks



Today's structured networks

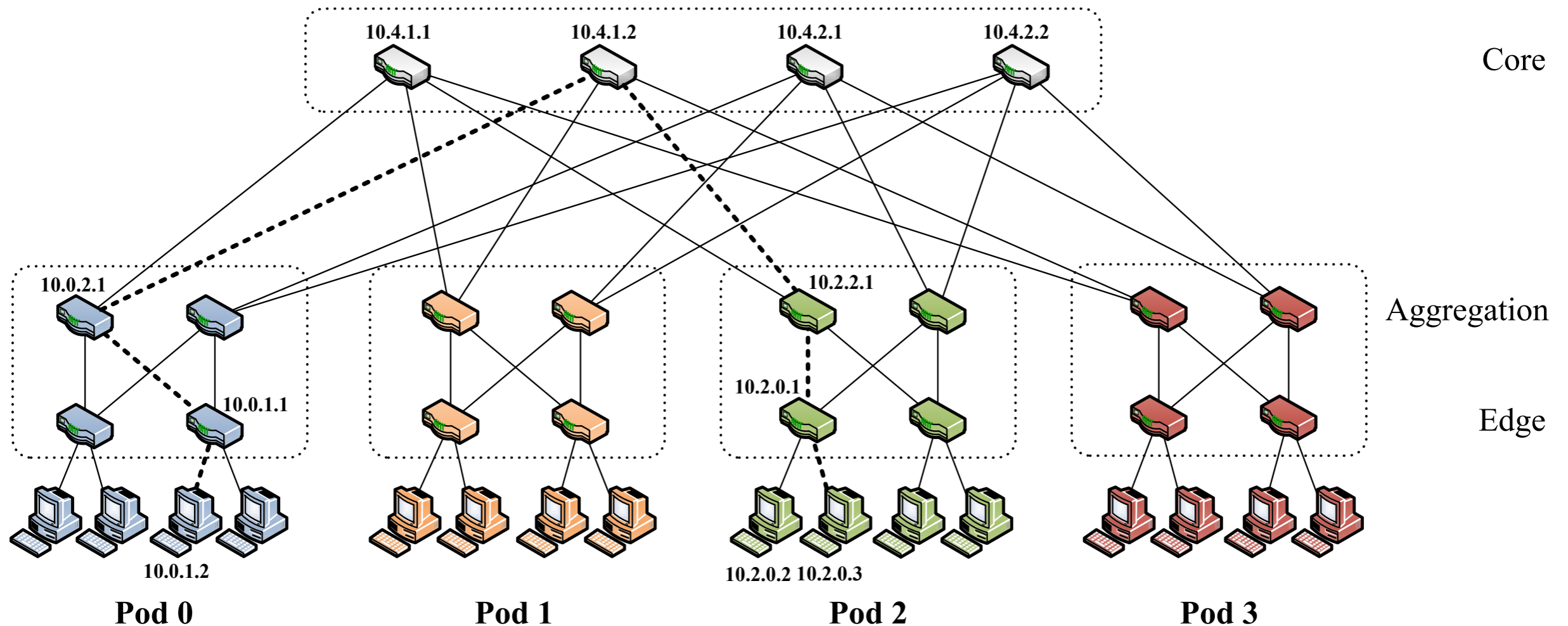
Fat tree



[Al-Fares,
Loukissas, Vahdat,
SIGCOMM '08]

Today's structured networks

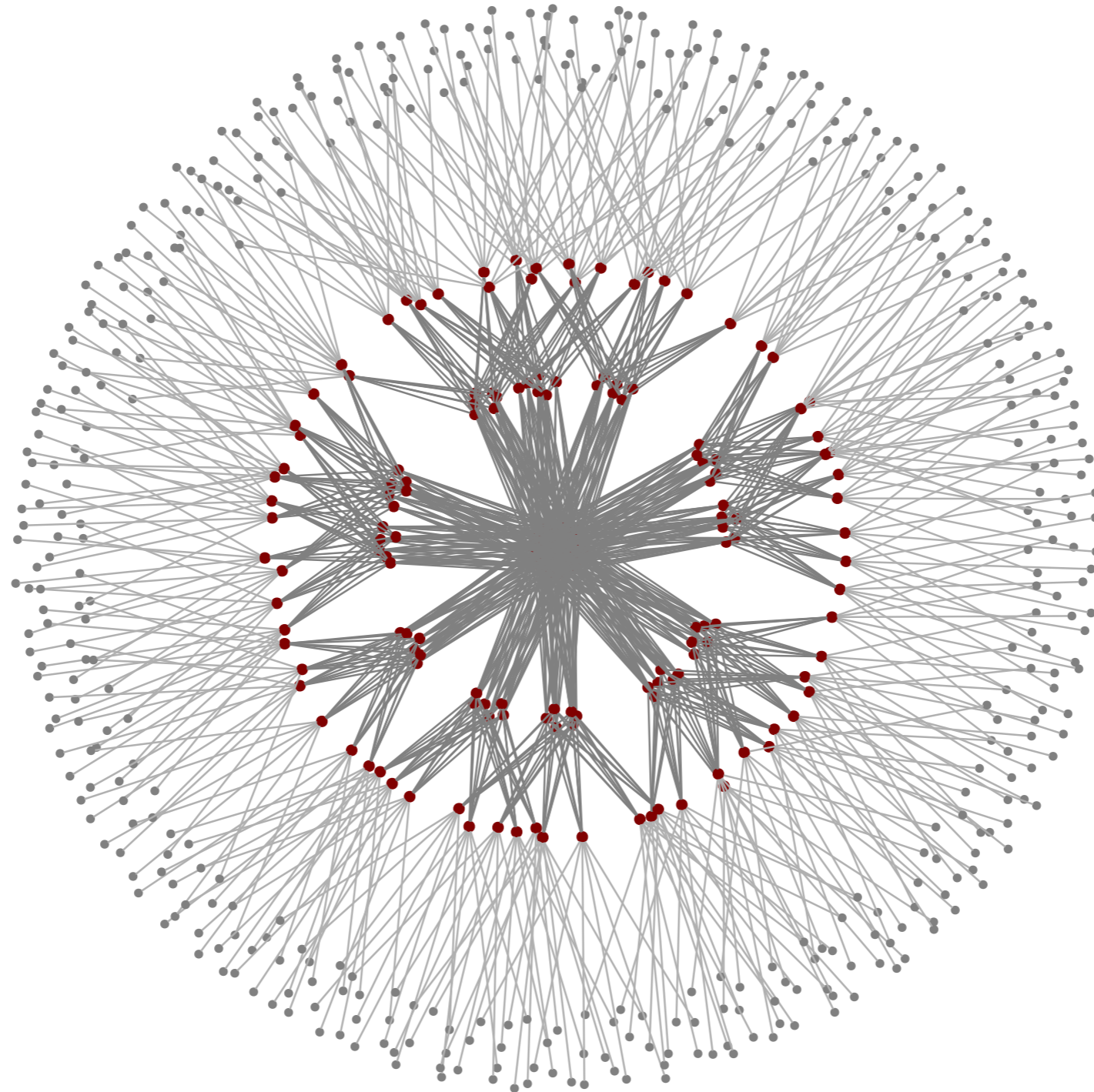
Fat tree



[Al-Fares,
Loukissas, Vahdat,
SIGCOMM '08]

Today's structured networks

Fat tree



Structure constrains expansion

Coarse design points

- Hypercube: 2^k switches
- de Bruijn-like: 3^k switches
- 3-level fat tree: $5k^2/4$ switches

Fat trees by the numbers:

- (3-level, with commodity 24, 32, 48, ... port switches)
- 3456 servers, 8192 servers, 27648 servers, ...

Unclear how to maintain structure incrementally

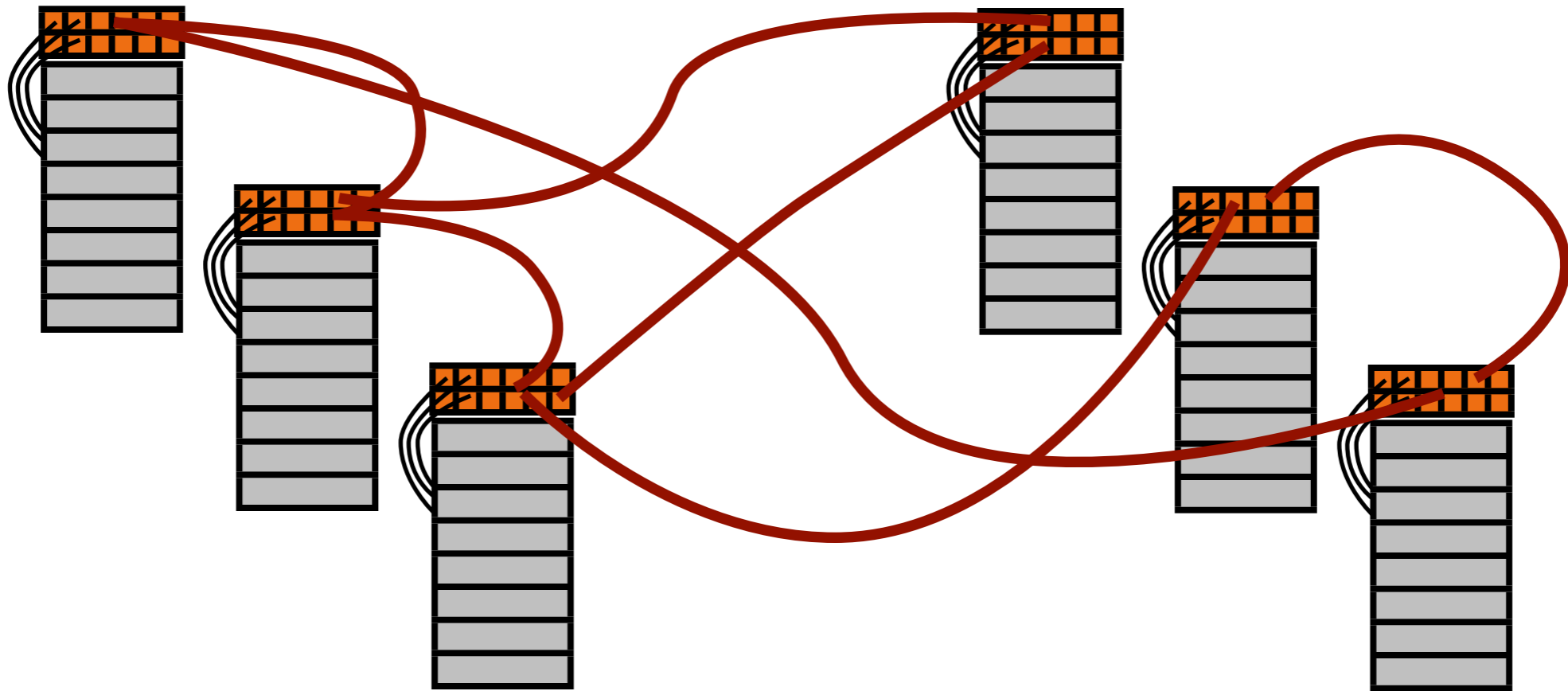
- Overutilize switches? Uneven / constrained bandwidth
- Leave ports free for later? Wasted investment

Our Solution

Forget about structure –
let's have **no structure at all!**

Jellyfish: The Topology

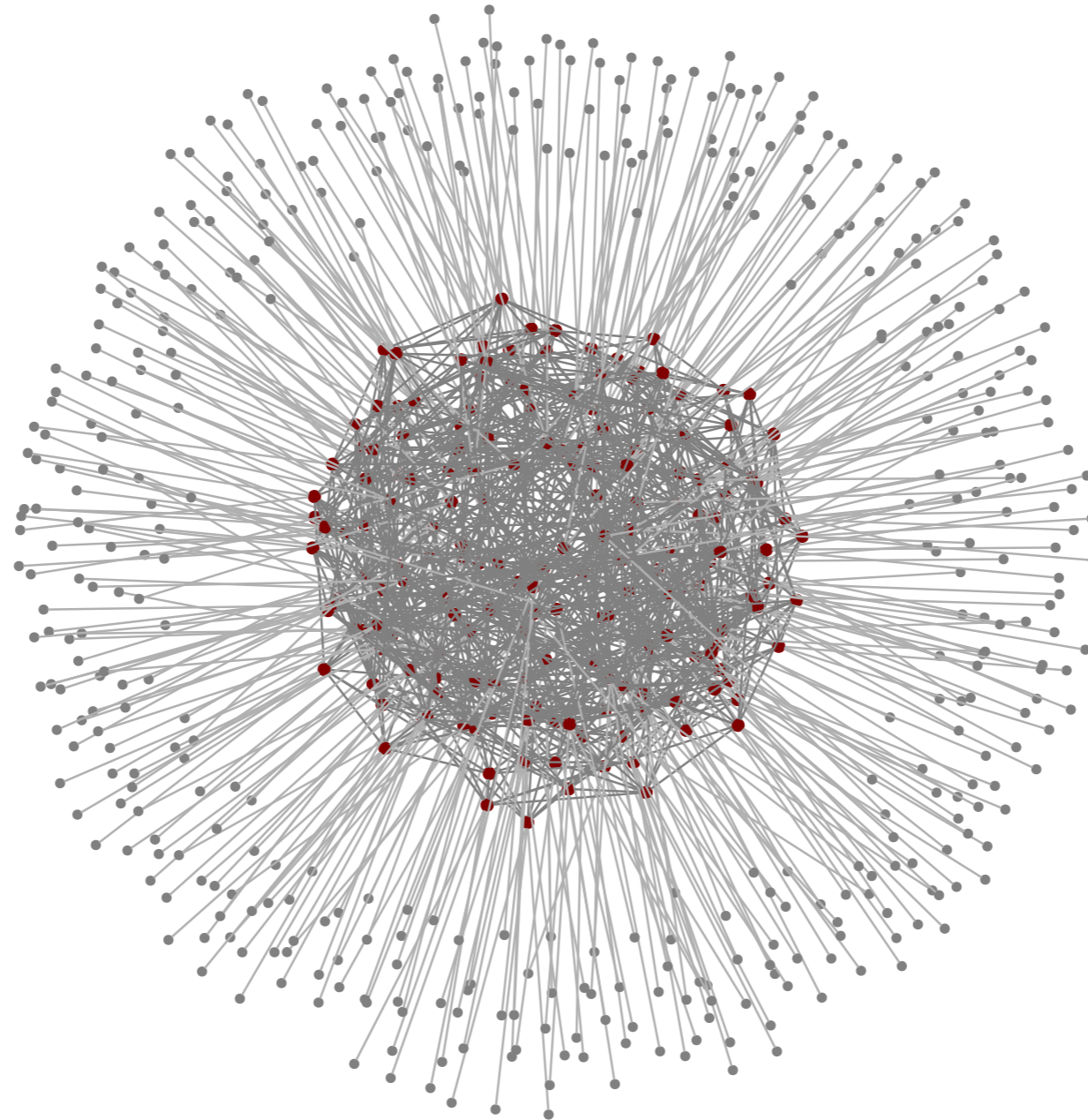
Jellyfish: The Topology



Servers connected to top-of-rack switch

Switches form uniform-random interconnections

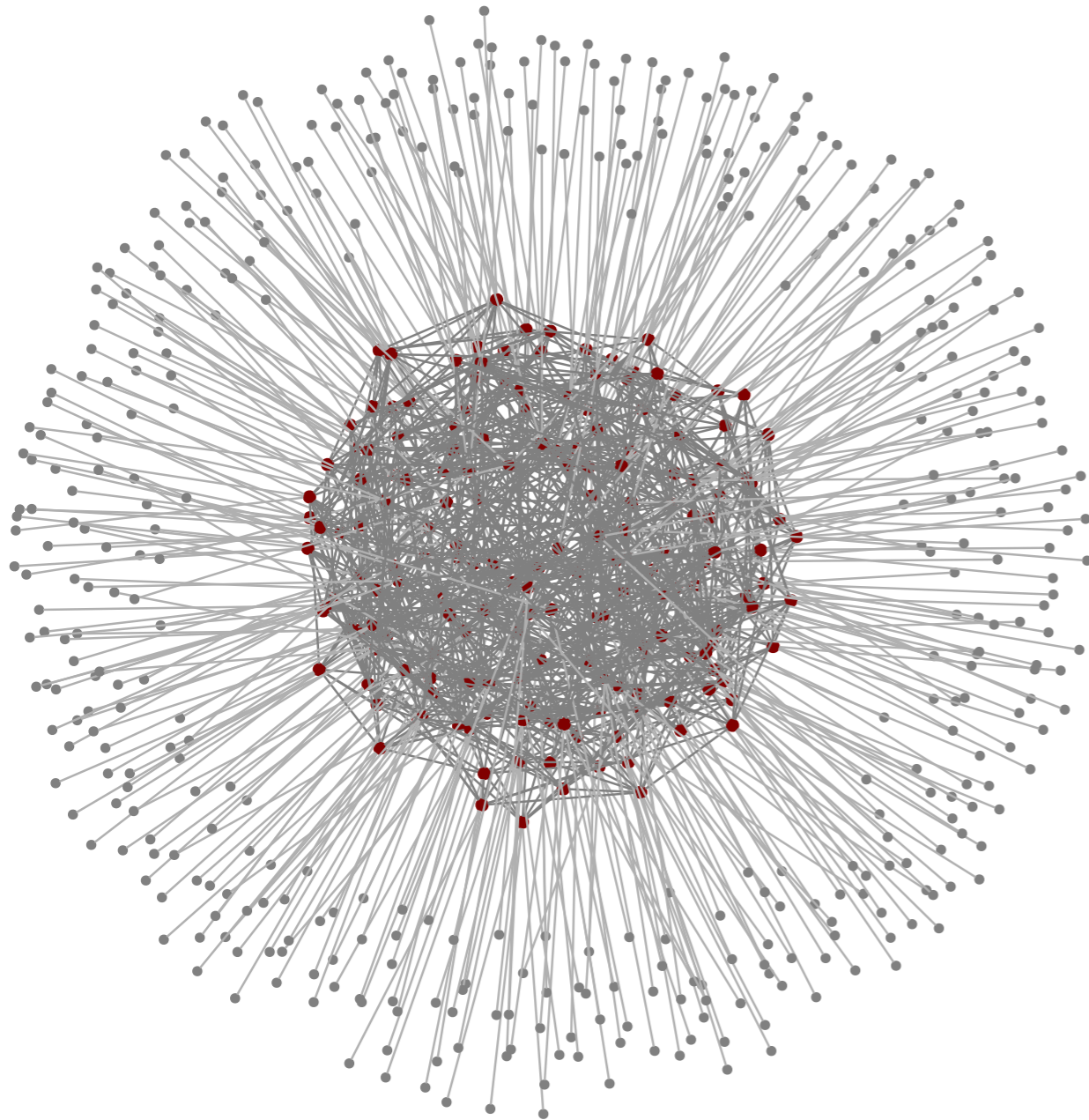
Capacity as a fluid



Jellyfish random graph

432 servers, 180 switches, degree 12

Capacity as a fluid



Jellyfish random graph

432 servers, 180 switches, degree 12

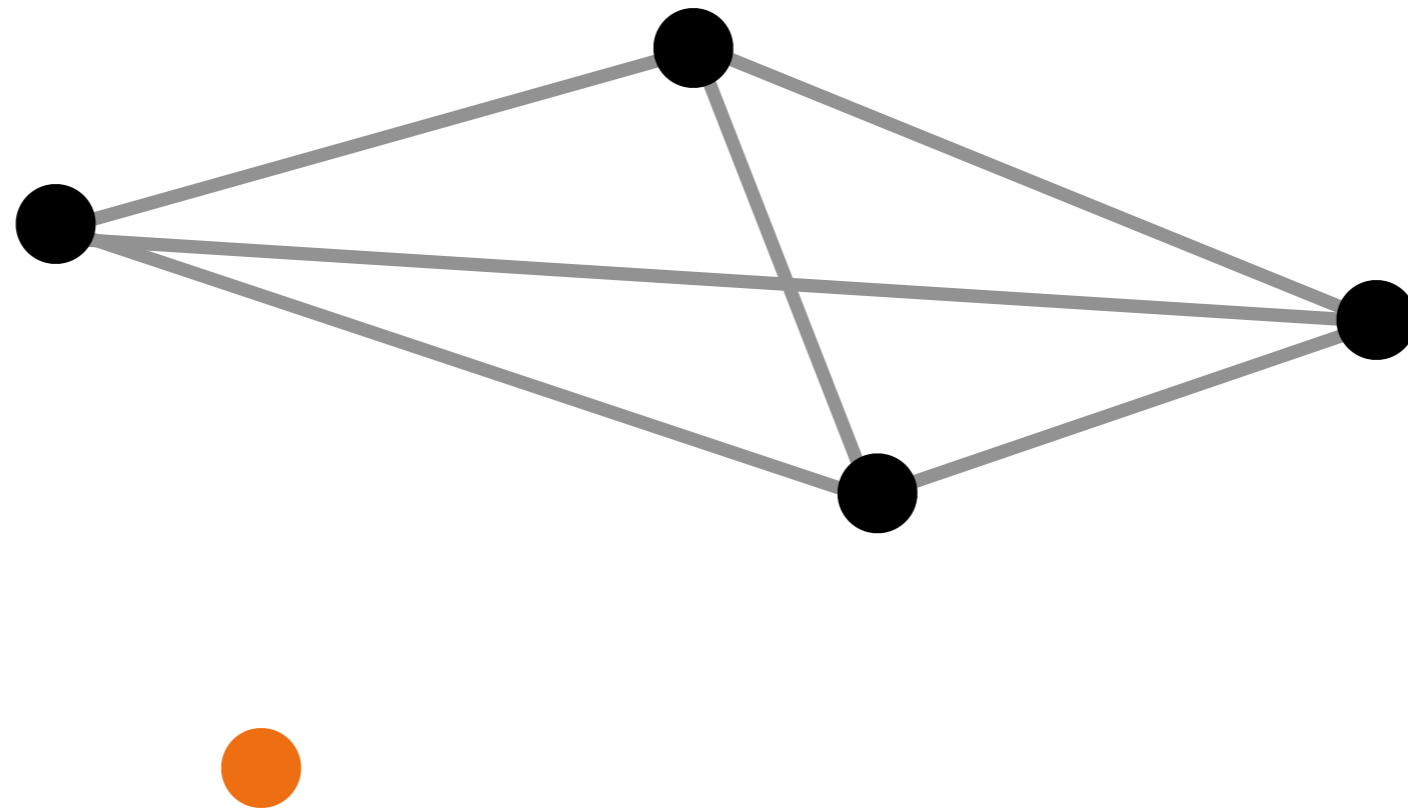


Jellyfish

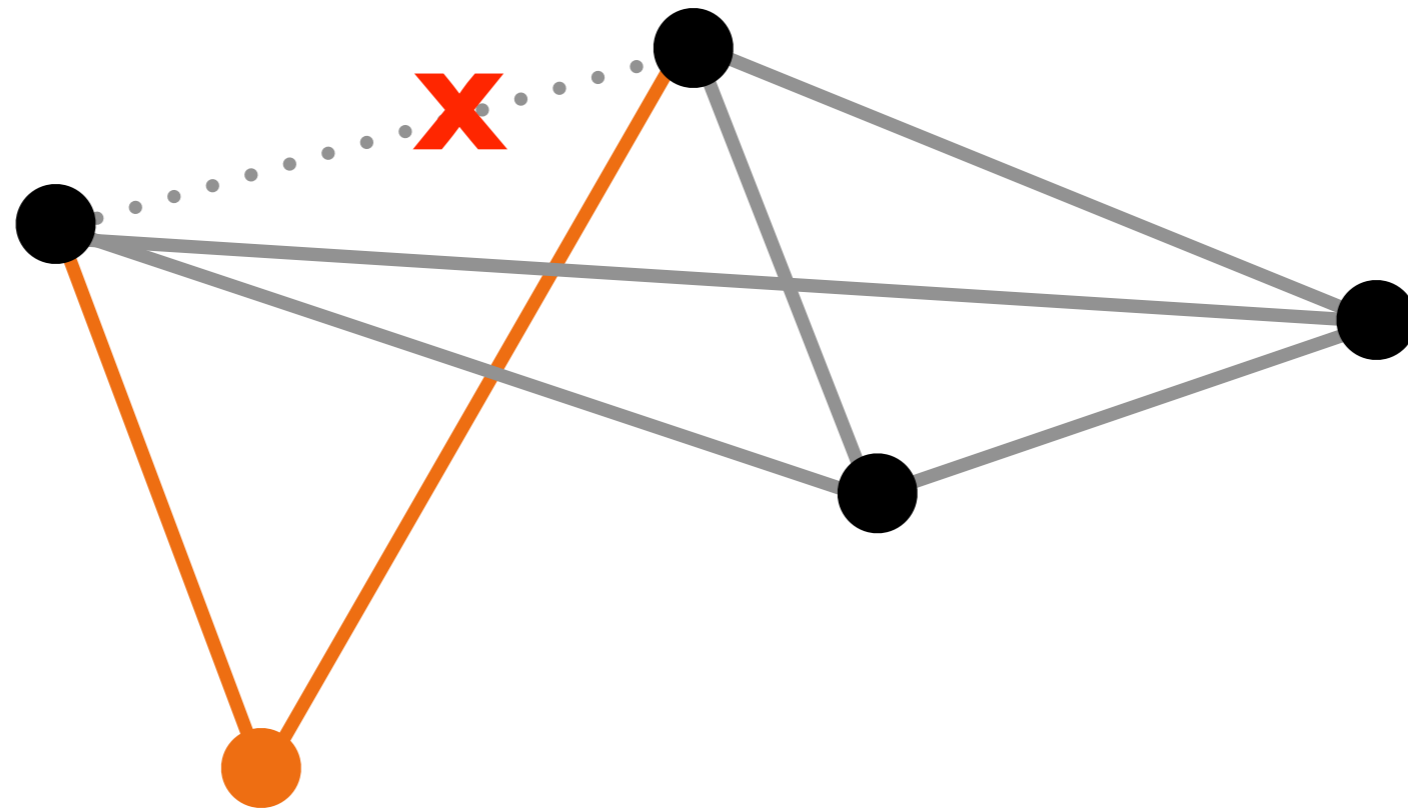
Crossota norvegica
Photo: Kevin Raskoff

Construction & Expansion

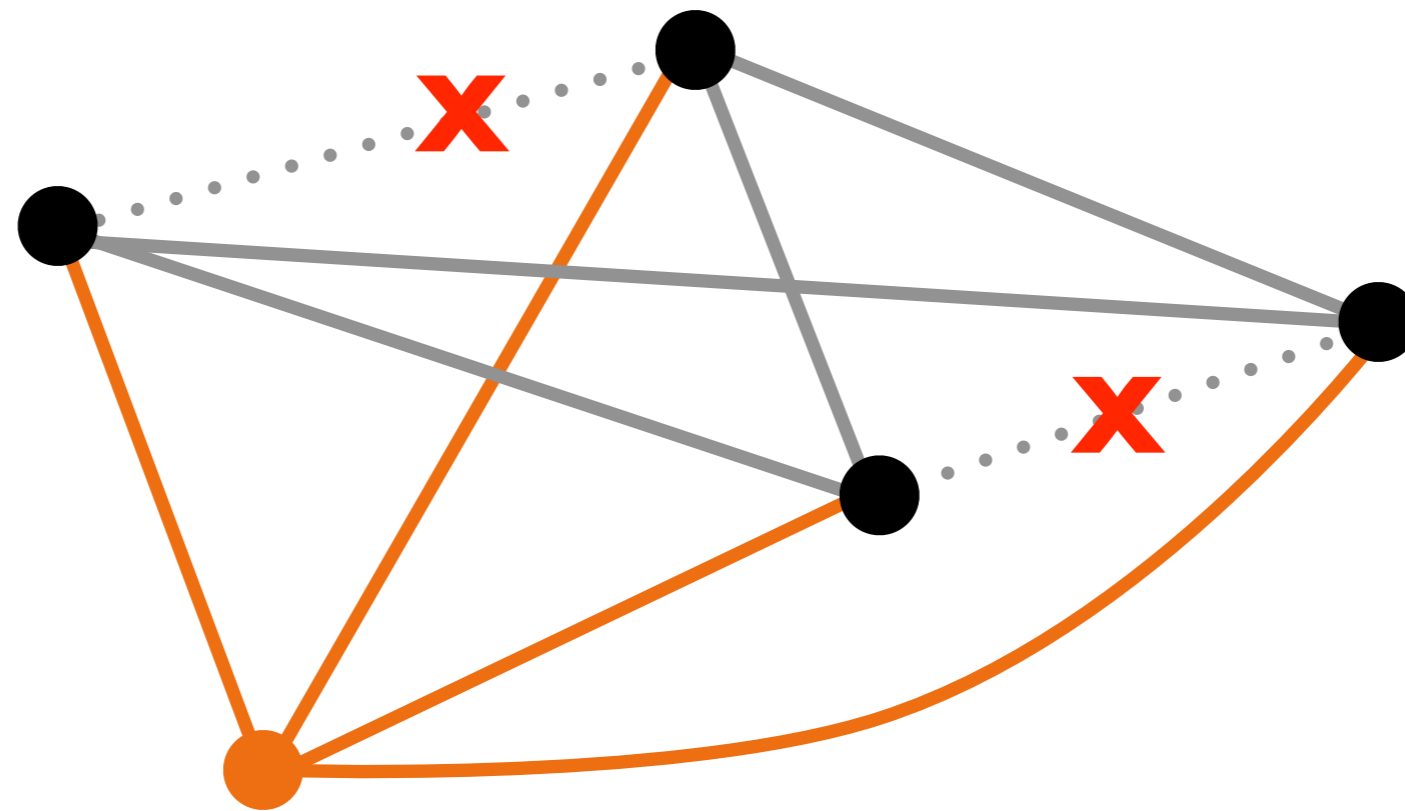
Building Jellyfish



Building Jellyfish



Building Jellyfish



Same procedure for initial construction
and incremental expansion

Can flexibly incorporate any type of equipment

Building Jellyfish

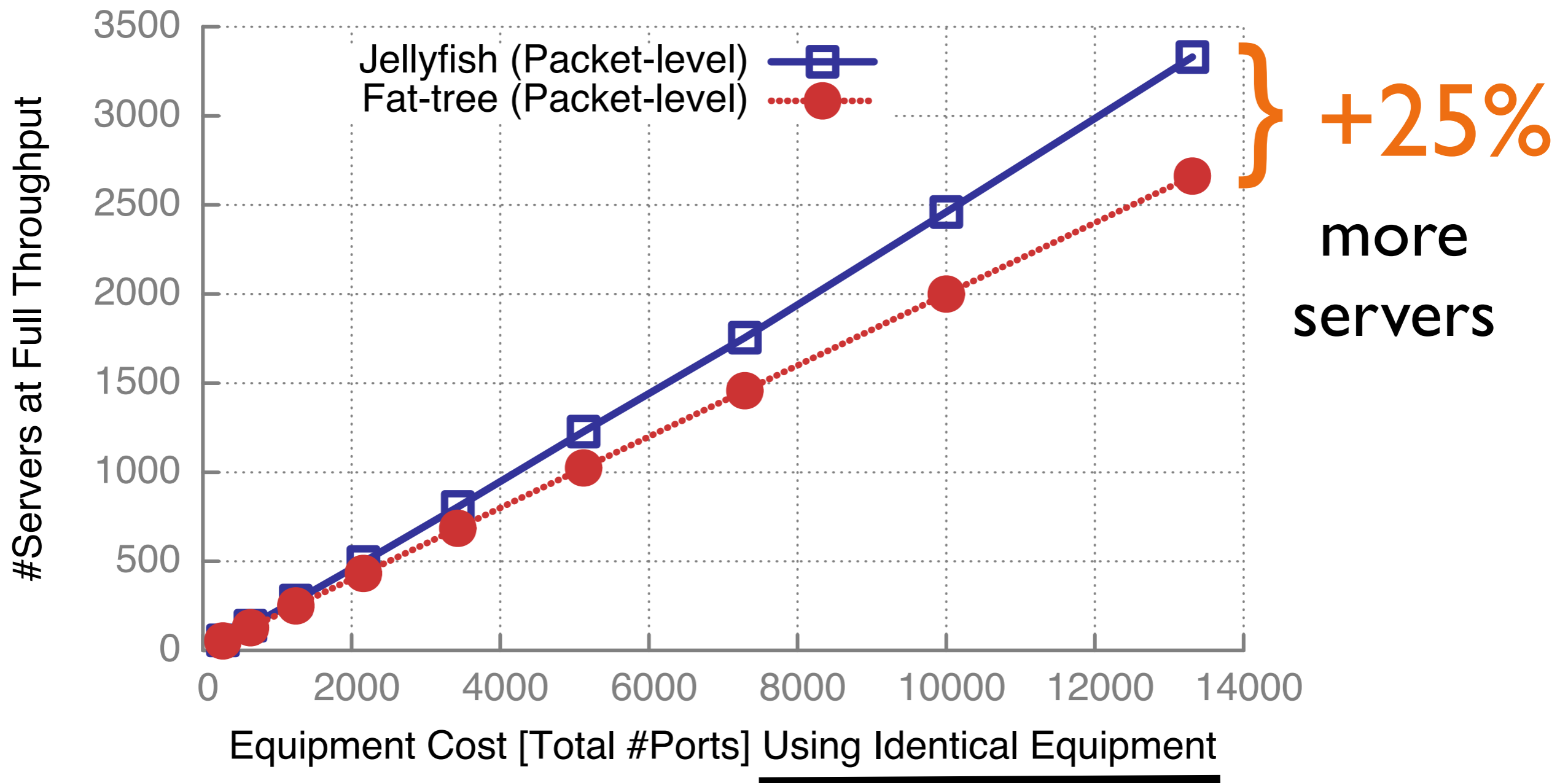
60% cheaper incremental expansion
compared with past technique for
traditional networks

LEGUP: [Curtis, Keshav, Lopez-Ortiz, CoNEXT'10]

Throughput

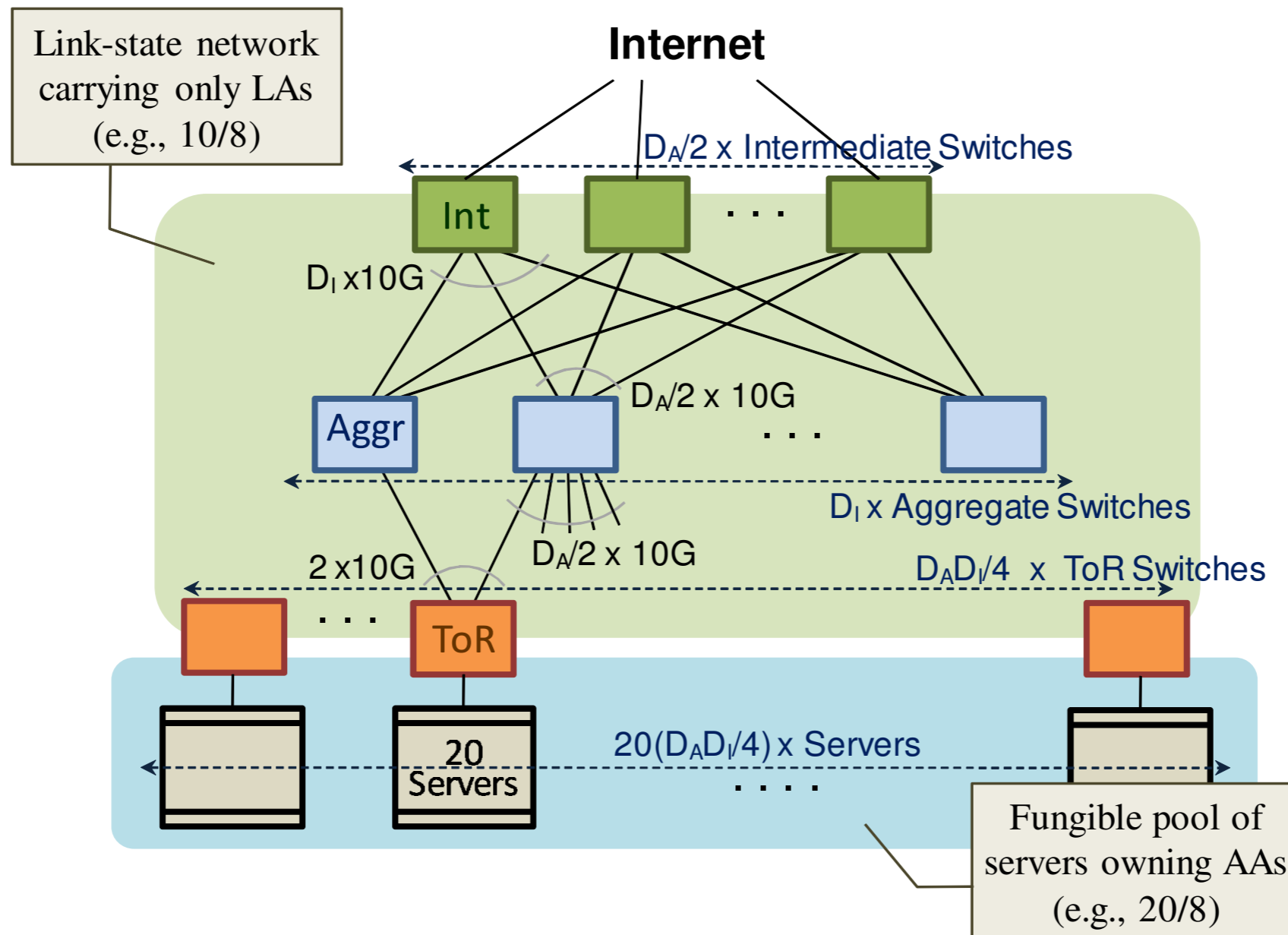
By giving up on structure,
do we take a hit on throughput?

