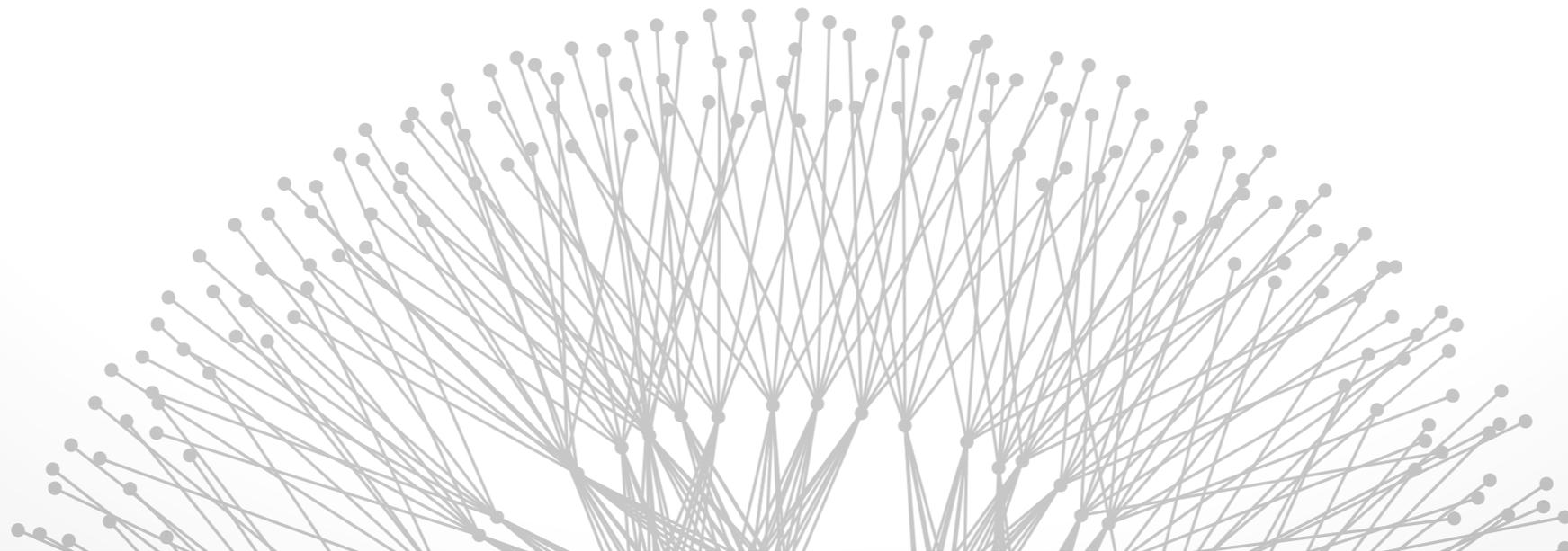


High Throughput Data Center Topology Design

Ankit Singla, P. Brighten Godfrey, Alexandra Kolla







***“How long must we wait
until our **pigeon system**
rivals those of the
Continental Powers?”***

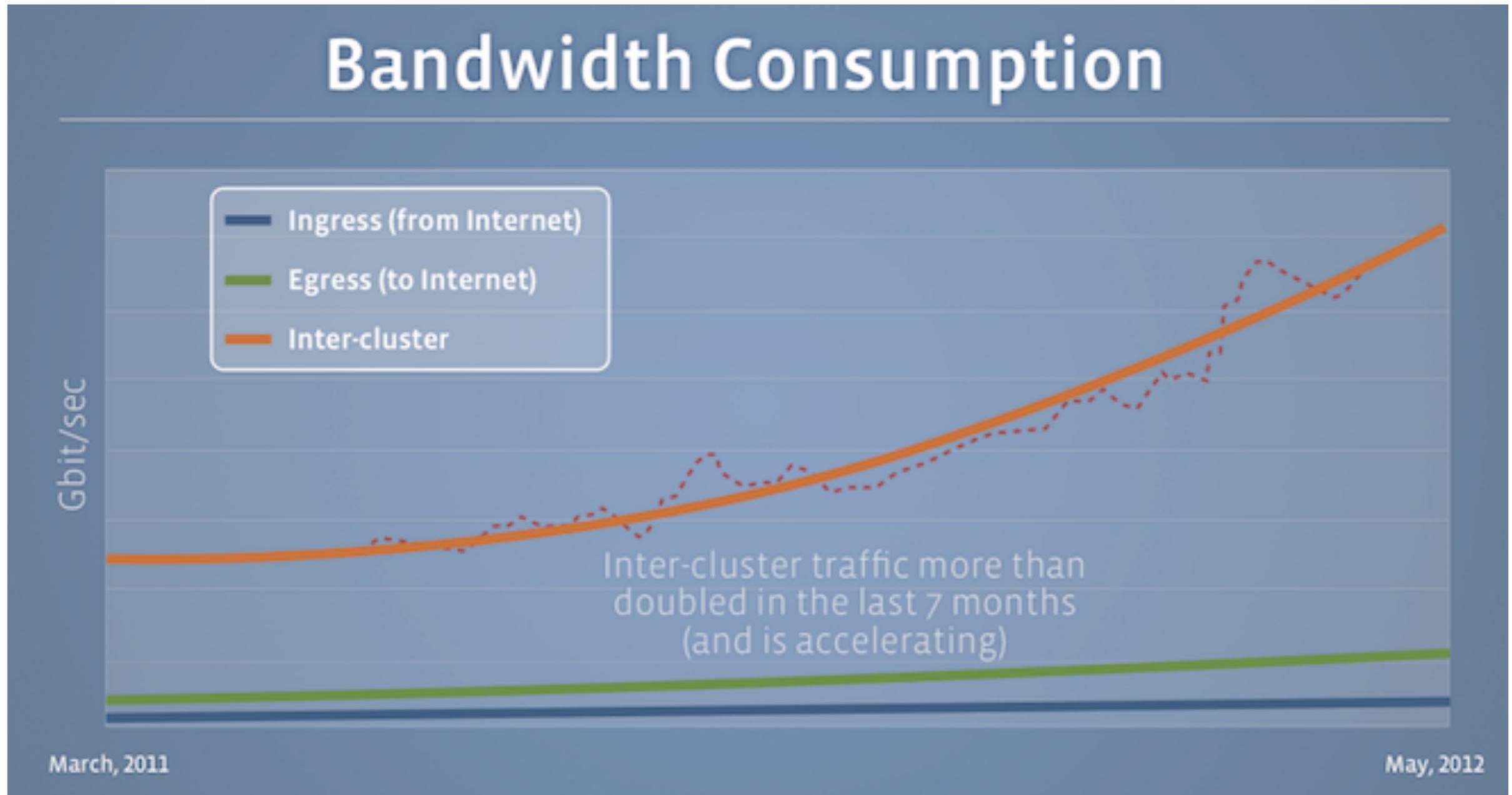
- The Nineteenth Century, 1899



Google

google.com/datacenter/

The need for throughput

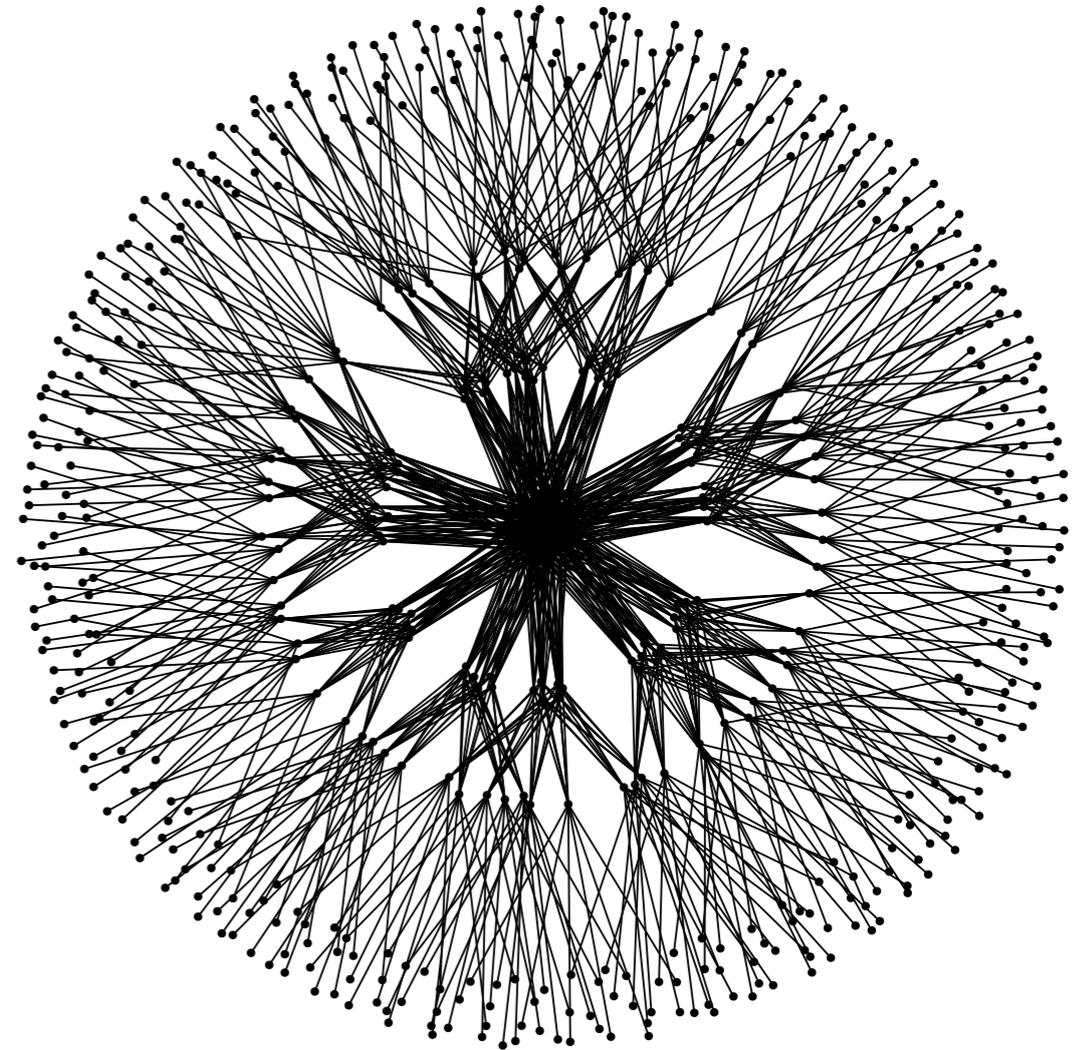
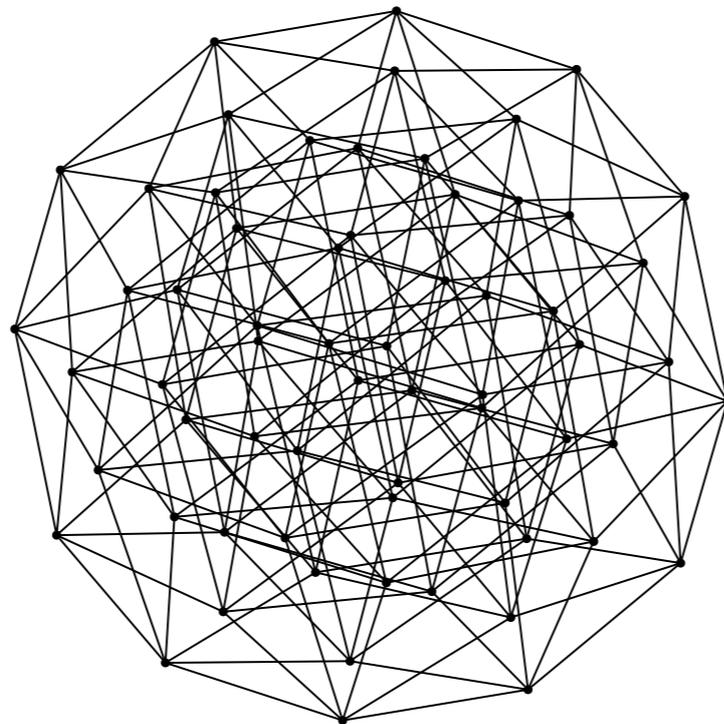
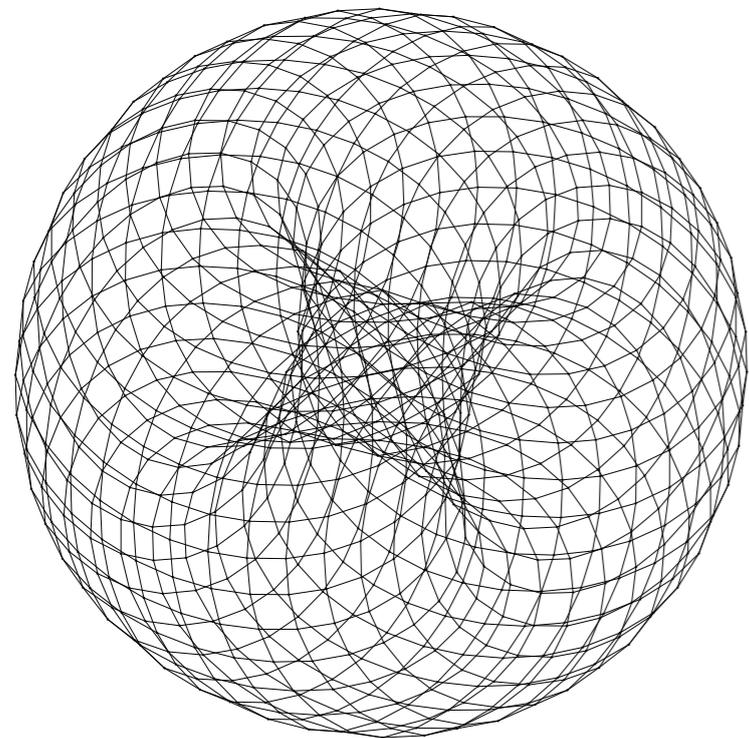


March
2011

May
2012

[Facebook, via Wired]

Many topology options ...



How do we design throughput
optimal network topologies?

How do we
design throughput
optimal network
topologies?

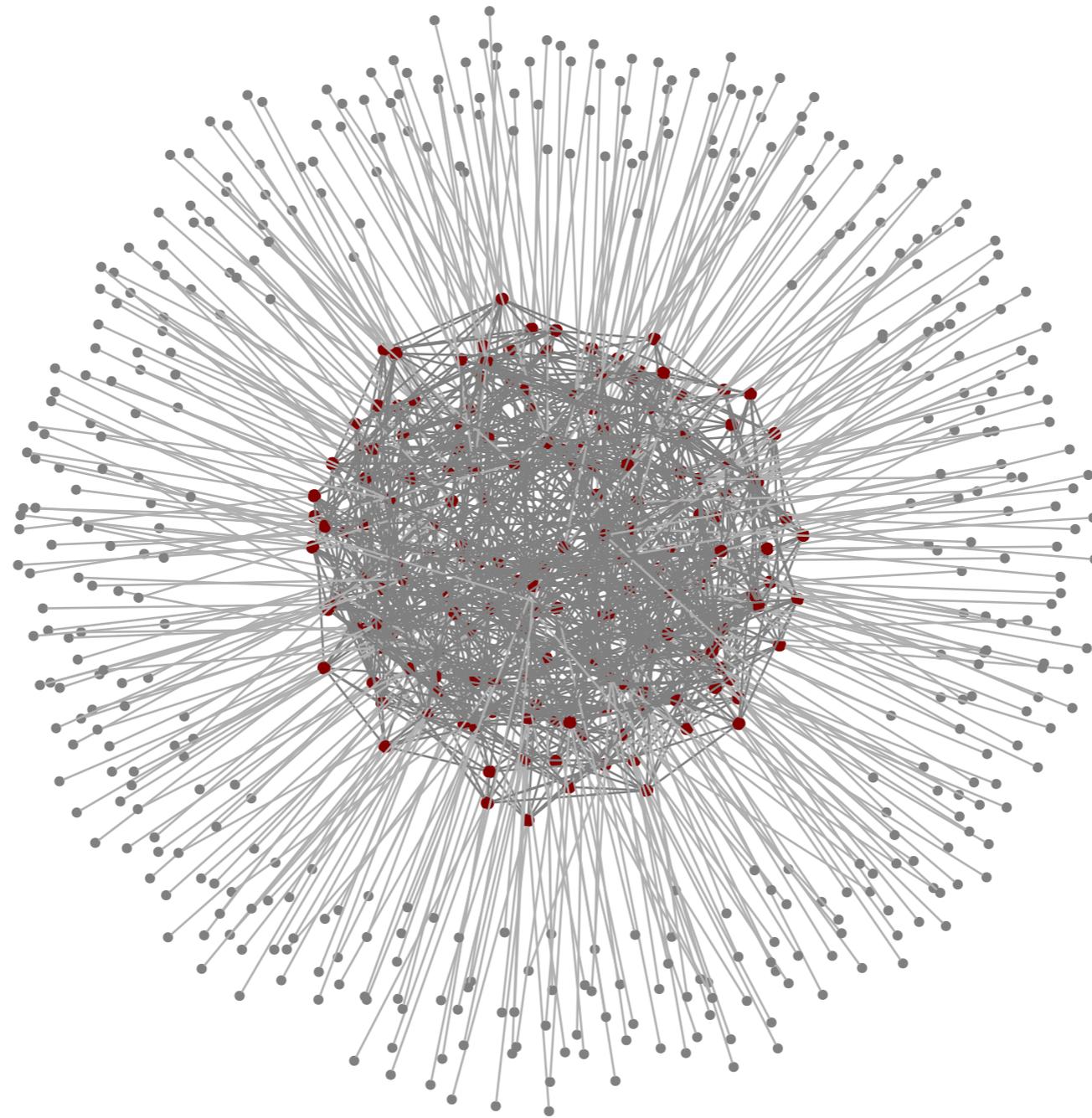


How close can we get to **optimal**
network capacity?

1 How close can we get to **optimal** network capacity?

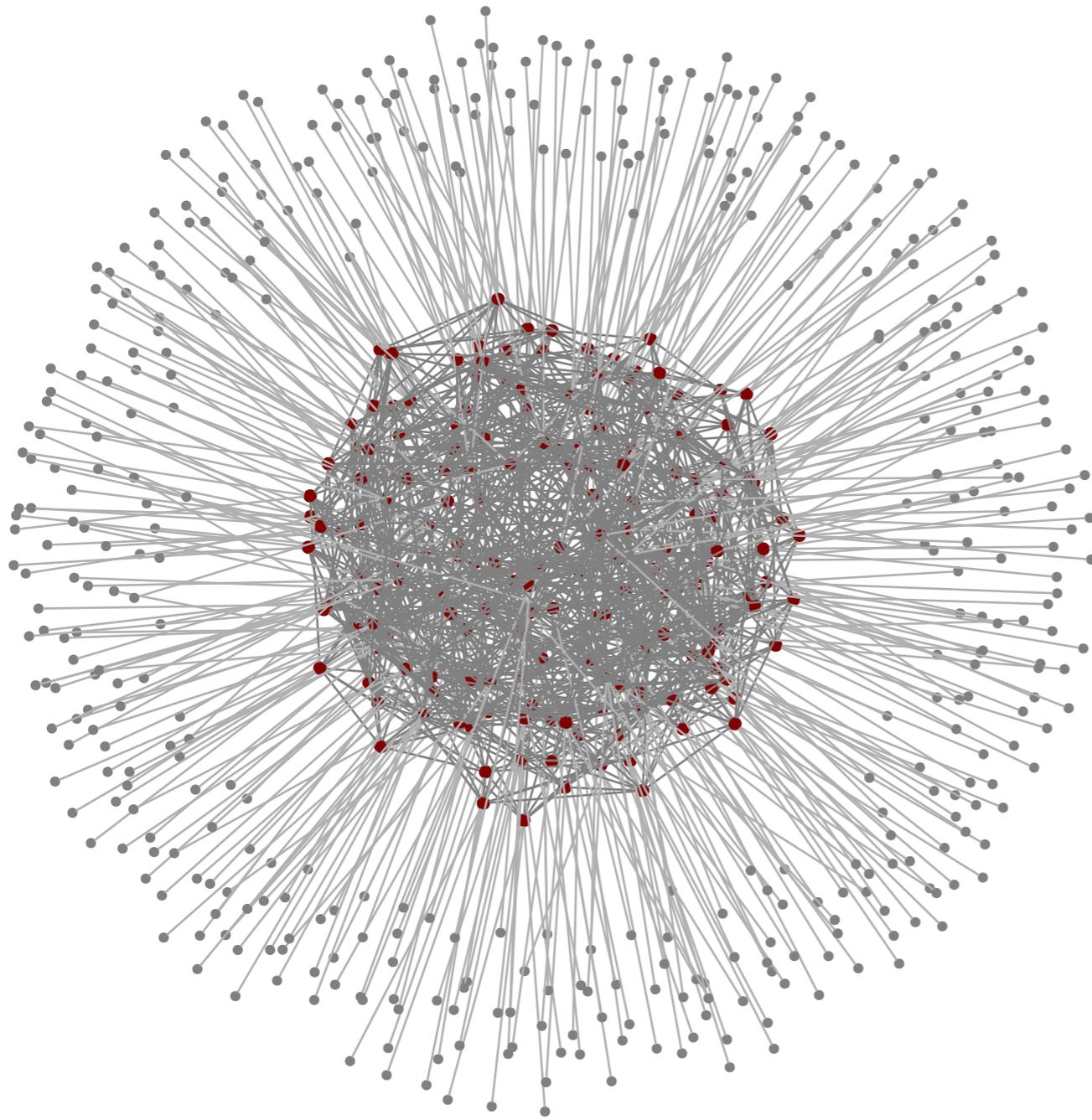
2 How do we handle **heterogeneity**?

Jellyfish: Networking Data Centers Randomly



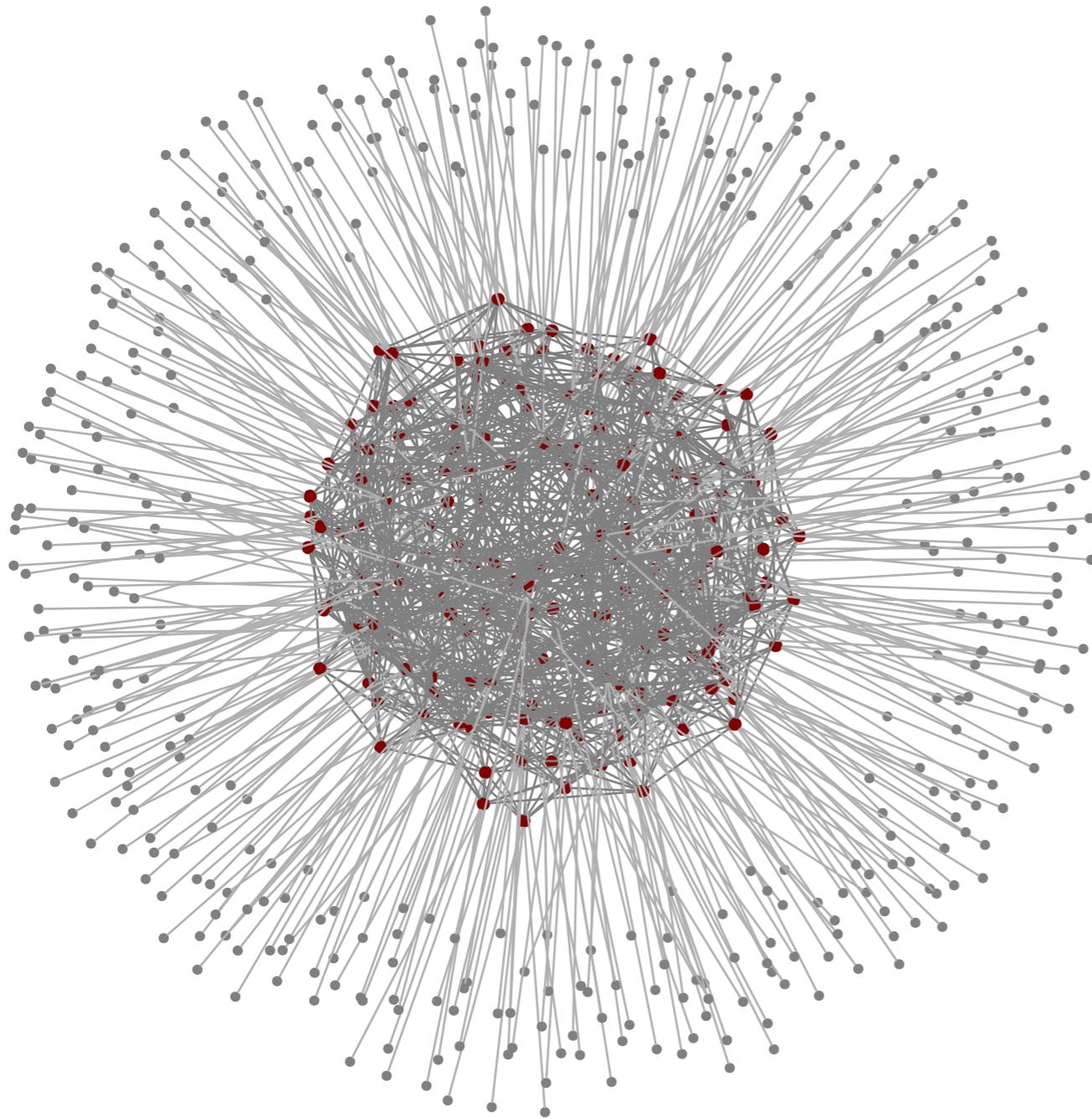
[NSDI 2012: Singla, Hong, Popa, Godfrey]

Jellyfish: Networking Data Centers Randomly



[NSDI 2012: Singla, Hong, Popa, Godfrey]