Throughput: Jellyfish vs. fat tree

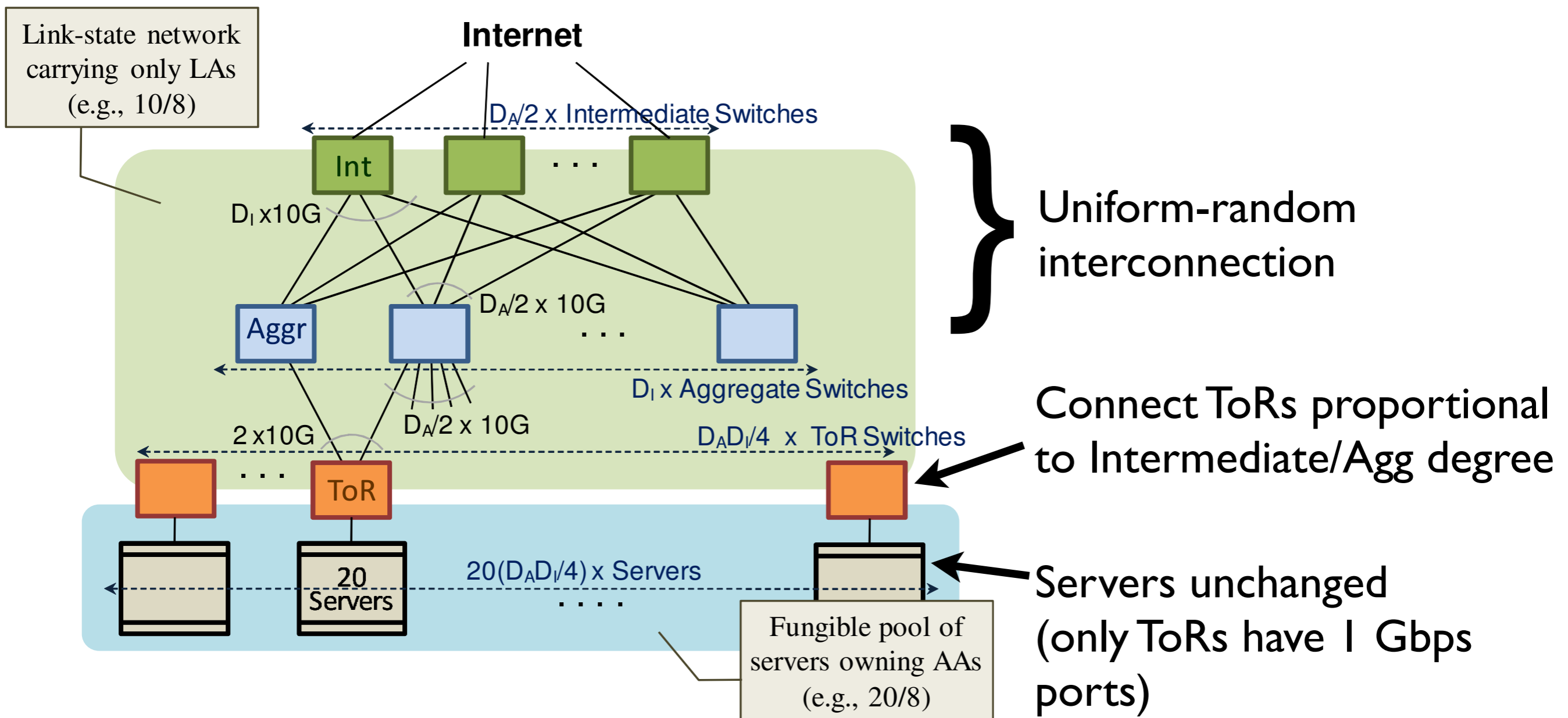


The VL2 topology

[Greenburg, Hamilton, Jain, Kandula, Kim, Lahiri, Maltz, Patel, Sengupta, SIGCOMM'09]

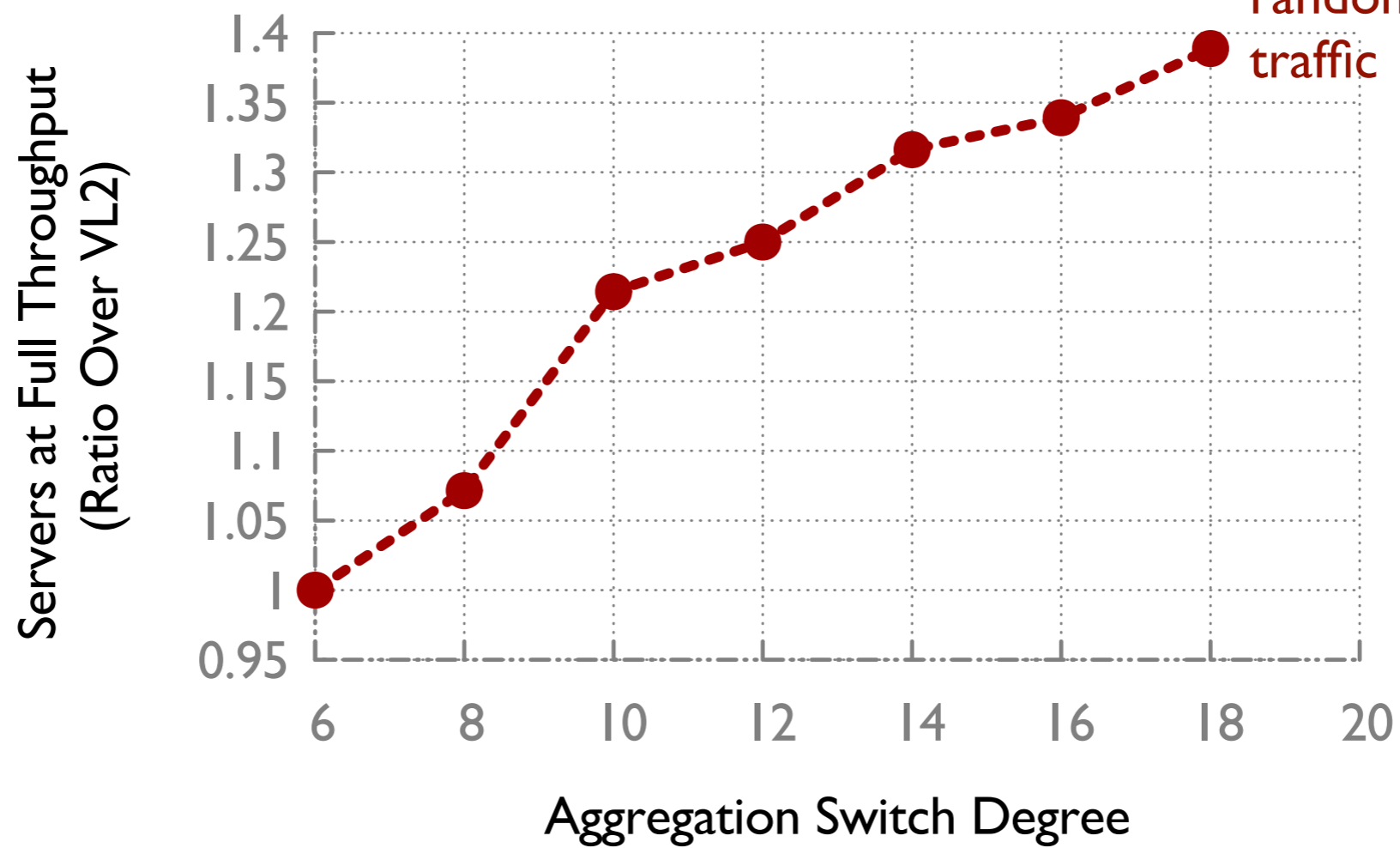


Rewiring VL2

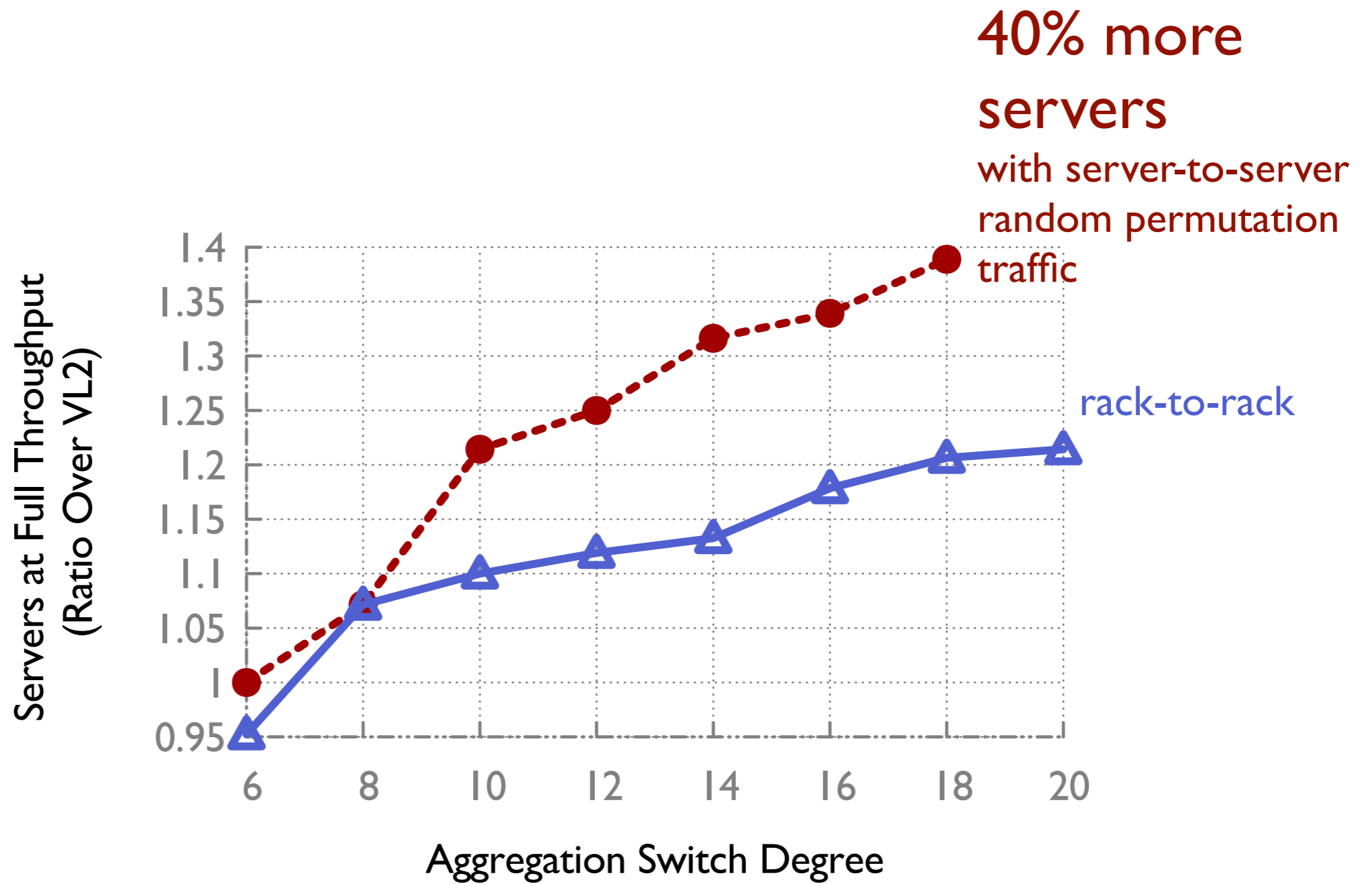


Rewiring VL2

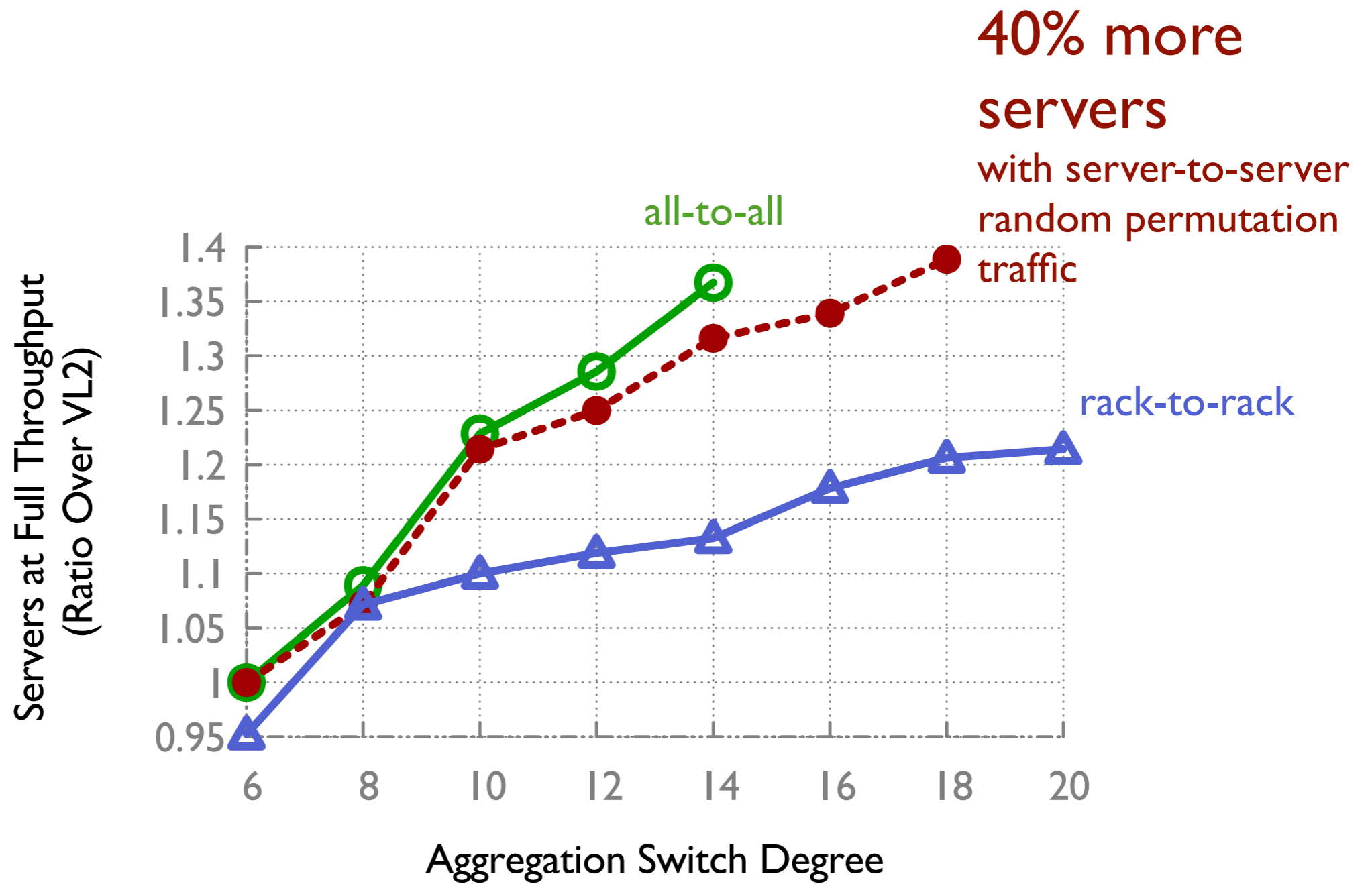
40% more servers
with server-to-server
random permutation
traffic



Rewiring VL2



Rewiring VL2



Just the beginning

Just the beginning

Topology design

- How close are random graphs to optimal?
- What if switches are heterogeneous?

System design (or: “*But what about...*”)

- Performance consistency?
- Cabling spaghetti?
- Routing and congestion control without structure?

Just the beginning

Topology design

- How close are random graphs to optimal?
- What if switches are heterogeneous?

System design (or: “*But what about...*”)

- Performance consistency?
- Cabling spaghetti?
- Routing and congestion control without structure?

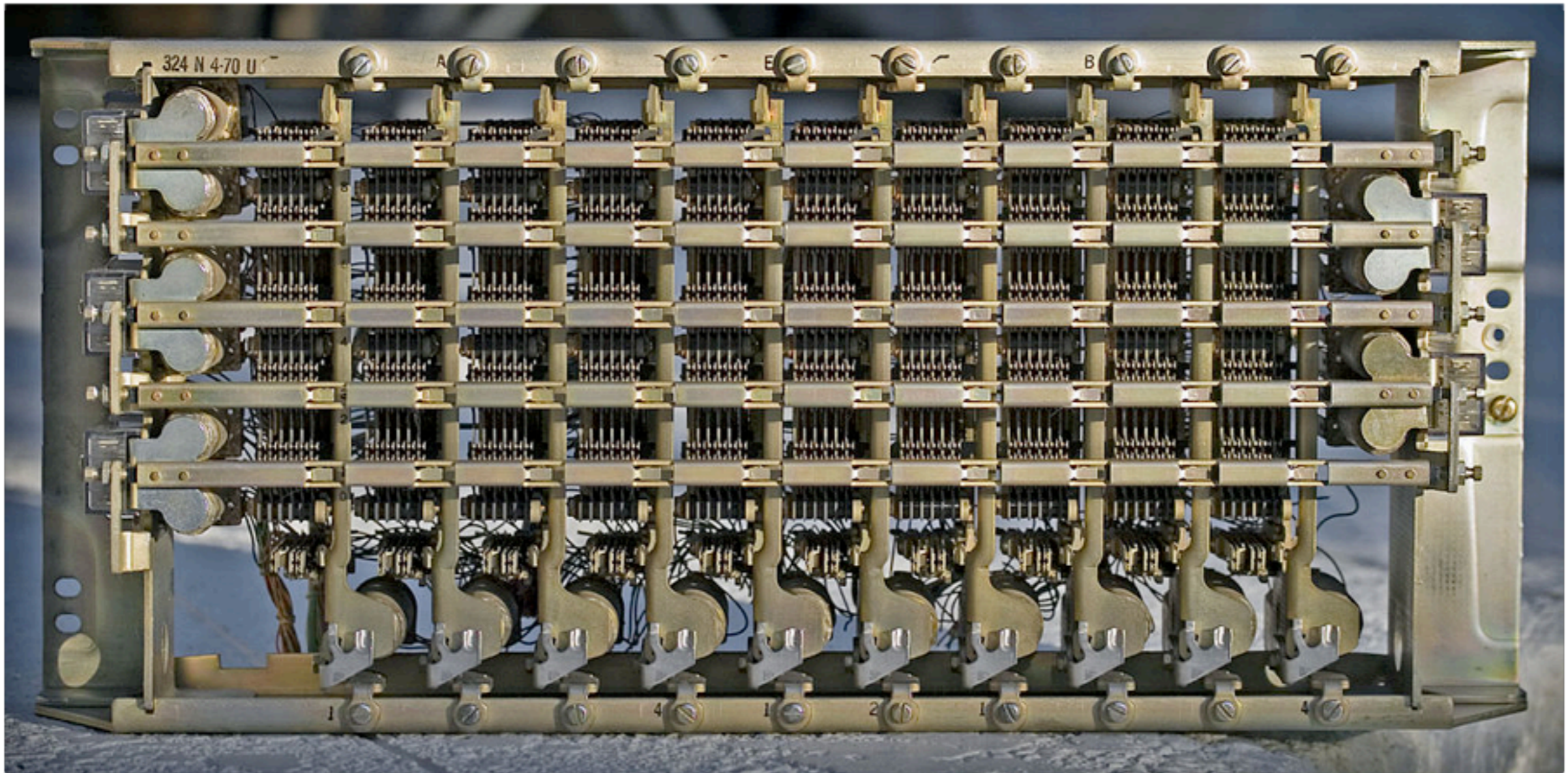
Topology Design in Context

“ “ *It is anticipated that the whole of the populous parts of the United States will, within two or three years, be covered with network like a spider's web.* ” ”

It is anticipated that the whole of the populous parts of the United States will, within two or three years, be covered with network like a spider's web.

— *The London Anecdotes,*
1848





Western Electric crossbar switch

[Photo: Wikipedia user Yeatesh]

A Study of Non-Blocking Switching Networks

By CHARLES CLOS

(Manuscript received October 30, 1952)

This paper describes a method of designing arrays of crosspoints for use in telephone switching systems in which it will always be possible to establish a connection from an idle inlet to an idle outlet regardless of the number of calls served by the system.

INTRODUCTION

The impact of recent discoveries and developments in the electronic art is being felt in the telephone switching field. This is evidenced by

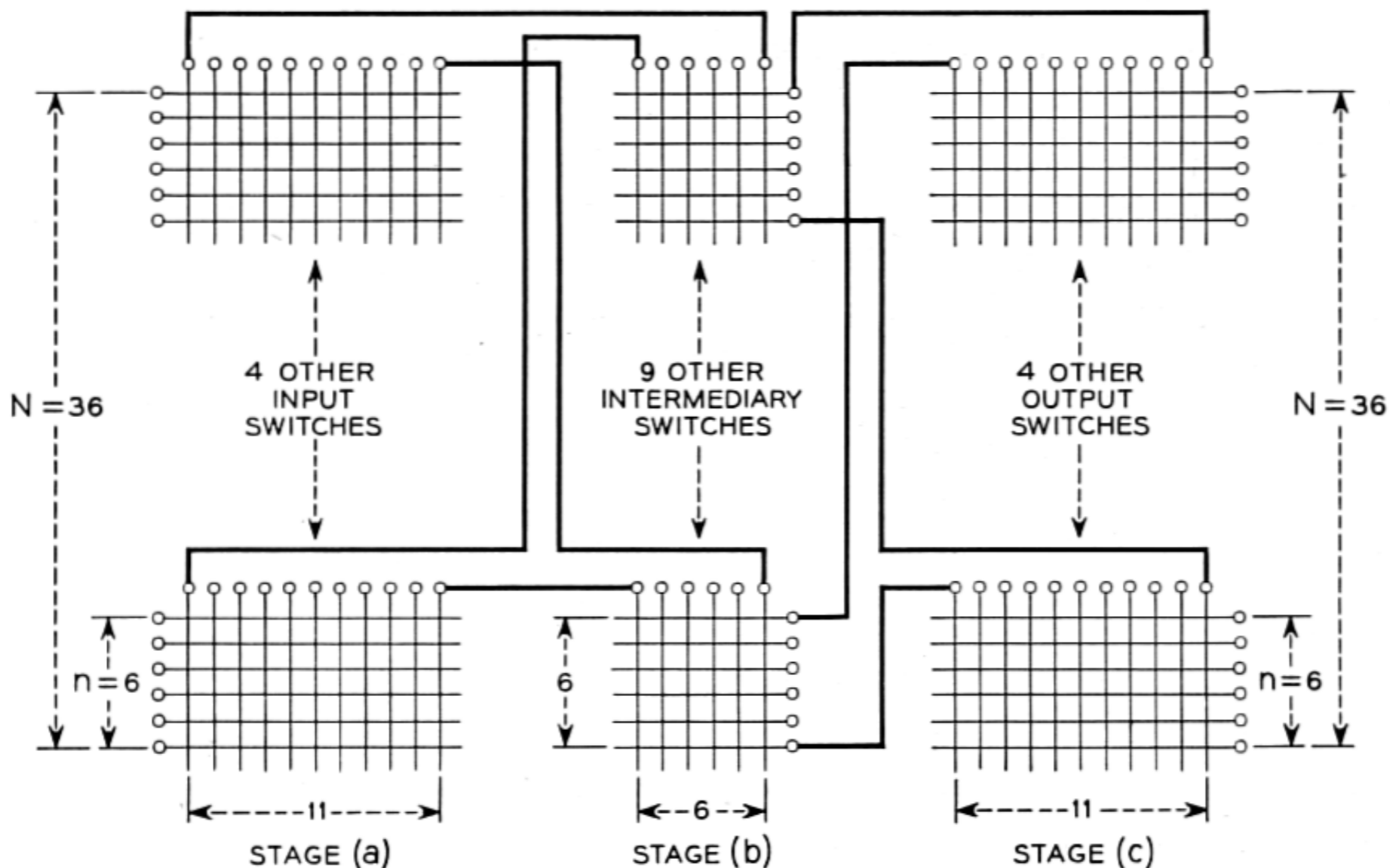
A Study of Non-Blocking Switching Networks

(Me

This paper describes use in telephone switching to establish a connection between a large number of calls served.

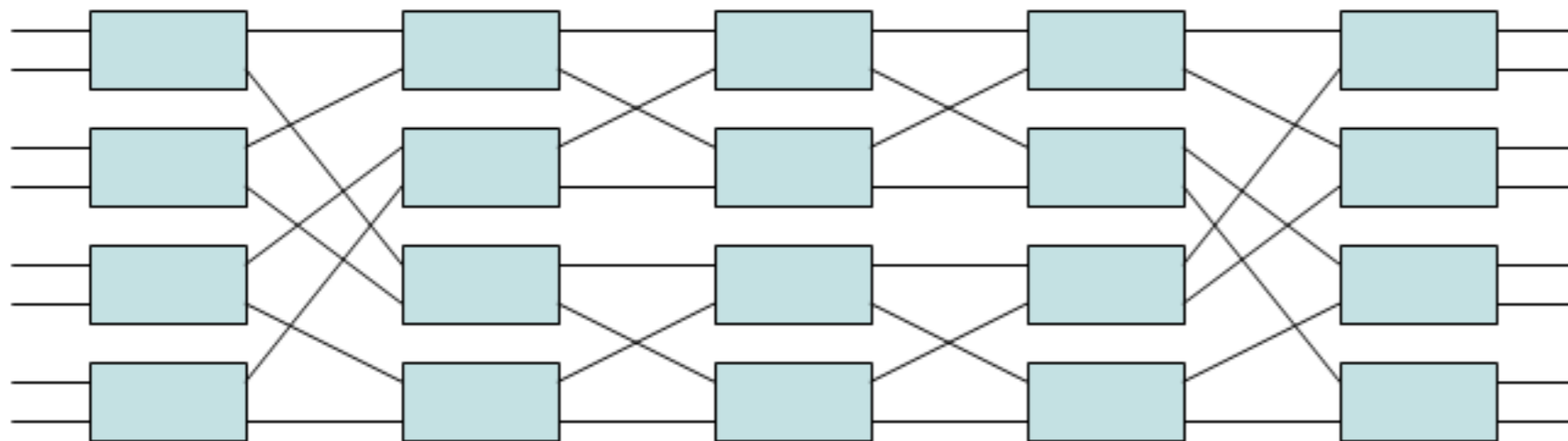
INTRODUCTION

The impact of r

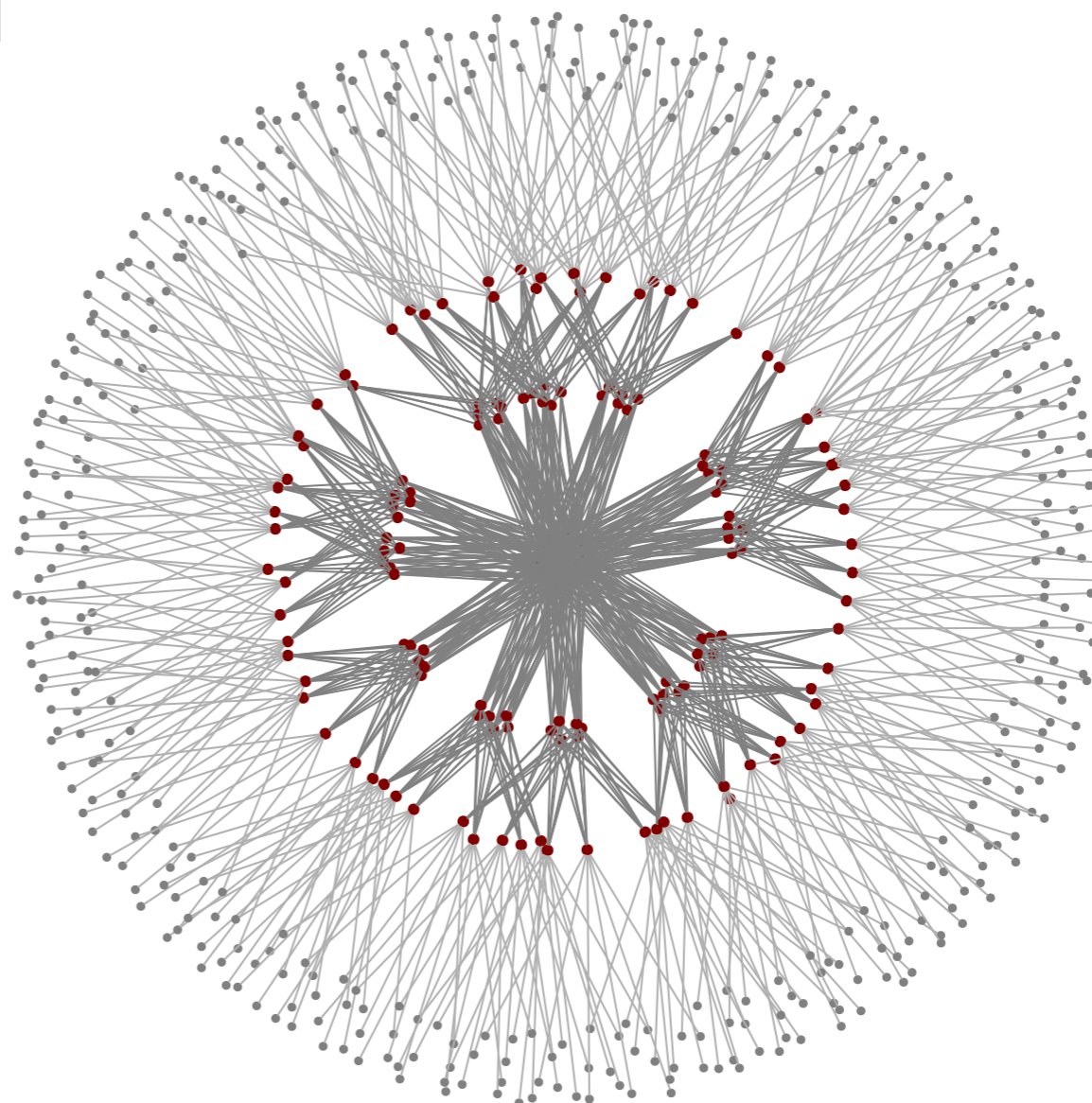


NUMBER OF CROSSPOINTS = $6N^{3/2} - 3N$ (1188 CROSSPOINTS WHEN $N = 36$)

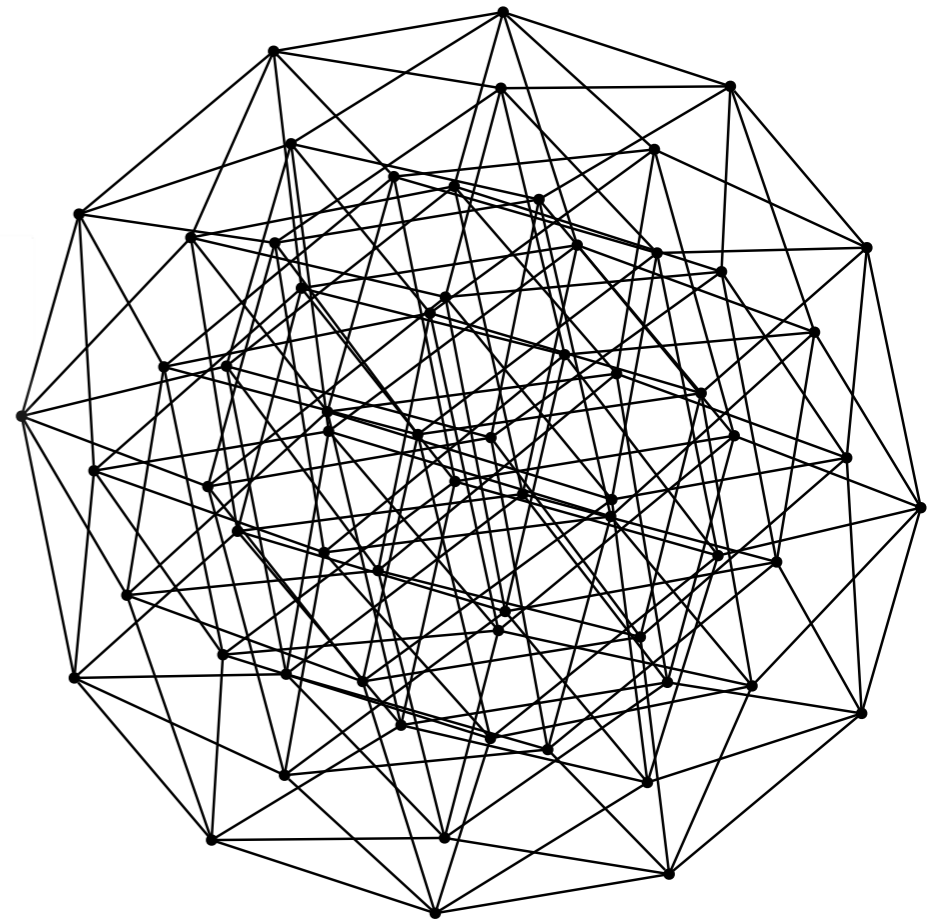
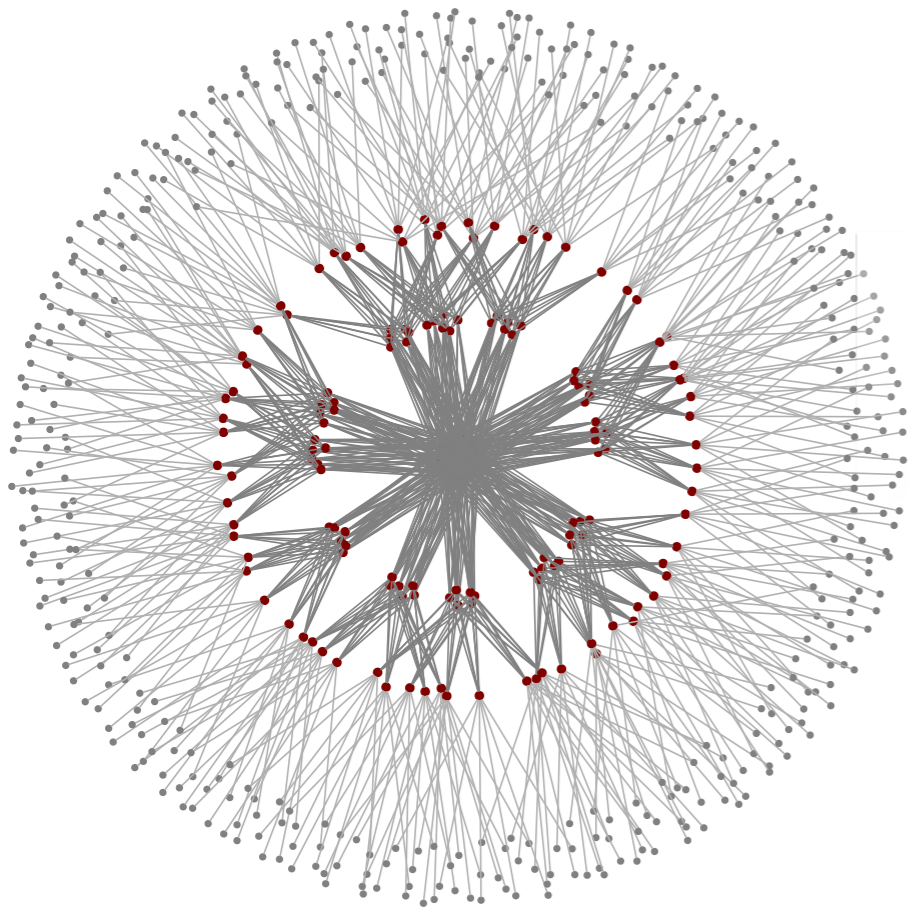
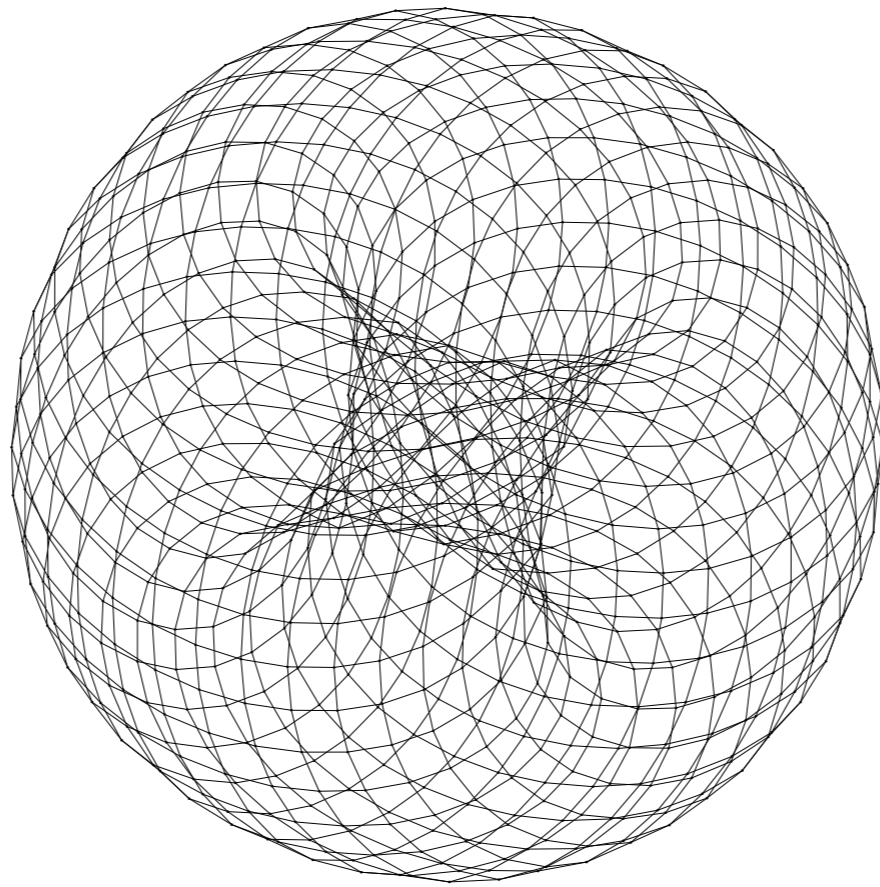
Fig. 2 — Three-stage switching array.



[Benes network: Wikipedia user Piggly]







What's different about data centers

Flexible forwarding

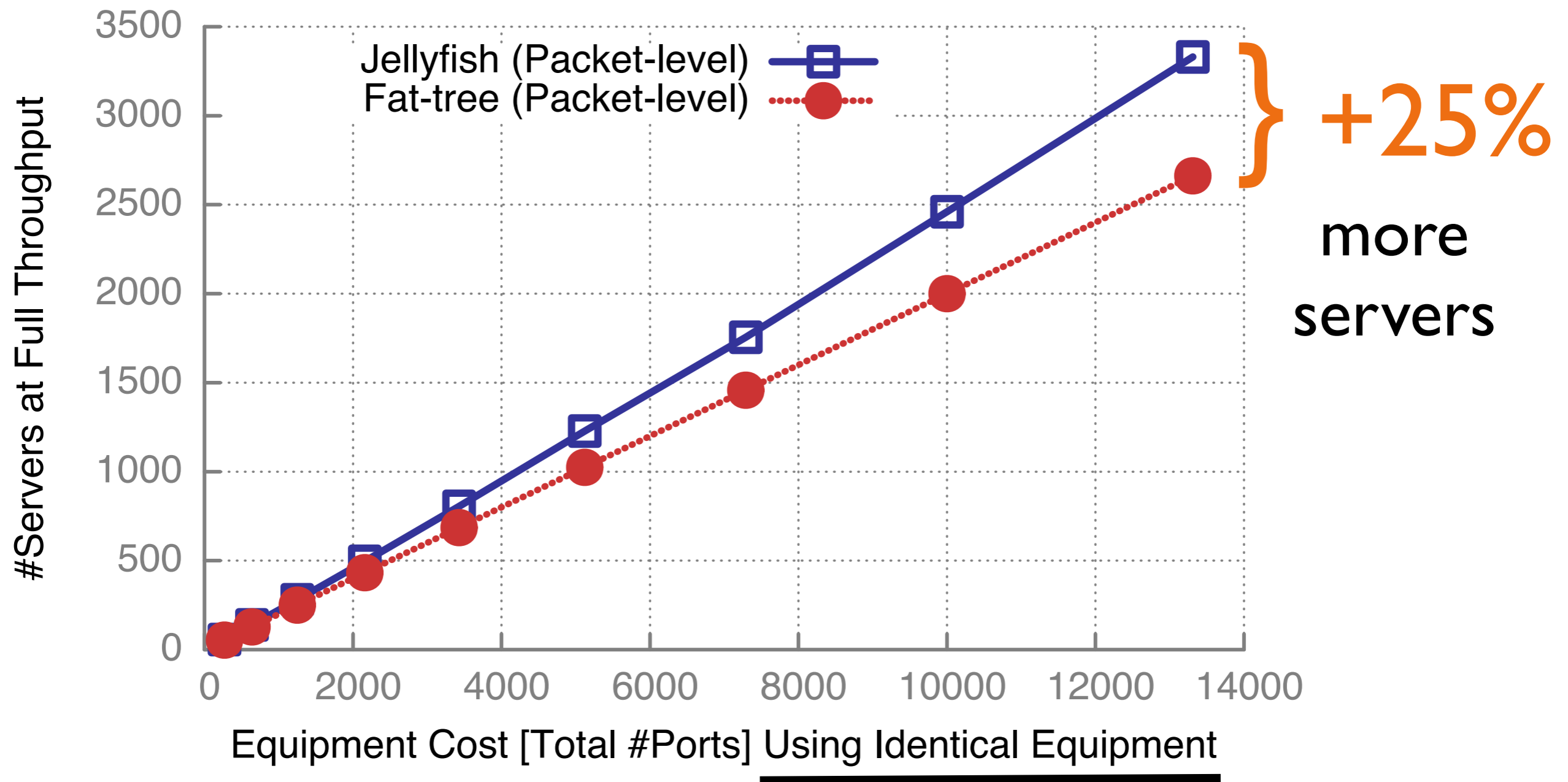
(compared with supercomputers)

Flexible routing & congestion control

(especially with software-defined networking)

Understanding Throughput

Throughput: Jellyfish vs. fat tree



Intuition

if we **fully utilize** all available capacity ...

$$\# \text{ 1 Gbps flows} = \frac{\text{total capacity}}{\text{used capacity per flow}}$$

Intuition

if we **fully utilize** all available capacity ...

$$\# \text{ 1 Gbps flows} = \frac{\sum_{\text{links}} \text{capacity}(\text{link})}{\text{used capacity per flow}}$$

Intuition

if we **fully utilize** all available capacity ...

$$\# \text{ 1 Gbps flows} = \frac{\sum_{\text{links}} \text{capacity}(\text{link})}{1 \text{ Gbps} \cdot \text{mean path length}}$$

Intuition

if we **fully utilize** all available capacity ...

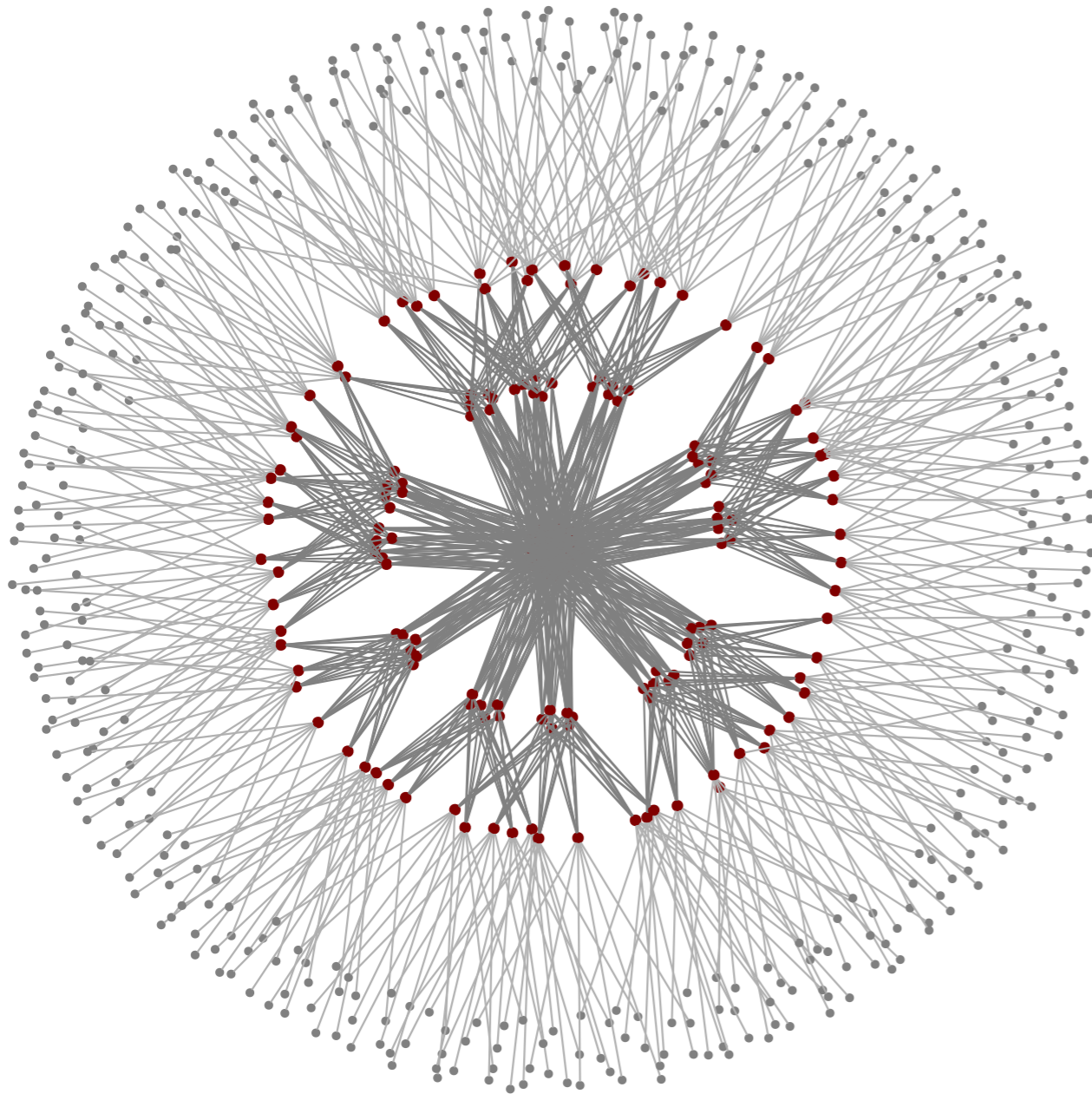
$$\# \text{ 1 Gbps flows} = \frac{\sum_{\text{links}} \text{capacity}(\text{link})}{1 \text{ Gbps} \cdot \text{mean path length}}$$

Mission:

minimize average path length

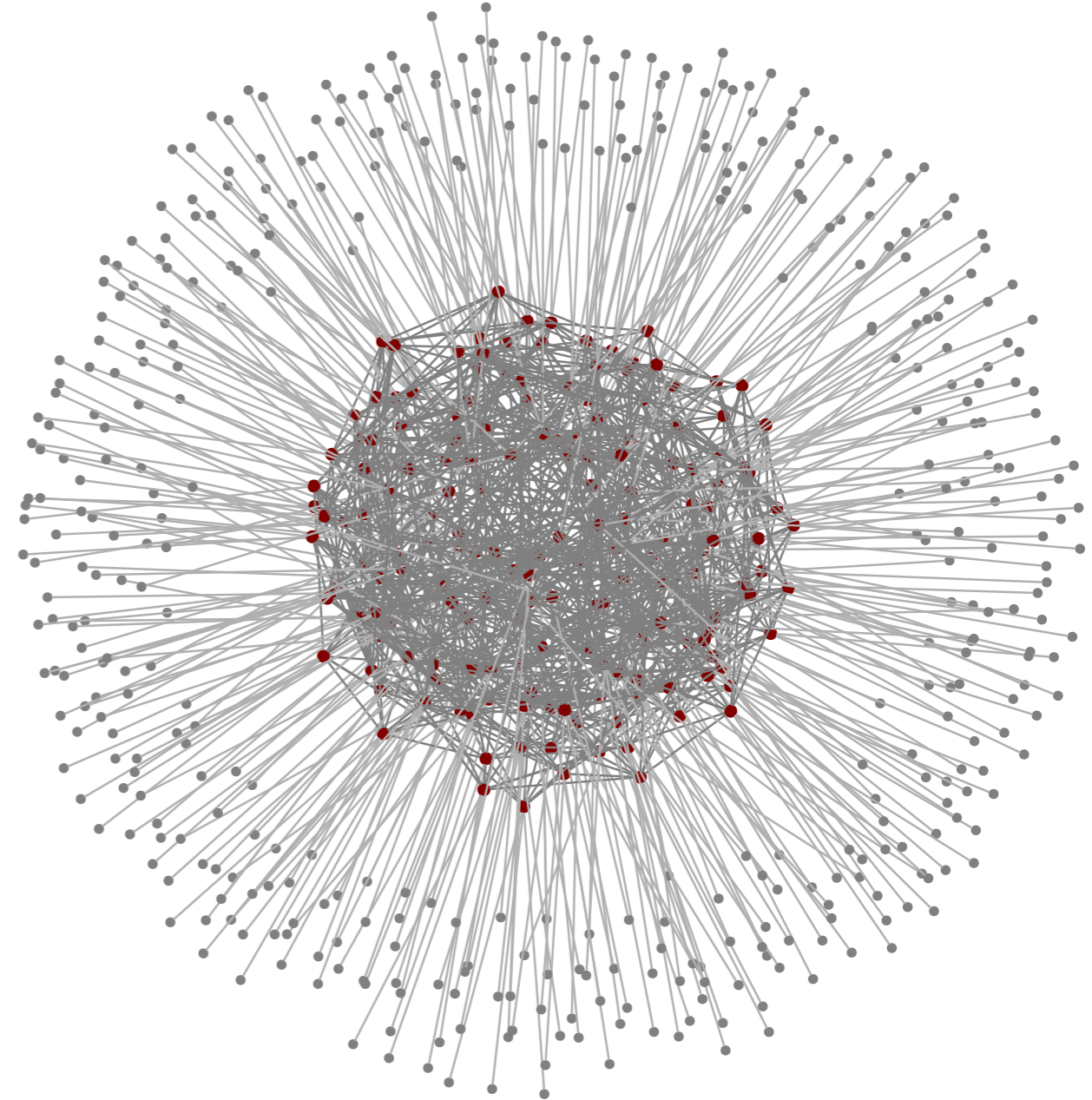
WWWIWSG SΛELS8G bscu I6U8cu

Example



Fat tree

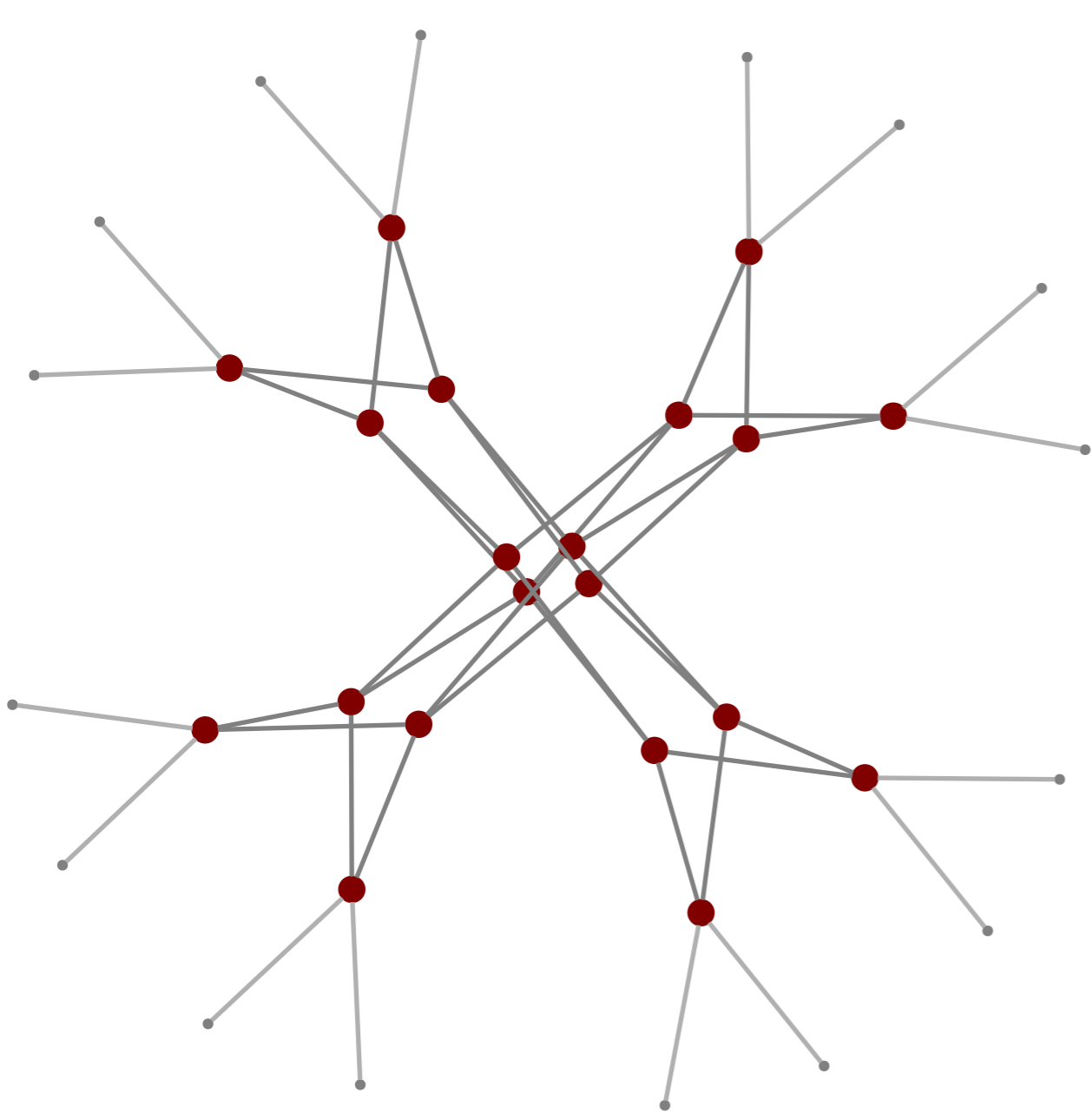
432 servers, 180 switches, degree 12



Jellyfish random graph

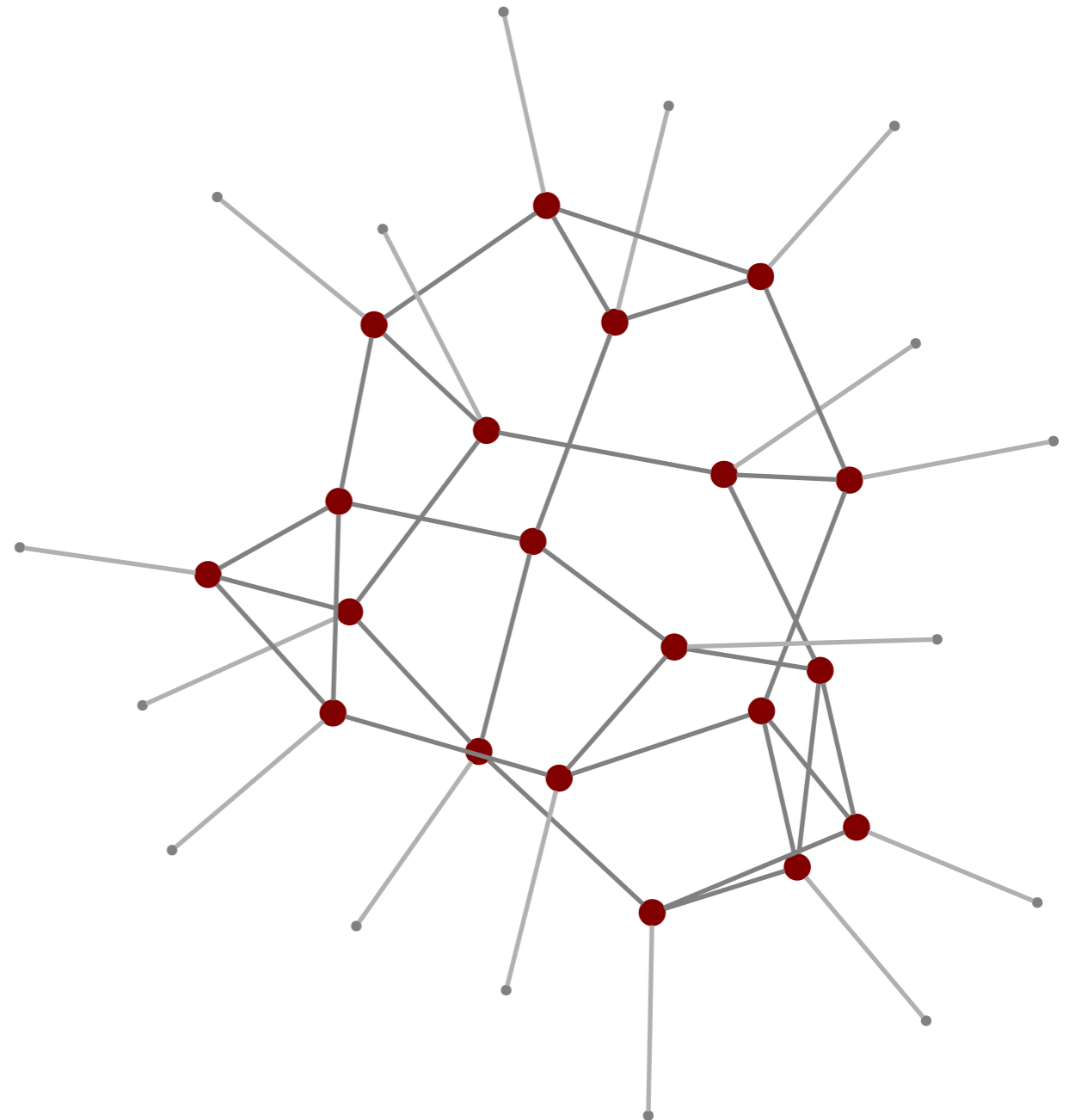
432 servers, 180 switches, degree 12

Example



Fat tree

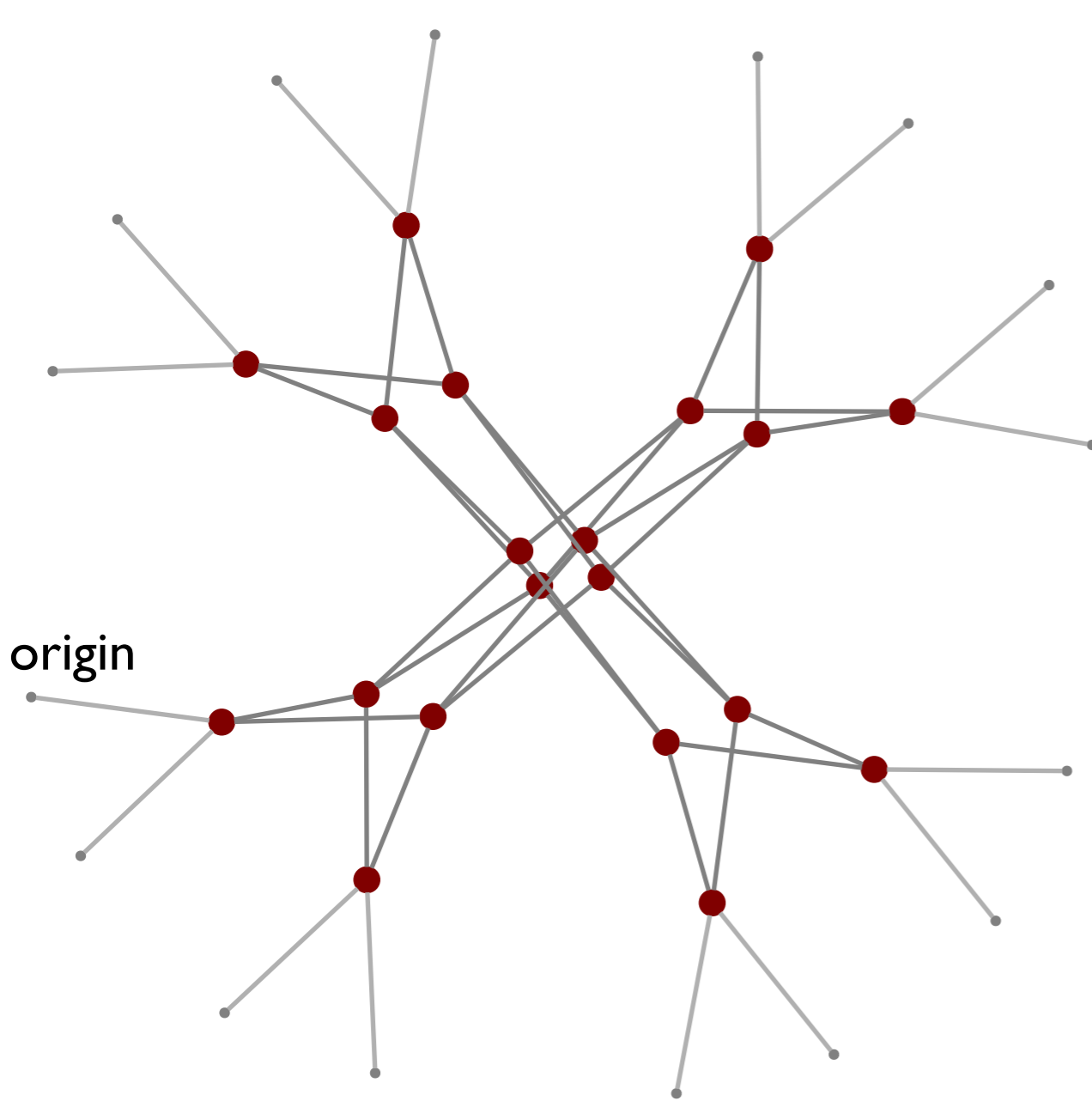
16 servers, 20 switches, degree 4



Jellyfish random graph

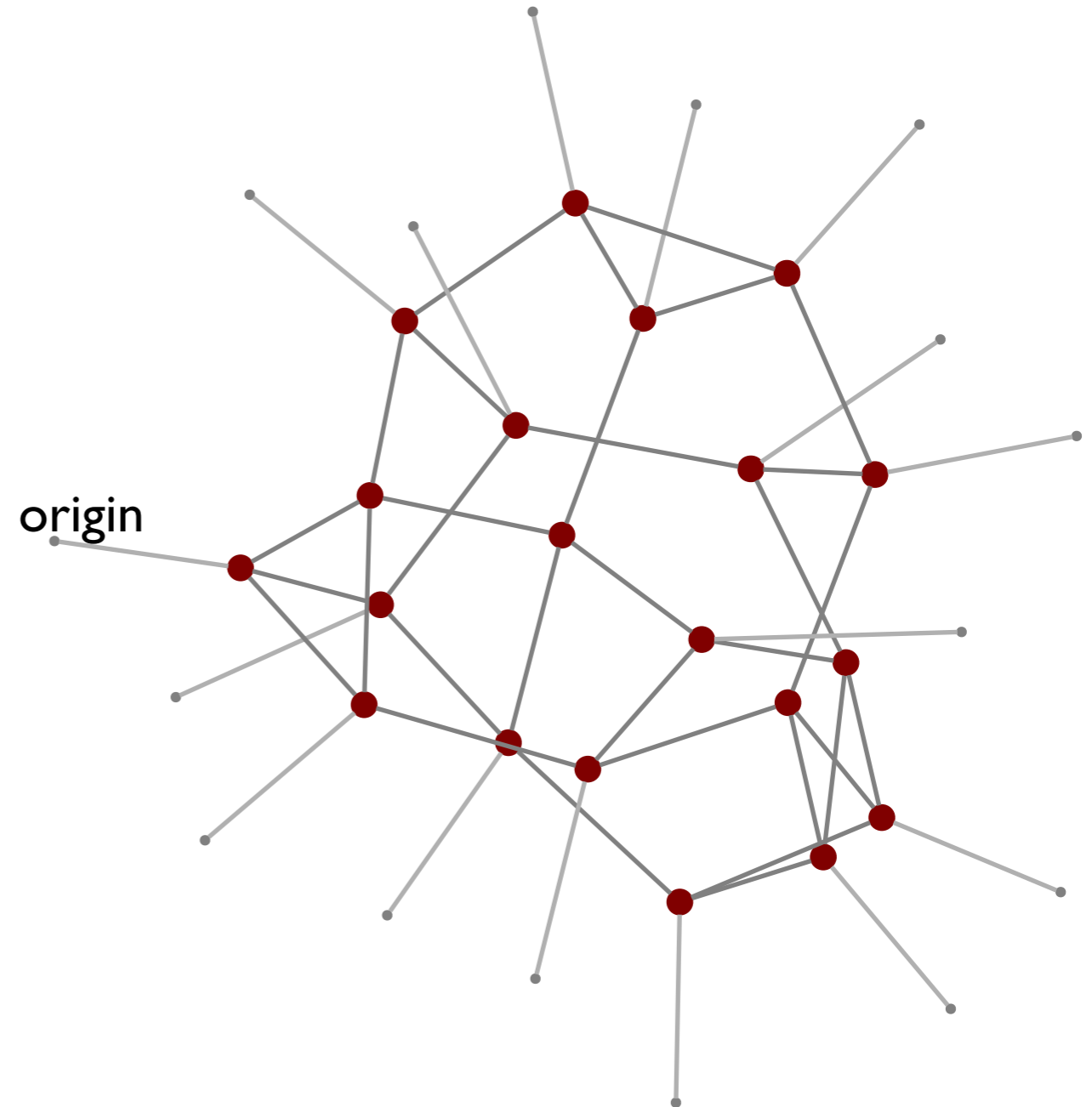
16 servers, 20 switches, degree 4

Example



Fat tree

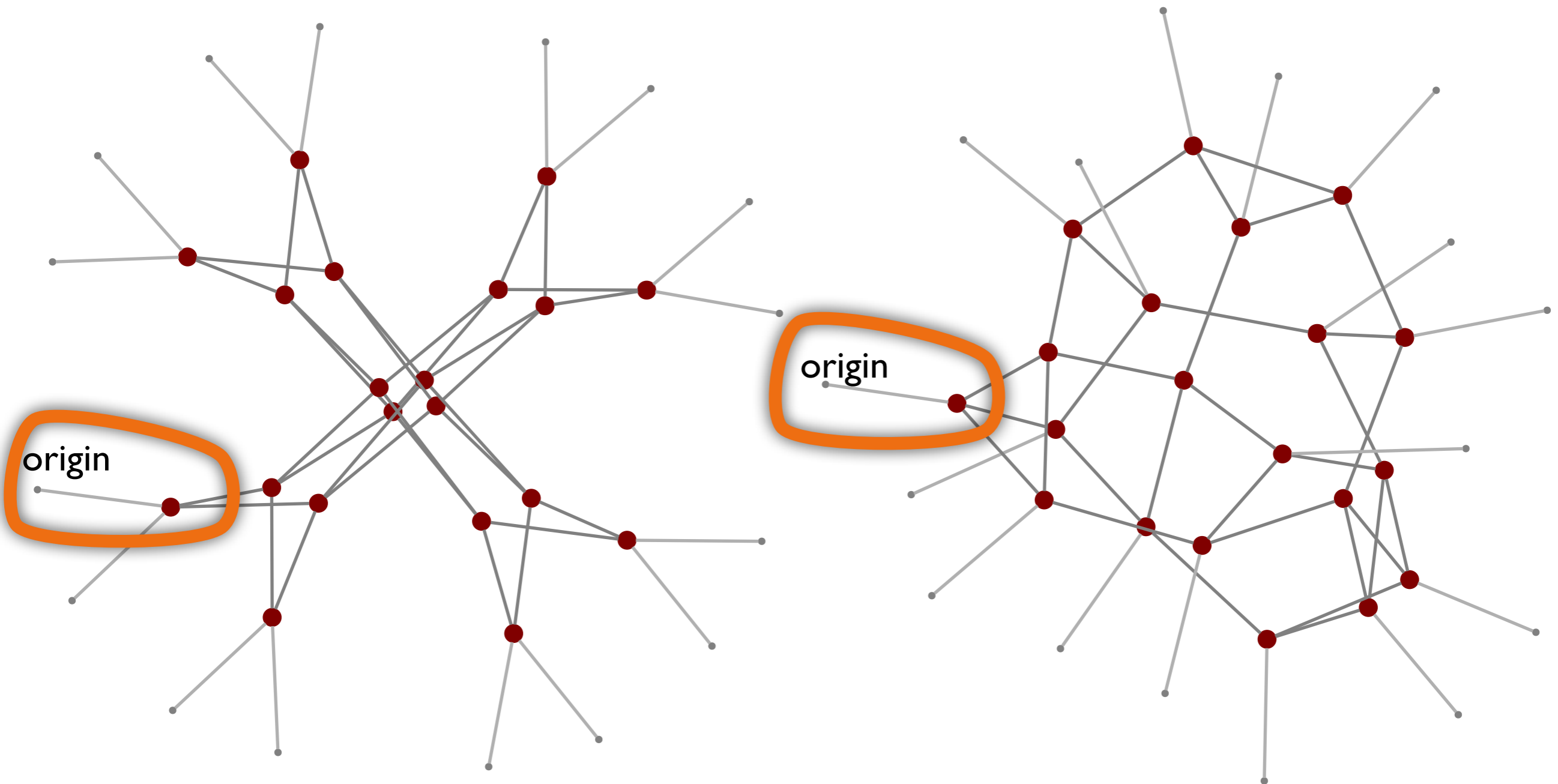
16 servers, 20 switches, degree 4



Jellyfish random graph

16 servers, 20 switches, degree 4

Example



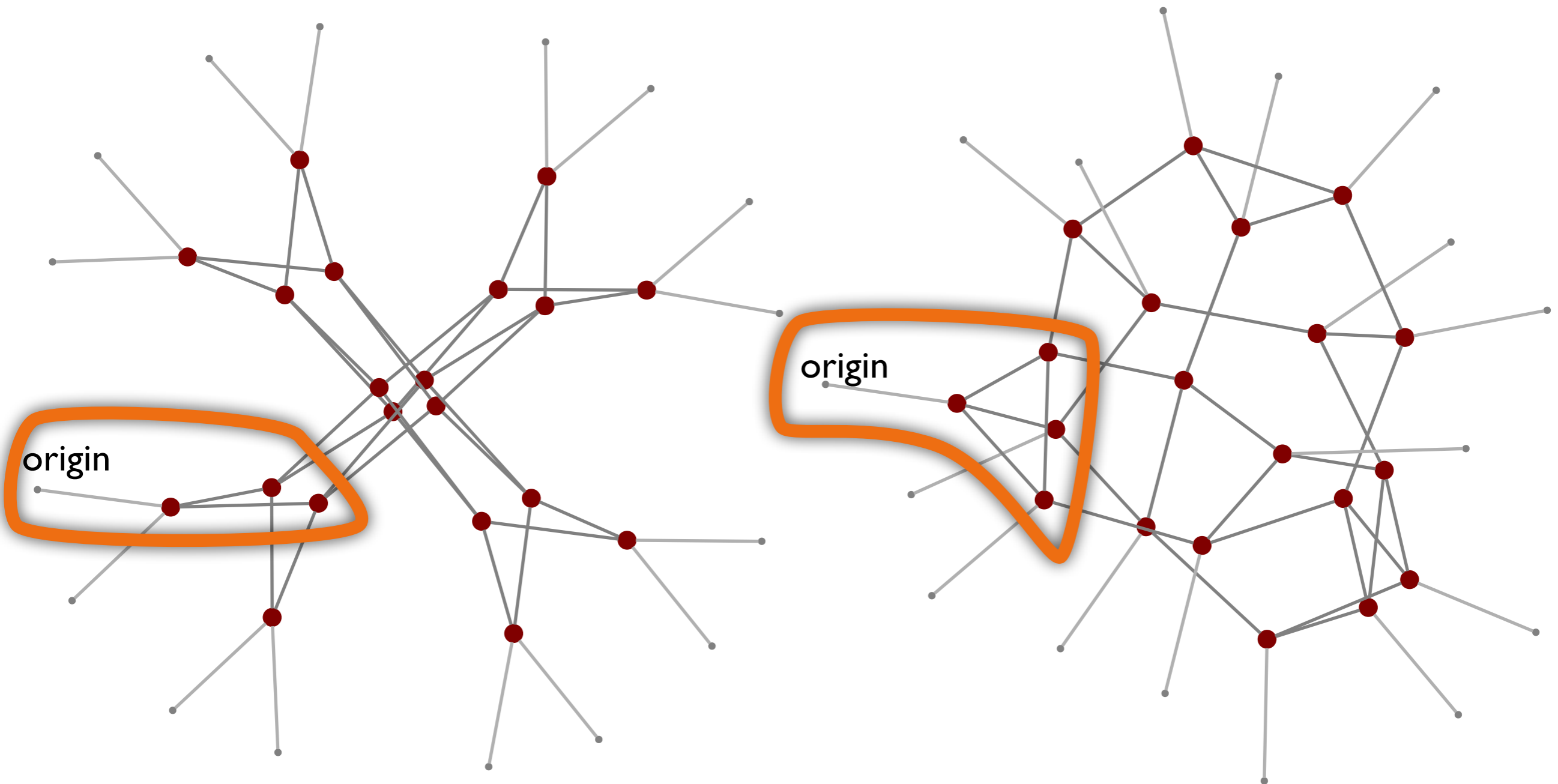
Fat tree

16 servers, 20 switches, degree 4

Jellyfish random graph

16 servers, 20 switches, degree 4

Example



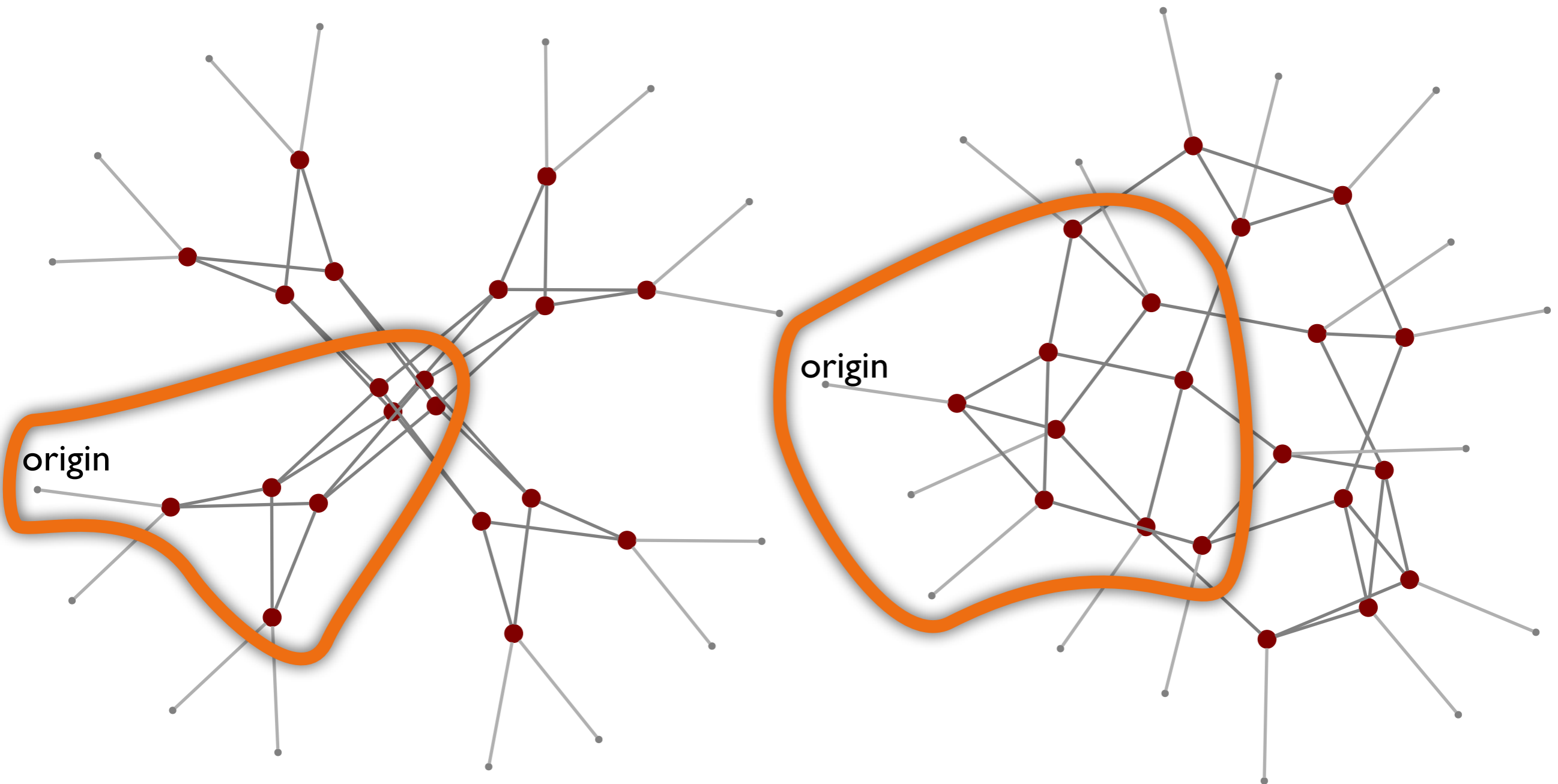
Fat tree

16 servers, 20 switches, degree 4

Jellyfish random graph

16 servers, 20 switches, degree 4

Example



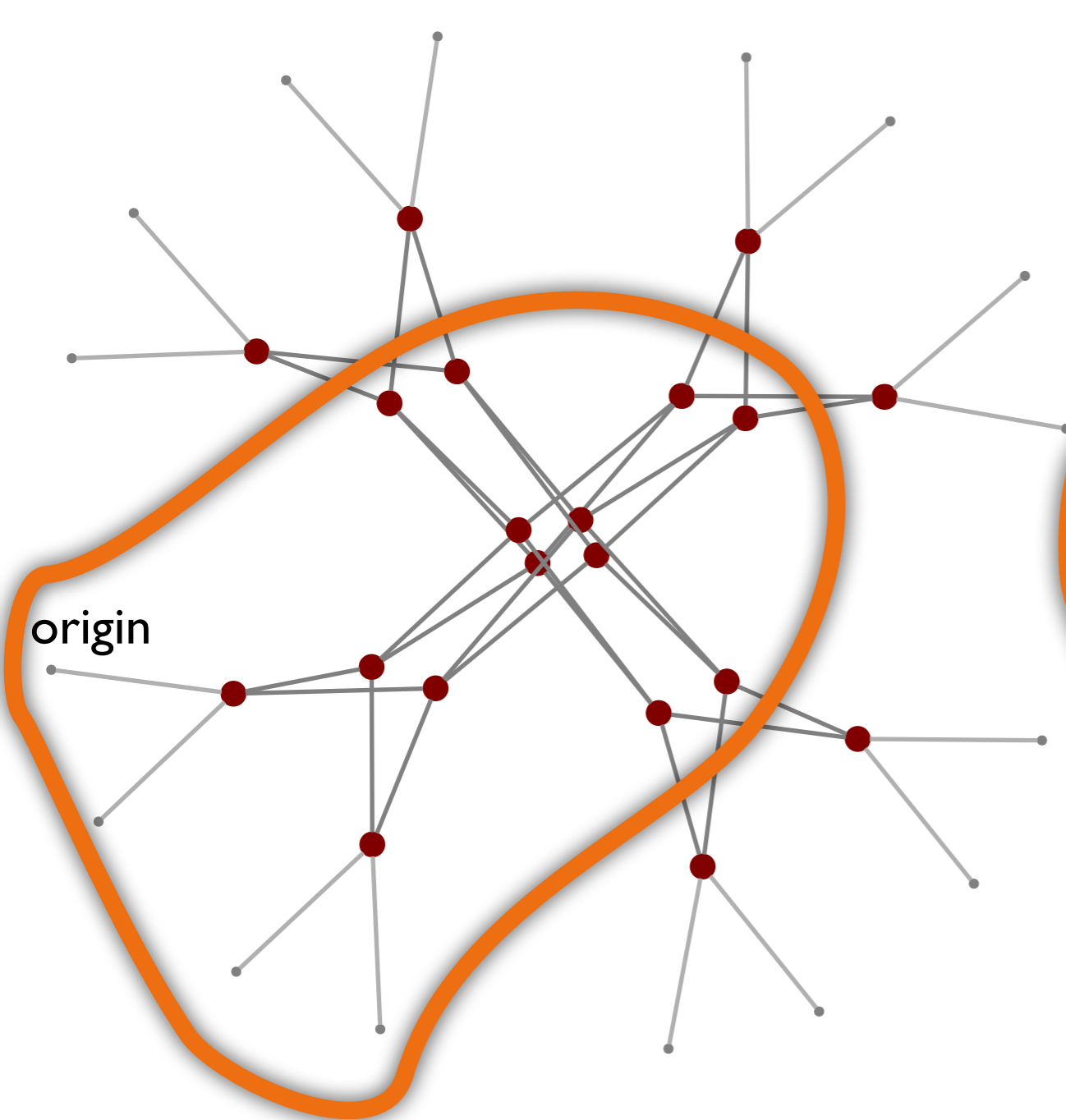
Fat tree

16 servers, 20 switches, degree 4

Jellyfish random graph

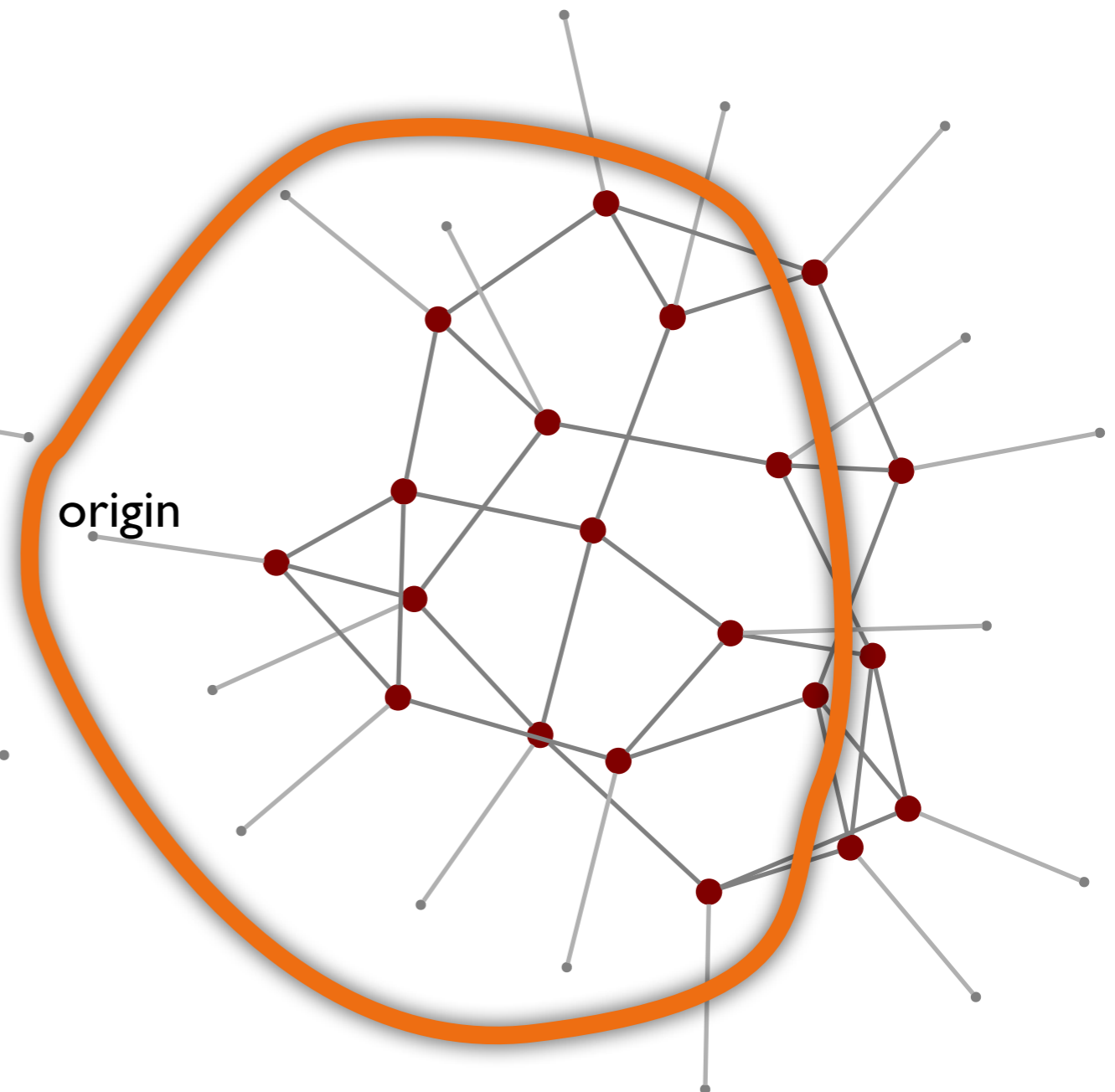
16 servers, 20 switches, degree 4

Example



Fat tree

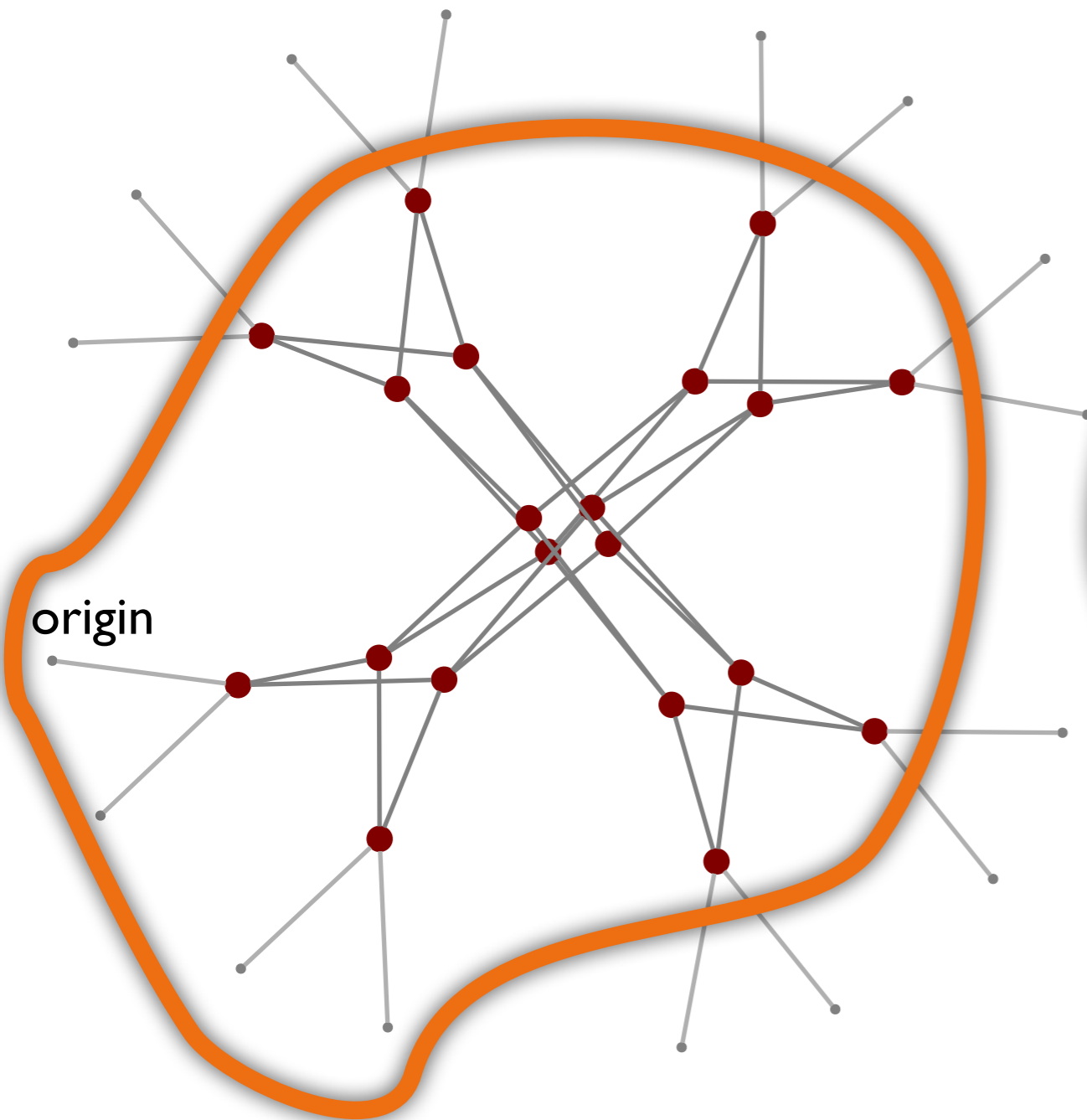
16 servers, 20 switches, degree 4



Jellyfish random graph

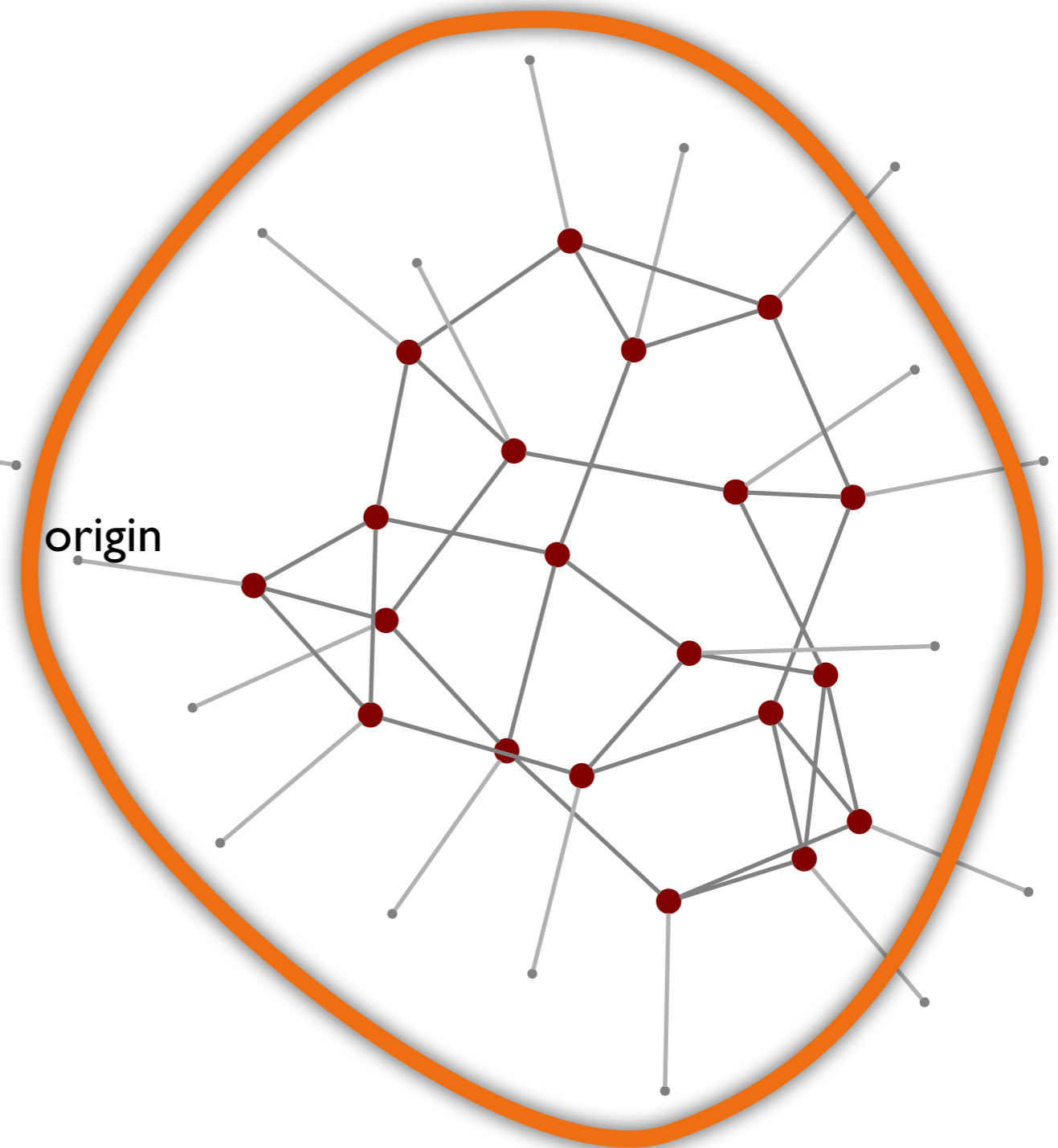
16 servers, 20 switches, degree 4

Example



Fat tree

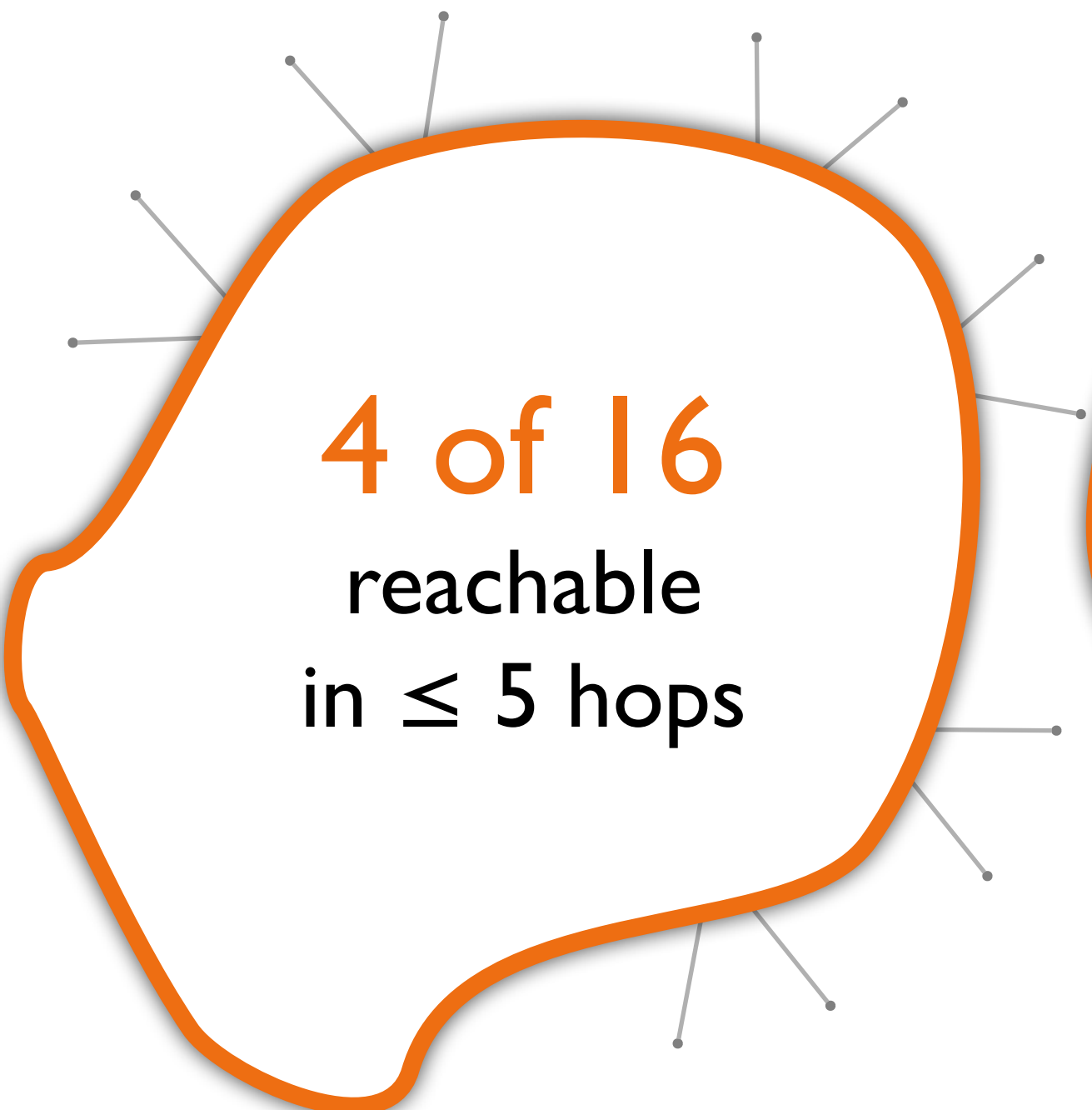
16 servers, 20 switches, degree 4



Jellyfish random graph

16 servers, 20 switches, degree 4

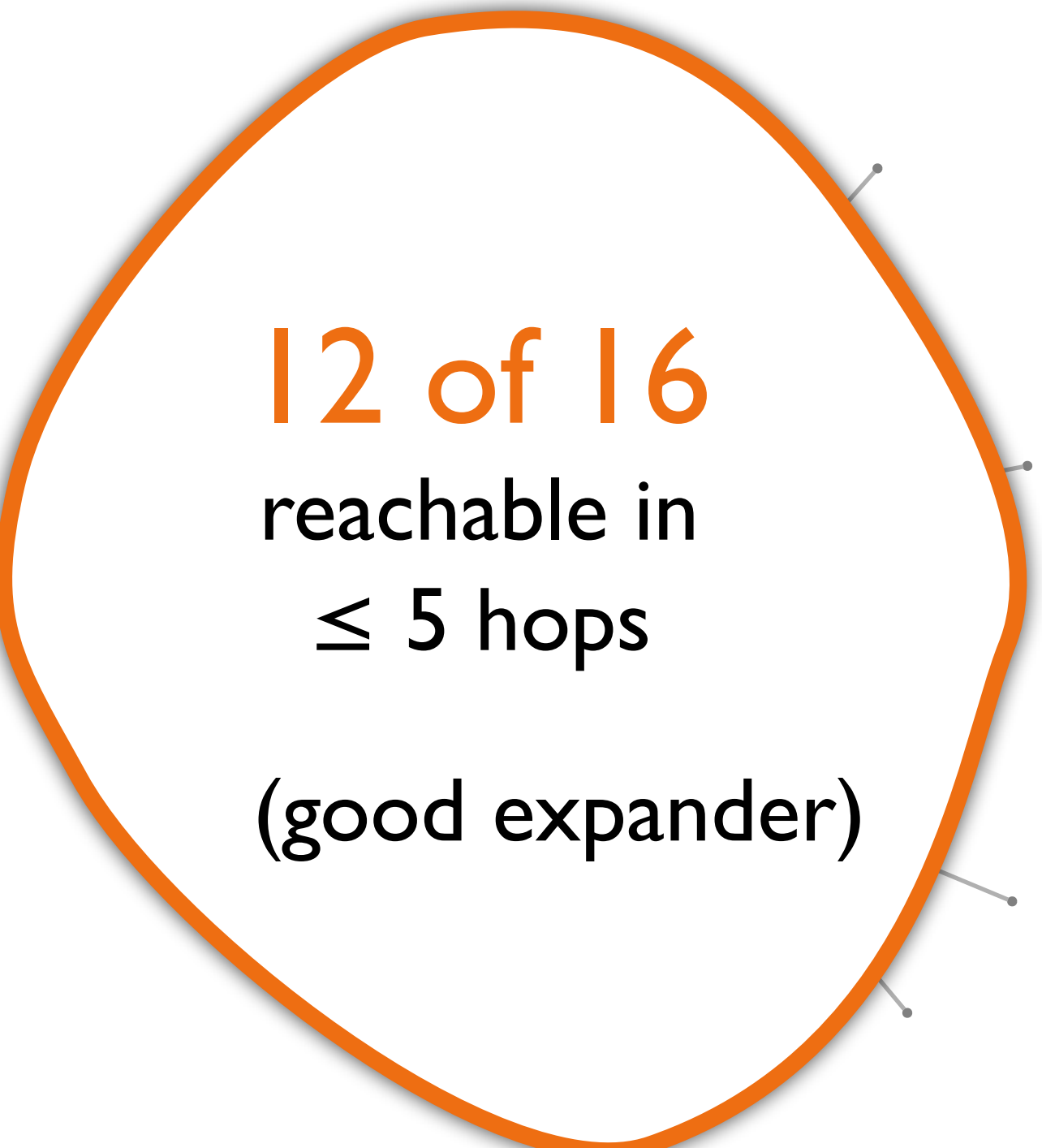
Example



4 of 16
reachable
in ≤ 5 hops

Fat tree

16 servers, 20 switches, degree 4

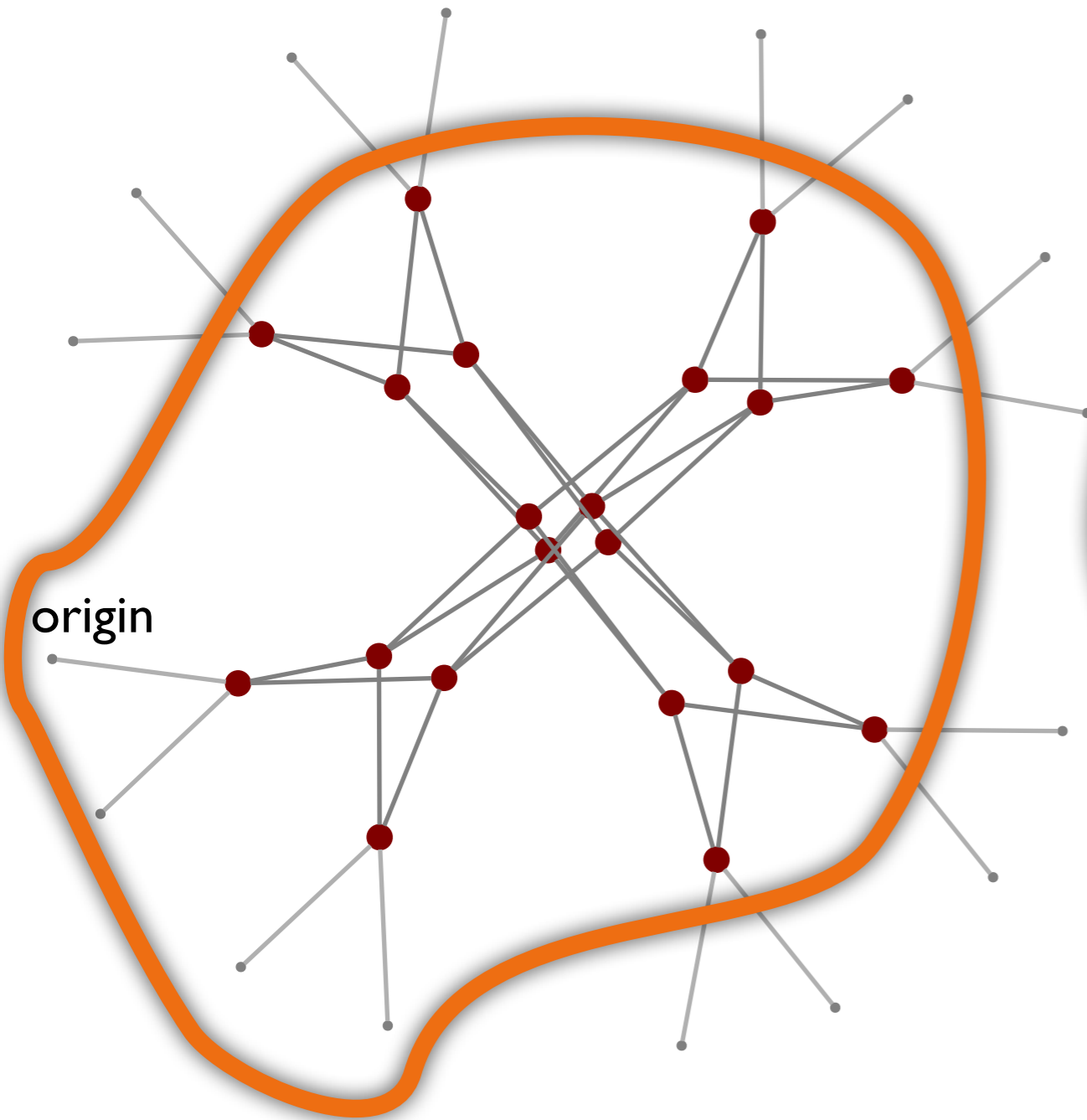


12 of 16
reachable in
 ≤ 5 hops
(good expander)

Jellyfish random graph

16 servers, 20 switches, degree 4

Example



Fat tree

16 servers, 20 switches, degree 4



12 of 16

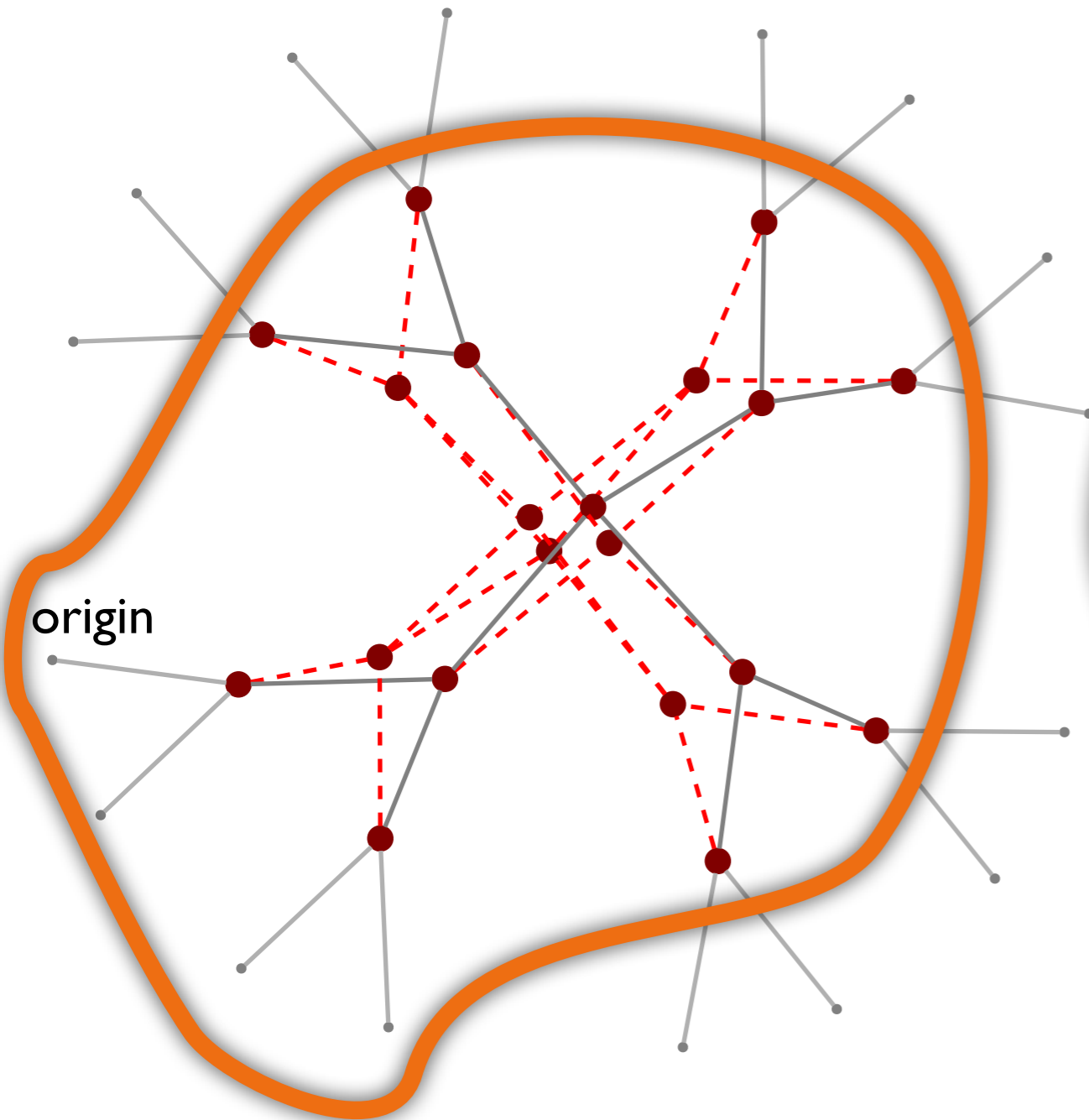
reachable in
 ≤ 5 hops

(good expander)

Jellyfish random graph

16 servers, 20 switches, degree 4

Example



Fat tree

16 servers, 20 switches, degree 4



12 of 16

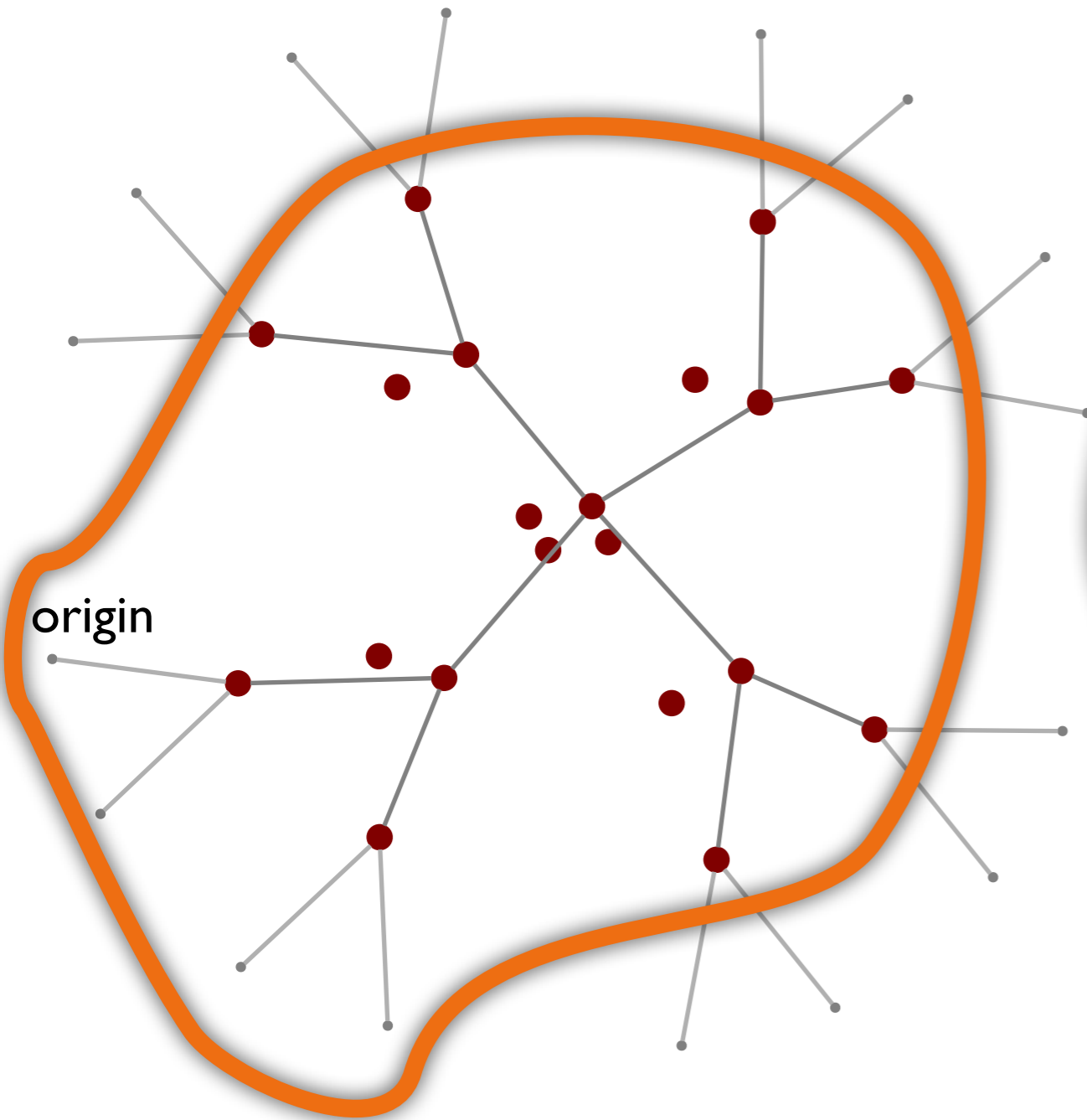
reachable in
 ≤ 5 hops

(good expander)

Jellyfish random graph

16 servers, 20 switches, degree 4

Example



Fat tree

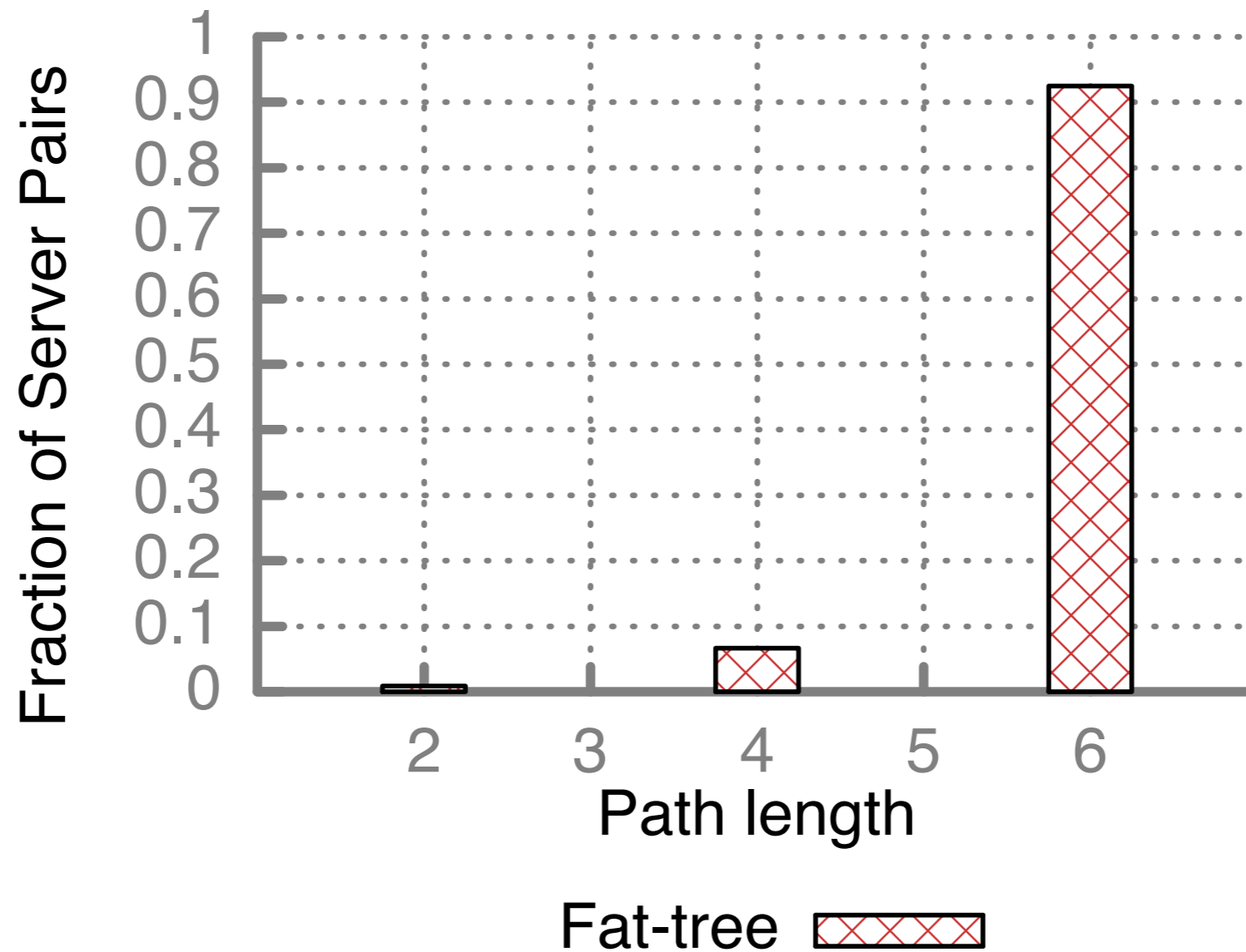
16 servers, 20 switches, degree 4



Jellyfish random graph

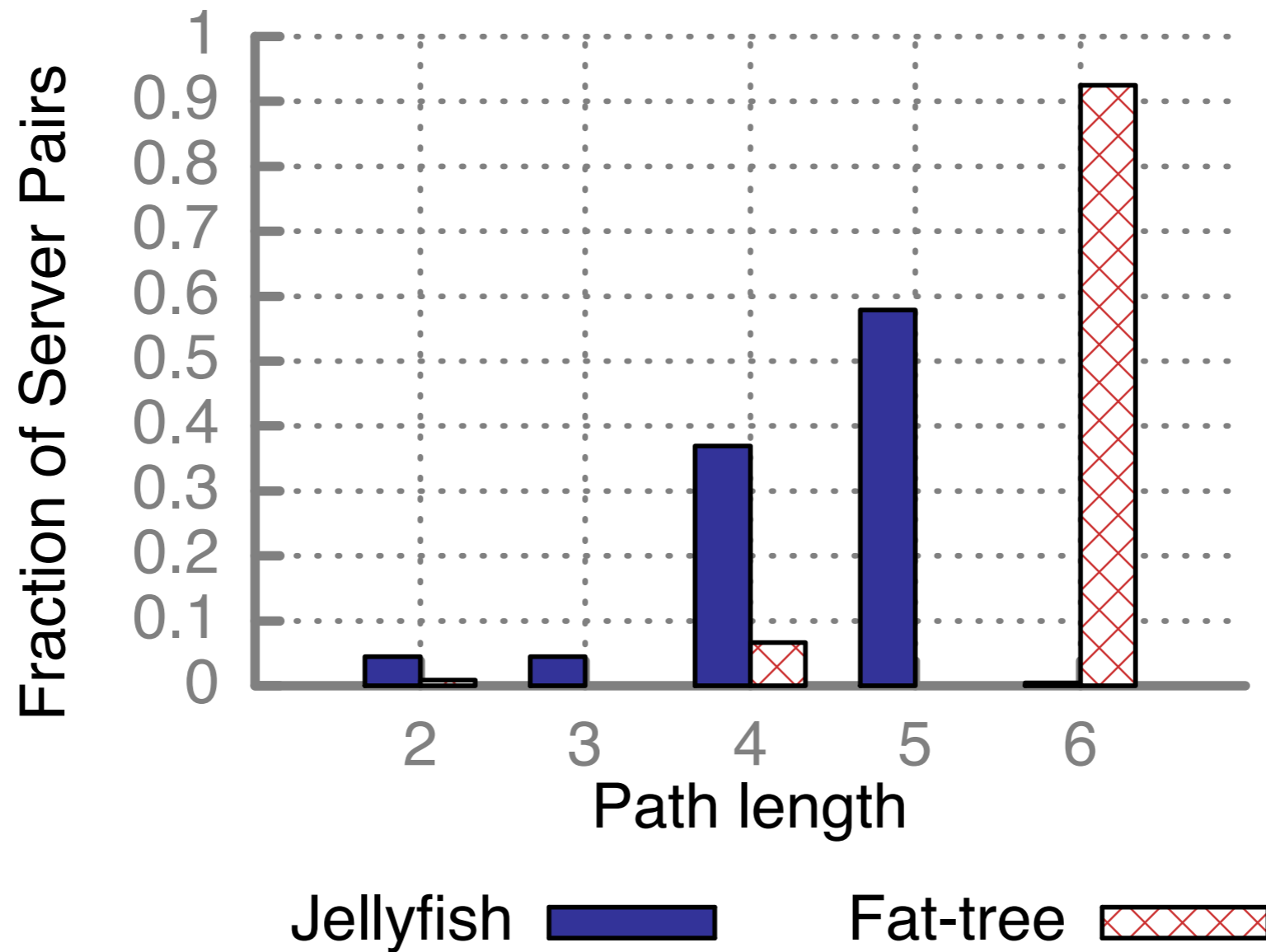
16 servers, 20 switches, degree 4

Jellyfish has short paths



Fat-tree with 686 servers

Jellyfish has short paths



Jellyfish, same equipment

System Design:

Performance Consistency

Is performance more variable?

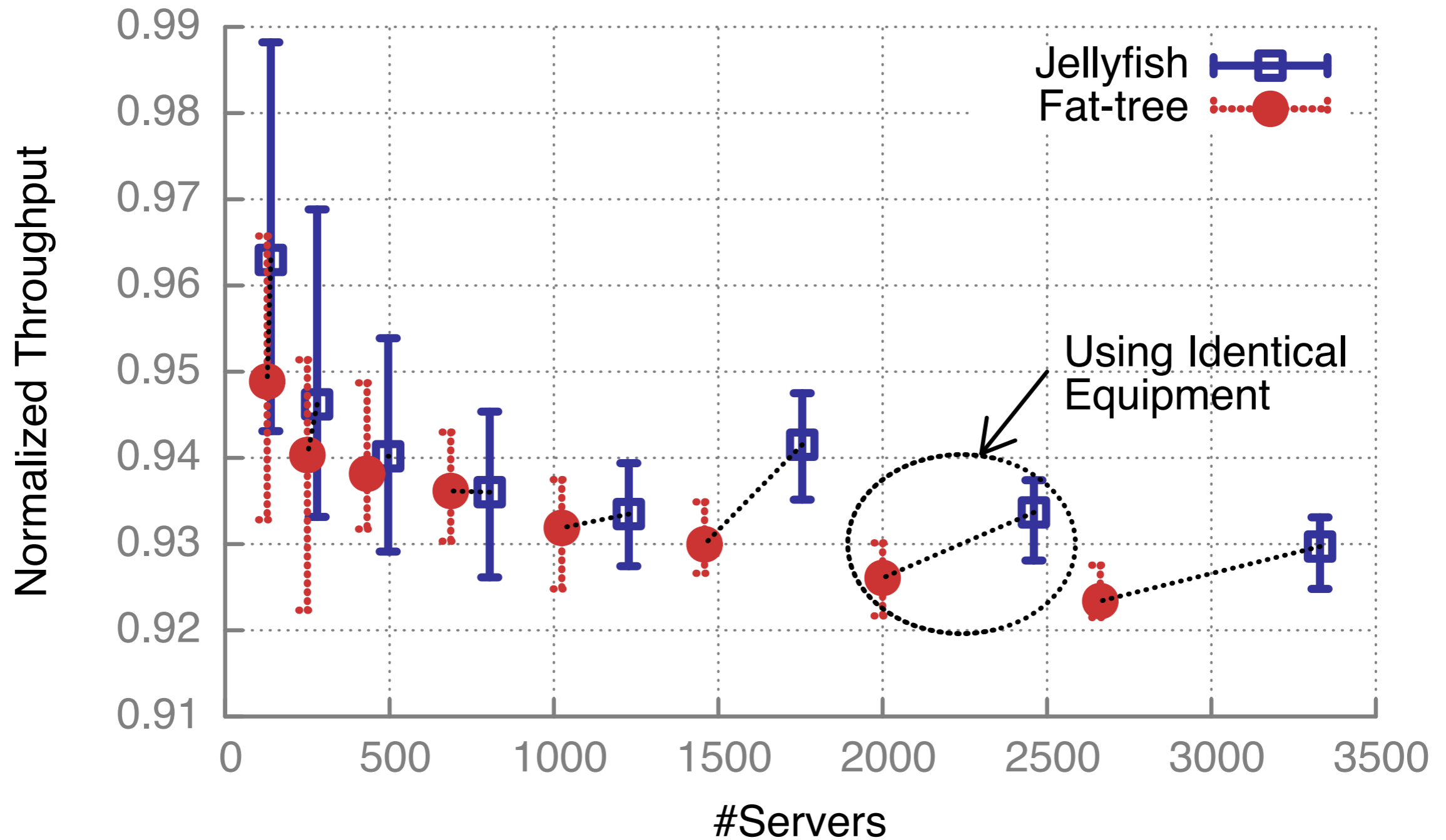
Performance depends on choice of random graph

- if you expand the network, would performance change dramatically?

Extreme case: graph could be disconnected!

- never happens, with high probability

Little variation if size is moderate



{min, avg, max} of 20 trials shown

System Design:

Routing

Routing

Intuition

if we fully utilize all available capacity ...

$$\# \text{ 1 Gbps flows} = \frac{\text{total capacity}}{\text{used capacity per flow}}$$

How do we effectively utilize capacity without structure?

Routing without structure

In theory, just a multicommodity flow (MCF) problem

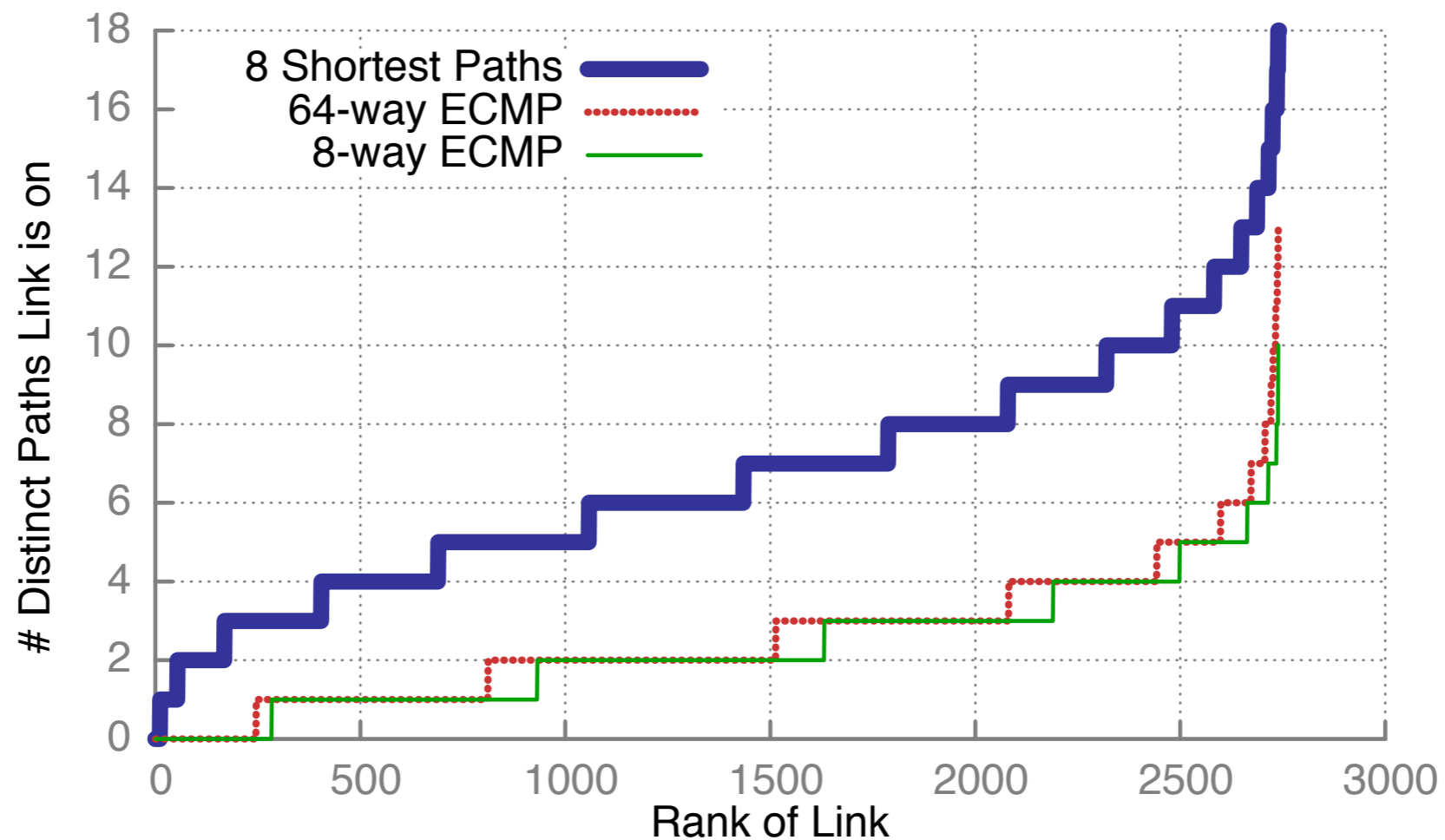
Potential issues:

- Solve MCF using a distributed protocol?
- Optimal solution could have too many small subflows

Routing

Does ECMP work?