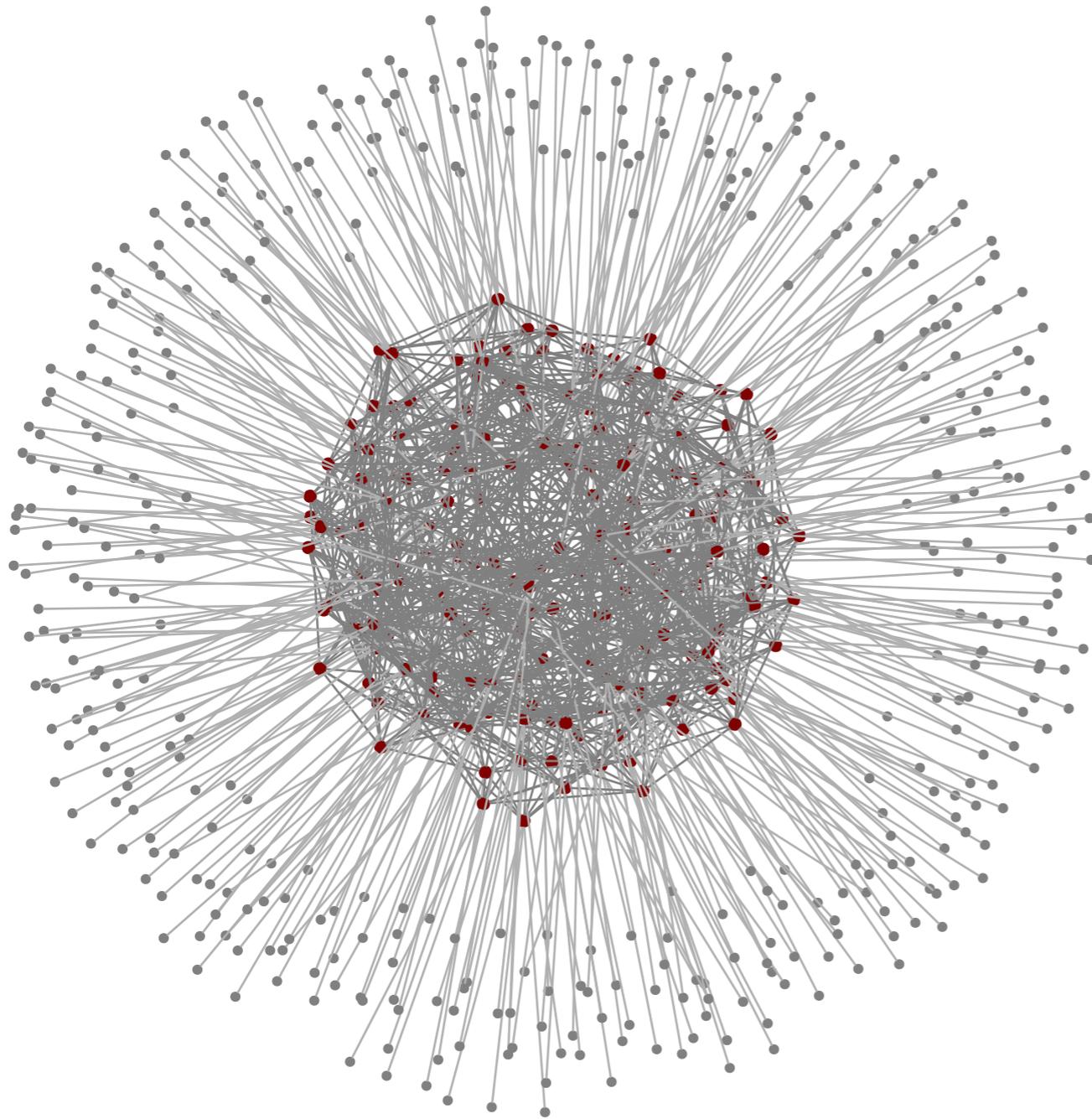
Jellyfish: Networking Data Centers Randomly



- High capacity
 - Beat fat-trees by 25%+

[NSDI 2012: Singla, Hong, Popa, Godfrey]

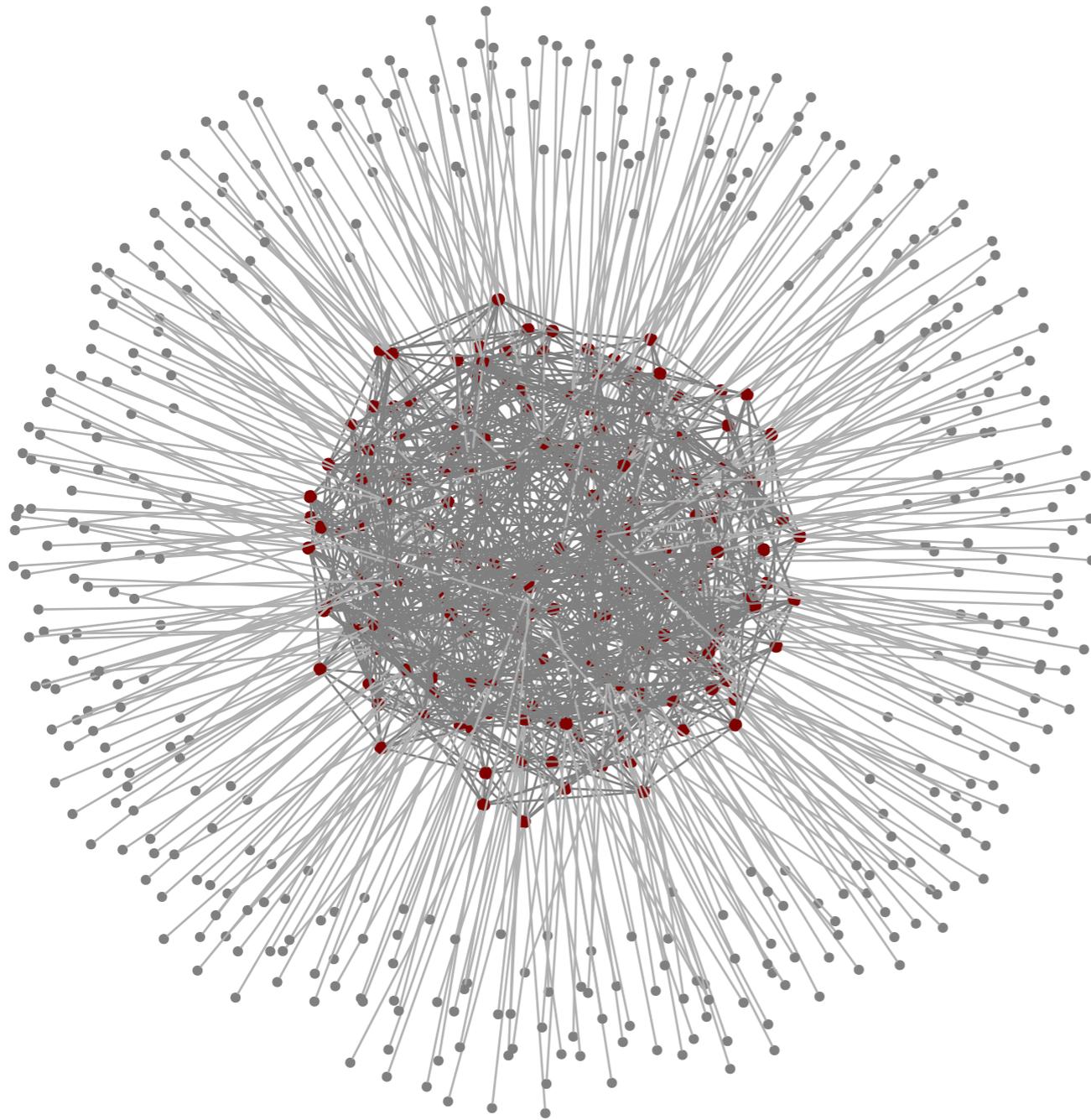
Jellyfish: Networking Data Centers Randomly



- High capacity
 - Beat fat-trees by 25%+
- Easier to expand
 - 60% cheaper expansion

[NSDI 2012: Singla, Hong, Popa, Godfrey]

Jellyfish: Networking Data Centers Randomly



- High capacity
 - Beat fat-trees by 25%+
- Easier to expand
 - 60% cheaper expansion
- Routing and cabling are solvable problems

[NSDI 2012: Singla, Hong, Popa, Godfrey]



How close can we get to **optimal**
network capacity?

- 1 How close can we get to **optimal** network capacity?
- 2 How do we handle **heterogeneity**?

How do we measure
throughput?

How do we measure throughput?

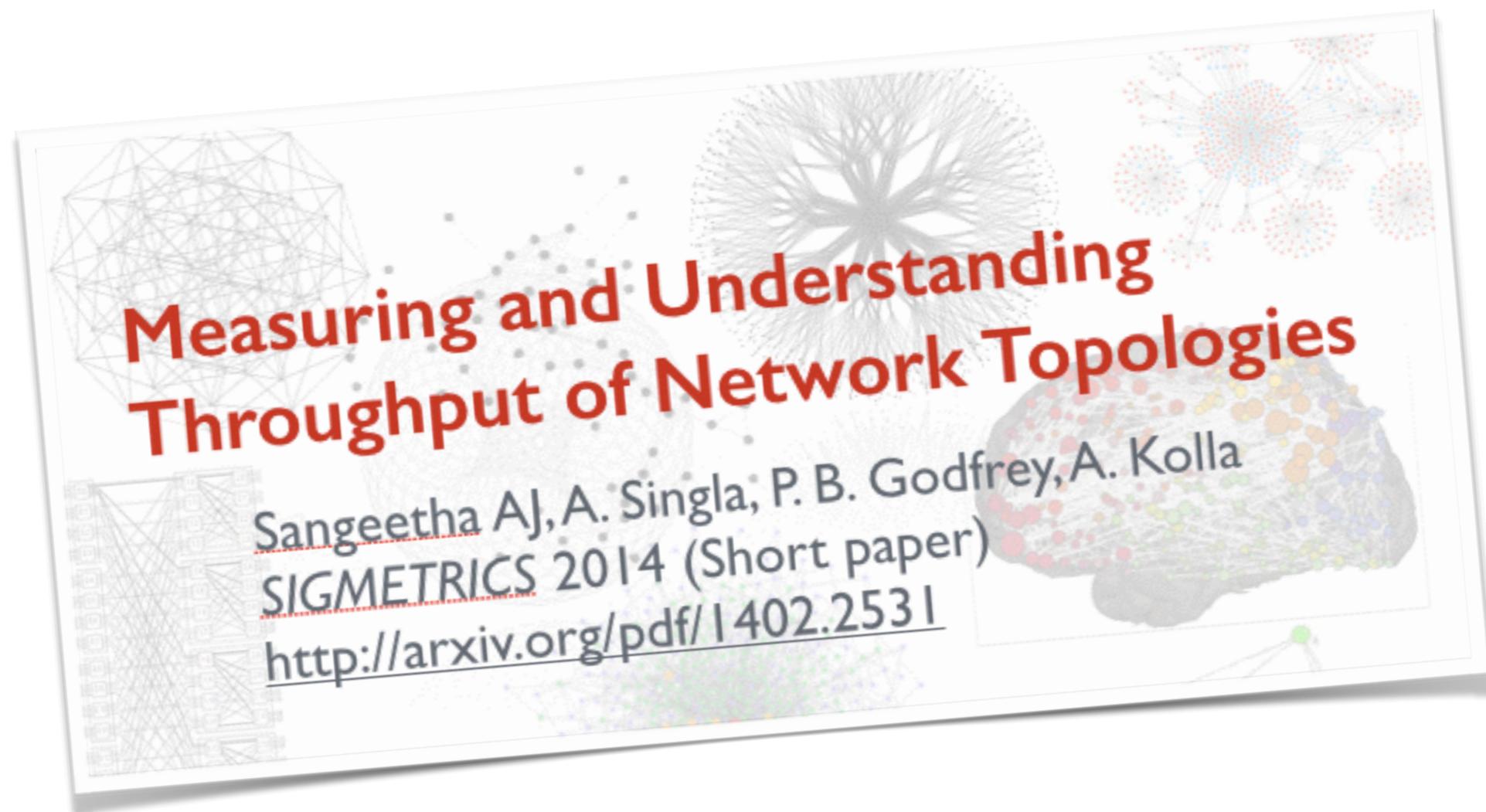
Maximize the minimum flow

How do we measure throughput?