- No
- ECMP doesn't use Jellyfish's path diversity

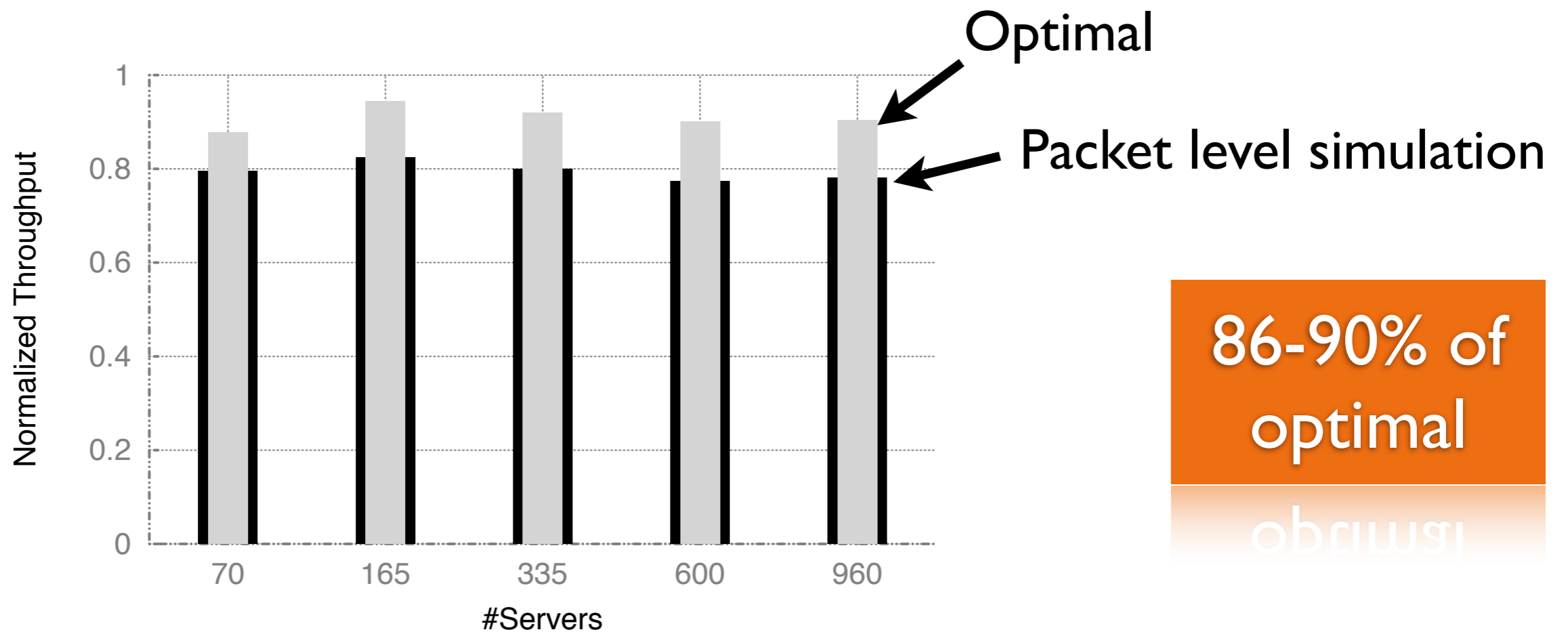


Routing: a simple solution

Find k shortest paths

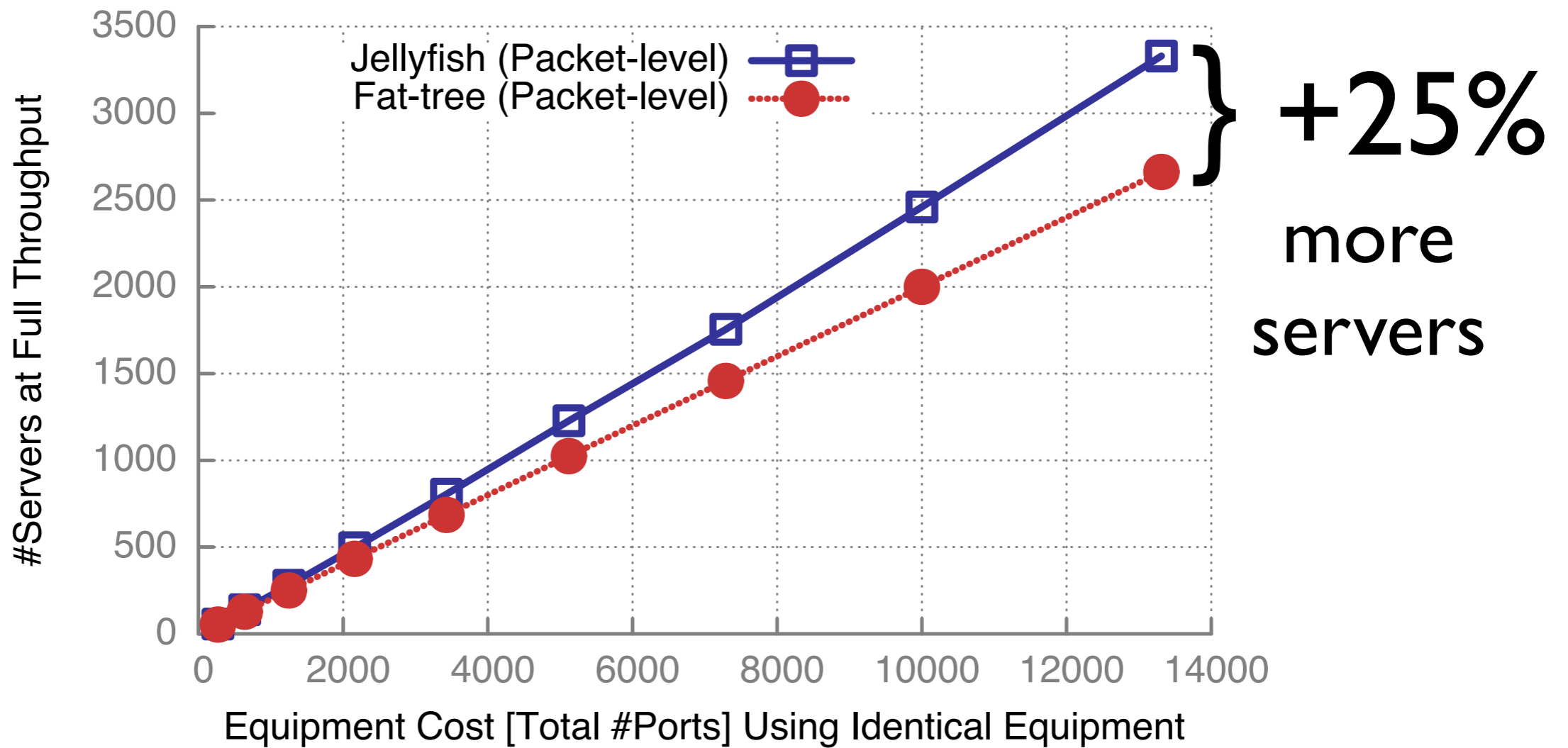
Let Multipath TCP do the rest

- [Wischik, Raiciu, Greenhalgh, Handley, NSDI'10]



(TCP is within 3 percentage points of MPTCP)

Throughput: Jellyfish vs. fat tree



8-shortest paths + MPTCP

Deploying k-shortest paths

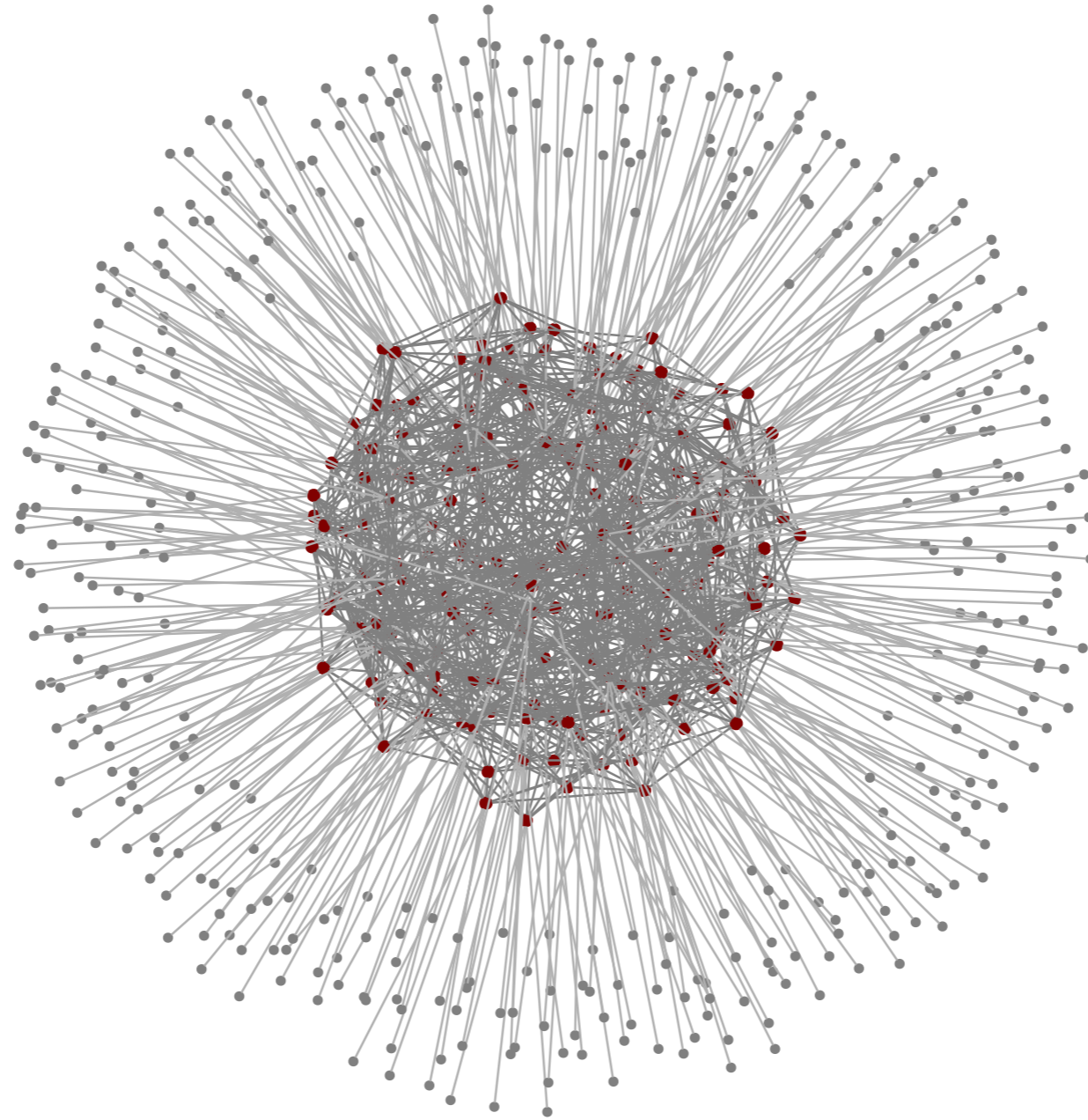
Multiple options:

- SPAIN [Mudigonda, Yalagandula, Al-Fares, Mogul, NSDI' 10]
- Equal-cost MPLS tunnels
- IBM Research's SPARTA [CoNEXT 2012]
- SDN controller based methods

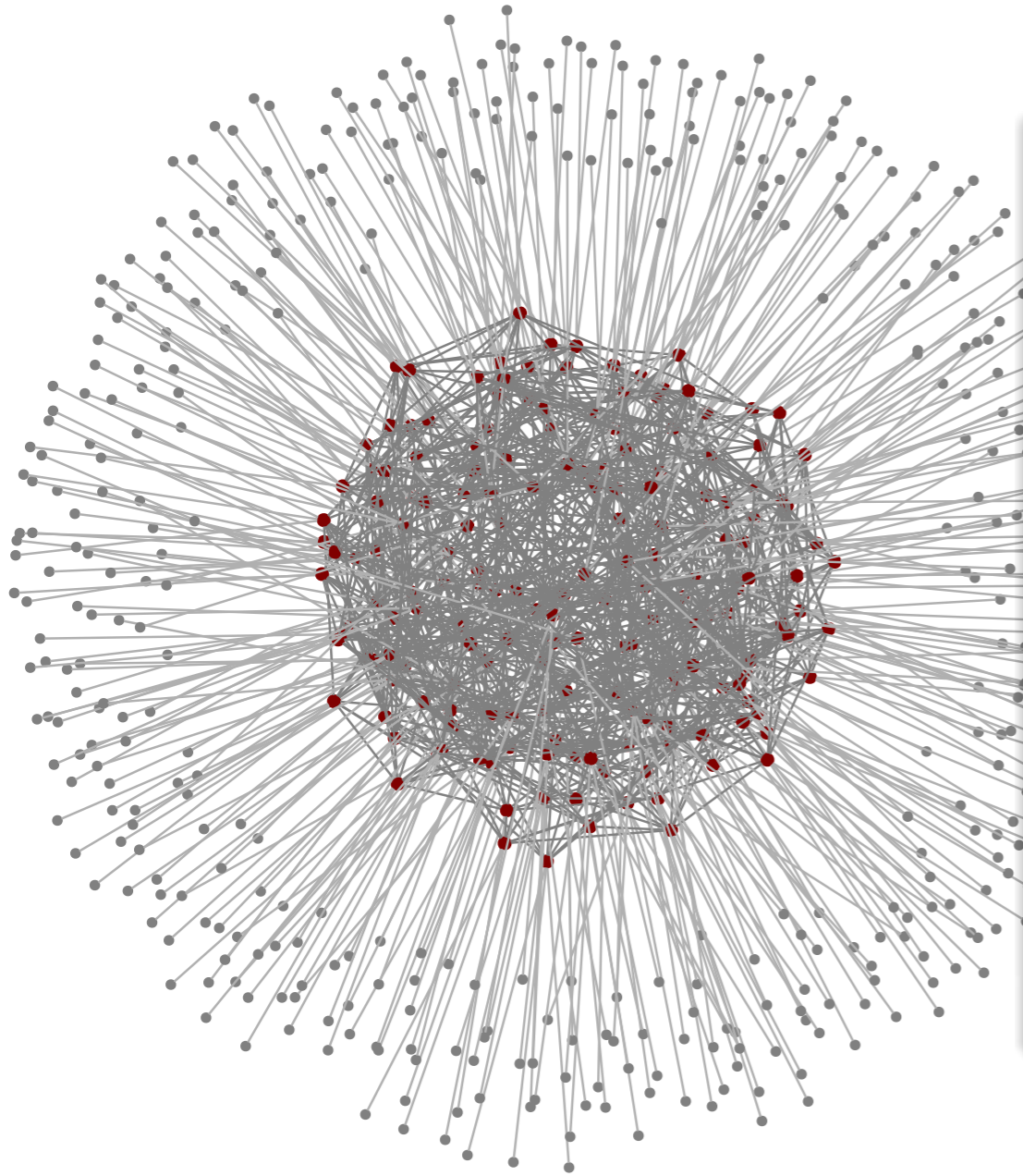
System Design:

Cabling

Cabling



Cabling

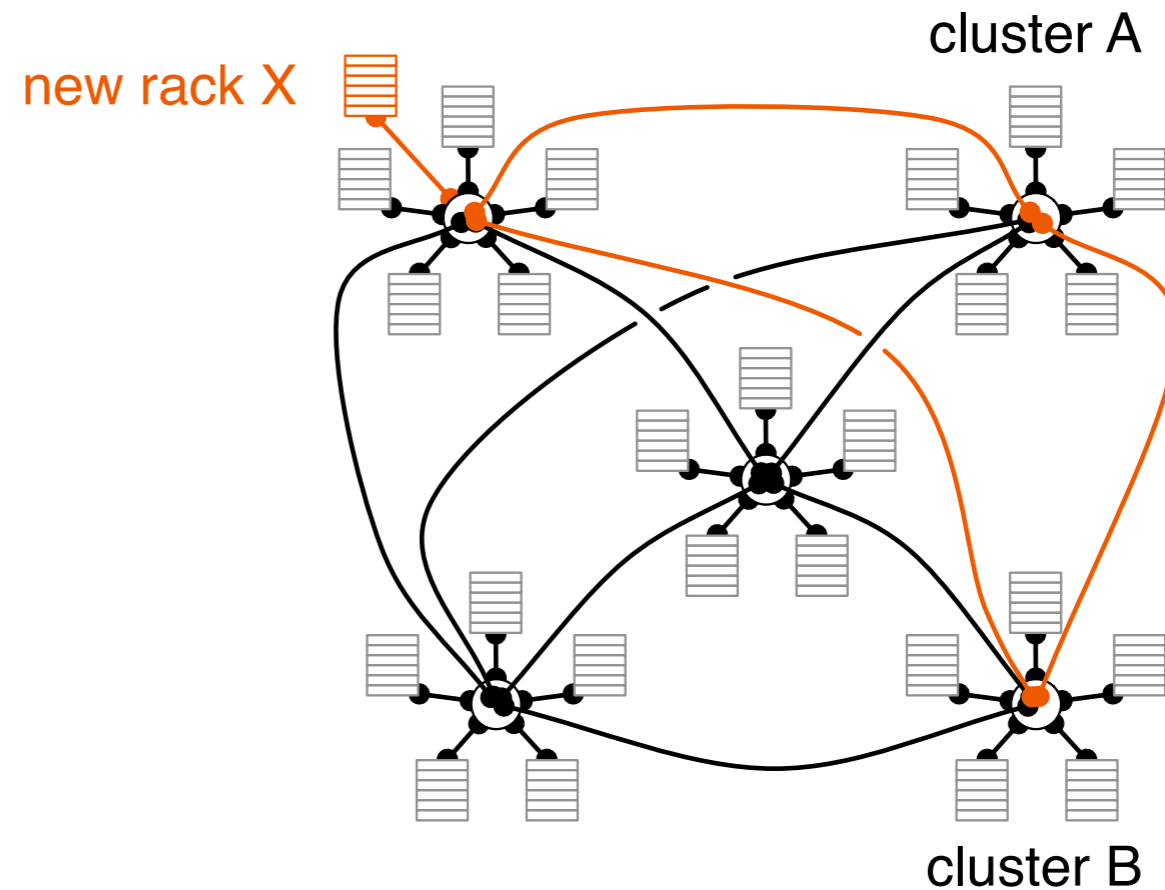


Cabling solutions

Fewer
cables

for same #
servers as
fat tree

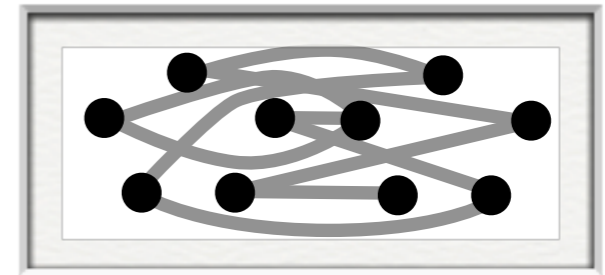
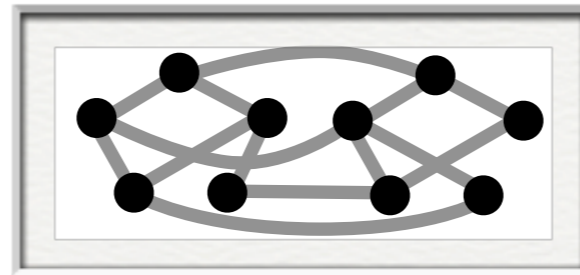
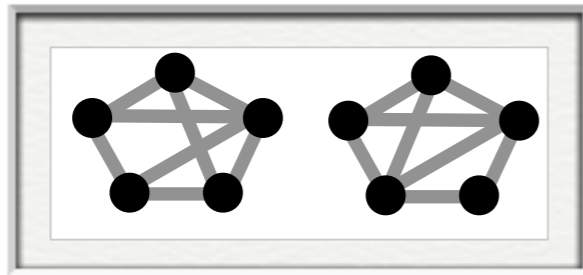
Aggregate
bundles



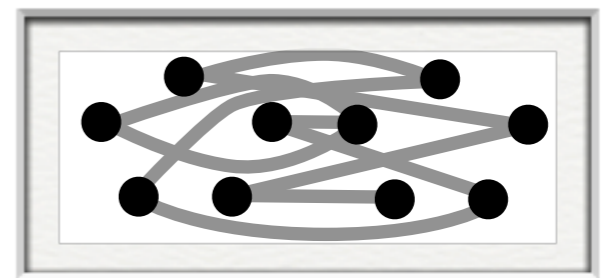
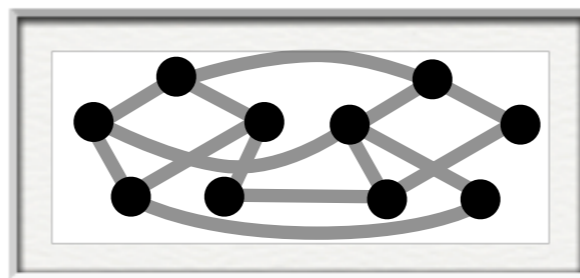
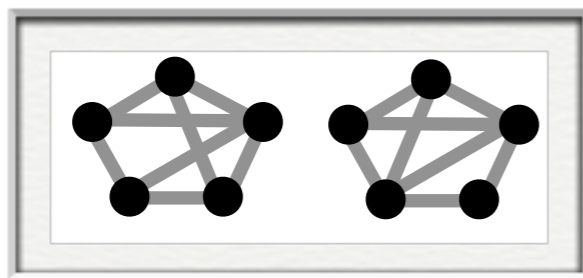
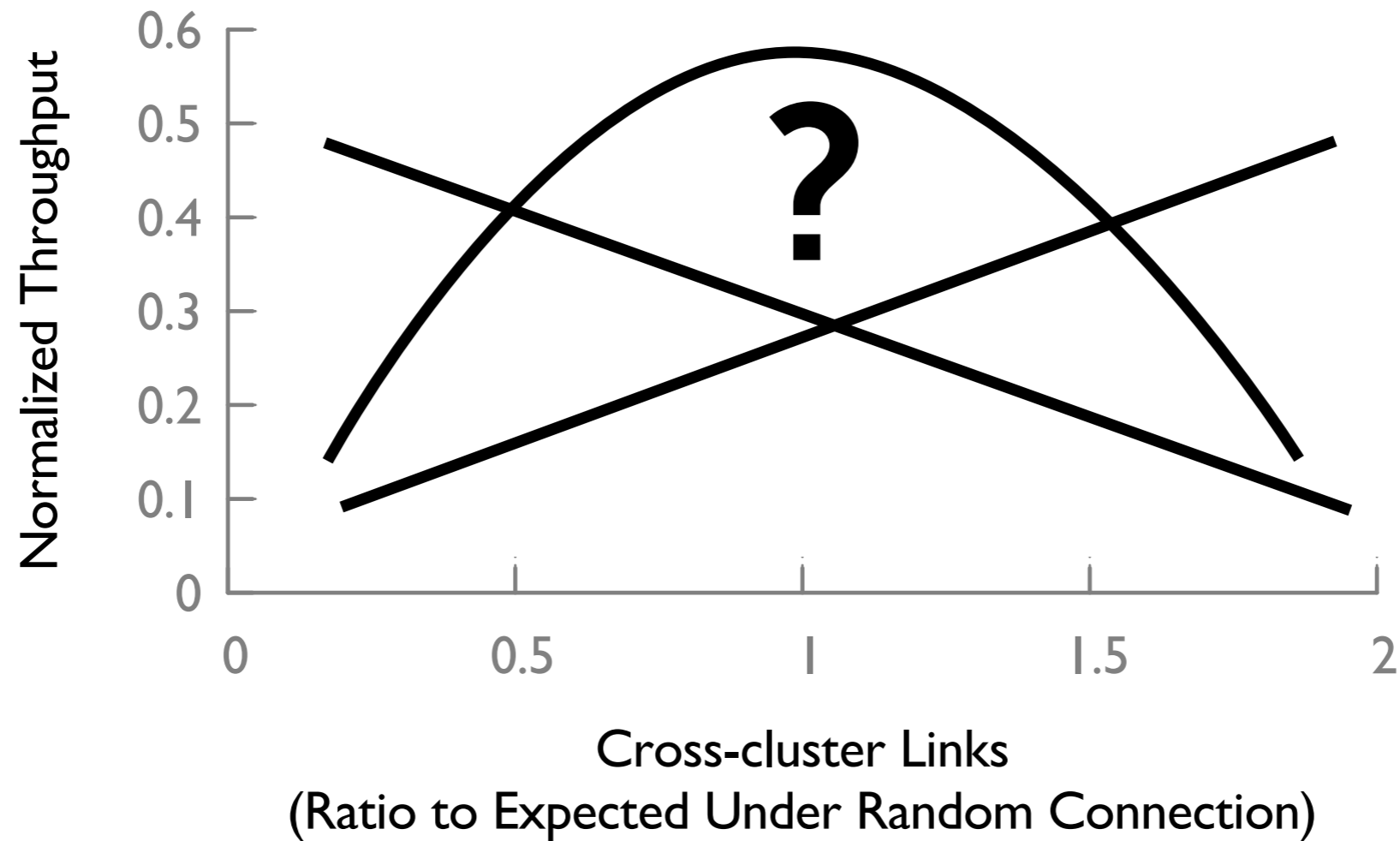
Generic optimization: Place all switches centrally

Interconnecting clusters

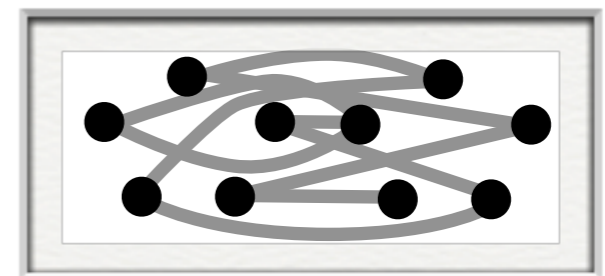
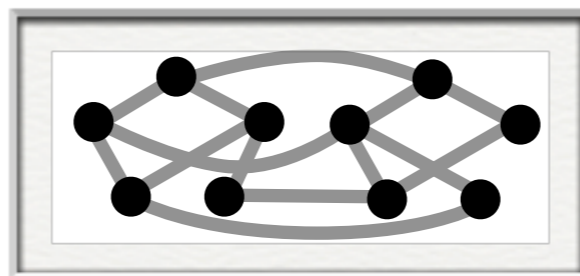
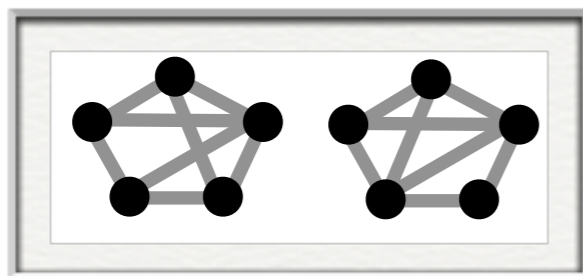
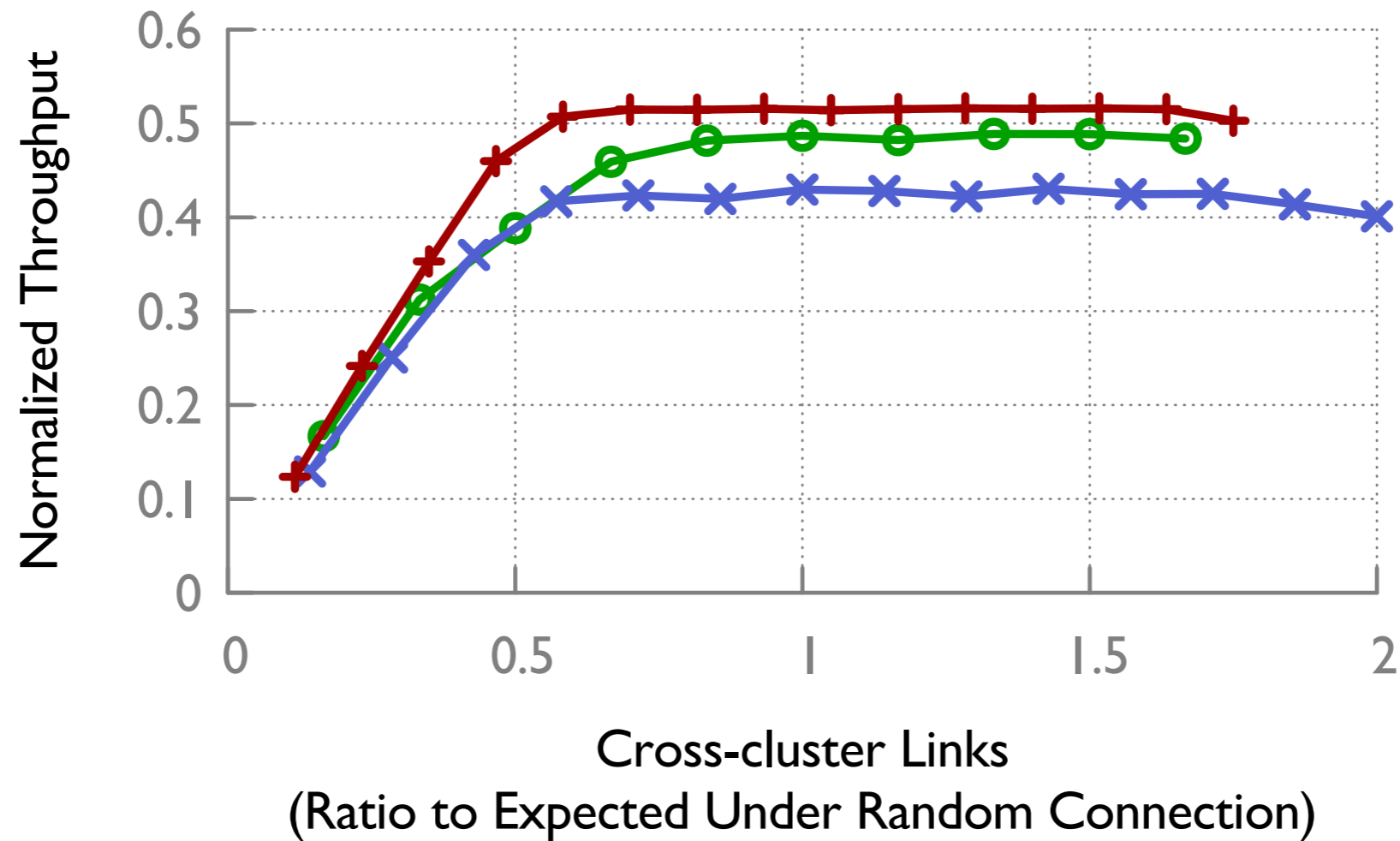
How many “long” cables do we need?



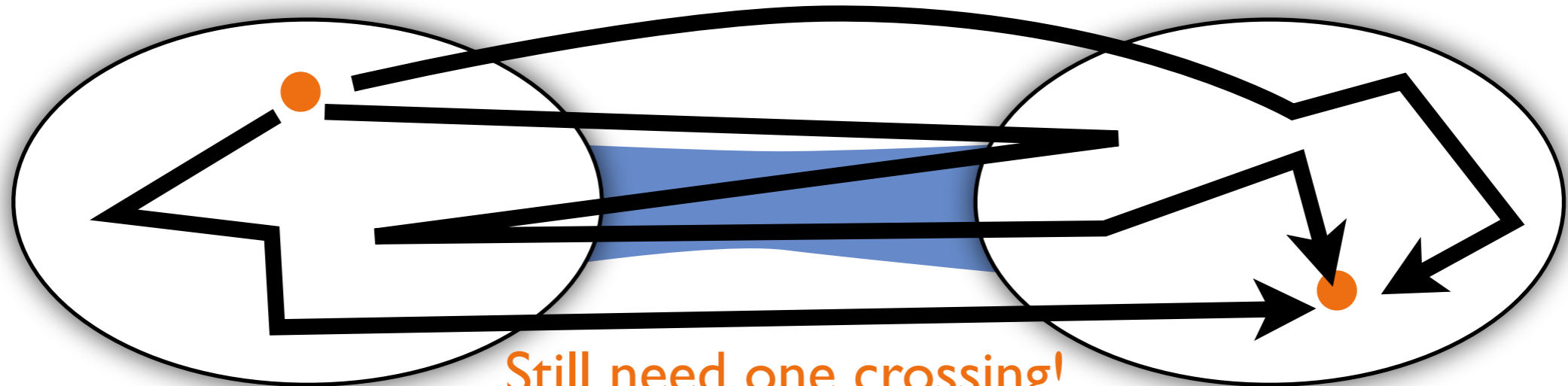
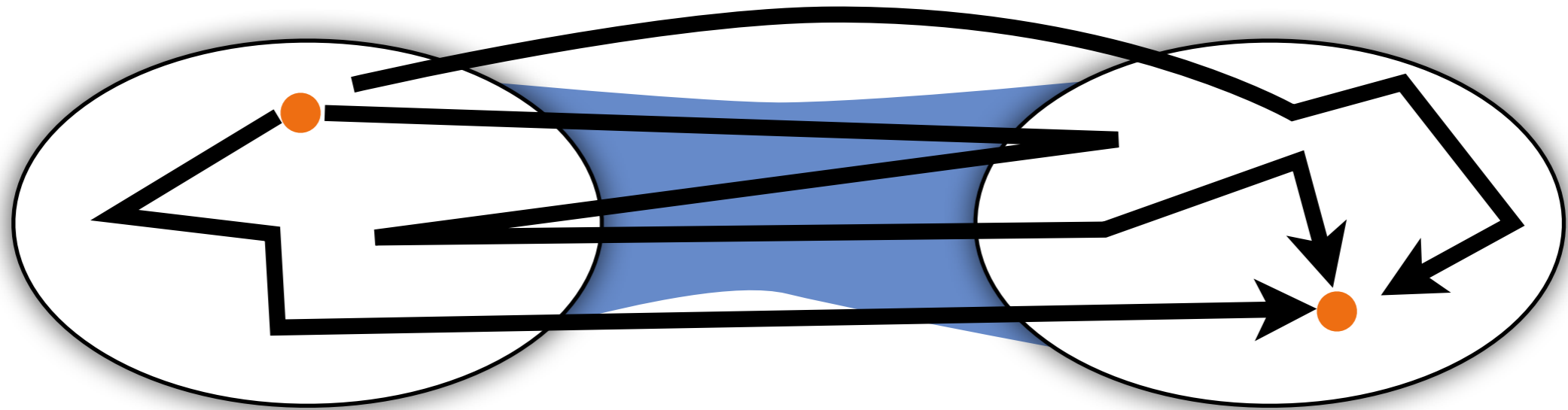
Interconnecting clusters



Interconnecting clusters



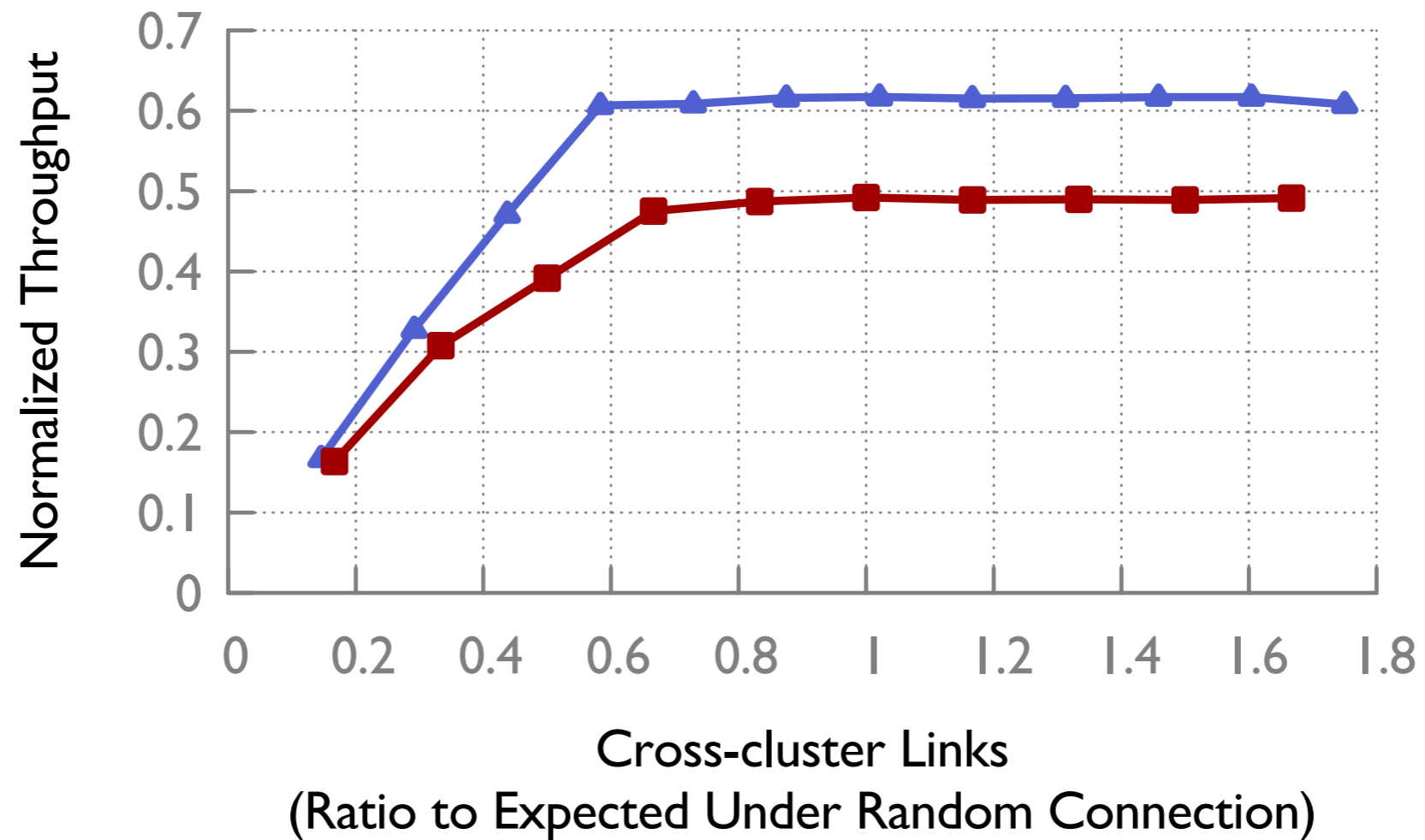
Intuition



Still need one crossing!

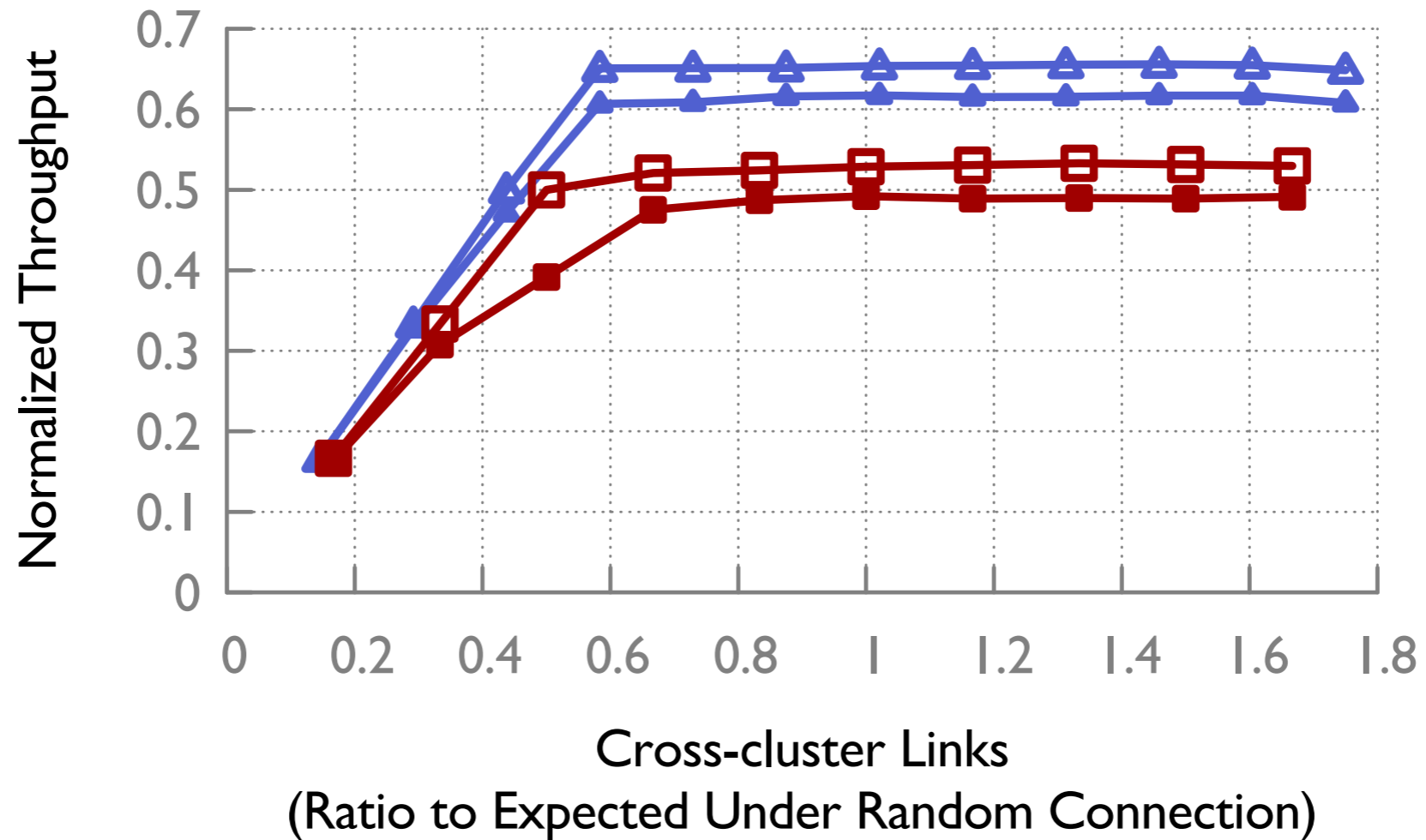
Throughput should drop when less than $\Theta\left(\frac{1}{APL}\right)$ of total capacity crosses the cut!

Explaining throughput



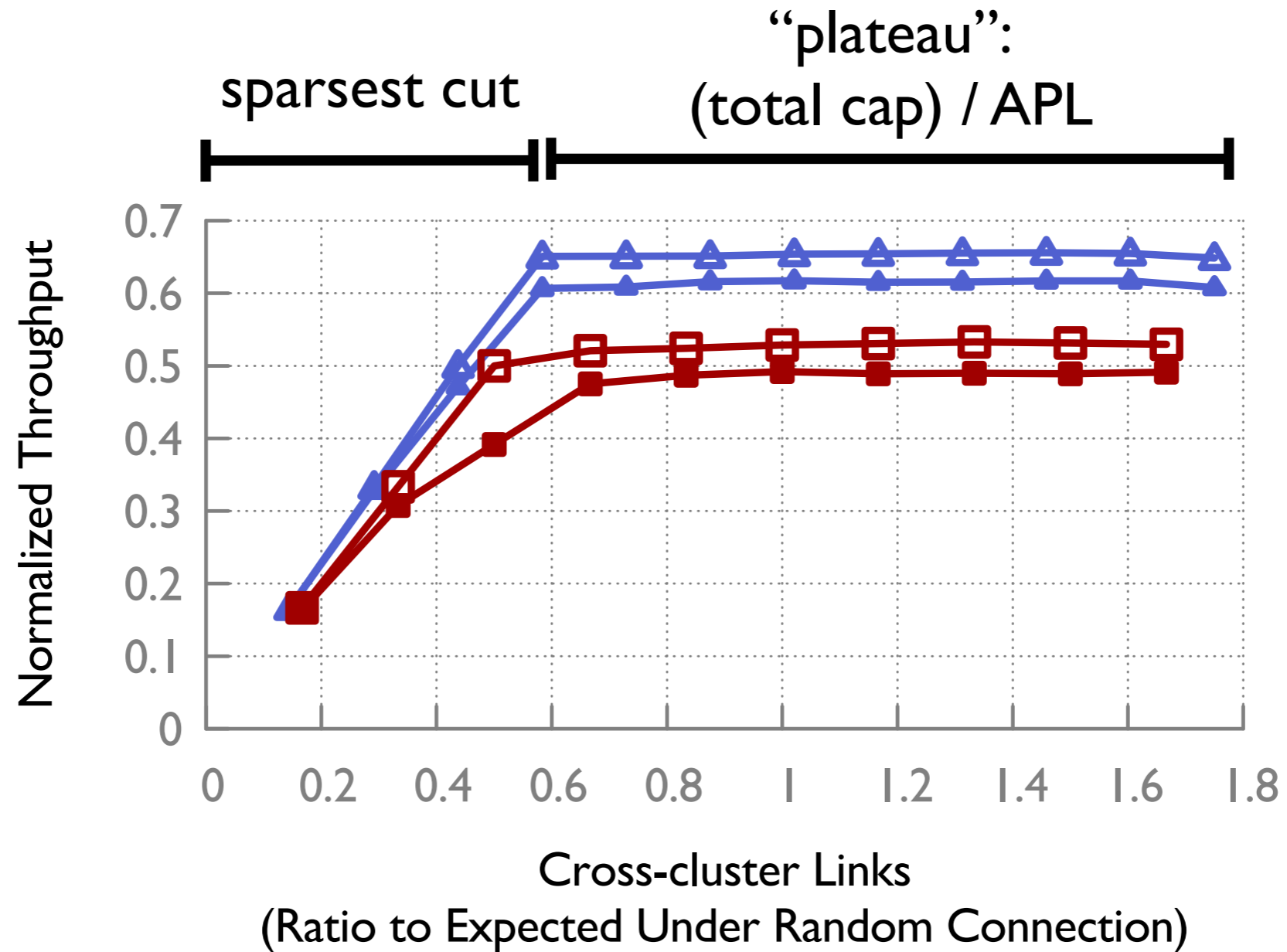
Explaining throughput

Upper bounds...

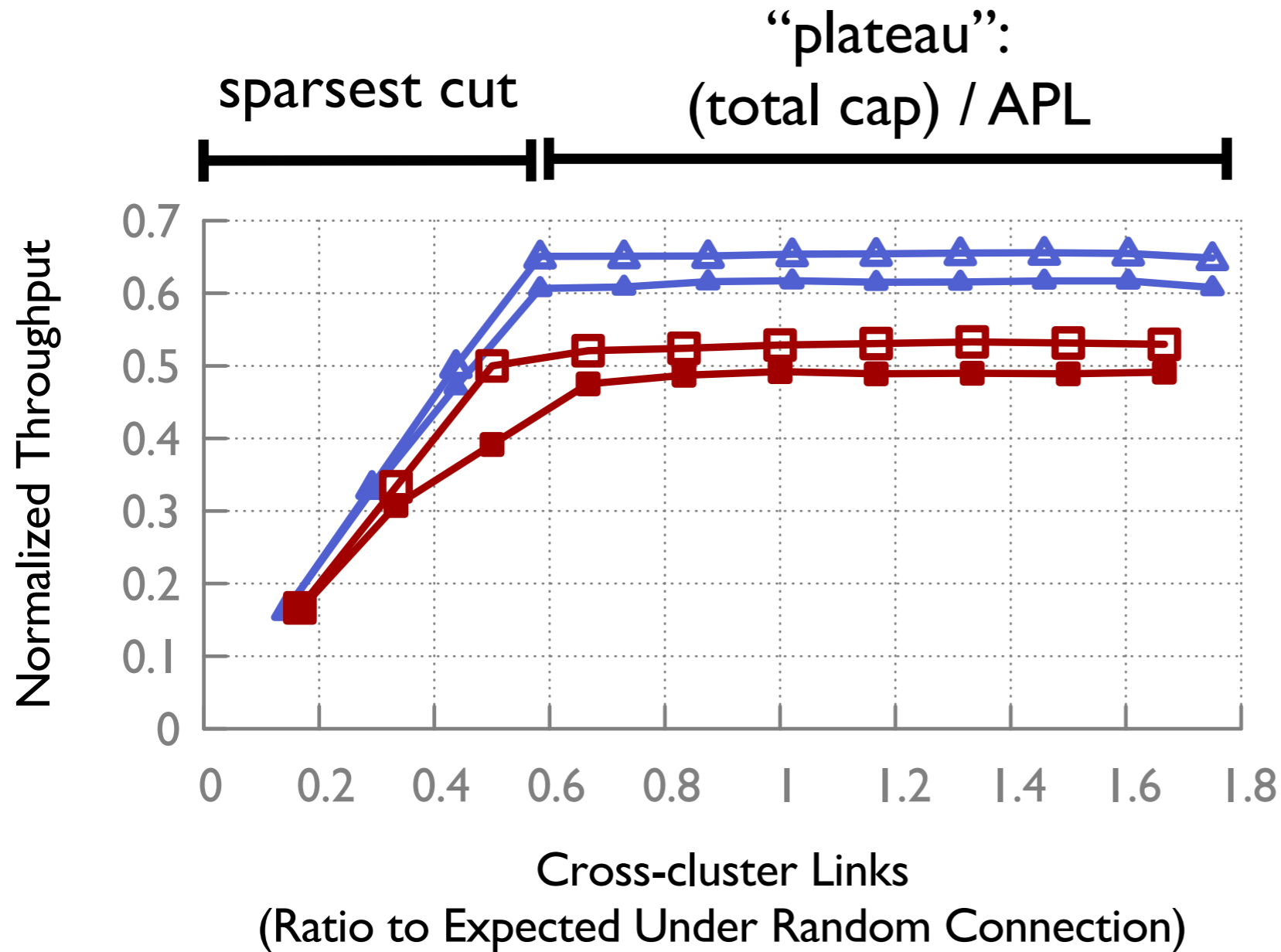


And constant-factor matching lower bounds in special case.

Two regimes of throughput



Two regimes of throughput



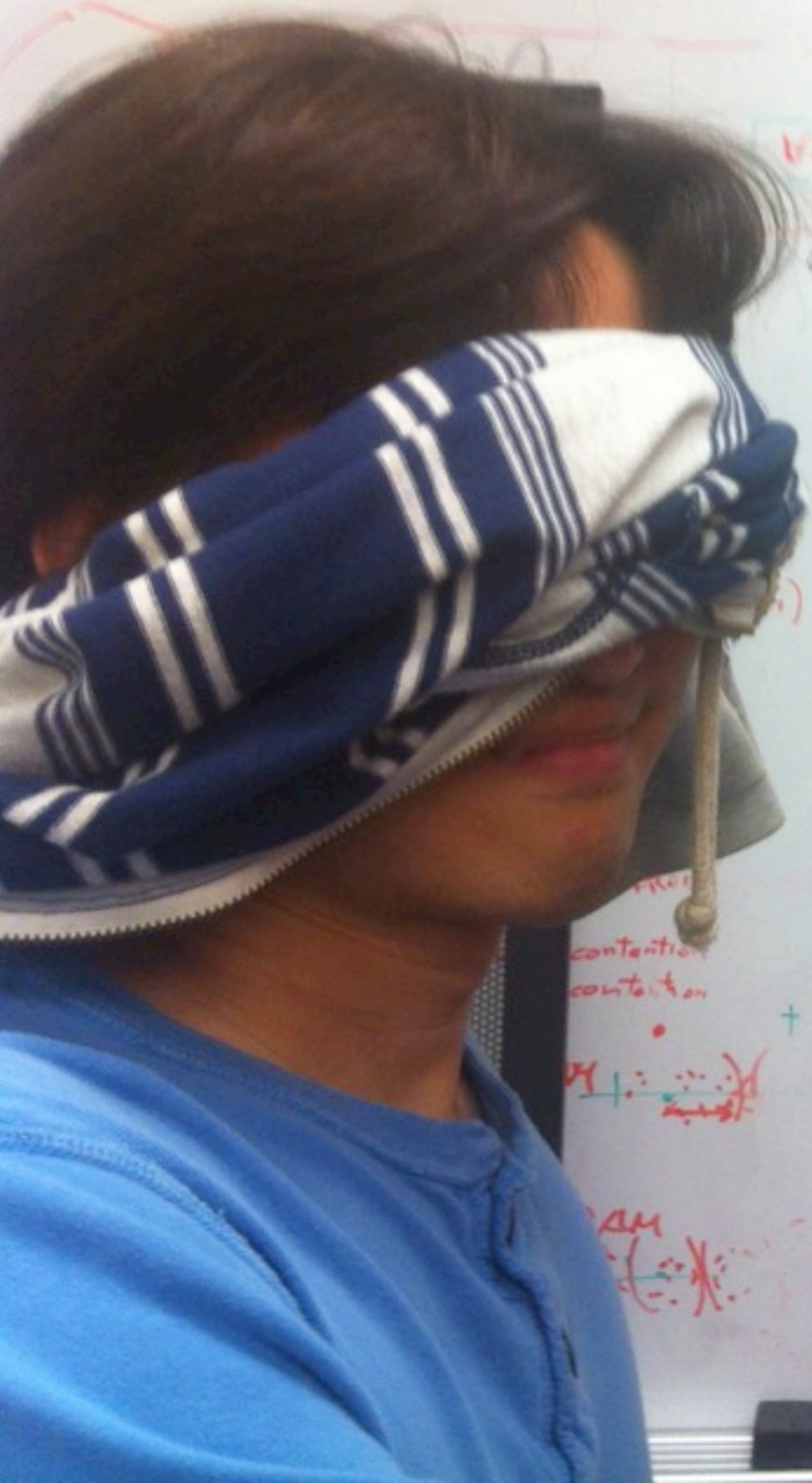
High-capacity switches needn't be clustered

Bisection bandwidth is poor predictor of performance!

Cables can be localized

What's Next





API

API

API

14

13

queues (empty)

min 12

2
1
0

0

189
338

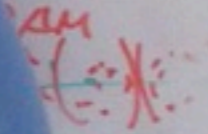
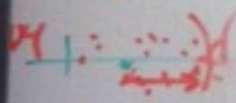
$13 \times 13 \times 2$

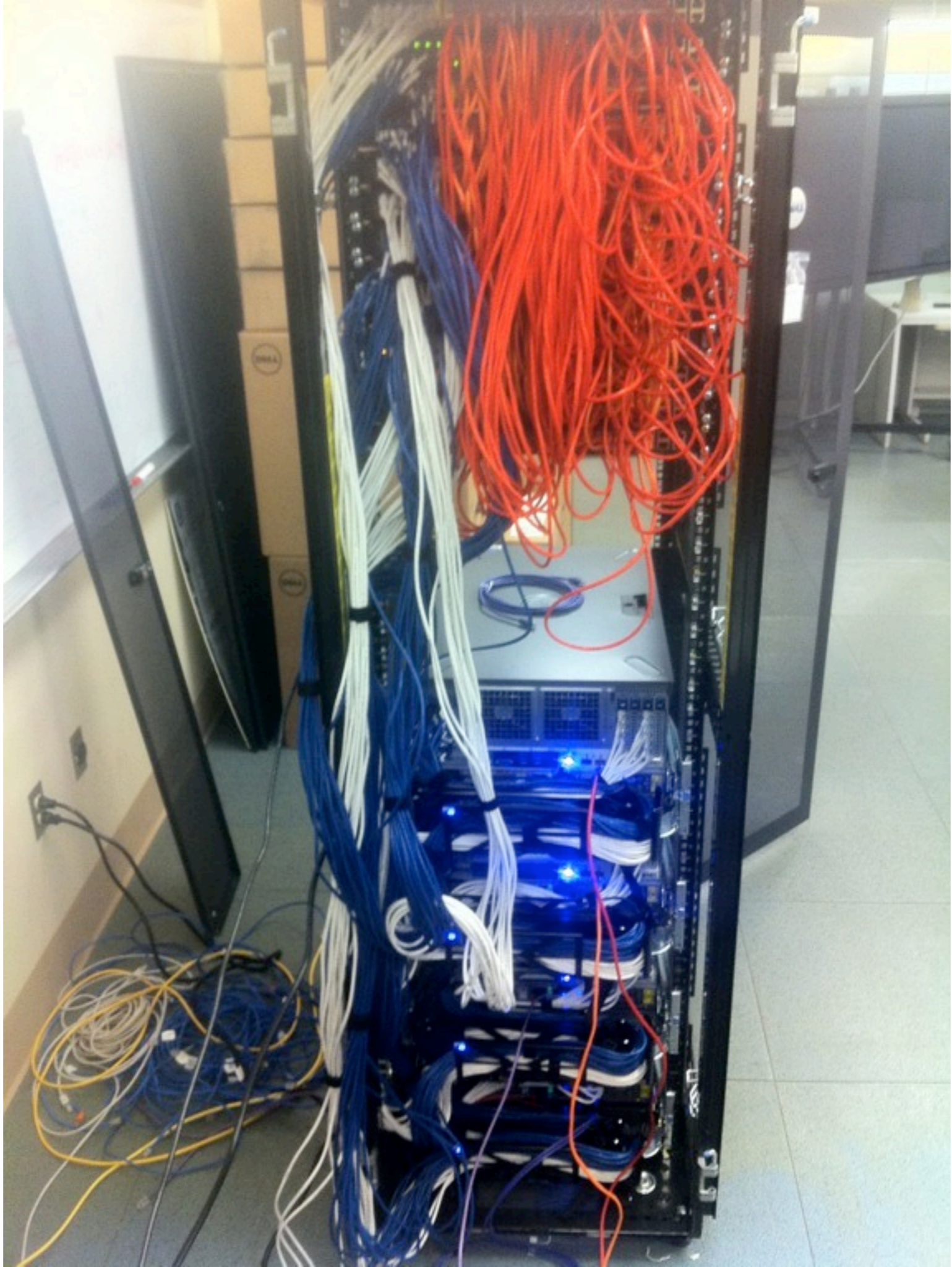
+ 144
12
156

312

450

contention
contention





Research agenda

Prototype in the lab

- High throughput routing even in unstructured networks
- New techniques for near-optimal TE applicable generally
- SDN-based implementation

Topology-aware application & VM placement

Tech transfer

For more...

“Networking Data Centers Randomly”

A. Singla, C. Hong, L. Popa, P. B. Godfrey
NSDI 2012

“High throughput data center topology design”

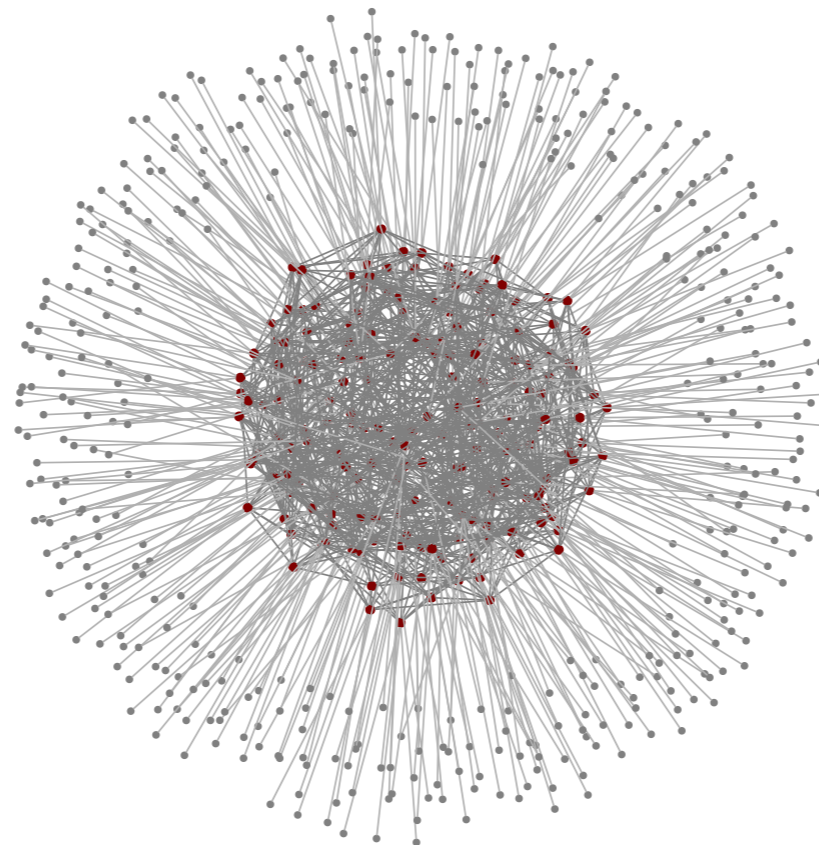
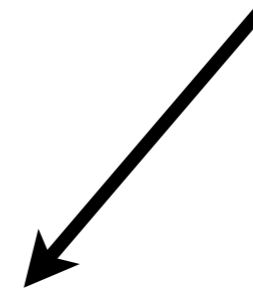
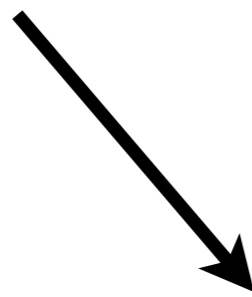
A. Singla, P. B. Godfrey, A. Kolla
Manuscript (check arxiv soon!)



Conclusion

High throughput

Expandability



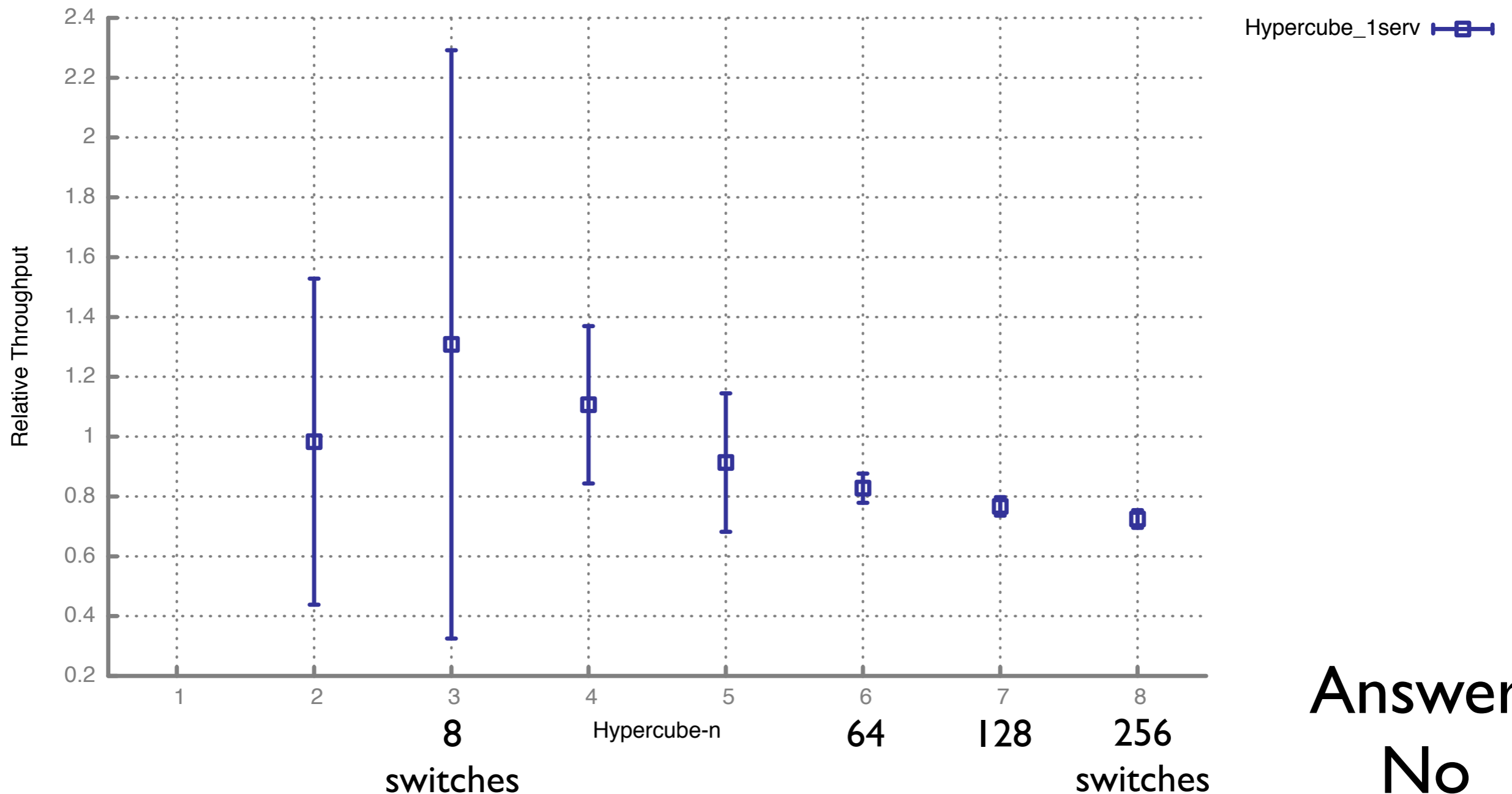


[Photo: Kevin Raskoff]

Backup Slides

Hypercube vs. Random Graph

Is Jellyfish's advantage just that it's a "direct" network?



Answer:
No

Are There Even
Better Topologies?

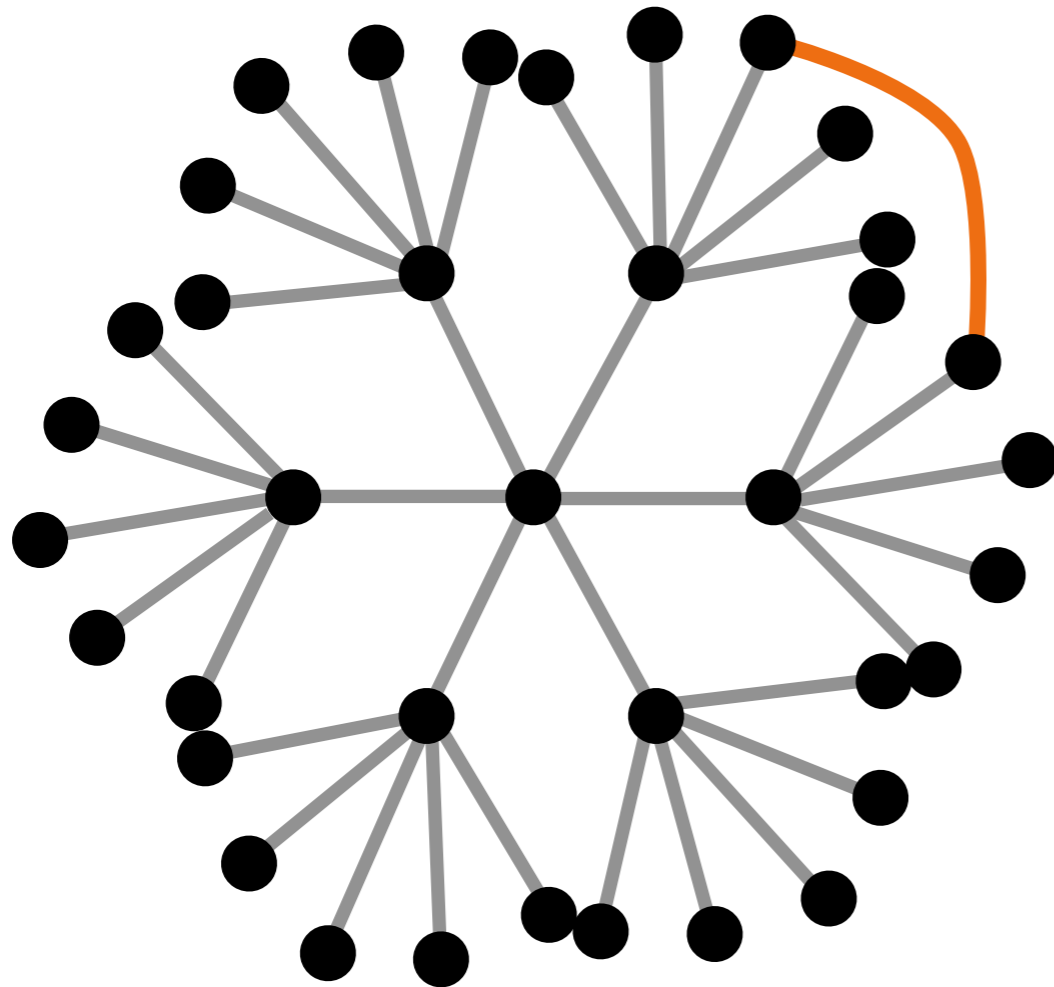
A simple upper bound

$$\text{Throughput per flow} \leq \frac{\sum_{\text{links}} \text{capacity}(\text{link})}{\# \text{ flows} \cdot \text{mean path length}}$$

Lower bound this!



Lower bound on mean path length

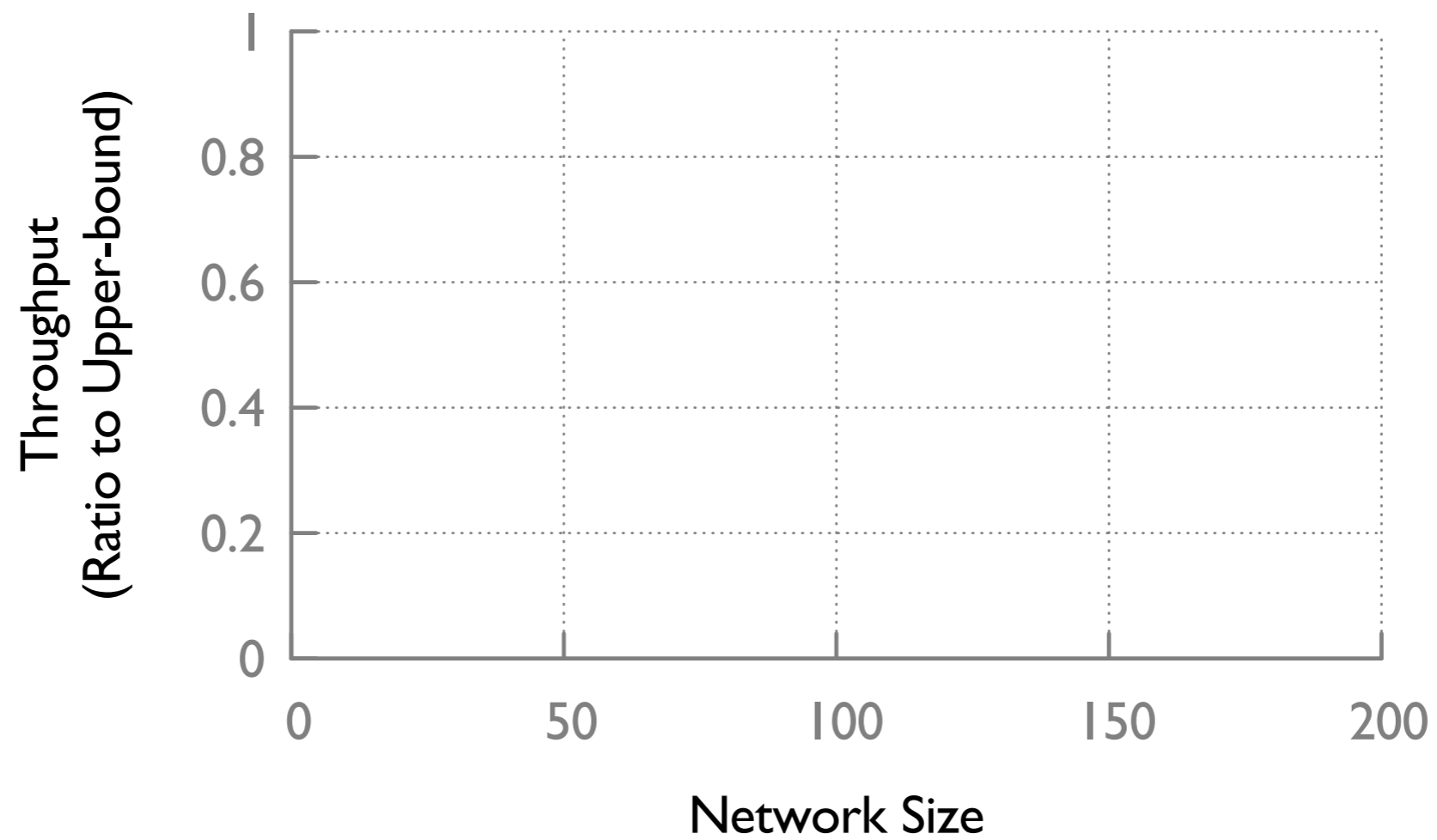


Distance	# Nodes
1	6
2	$6^2 - 6$

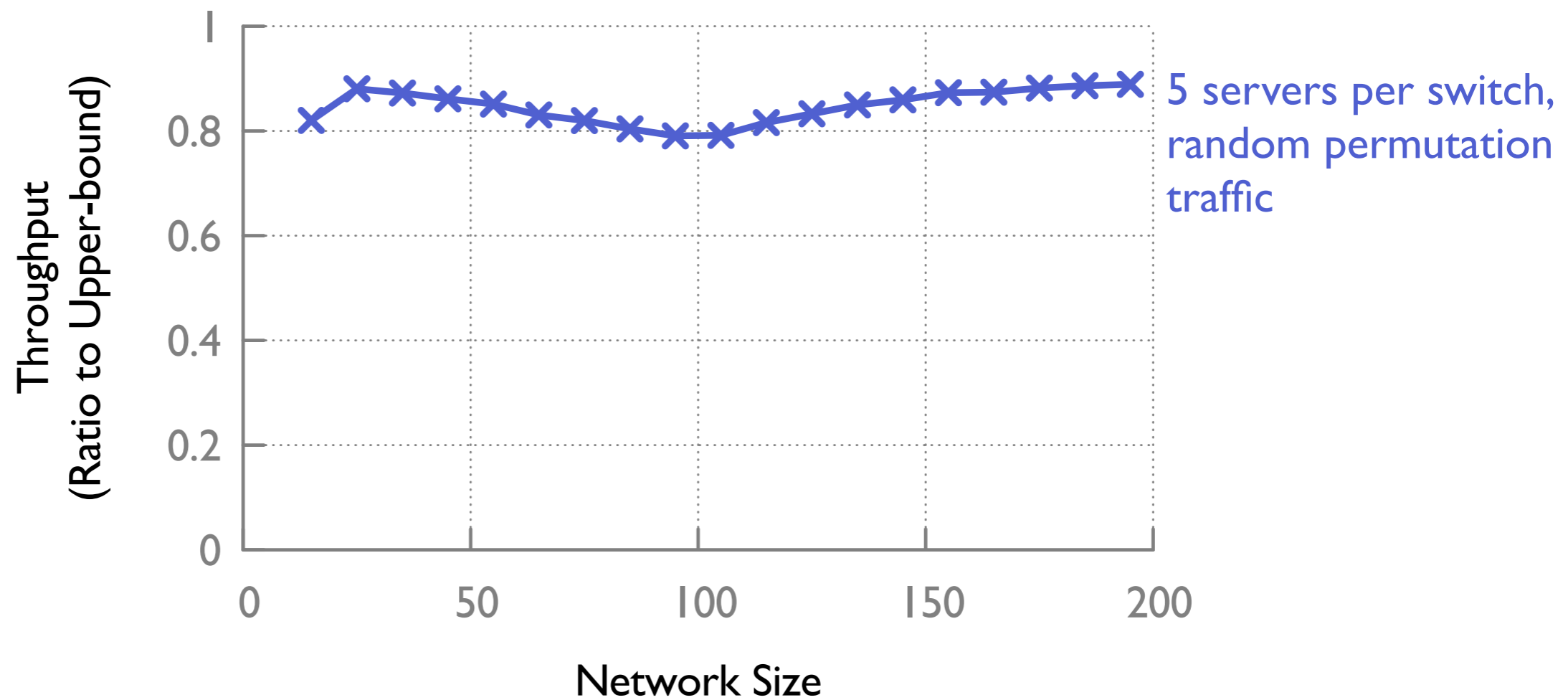
Ugliness omitted!

[Cerf et al., "A lower bound on the average shortest path length in regular graphs", 1974]

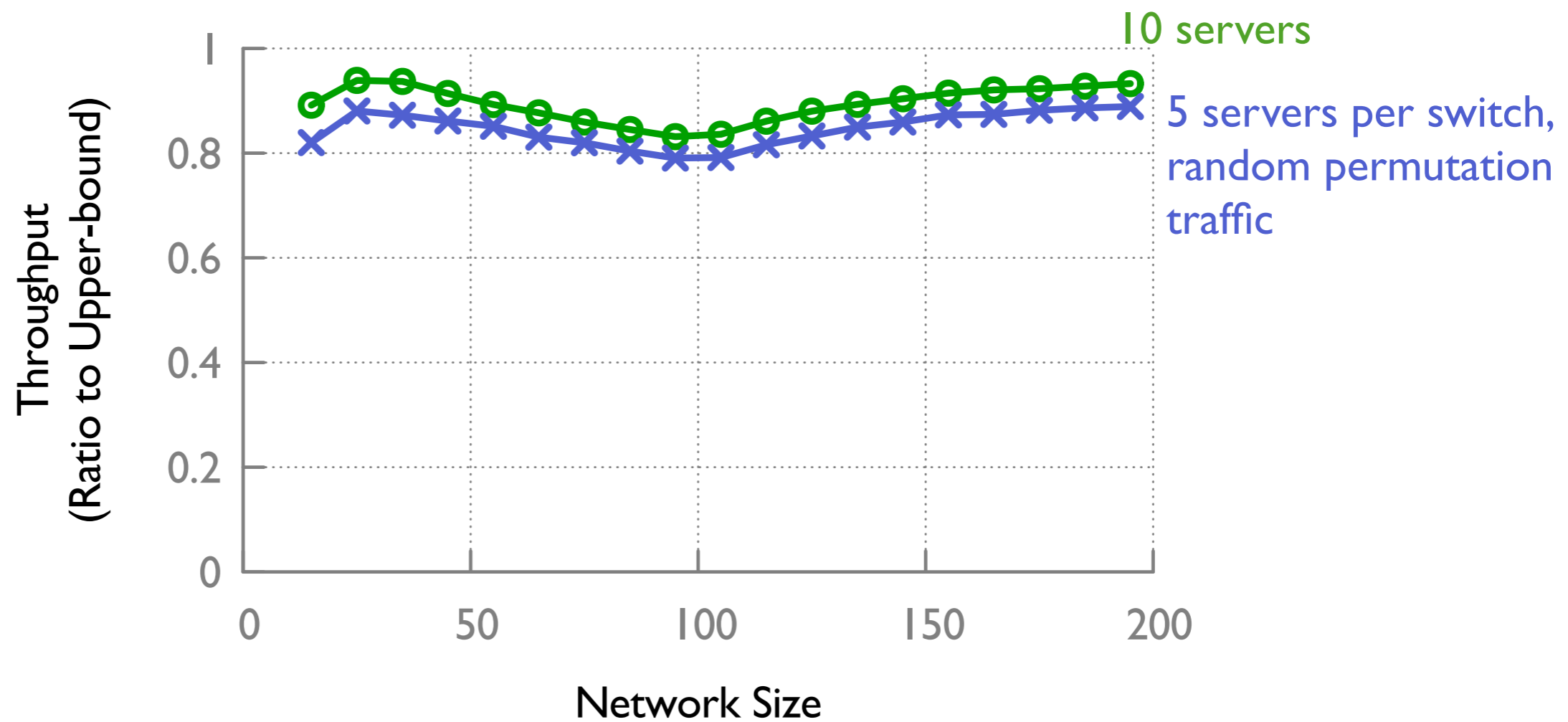
Random graphs vs. upper bound



Random graphs vs. upper bound

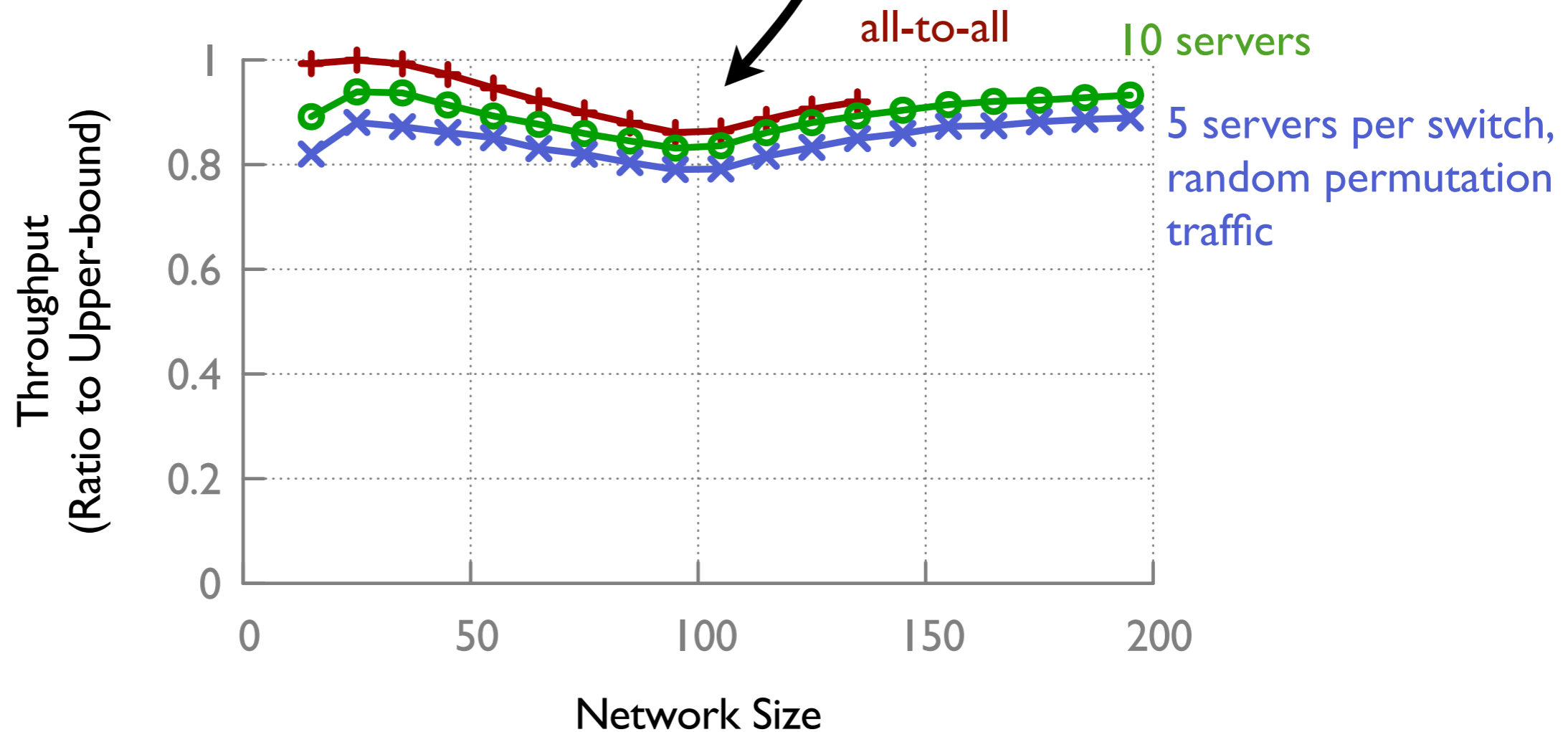


Random graphs vs. upper bound



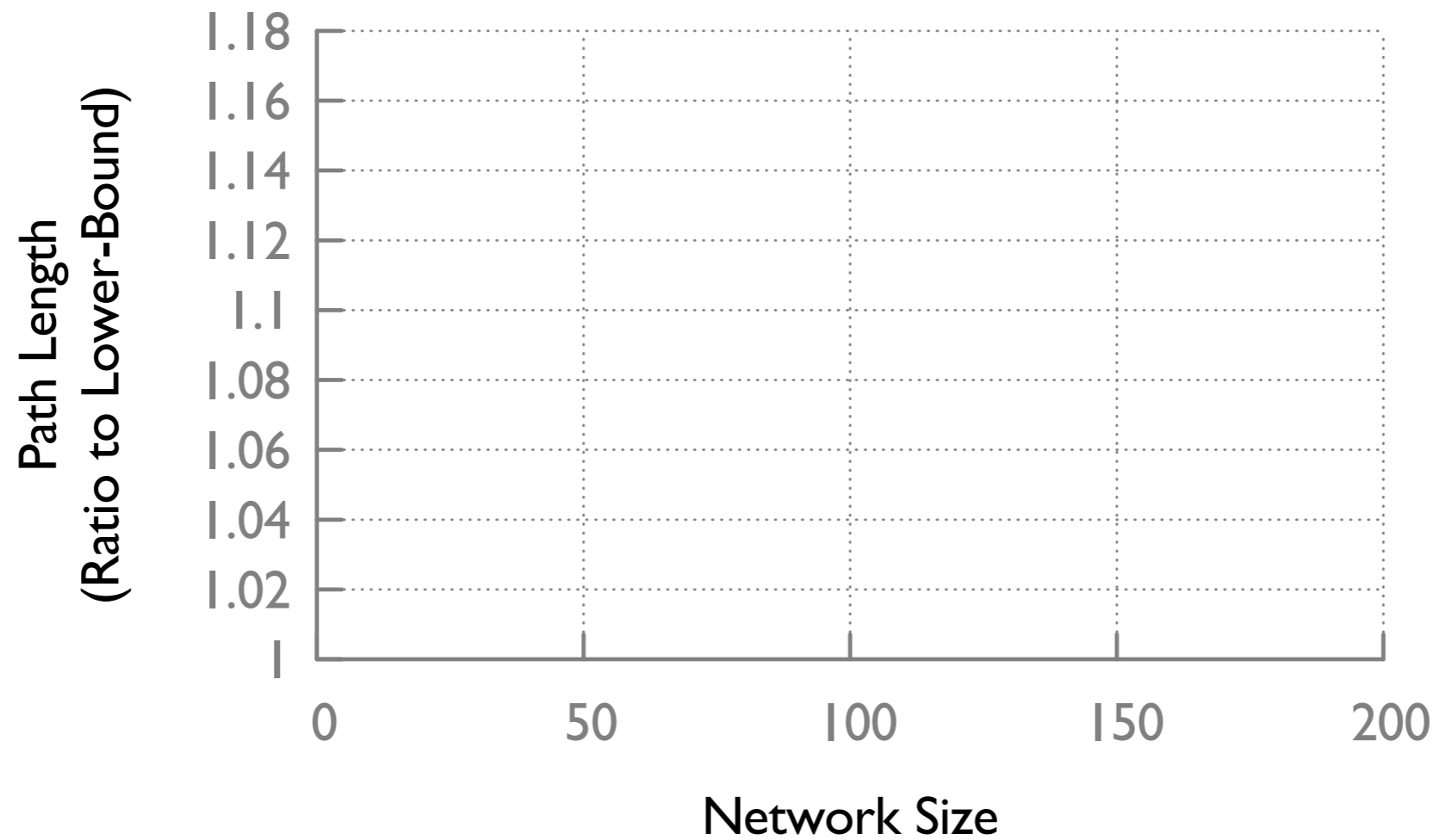
Random graphs vs. upper bound

(Aside: is any topology closer to the bound?)

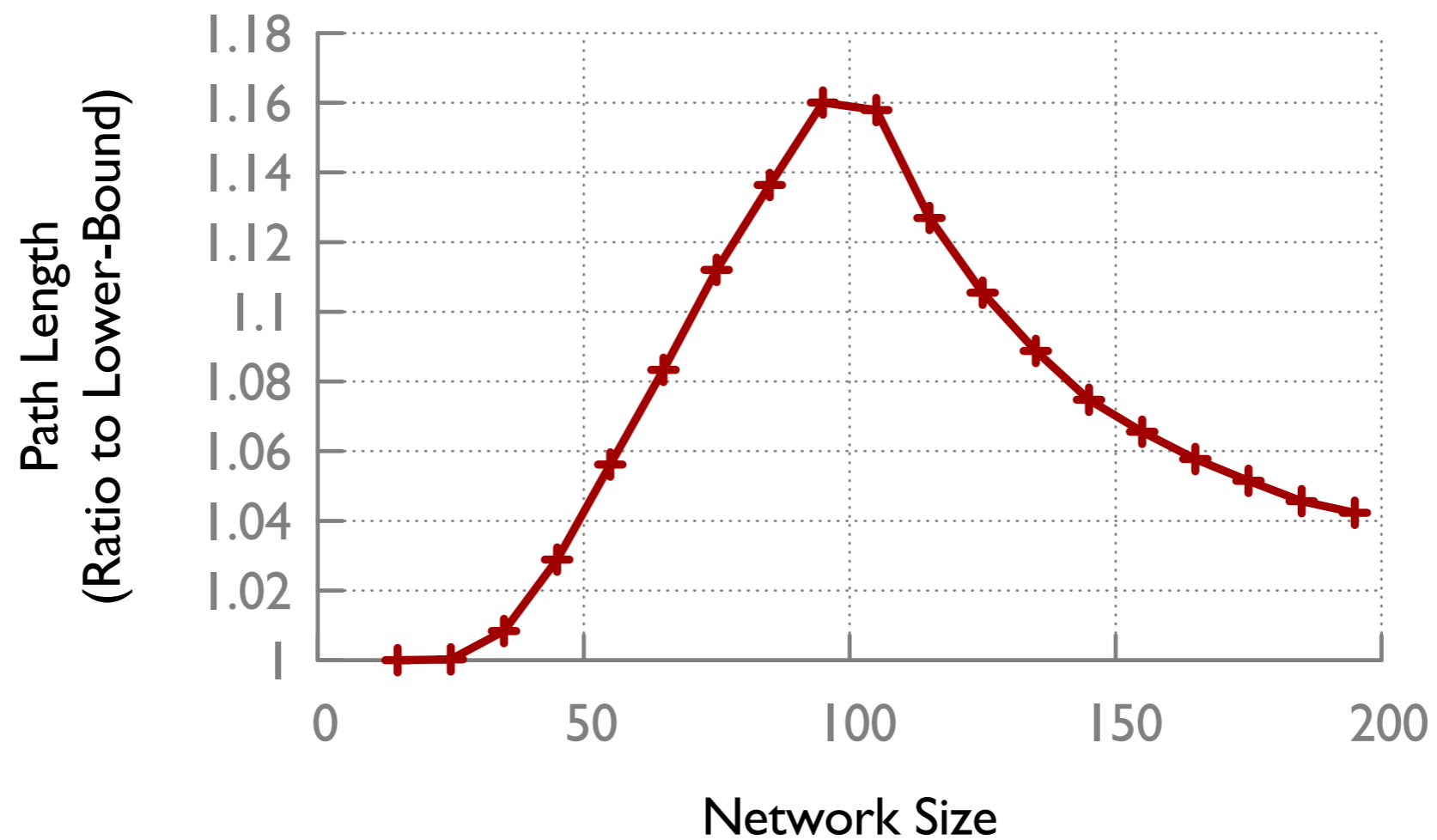


Random graphs within a few percent of optimal!