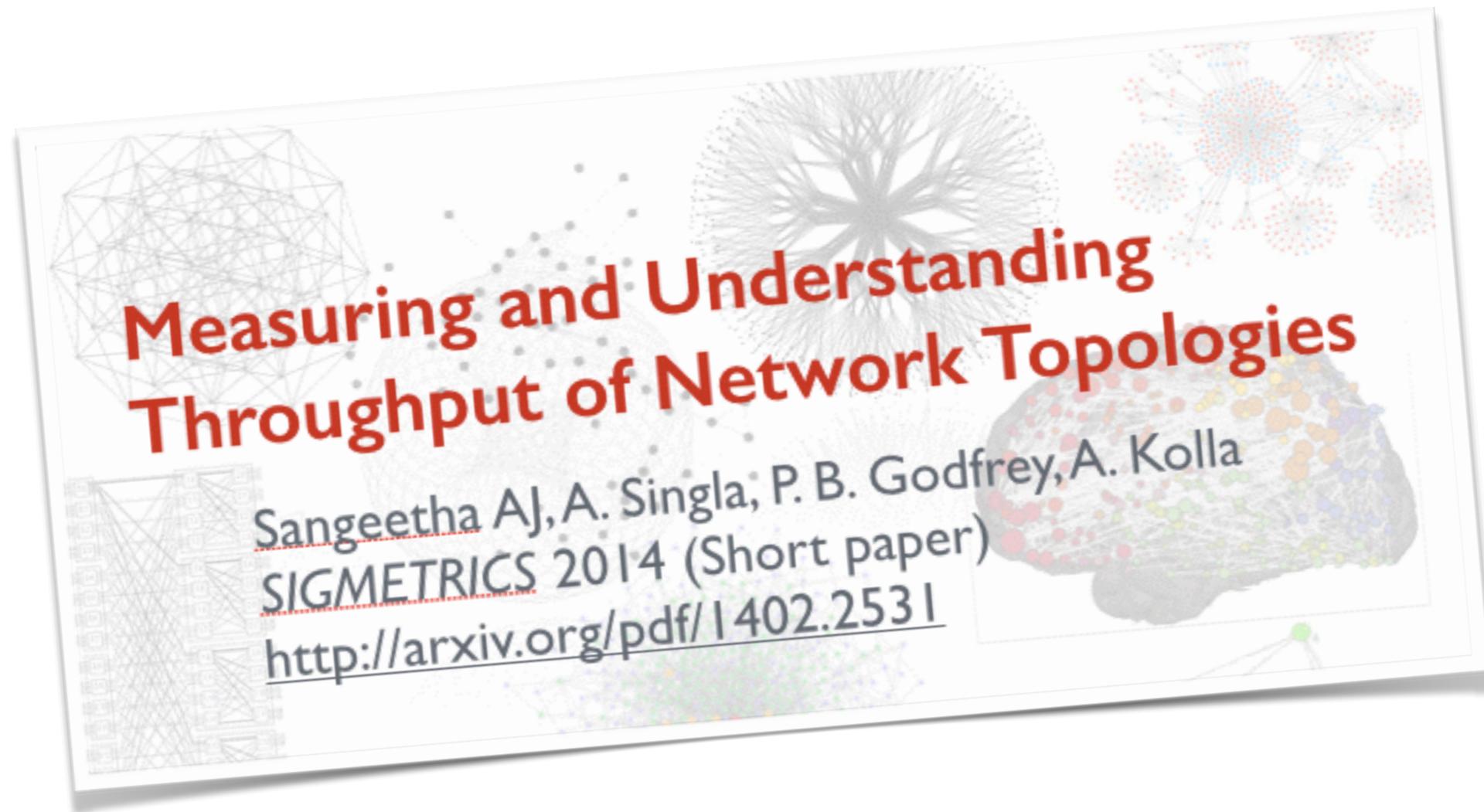
Maximize the minimum flow
under random permutation traffic

How do we measure throughput?

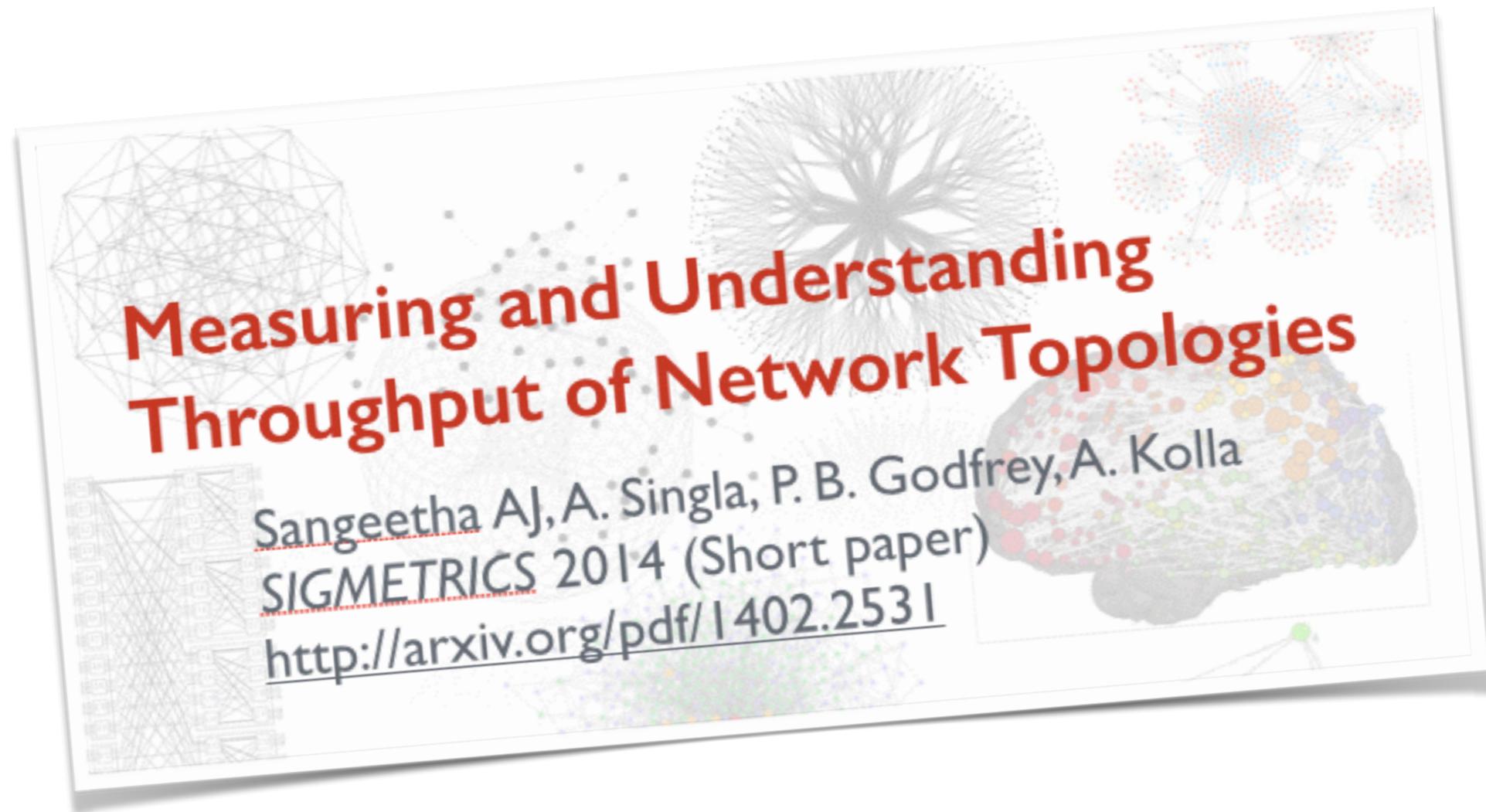
Maximize the minimum flow
under random permutation traffic



How do we measure throughput?

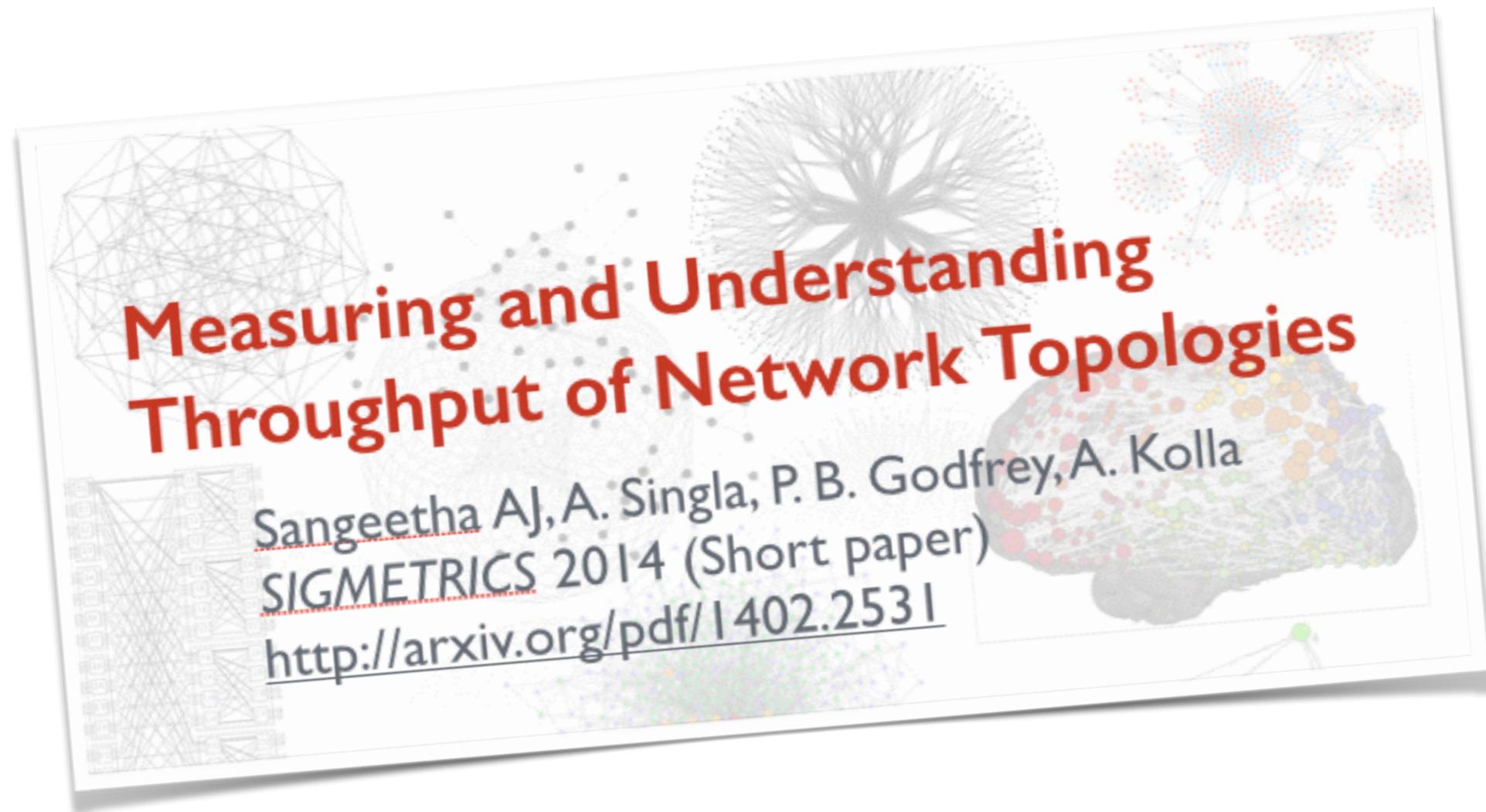


How do we measure throughput?



- Bisection bandwidth \neq throughput

How do we measure throughput?



- Bisection bandwidth \neq throughput
- Near-worst case traffic patterns

How close can we get to
optimal network capacity?

A simple upper bound

A simple upper bound

flows

A simple upper bound

flows • capacity used per flow

A simple upper bound

flows • capacity used per flow

\leq total capacity

A simple upper bound

flows • capacity used per flow

\leq total capacity

A simple upper bound

flows • throughput per flow • mean path length

\leq total capacity

A simple upper bound

$$\text{throughput per flow} \leq \frac{\text{total capacity}}{\# \text{ flows} \cdot \text{mean path length}}$$

A simple upper bound

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

A simple upper bound

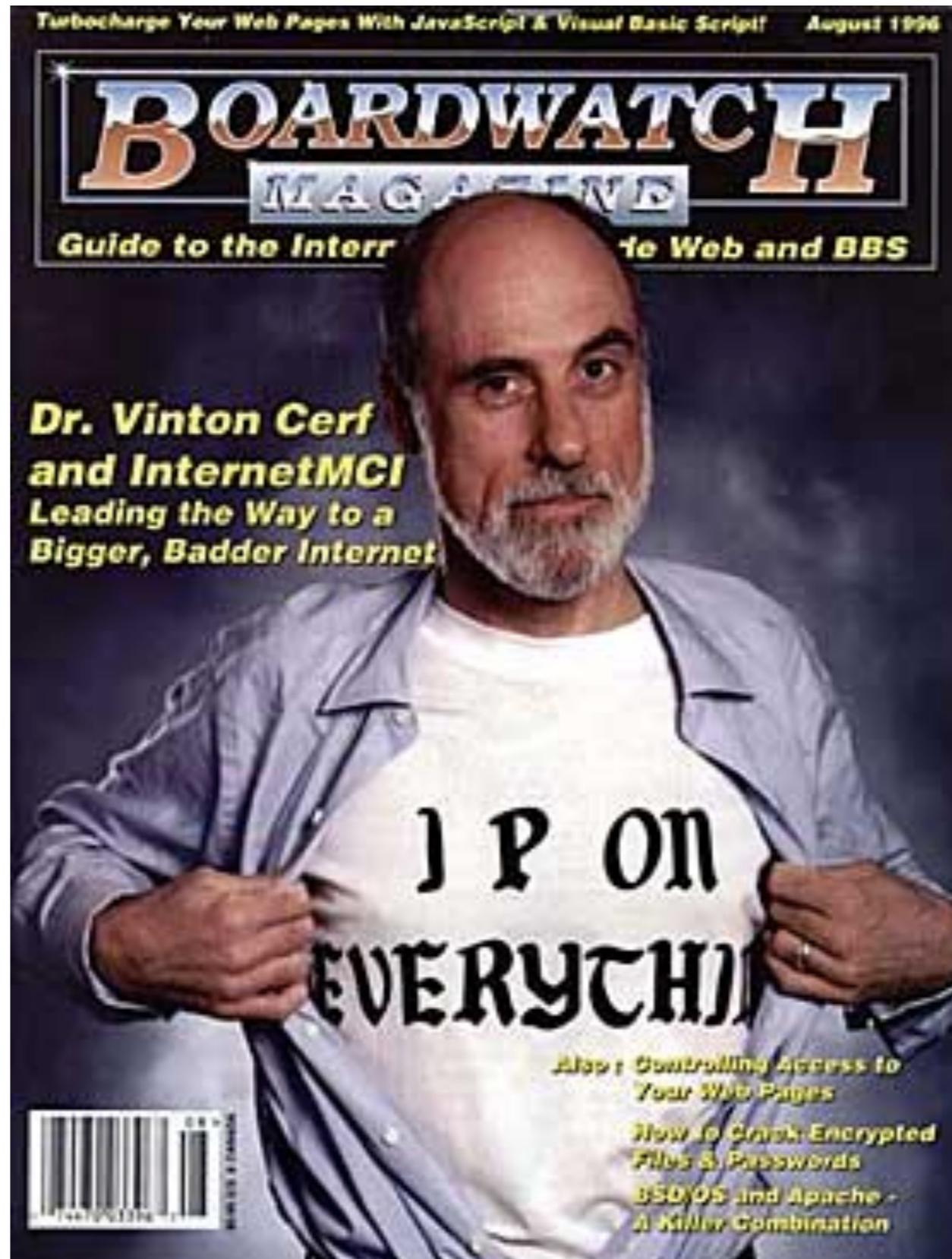
$$\text{throughput per flow} \leq \frac{\sum_{\text{links}} \text{capacity}(\text{link})}{\# \text{ flows} \cdot \text{mean path length}}$$

Lower bound this!