Random graphs vs. upper bound



Random graphs vs. upper bound



Designing Heterogeneous Networks

Random graphs as a building block

2 How should we interconnect switches?

Low-degree switches

High-degree switches

?

?

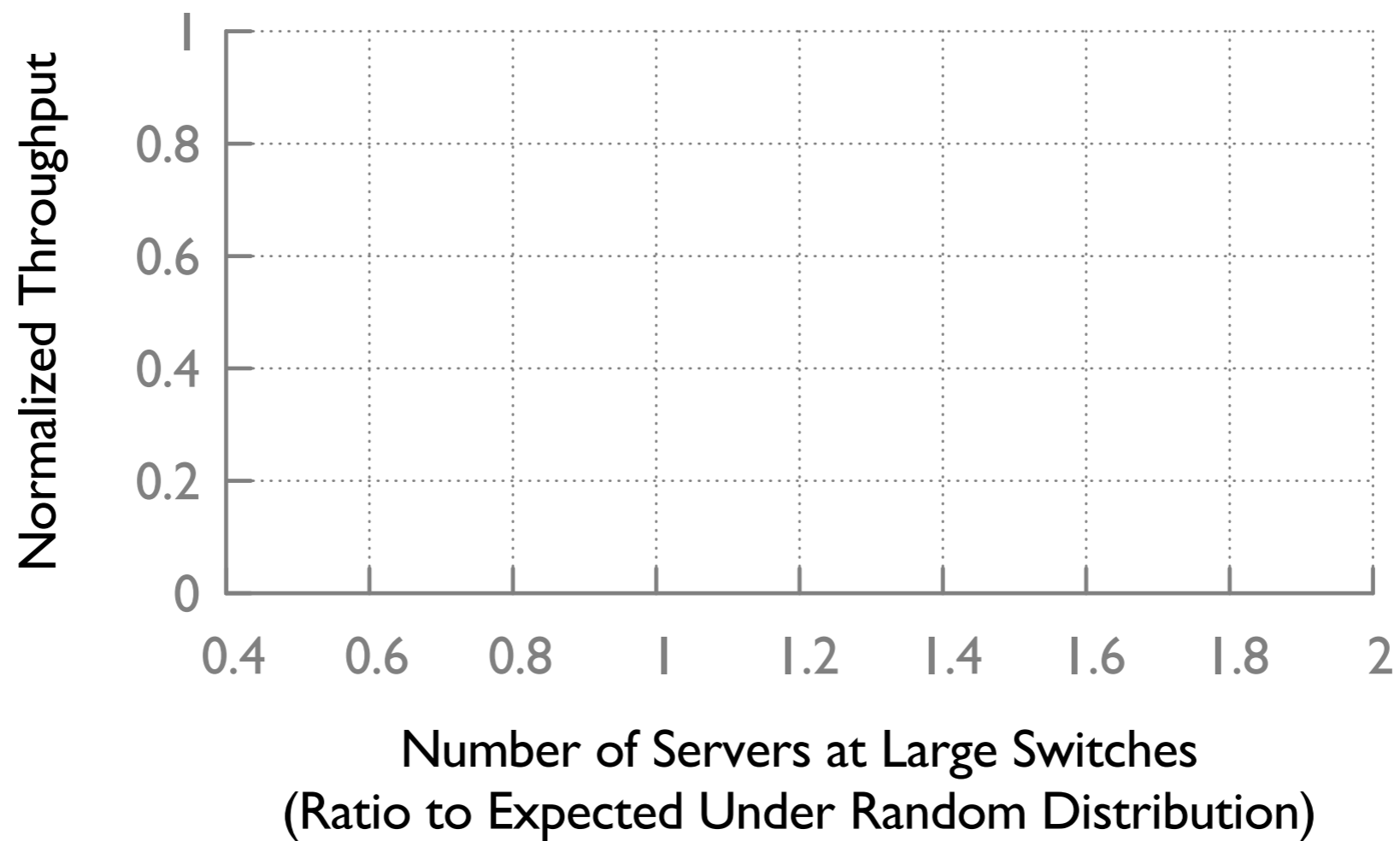
?

1 How should we distribute servers?

Servers

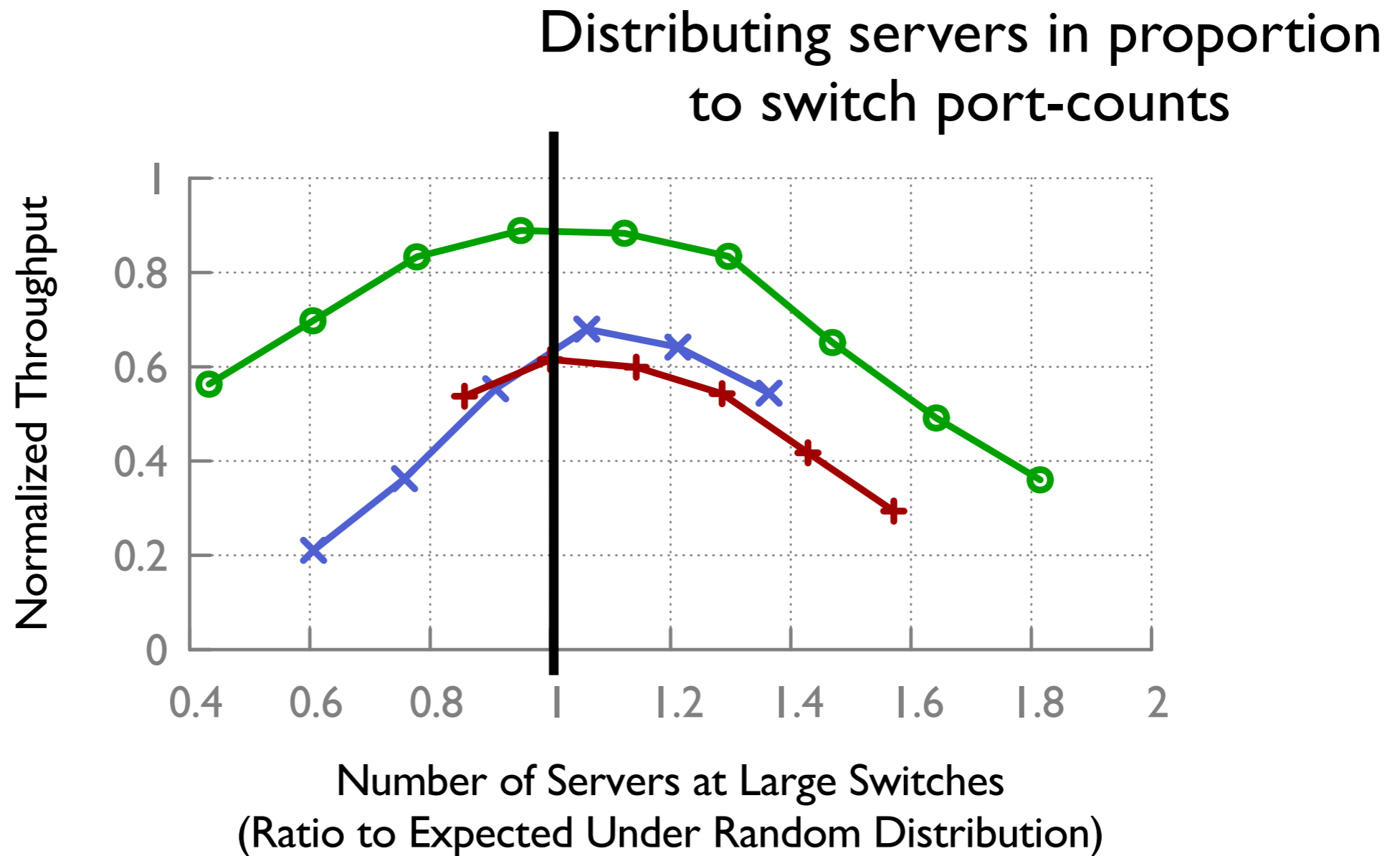
What would you do?

Distributing servers



(The switch interconnect being vanilla random)

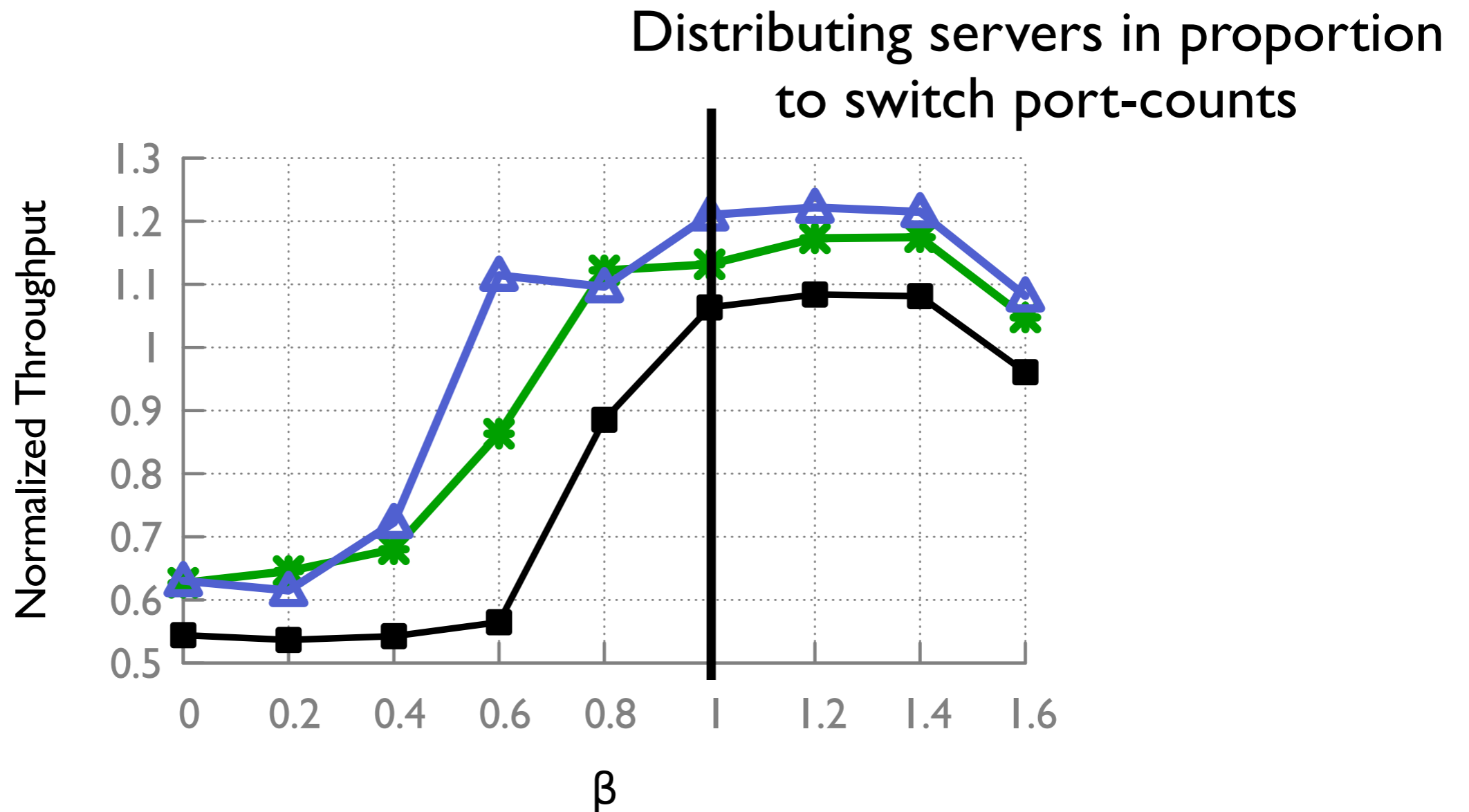
Distributing servers



(The switch interconnect being vanilla random)

Distributing servers

#Servers on switch $i \propto (\text{port-count of } i)^\beta$



Random graphs as a building block

2 How should we interconnect switches?

Low-degree switches

High-degree switches

?

?

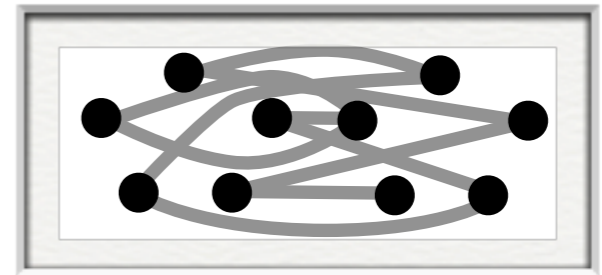
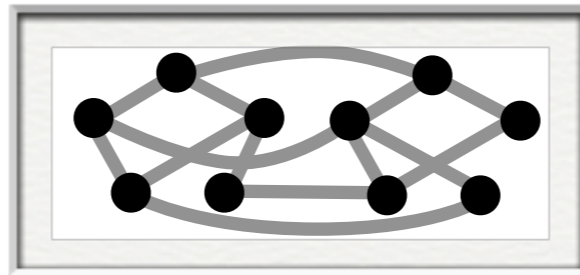
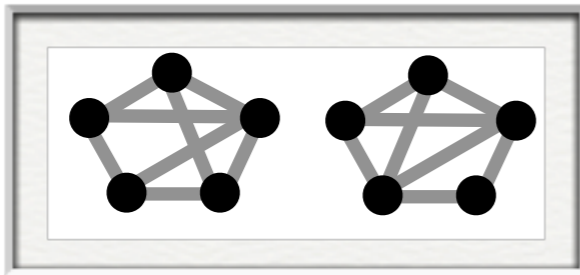
?

1 How should we distribute servers?

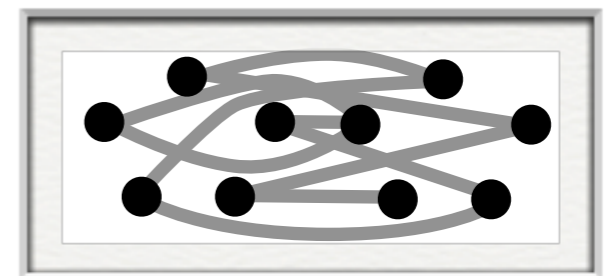
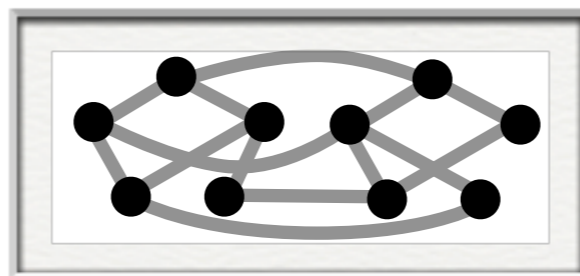
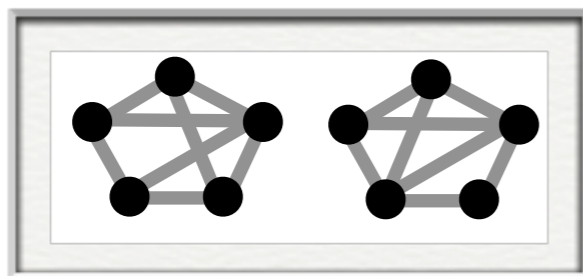
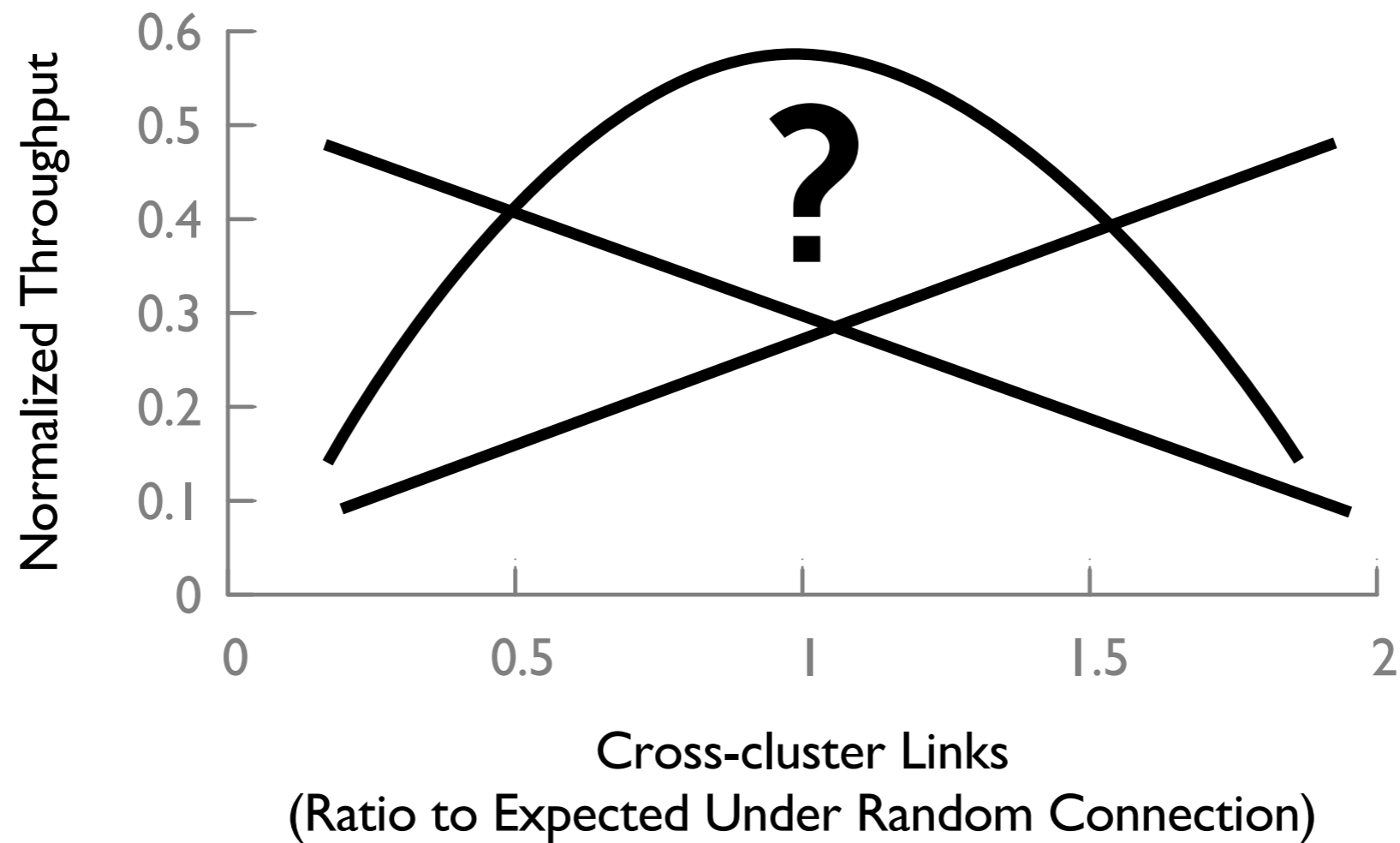
Servers

What would you do?

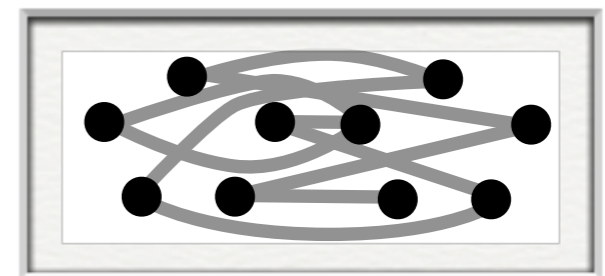
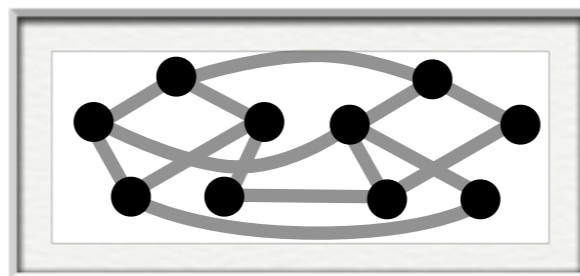
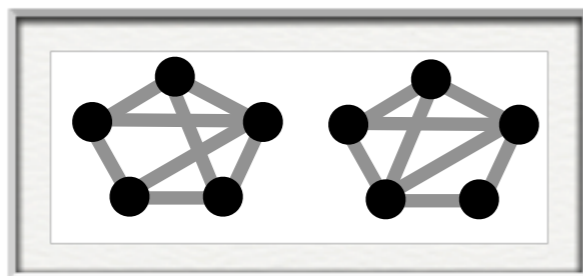
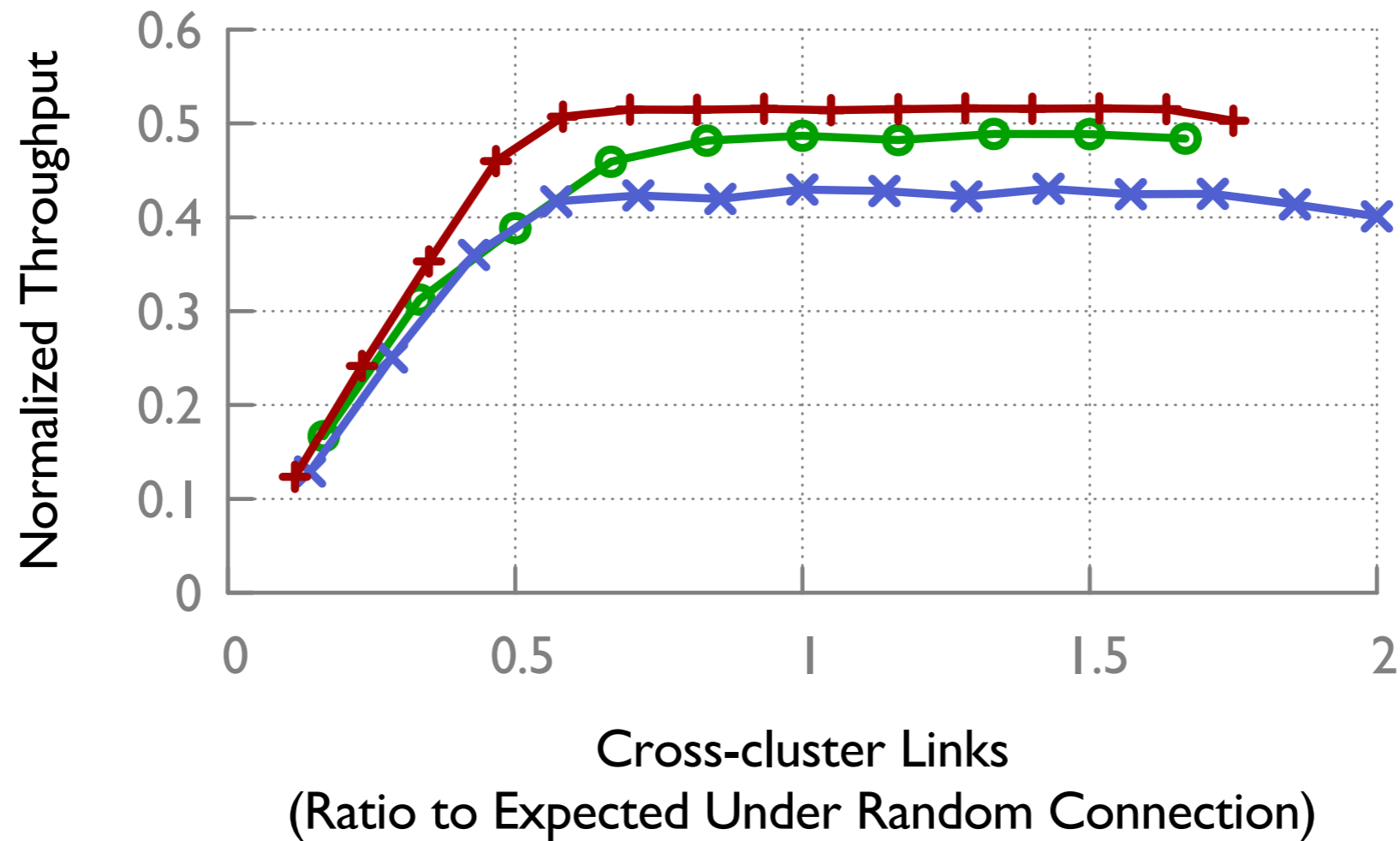
Interconnecting switches



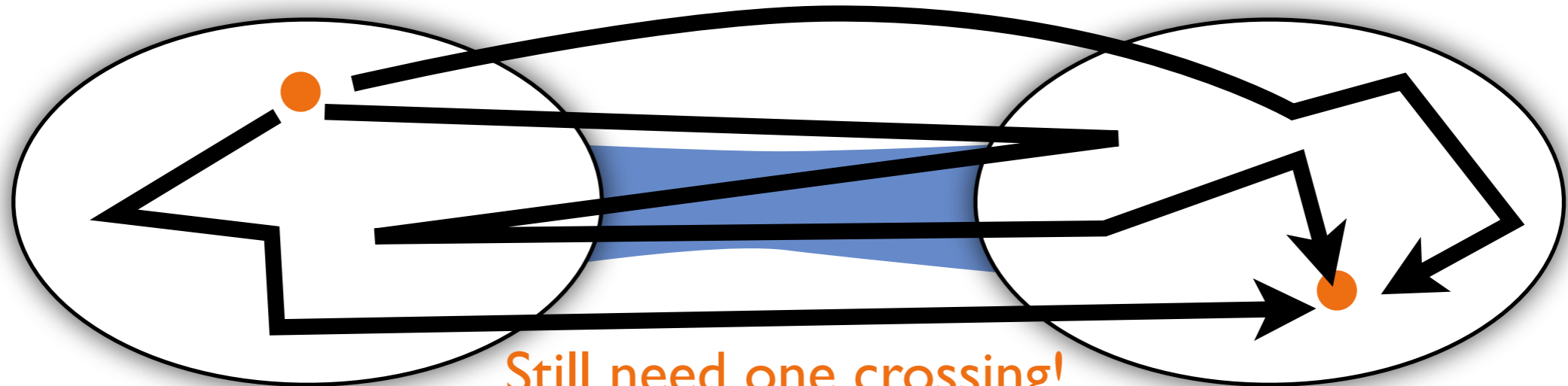
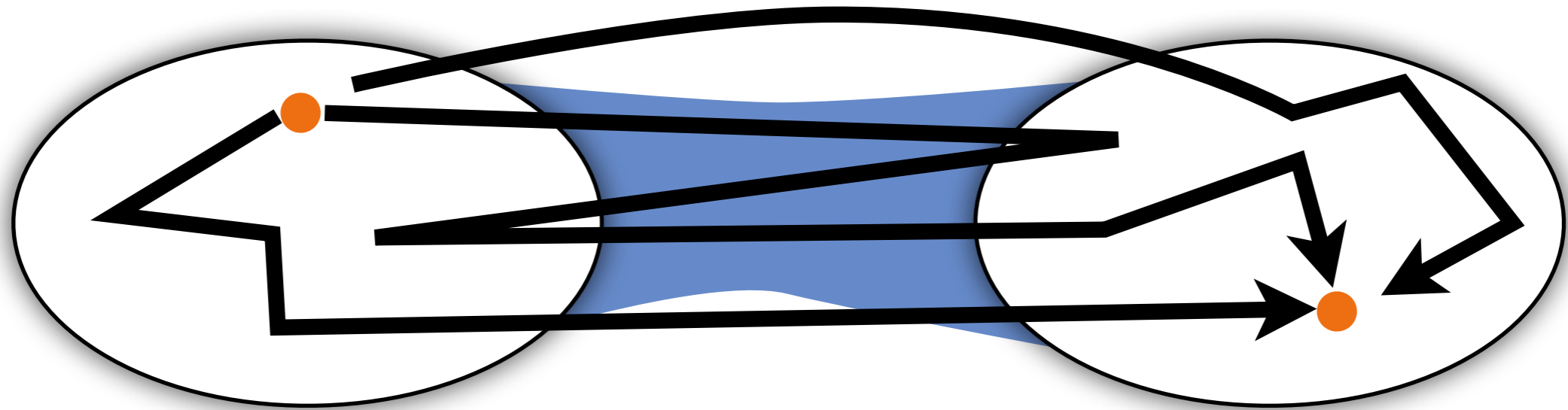
Interconnecting switches



Interconnecting switches



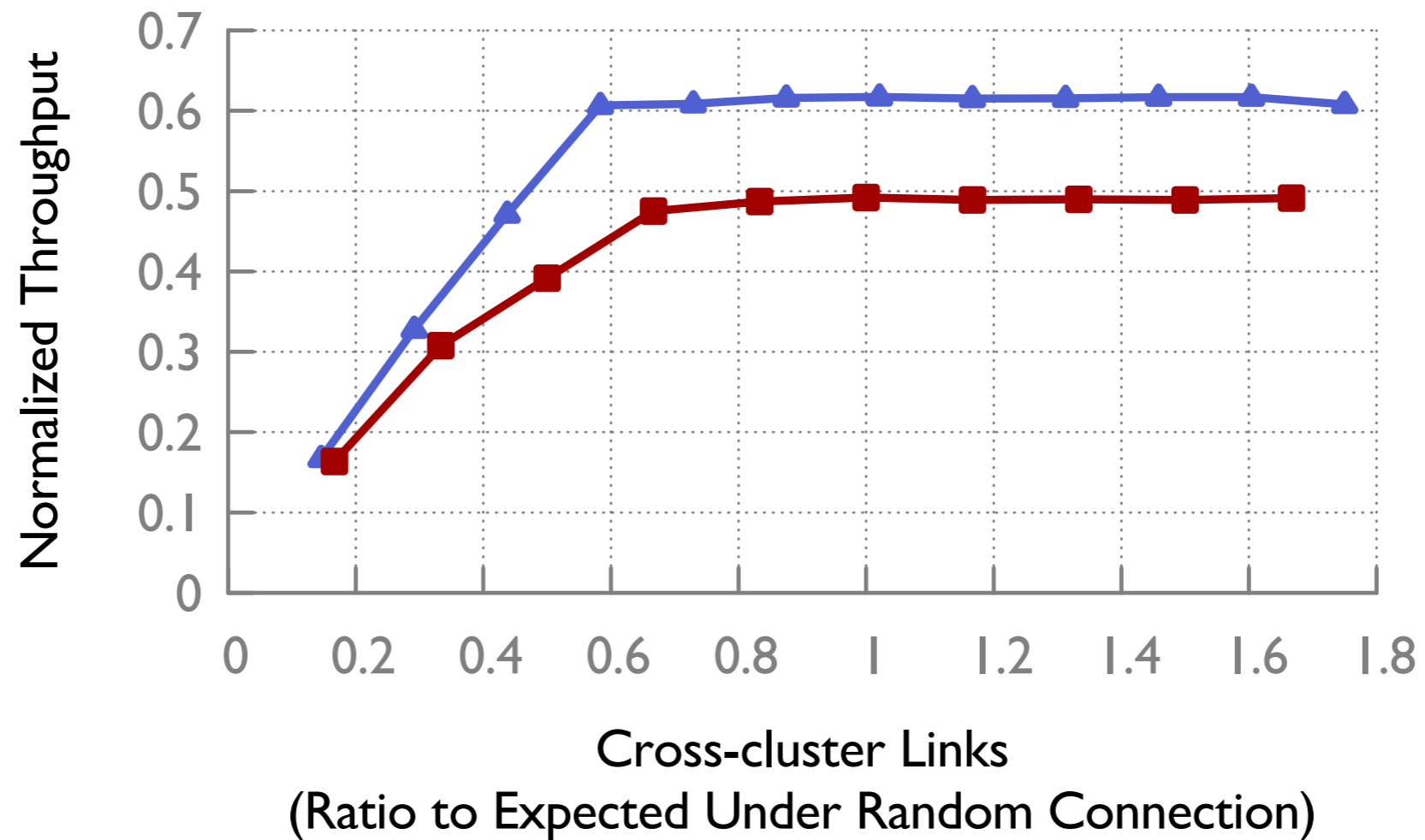
Intuition



Still need one crossing!

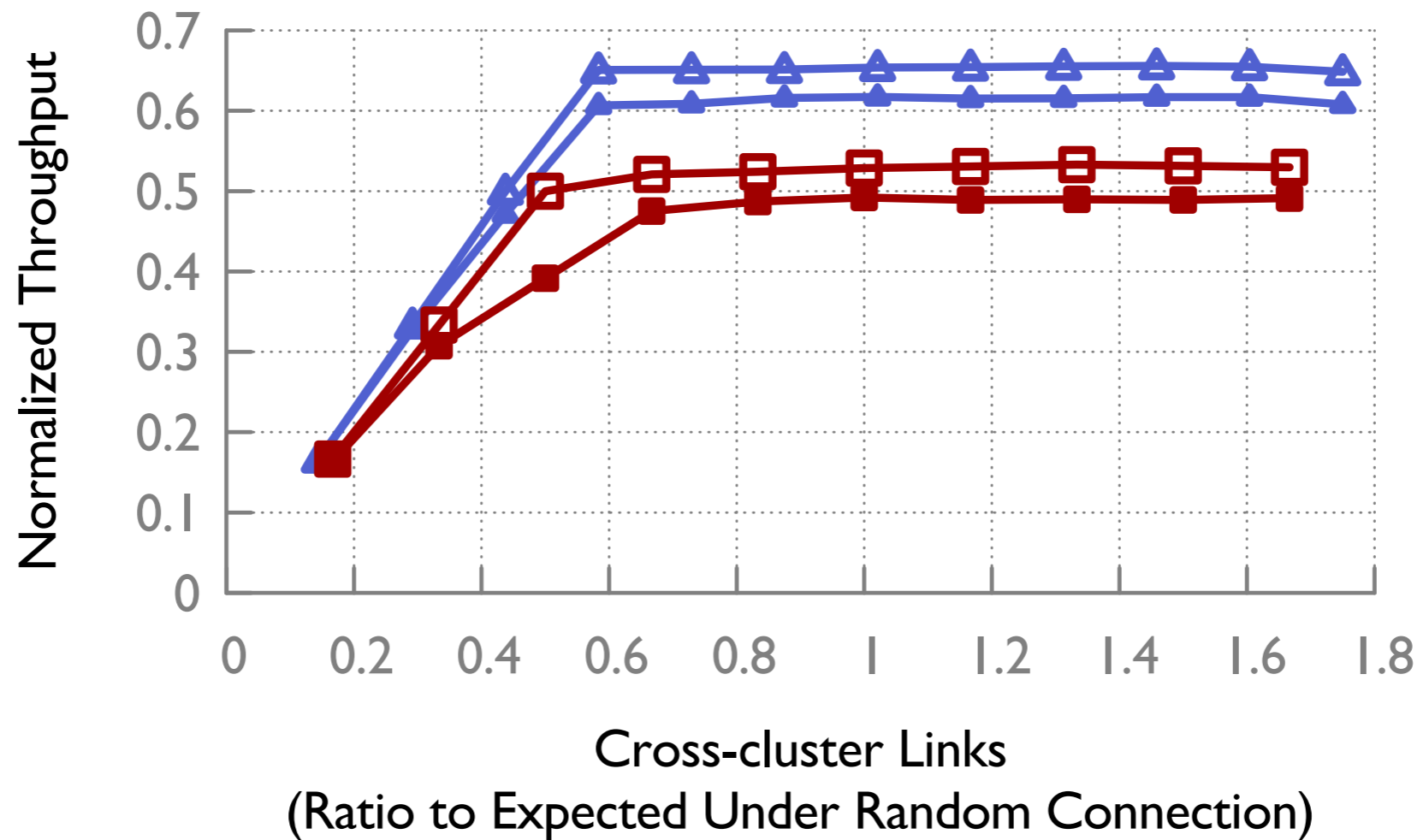
Throughput should drop when less than $\Theta\left(\frac{1}{APL}\right)$ of total capacity crosses the cut!

Explaining throughput



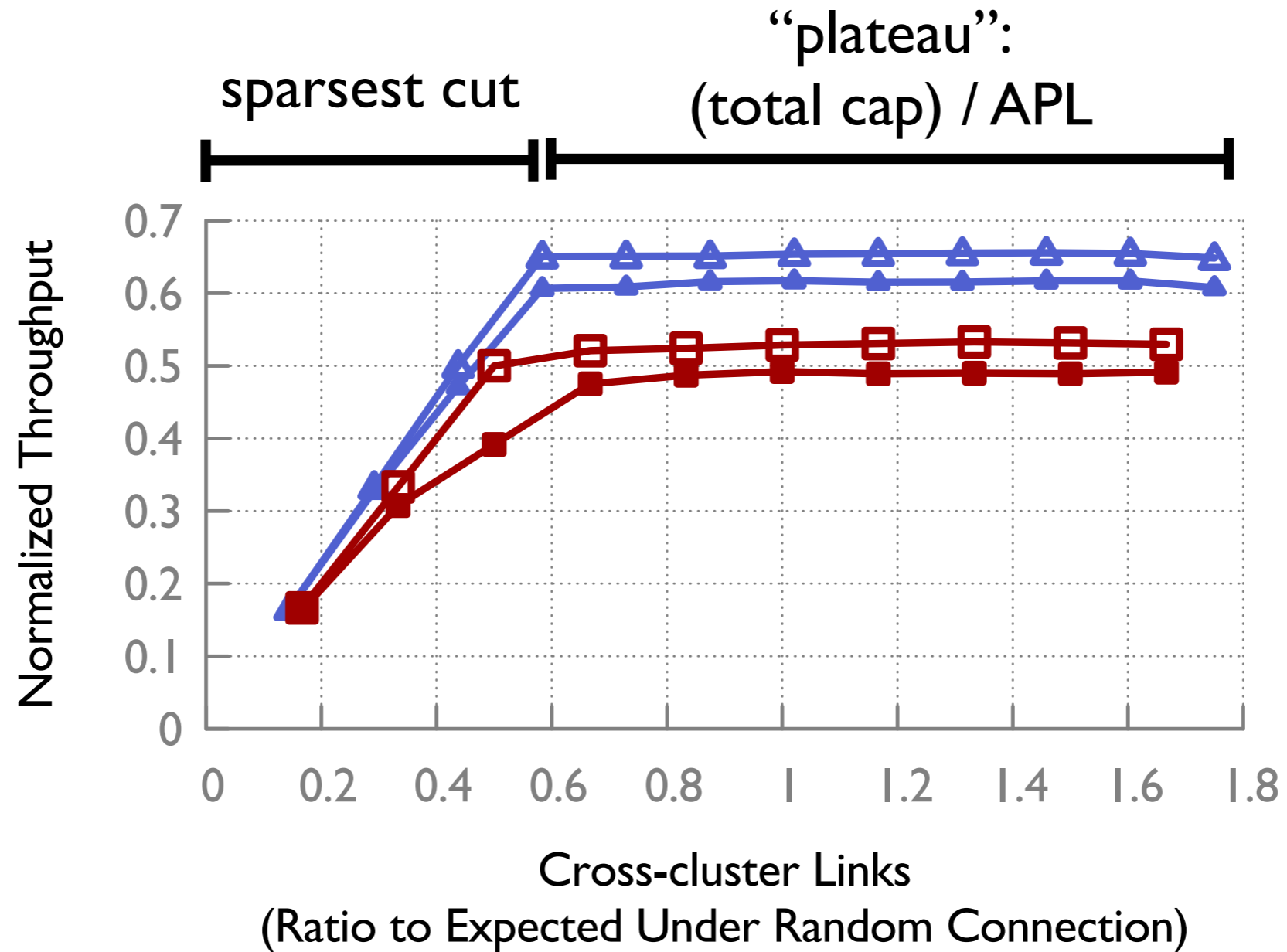
Explaining throughput

Upper bounds...