Lower bound on mean path length

Lower bound on mean path length



Lower bound on mean path length

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

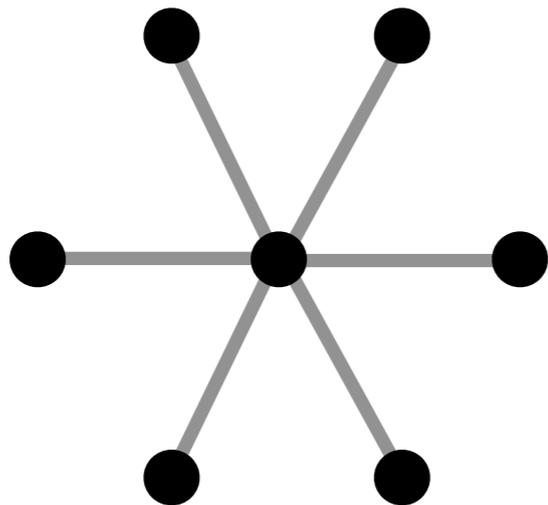
Lower bound on mean path length

Distance	# Nodes
----------	---------



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

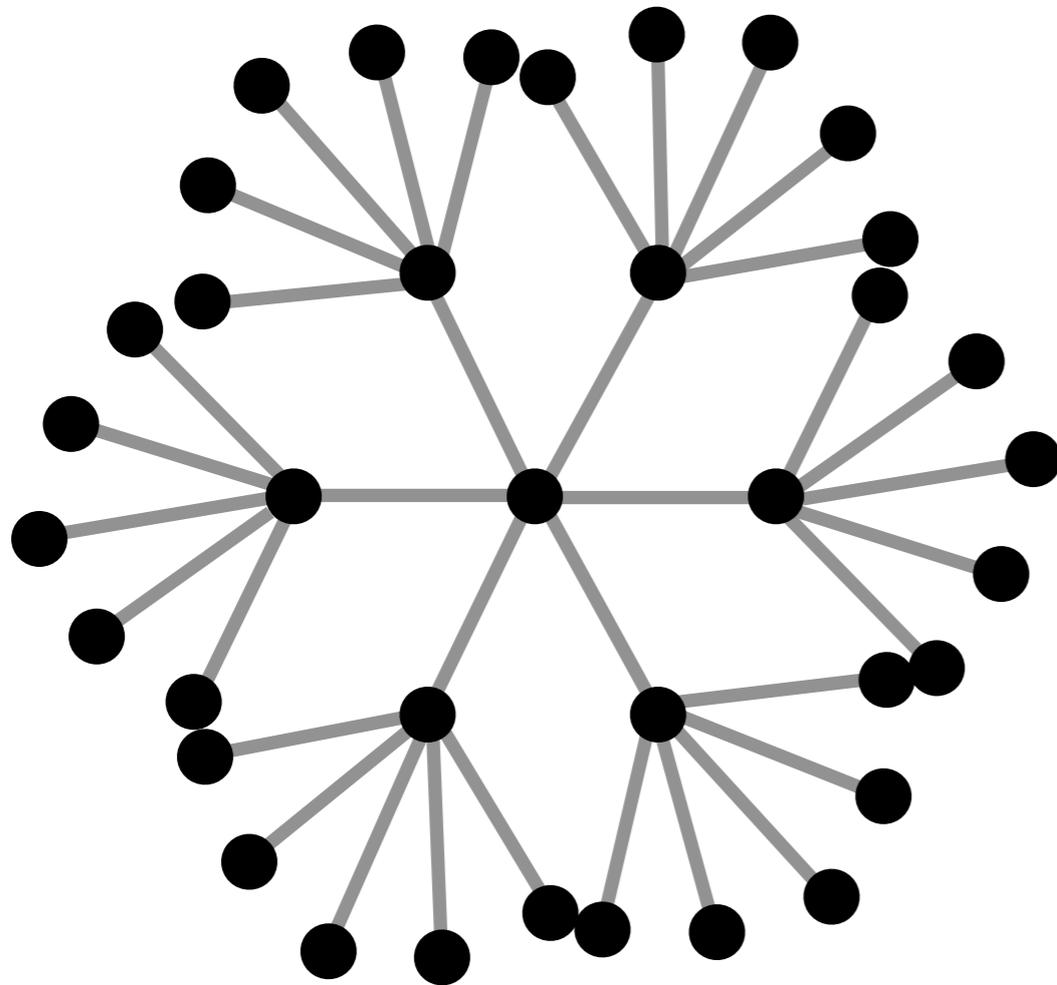
Lower bound on mean path length



Distance	# Nodes
1	6

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

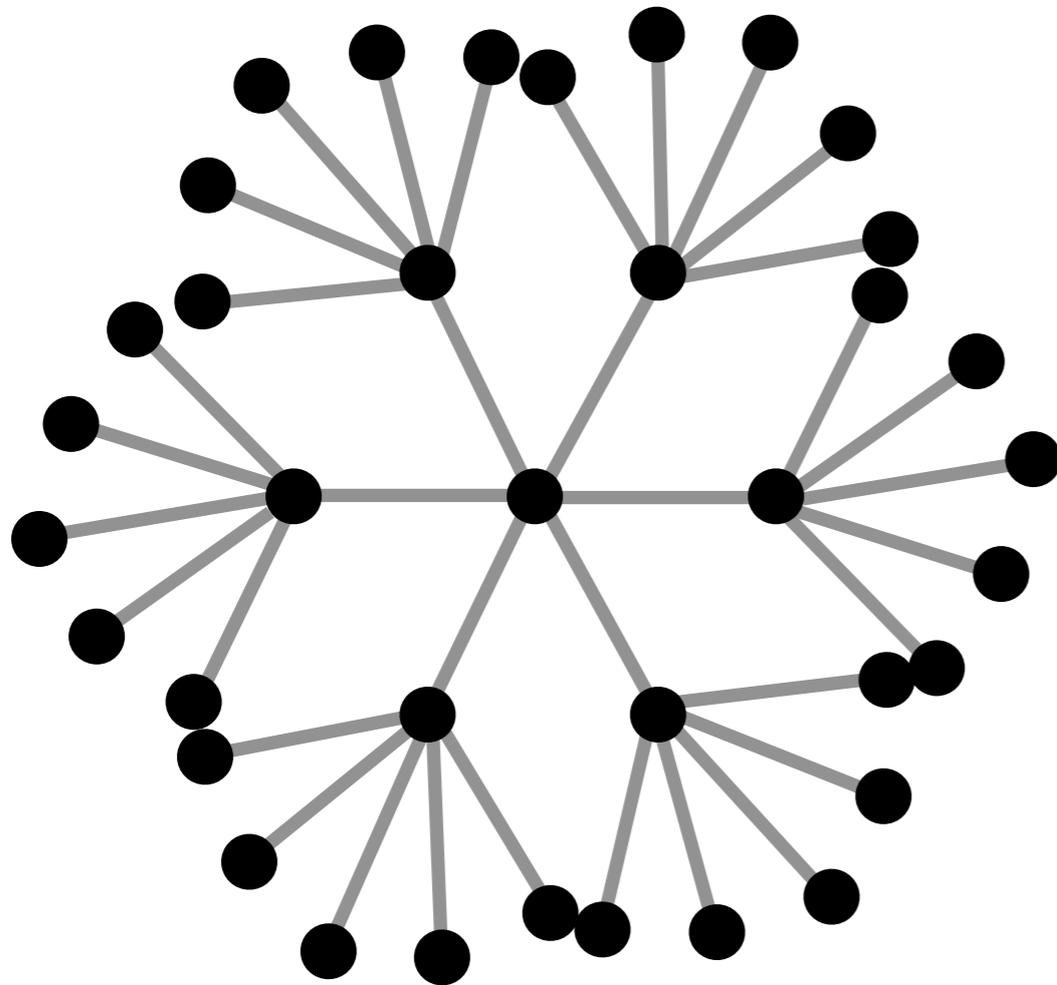
Lower bound on mean path length



Distance	# Nodes
1	6
2	6

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

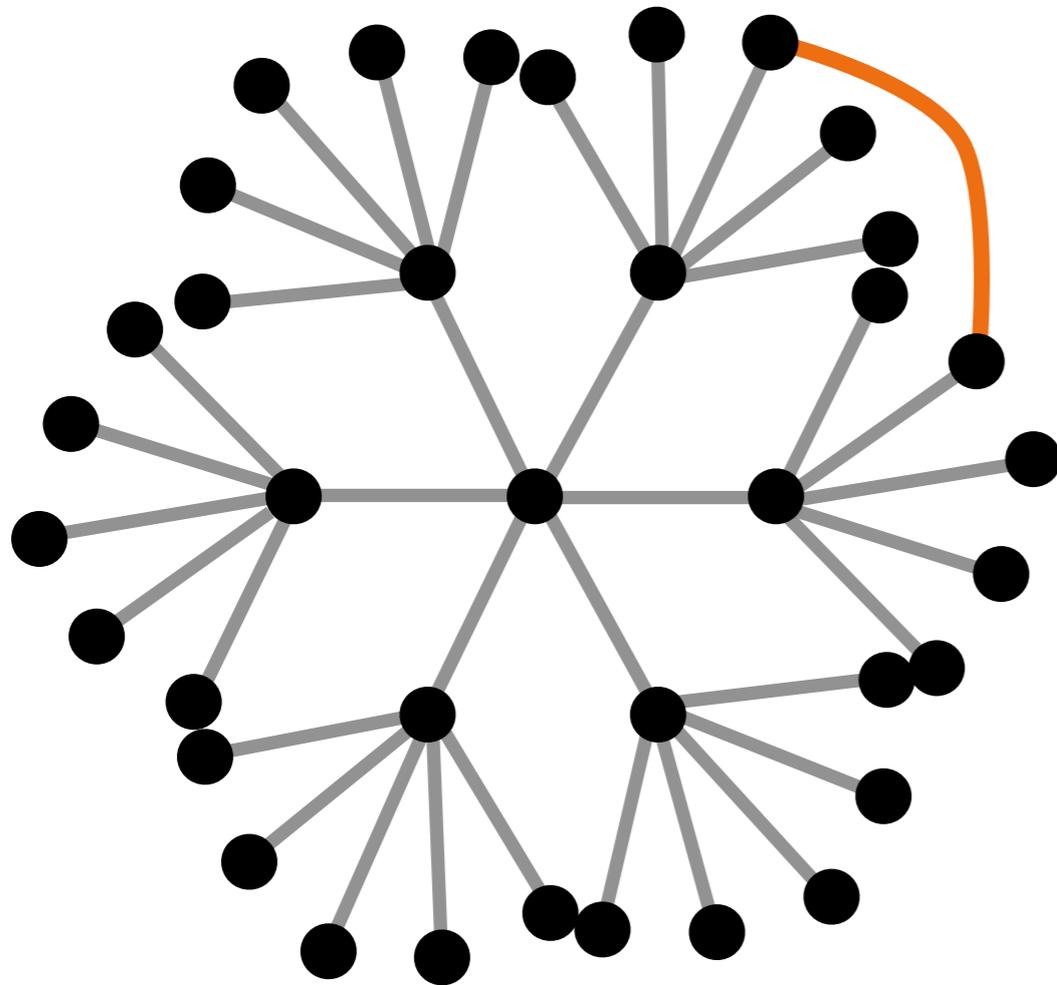
Lower bound on mean path length



Distance	# Nodes
1	6
2	6

(Ugliness omitted)

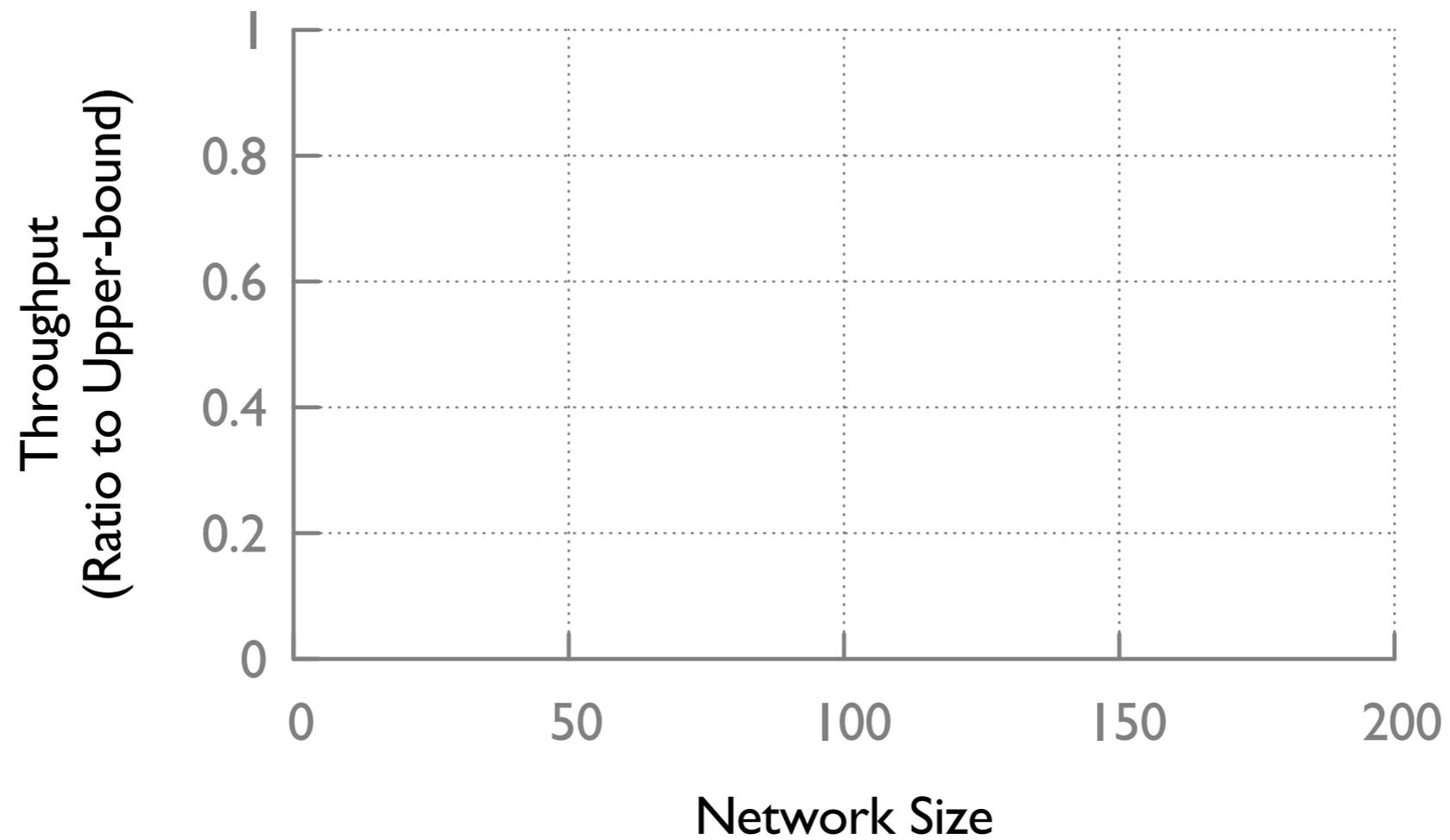
Lower bound on mean path length



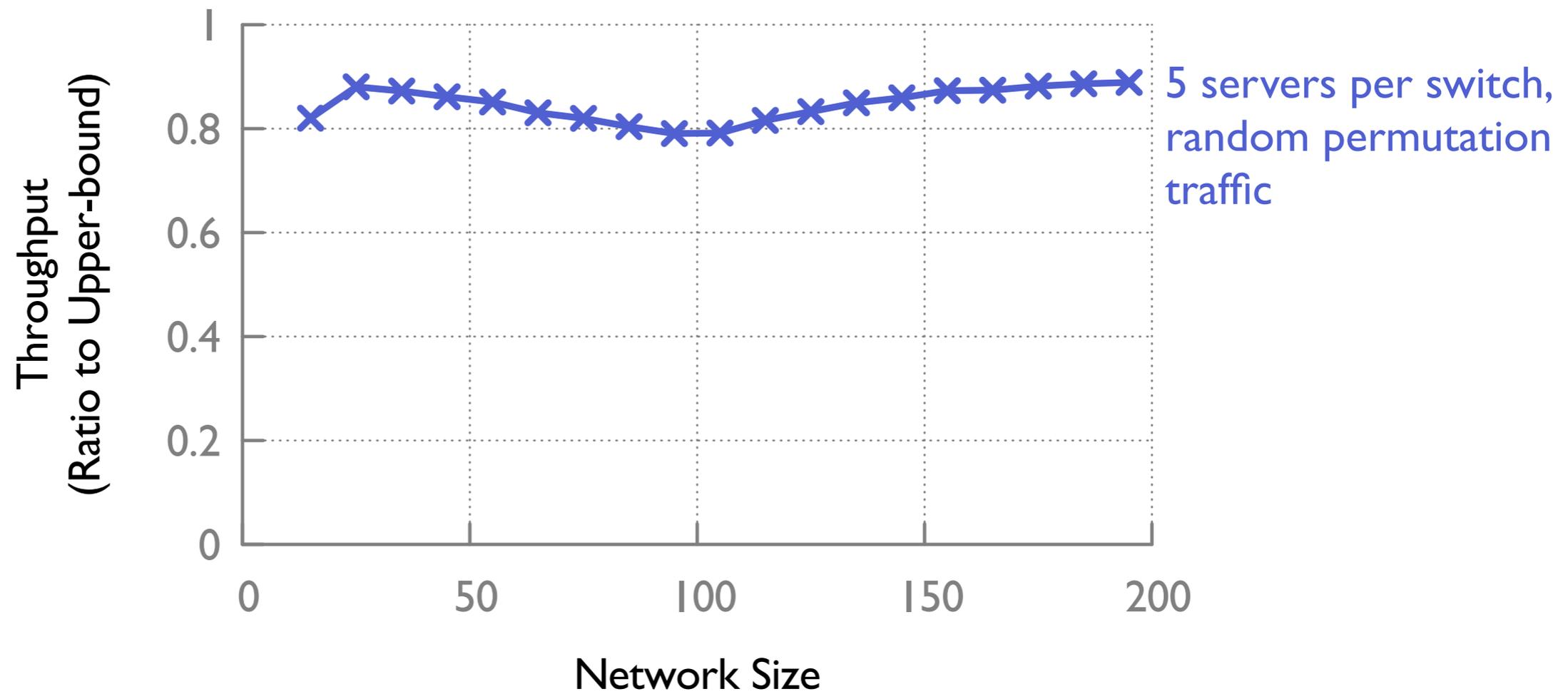
Distance	# Nodes
1	6
2	6

(Ugliness omitted)