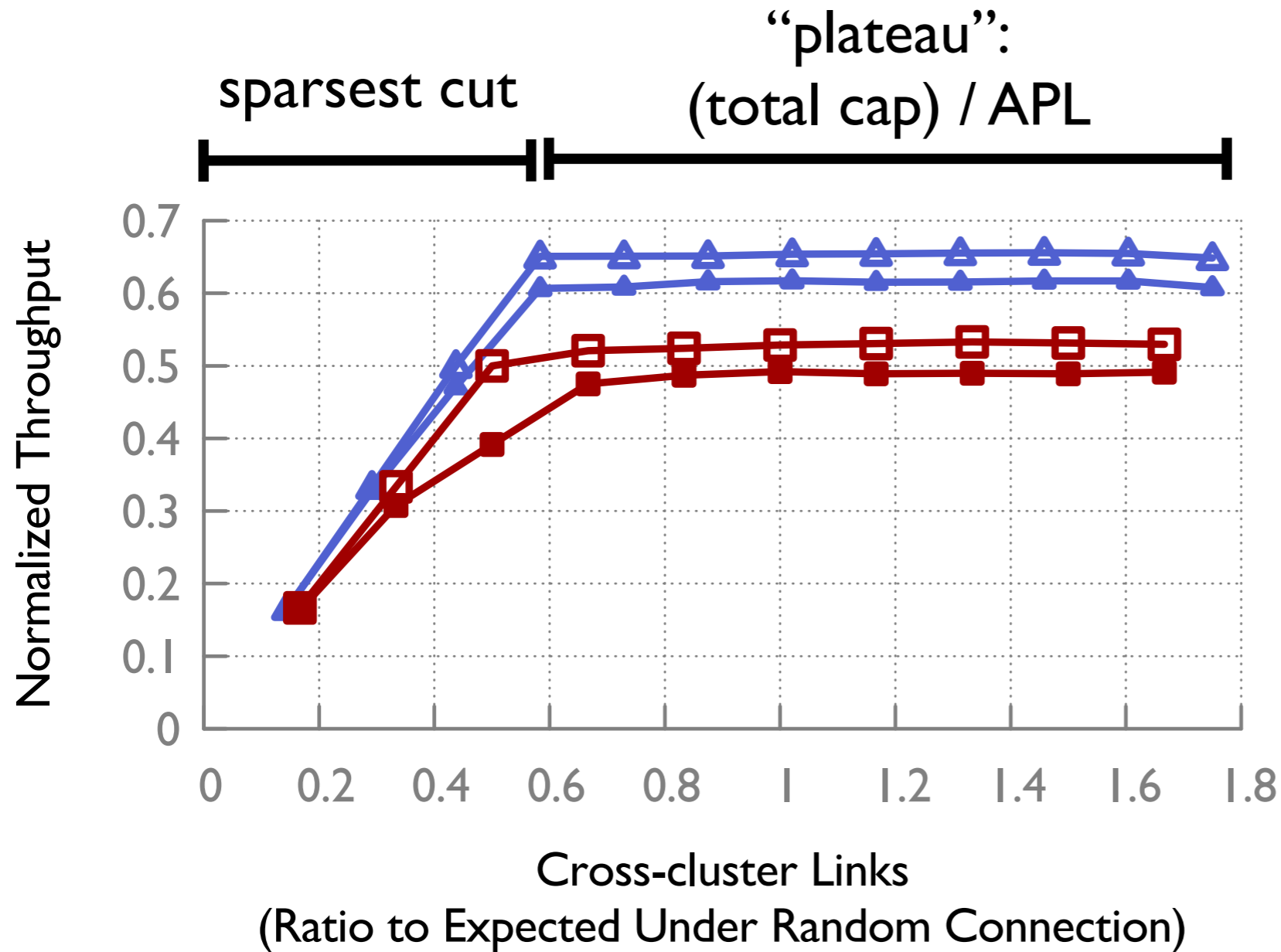


And constant-factor matching lower bounds in special case.

Two regimes of throughput



Two regimes of throughput



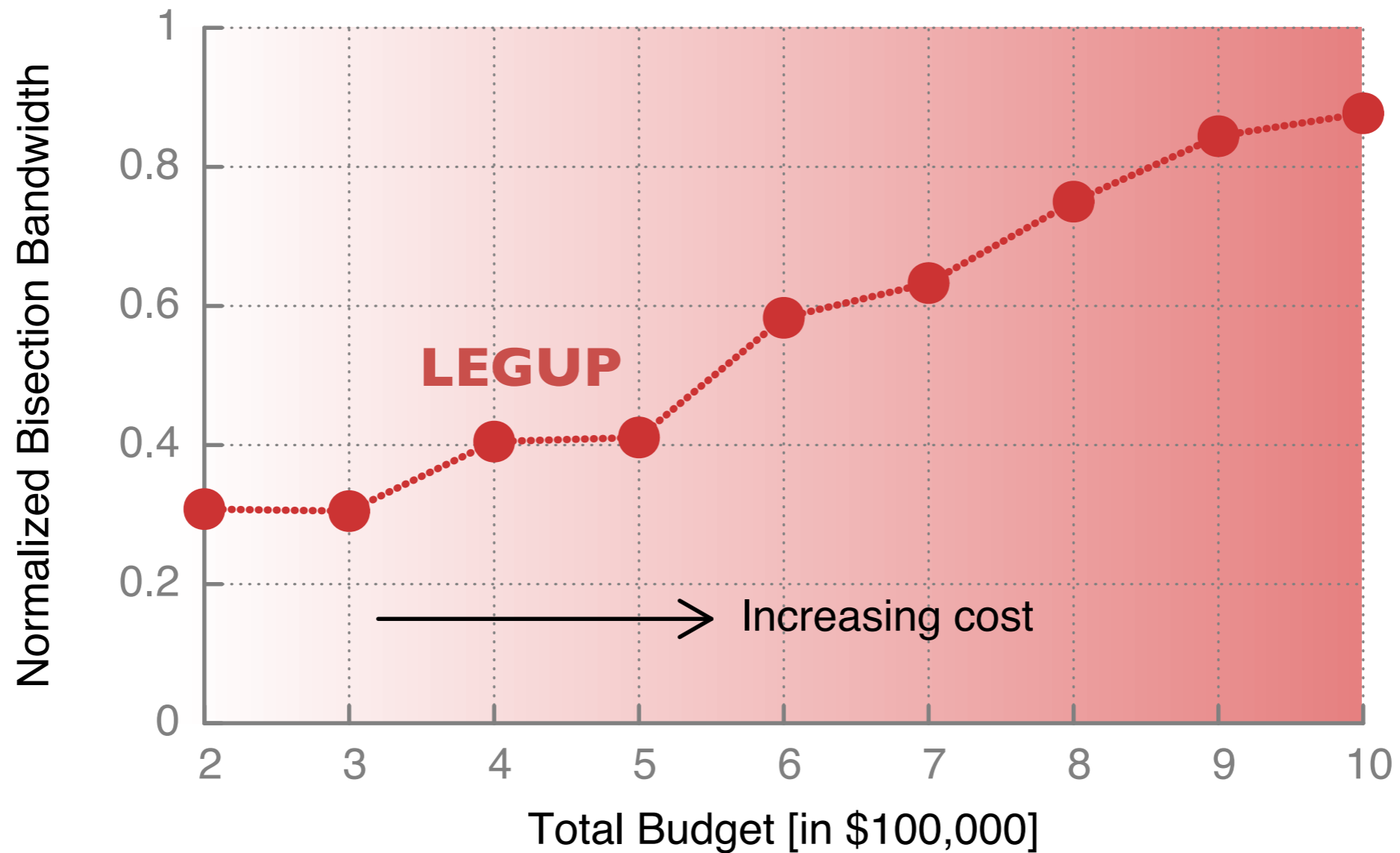
High-capacity switches needn't be clustered

Bisection bandwidth is poor predictor of performance!

Cables can be localized

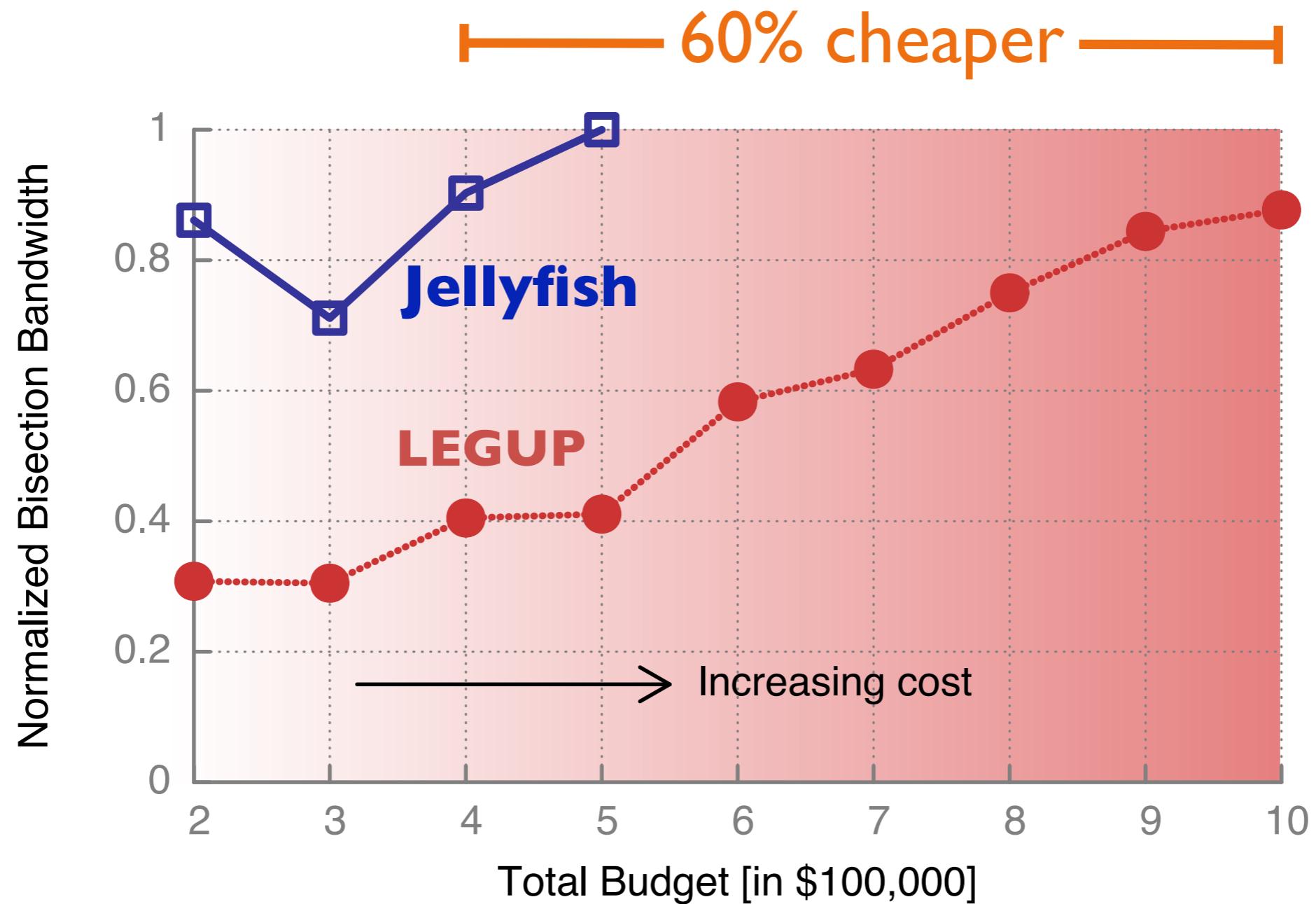
Quantifying Expandability

Quantifying expandability



LEGUP: [Curtis, Keshav, Lopez-Ortiz, CoNEXT'10]

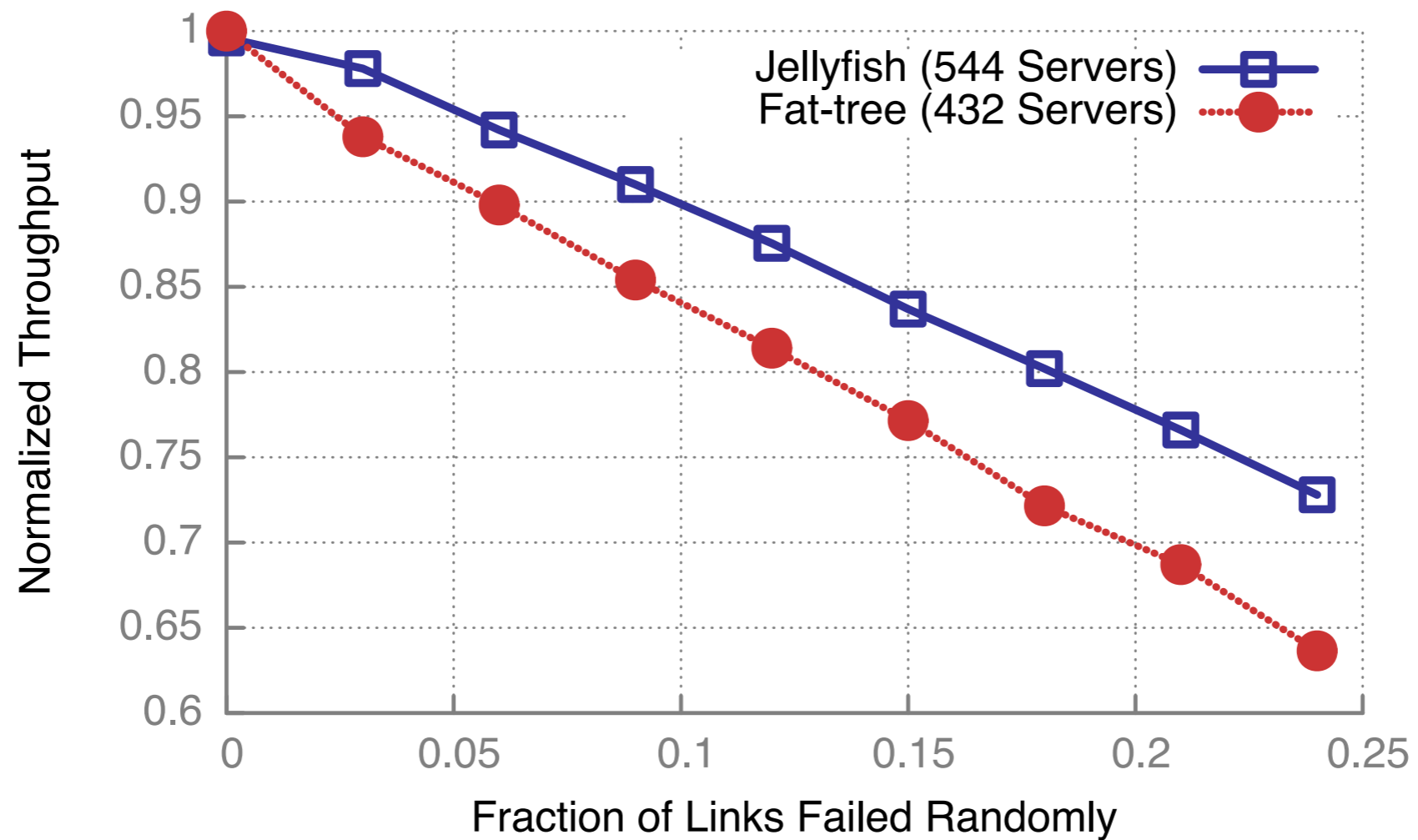
Quantifying expandability



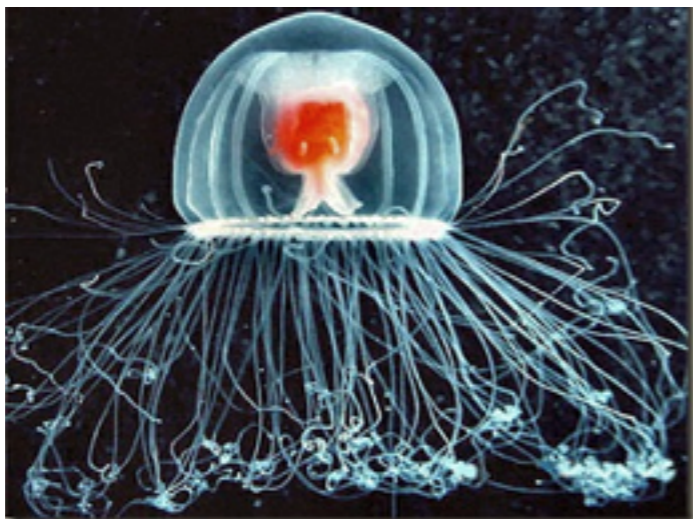
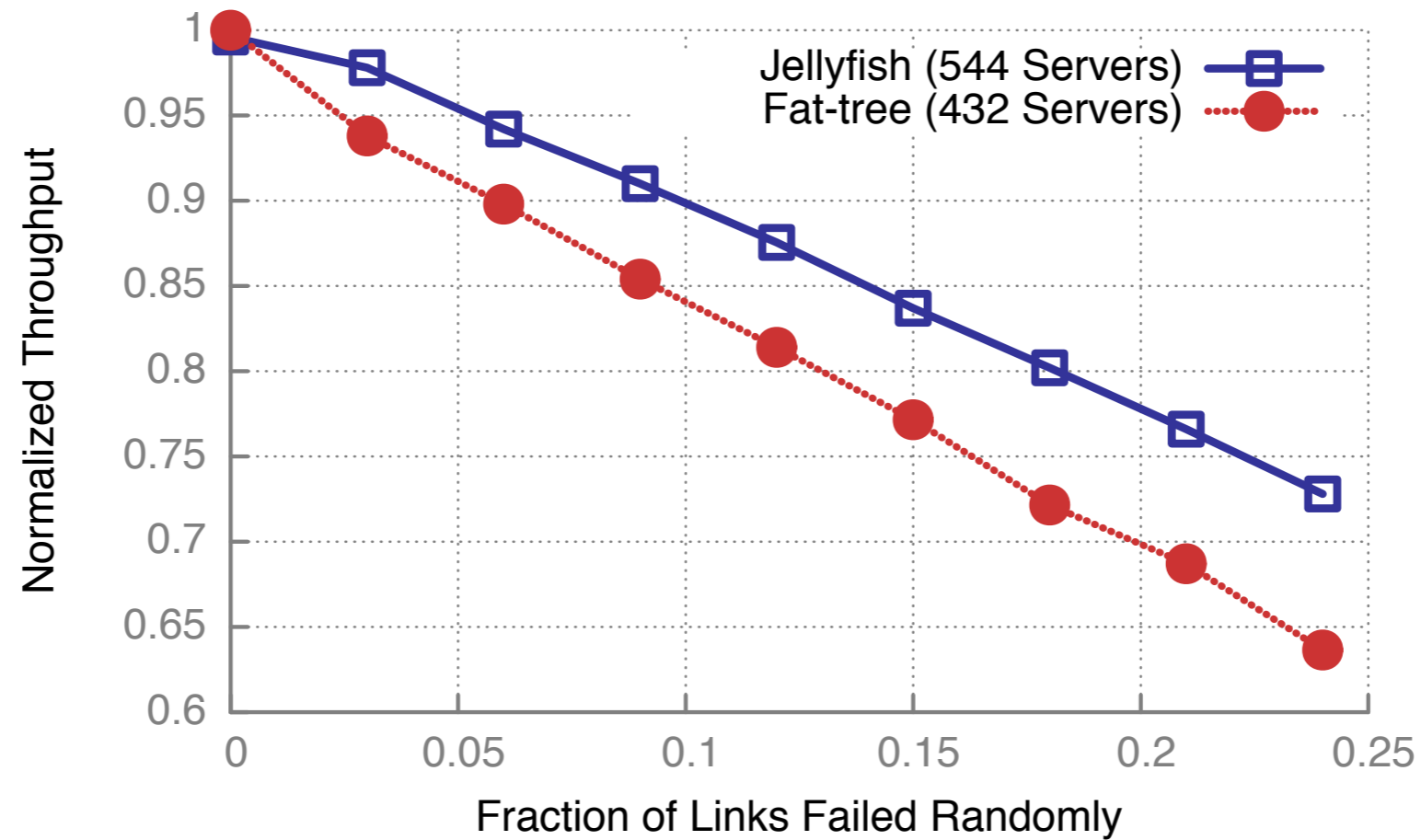
LEGUP: [Curtis, Keshav, Lopez-Ortiz, CoNEXT'10]

Failure Resilience

Throughput under link failures



Throughput under link failures



Turritopsis Nutricula?

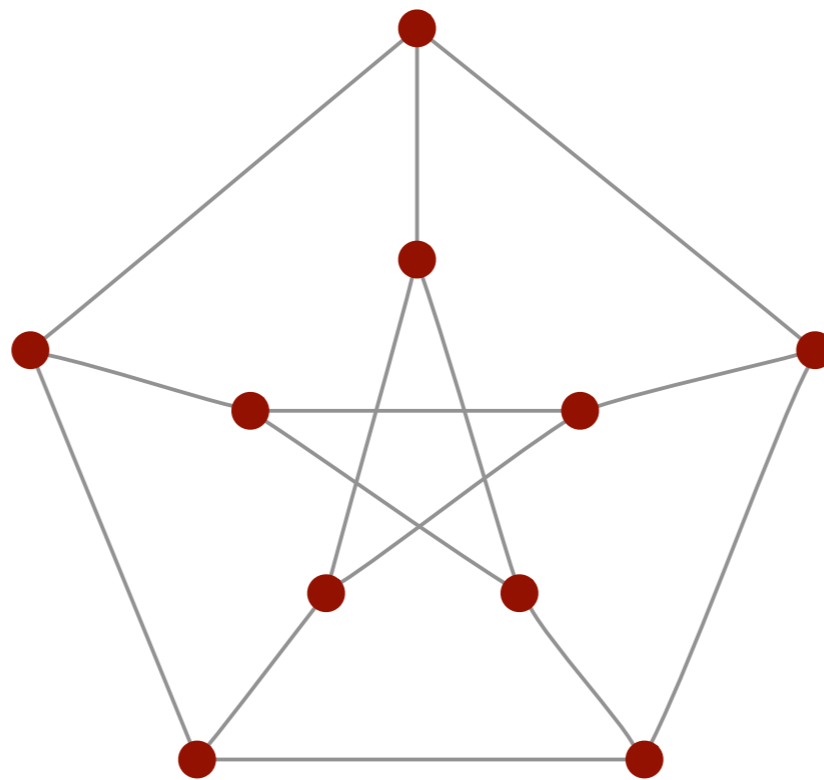
Beyond Random Graphs

Can we do even better?

What is the maximum number of nodes in any graph with degree \hat{d} and diameter d ?

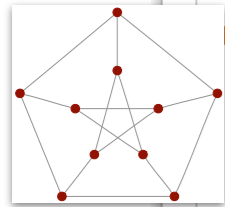
Can we do even better?

What is the maximum number of nodes in any graph with **degree 3** and **diameter 2**?



Peterson graph

Degree-diameter problem



LARGEST KNOWN (Δ, D) -GRAPHS. June 2010.

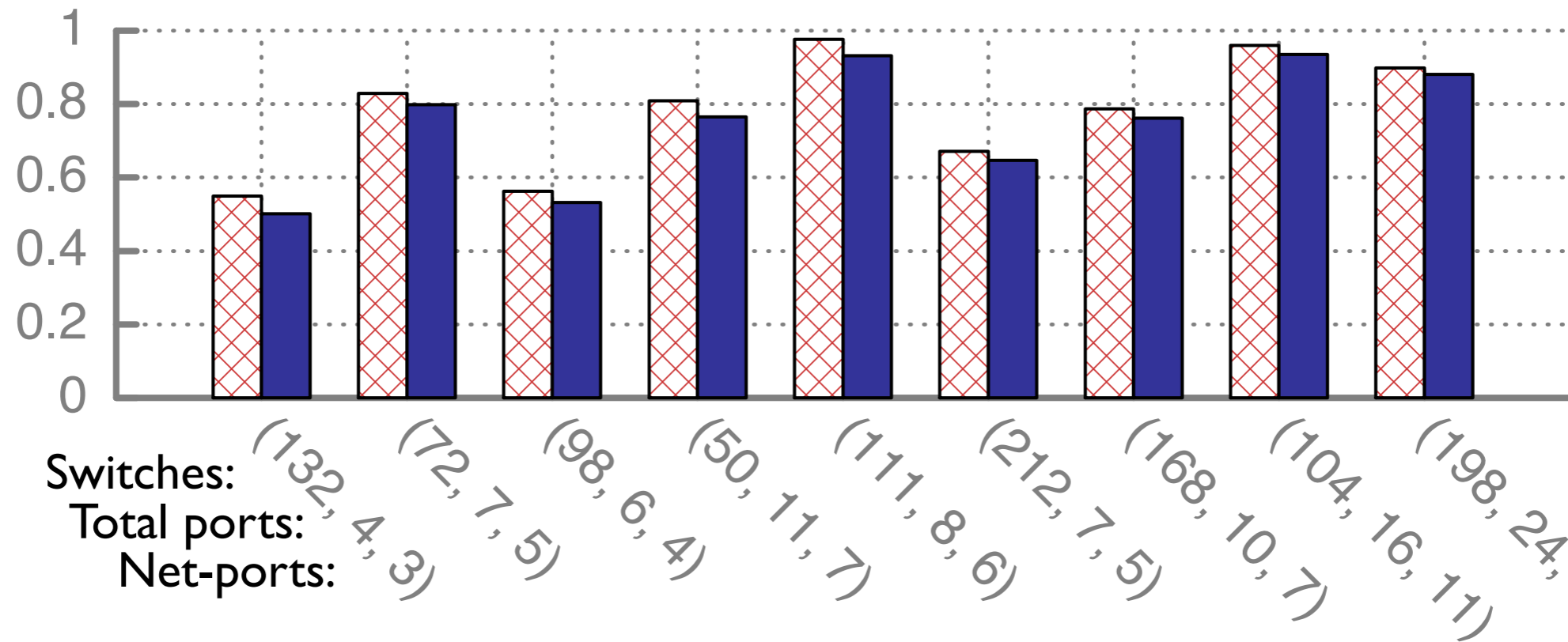
		Diameter								
		2	3	4	5	6	7	8	9	10
Degree	D \ D	<u>10</u>	<u>20</u>	<u>38</u>	<u>70</u>	<u>132</u>	<u>196</u>	<u>336</u>	<u>600</u>	<u>1 250</u>
	3	<u>15</u>	<u>41</u>	<u>98</u>	364	<u>740</u>	<u>1 320</u>	<u>3 243</u>	<u>7 575</u>	<u>17 703</u>
	4	<u>24</u>	<u>72</u>	<u>212</u>	<u>624</u>	<u>2 772</u>	<u>5 516</u>	<u>17 030</u>	<u>53 352</u>	<u>164 720</u>
	5	<u>32</u>	<u>111</u>	<u>390</u>	<u>1 404</u>	<u>7 917</u>	<u>19 282</u>	<u>75 157</u>	<u>295 025</u>	<u>1 212 117</u>
	6	<u>50</u>	<u>168</u>	<u>672</u>	<u>2 756</u>	<u>11 988</u>	<u>52 768</u>	<u>233 700</u>	<u>1 124 990</u>	<u>5 311 572</u>
	7	57	<u>253</u>	<u>1 100</u>	<u>5 060</u>	<u>39 672</u>	<u>130 017</u>	<u>714 010</u>	<u>4 039 704</u>	<u>17 823 532</u>
	8	74	585	<u>1 550</u>	<u>8 200</u>	<u>75 893</u>	<u>270 192</u>	<u>1 485 498</u>	<u>10 423 212</u>	<u>31 466 244</u>
	9	91	650	<u>2 223</u>	<u>13 140</u>	<u>134 690</u>	<u>561 957</u>	<u>4 019 736</u>	<u>17 304 400</u>	<u>104 058 822</u>
	10	<u>104</u>	715	3 200	<u>18 700</u>	156 864	<u>971 028</u>	<u>5 941 864</u>	<u>62 932 488</u>	<u>250 108 668</u>
	11	133	<u>786</u>	4 680	<u>29 470</u>	<u>359 772</u>	<u>1 900 464</u>	<u>10 423 212</u>	<u>104 058 822</u>	<u>600 105 100</u>
	12	<u>162</u>	<u>851</u>	6 560	<u>39 576</u>	531 440	<u>2 901 404</u>	<u>17 823 532</u>	<u>180 002 472</u>	<u>1 050 104 118</u>
	13	183	<u>916</u>	8 200	<u>56 790</u>	<u>816 294</u>	6 200 460	<u>41 894 424</u>	<u>450 103 771</u>	<u>2 050 103 984</u>
	14	186	1 215	11 712	<u>74 298</u>	1 417 248	<u>8 079 298</u>	<u>90 001 236</u>	<u>900 207 542</u>	<u>4 149 702 144</u>
	15	<u>198</u>	1 600	14 640	132 496	1 771 560	14 882 658	<u>104 518 518</u>	<u>1 400 103 920</u>	7 394 669 856

Degree-diameter problem

Do the best known degree-diameter graphs also work well for high throughput?

Degree-diameter vs. Jellyfish

Normalized Throughput



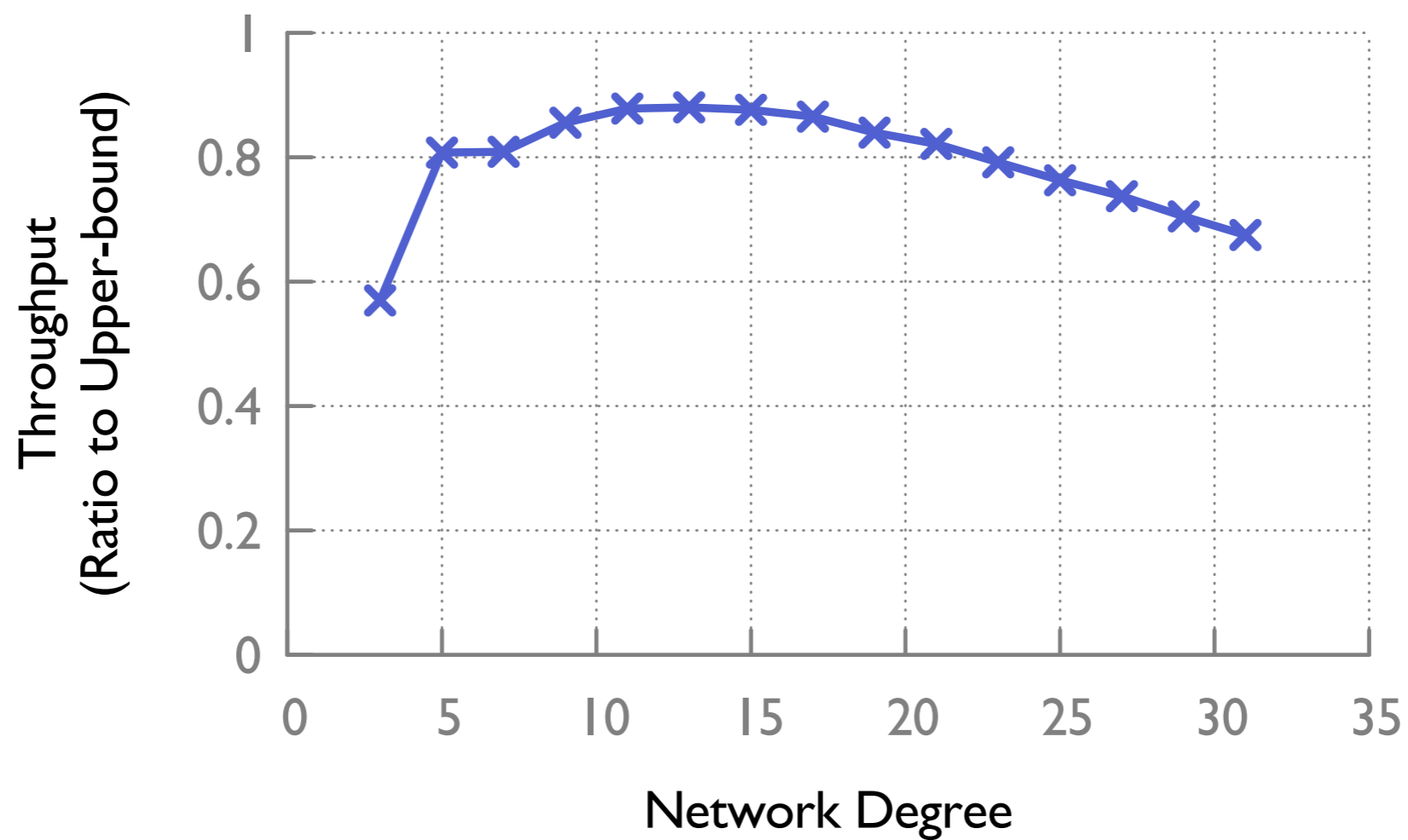
Best-known Degree-Diameter Graph 
 Jellyfish 

D-D graphs **do** have high throughput

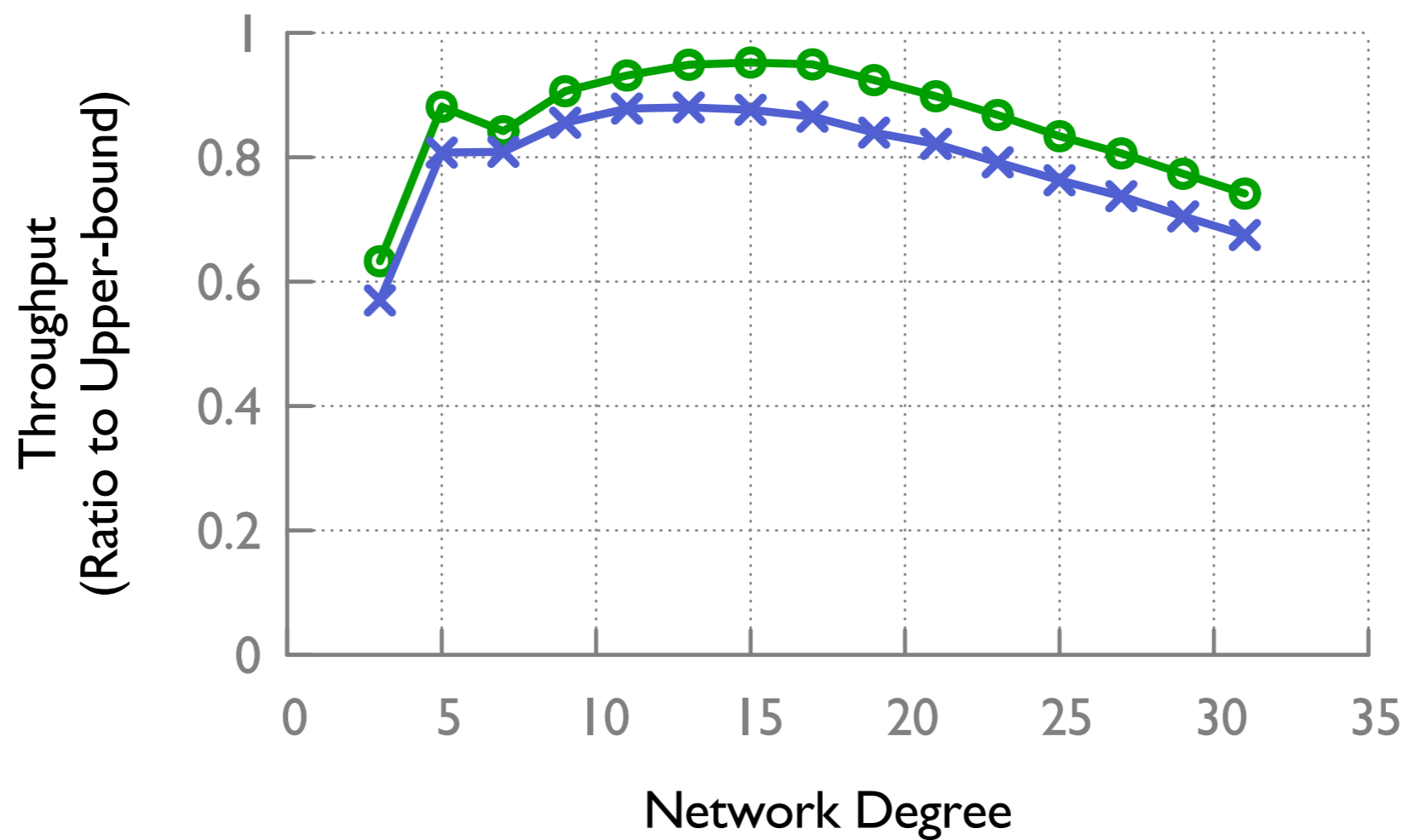
Jellyfish within 9%!

Random graphs vs. upper bound
for fixed size and increasing degree

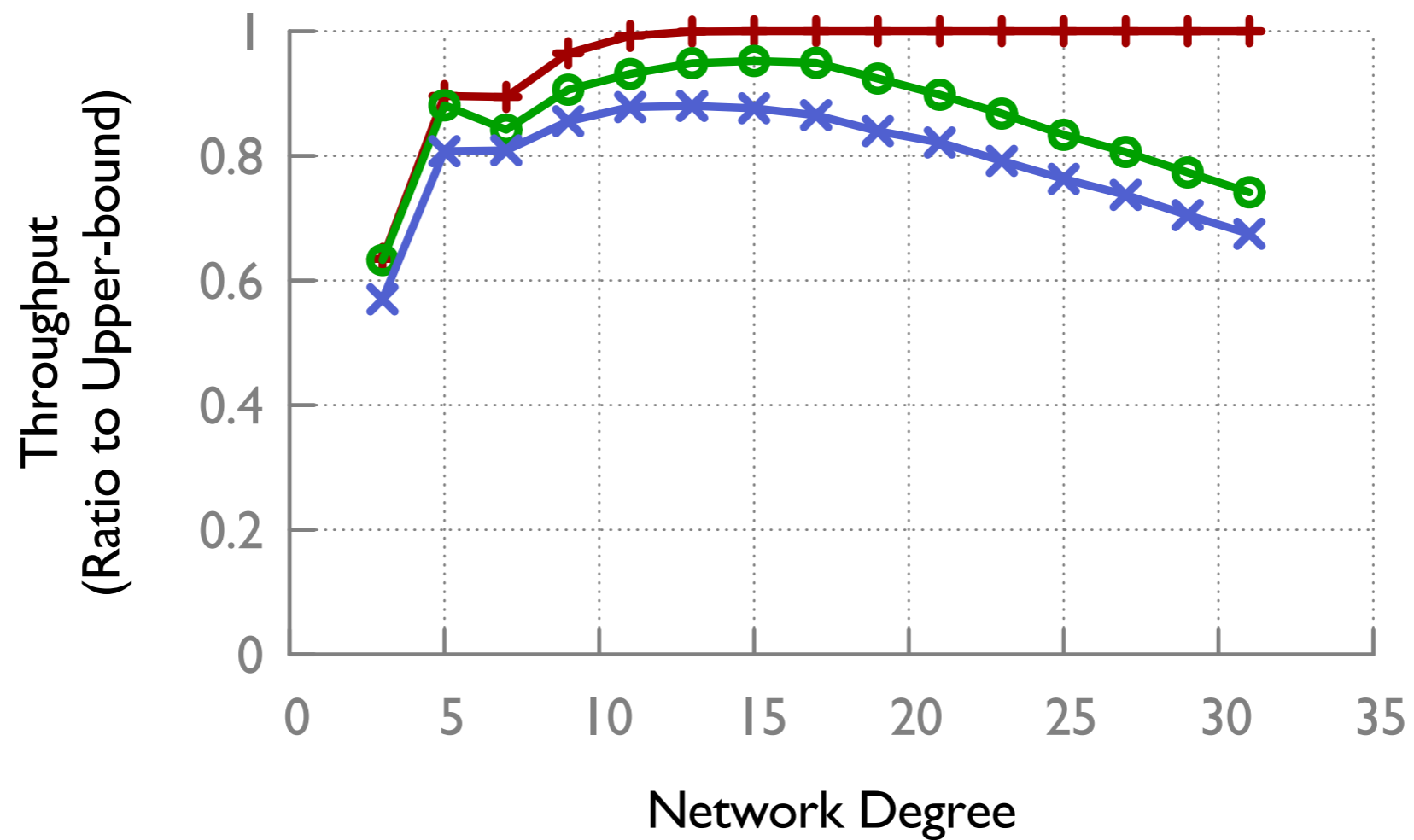
Random graphs vs. upper bound



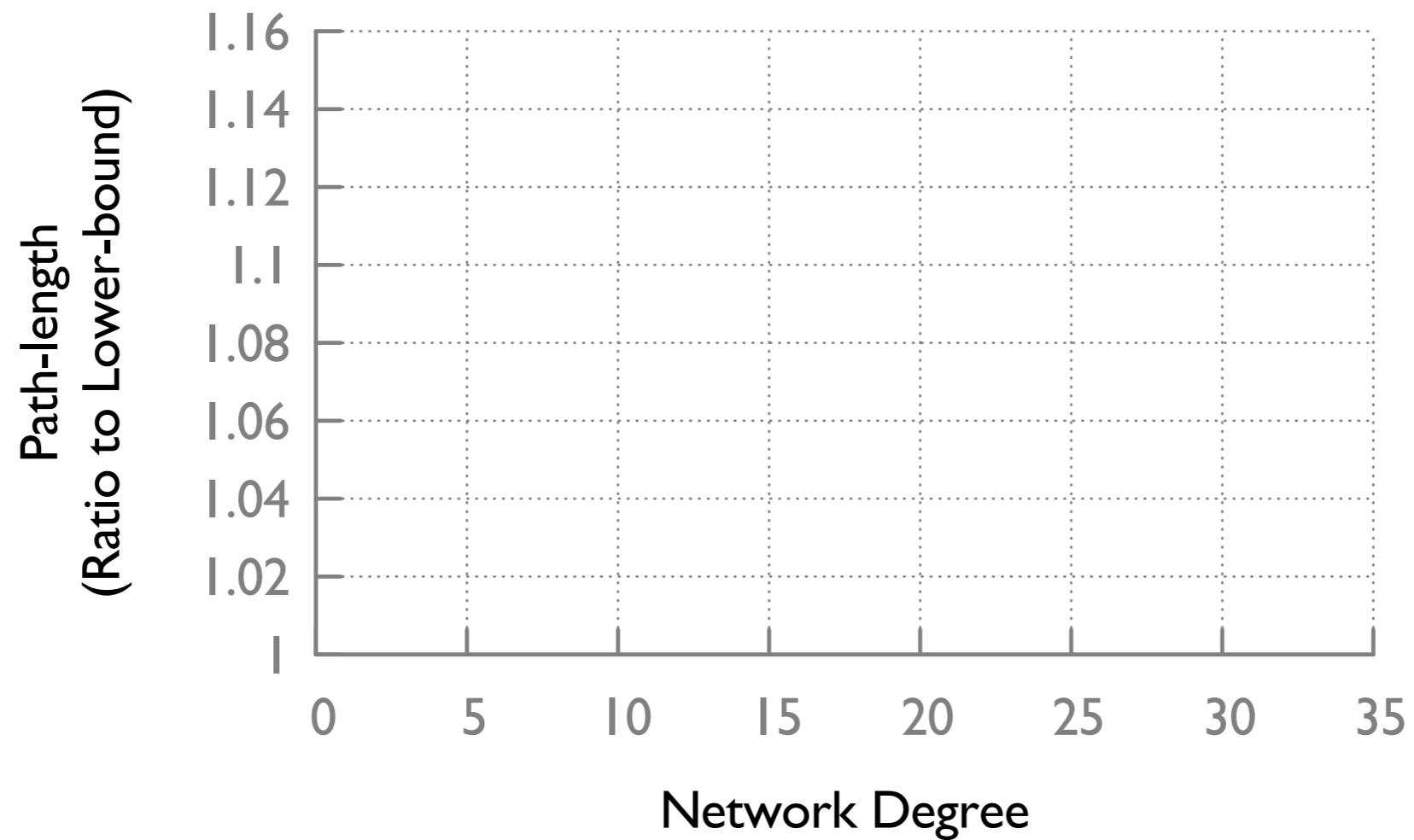
Random graphs vs. upper bound



Random graphs vs. upper bound



Random graphs vs. upper bound



Random graphs vs. upper bound