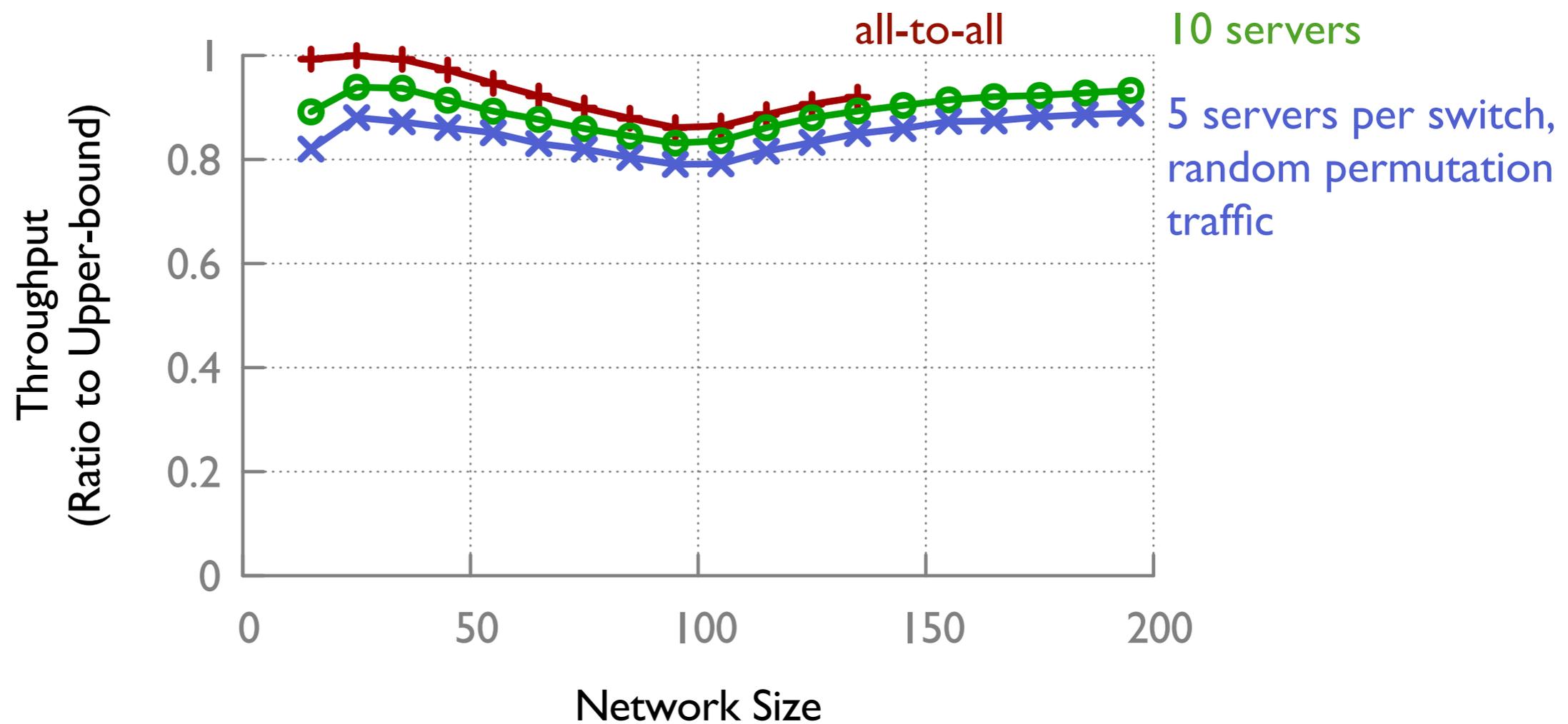
Random graphs vs. bound



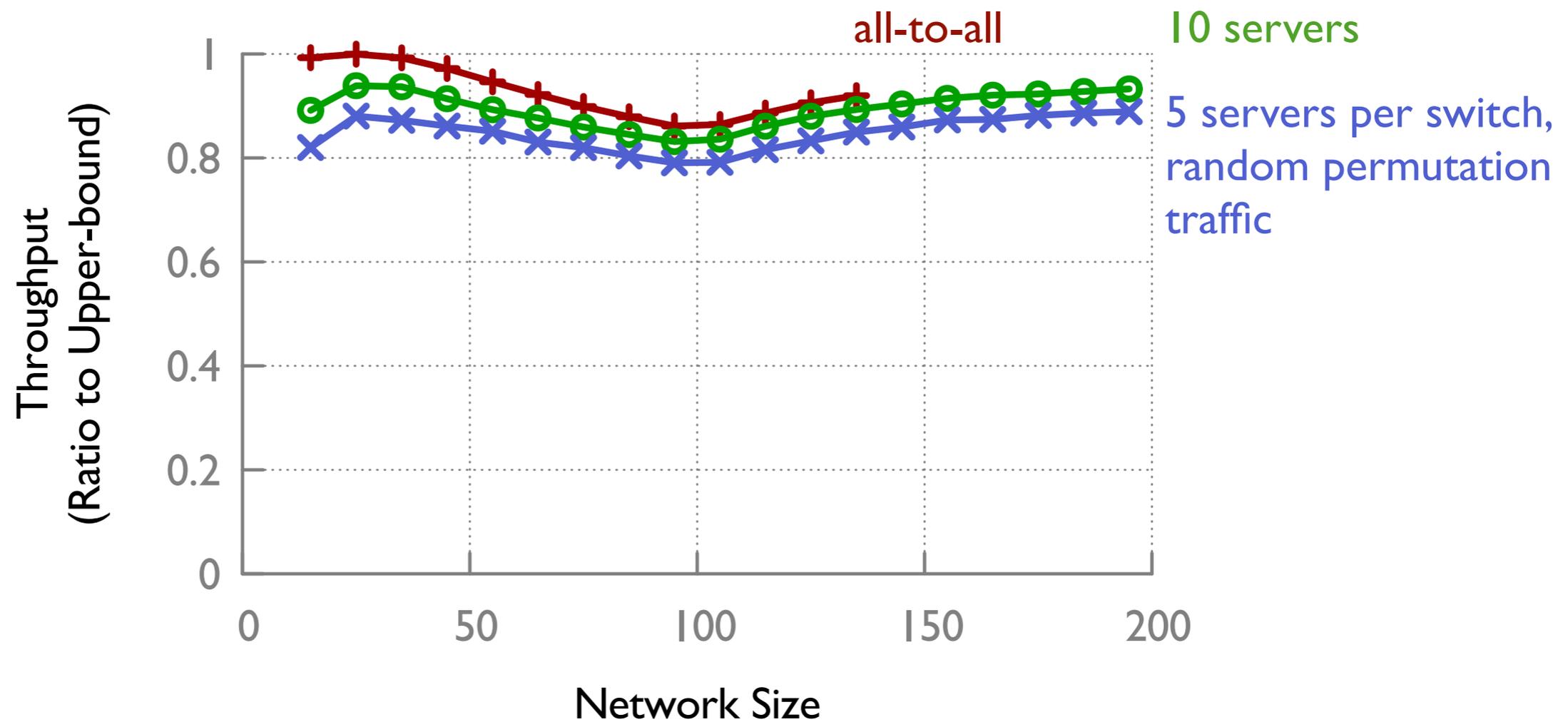
Random graphs vs. bound



Random graphs vs. bound

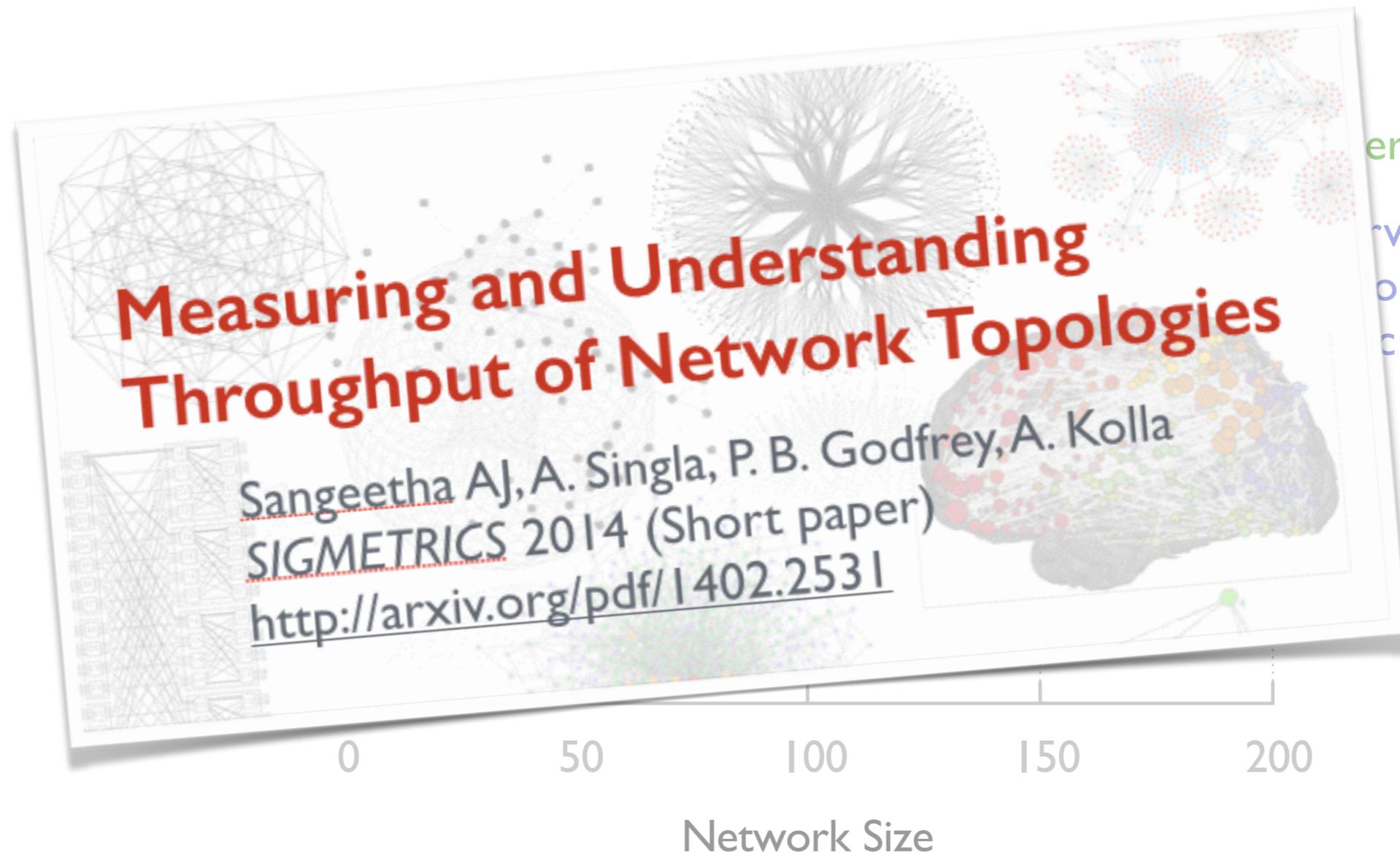


Random graphs vs. bound



Random graphs within a few percent of optimal!

Random graphs vs. bound



Random graphs within a few percent of optimal!

Random graphs exceed throughput of other topologies

How close can we get to
optimal network capacity?

Very close!!

How do we handle heterogeneity?

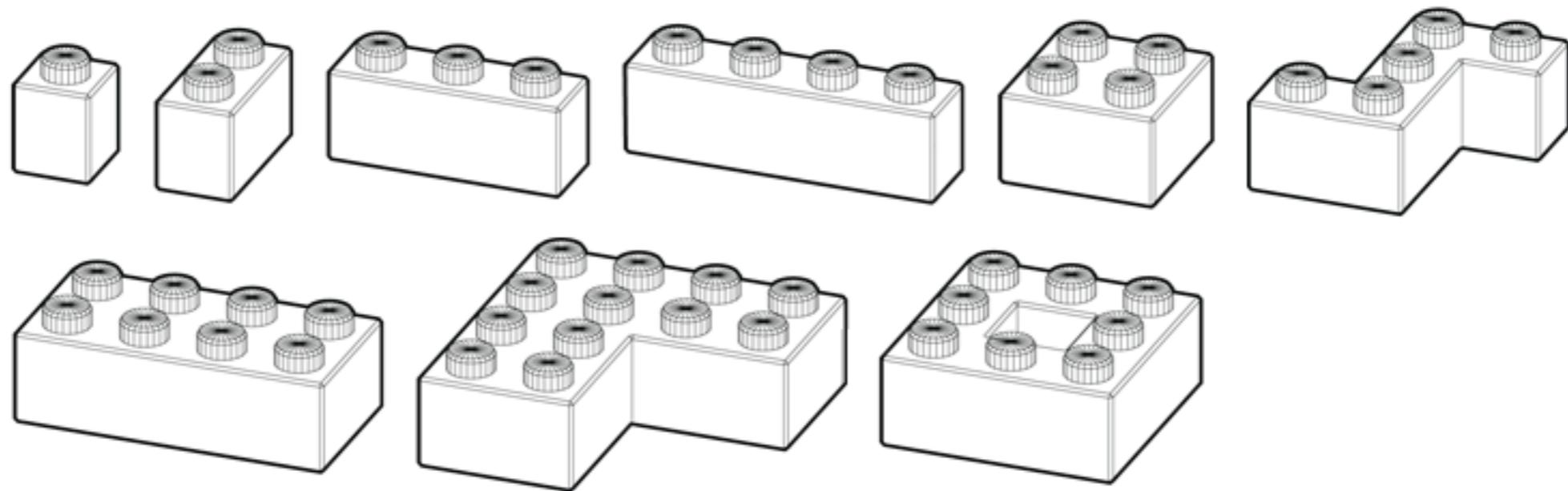
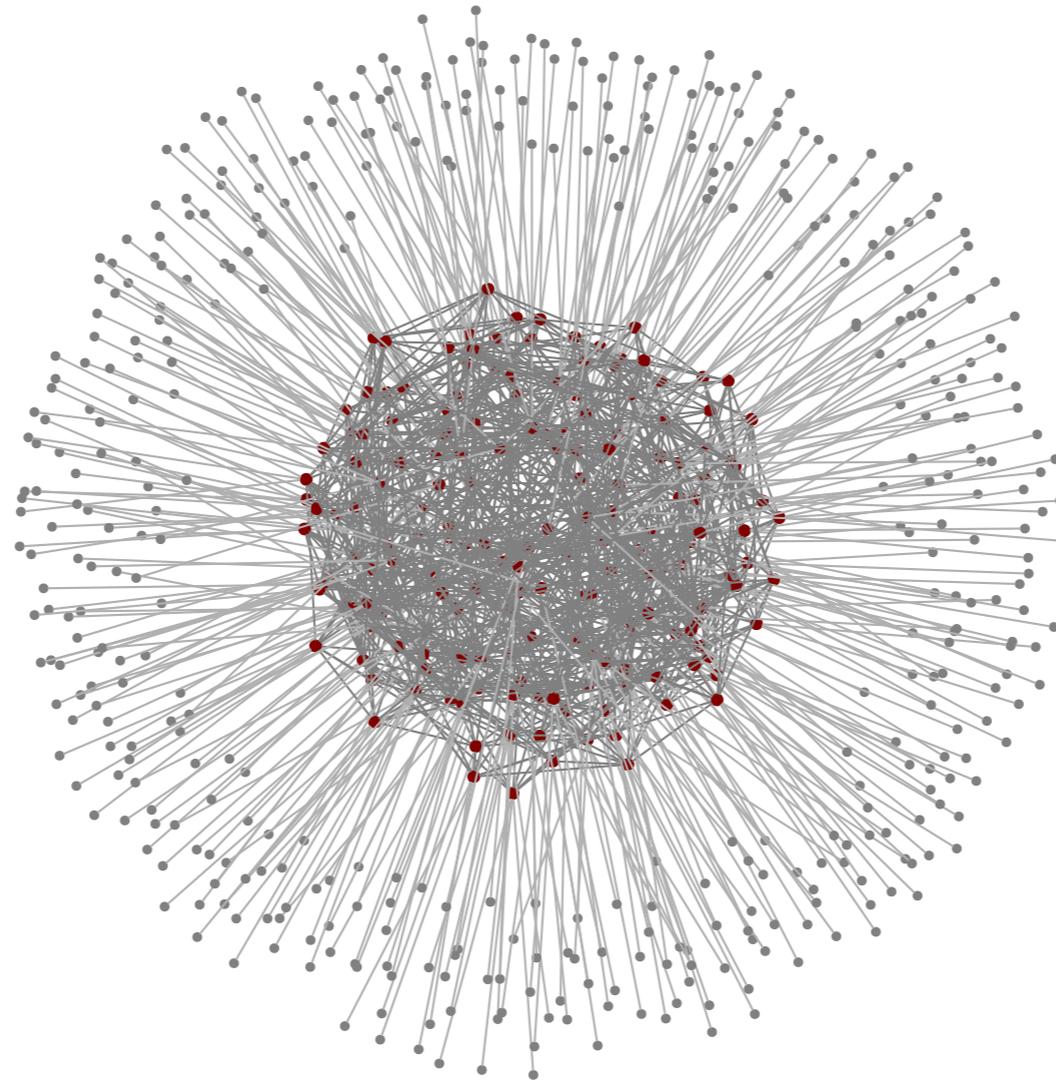
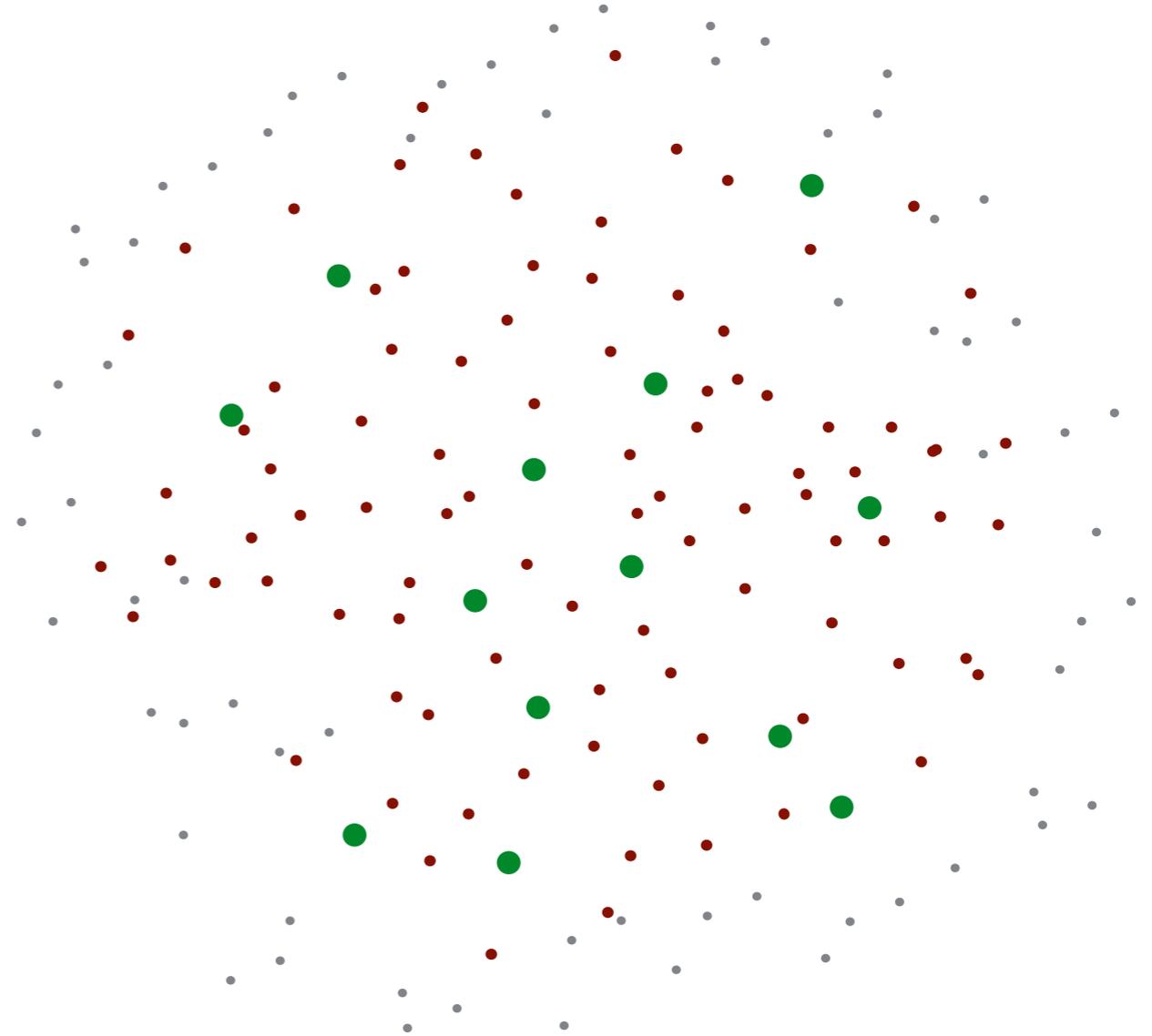
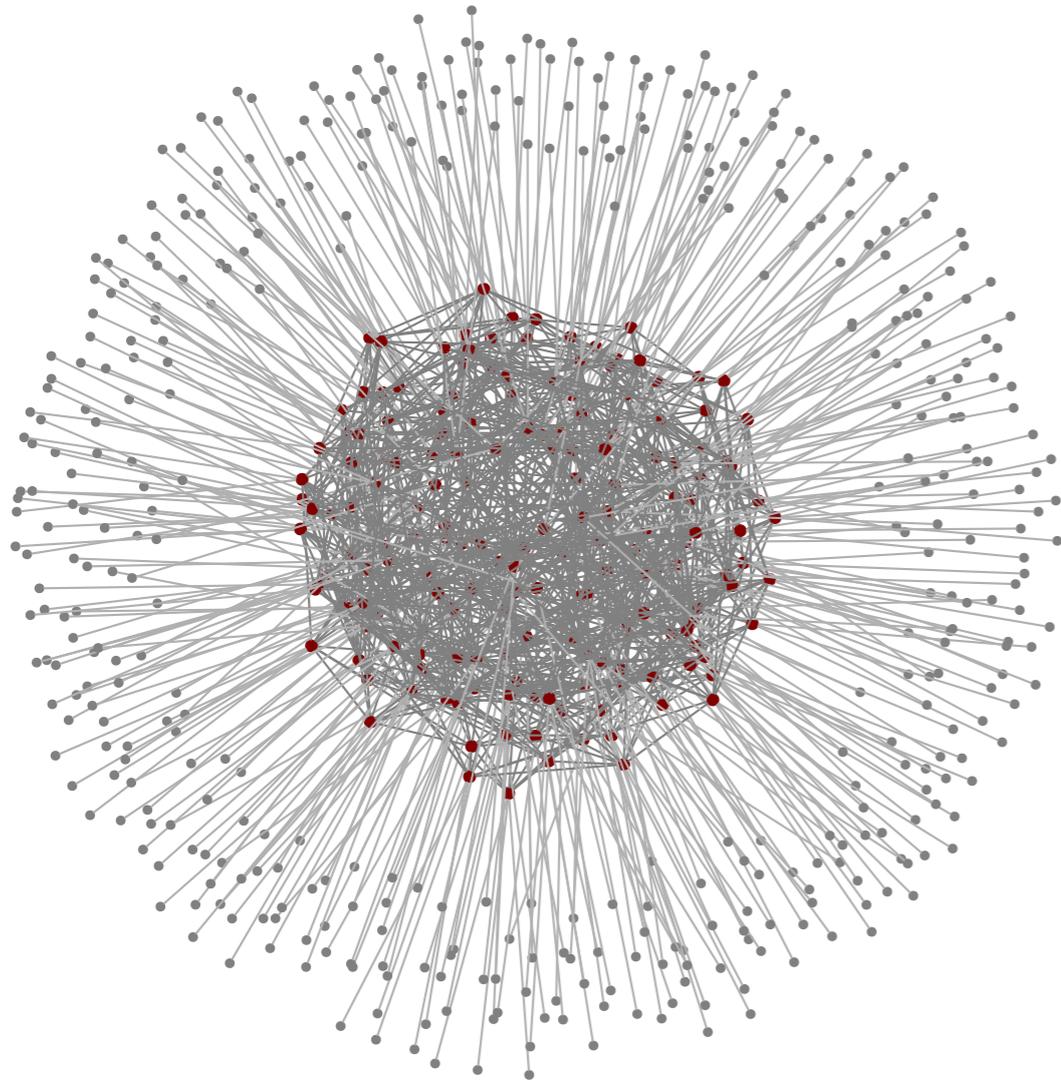


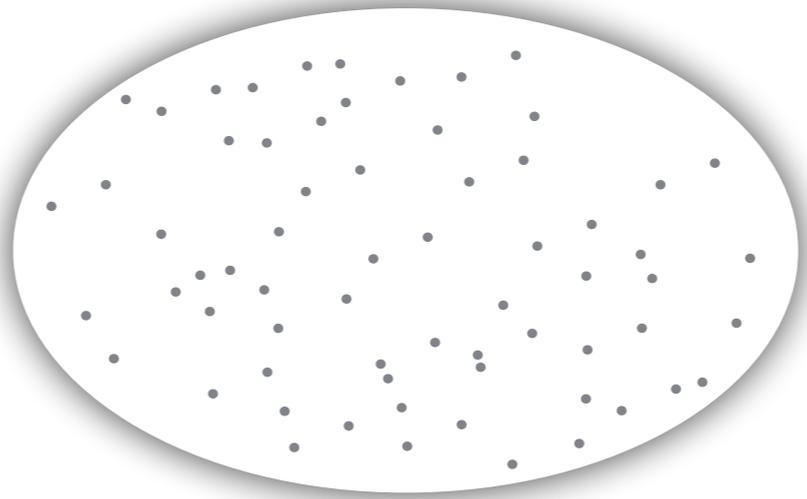
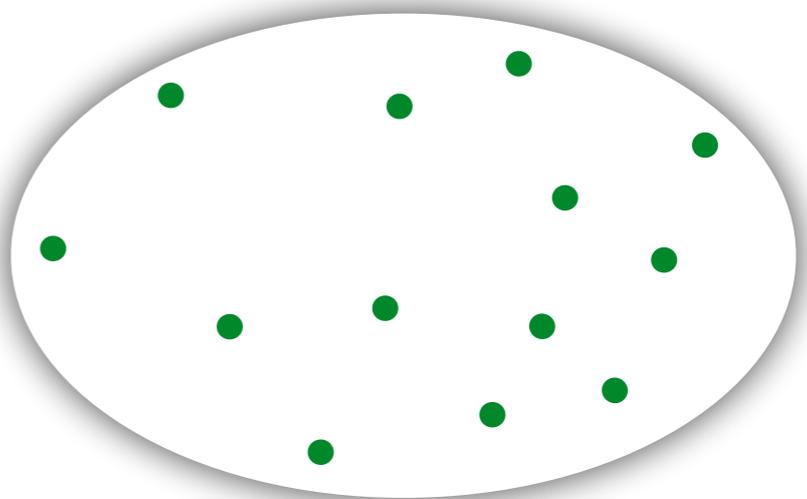
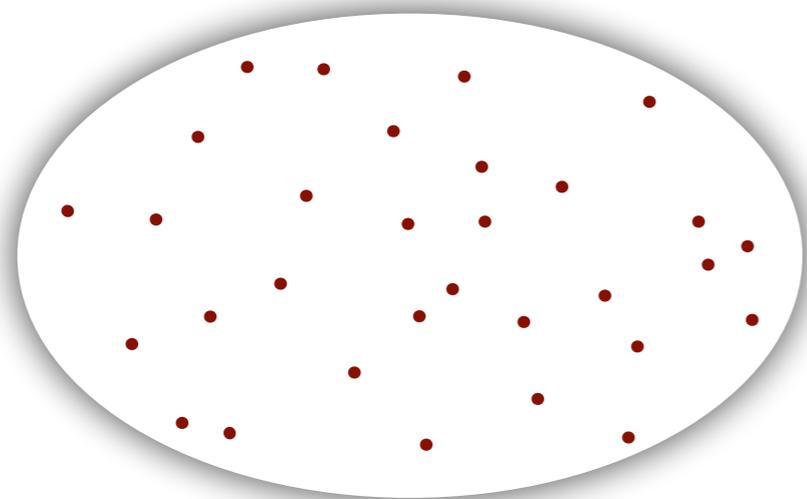
Image credit: Legolizer (www.drububu.com)

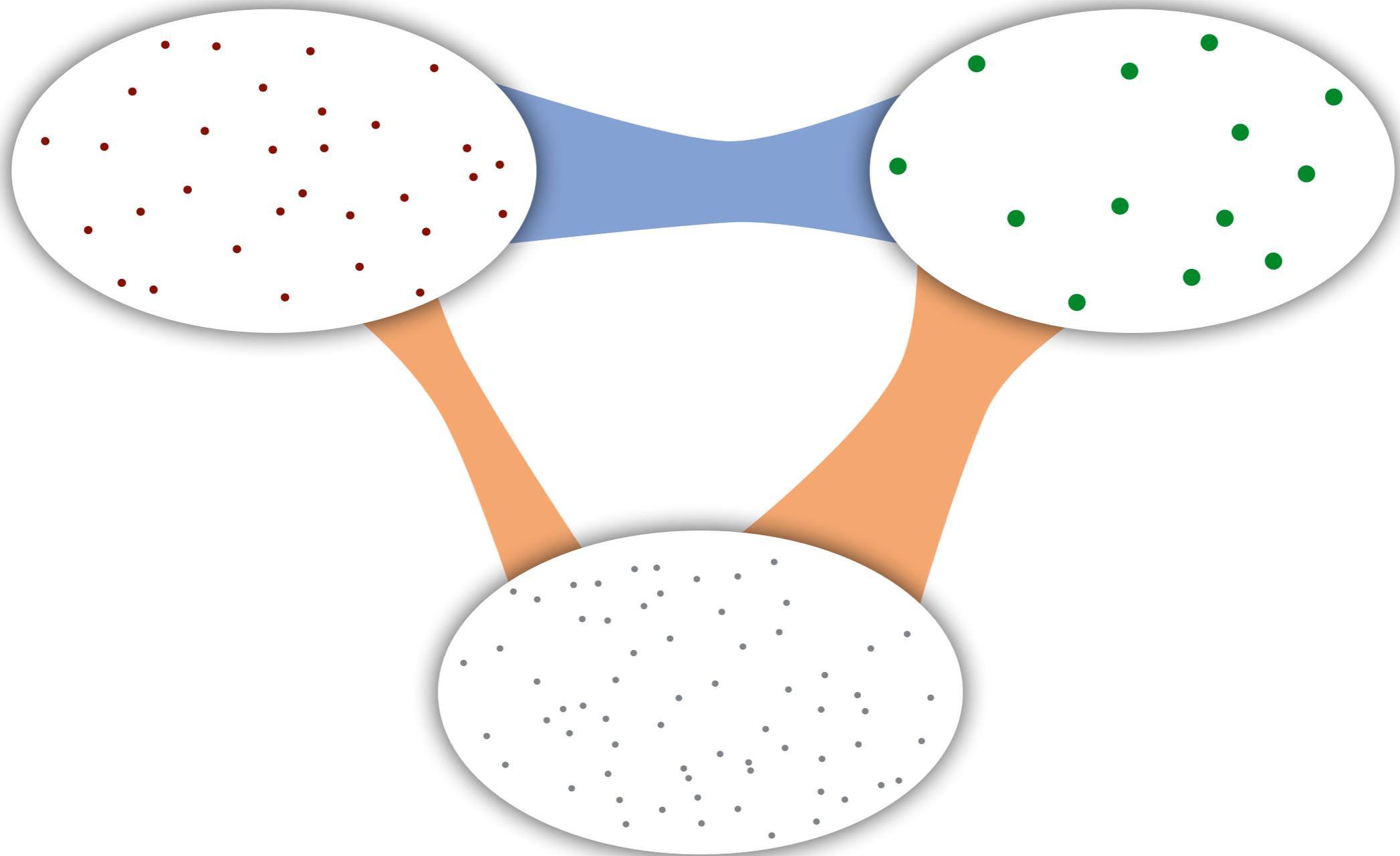
Heterogeneity



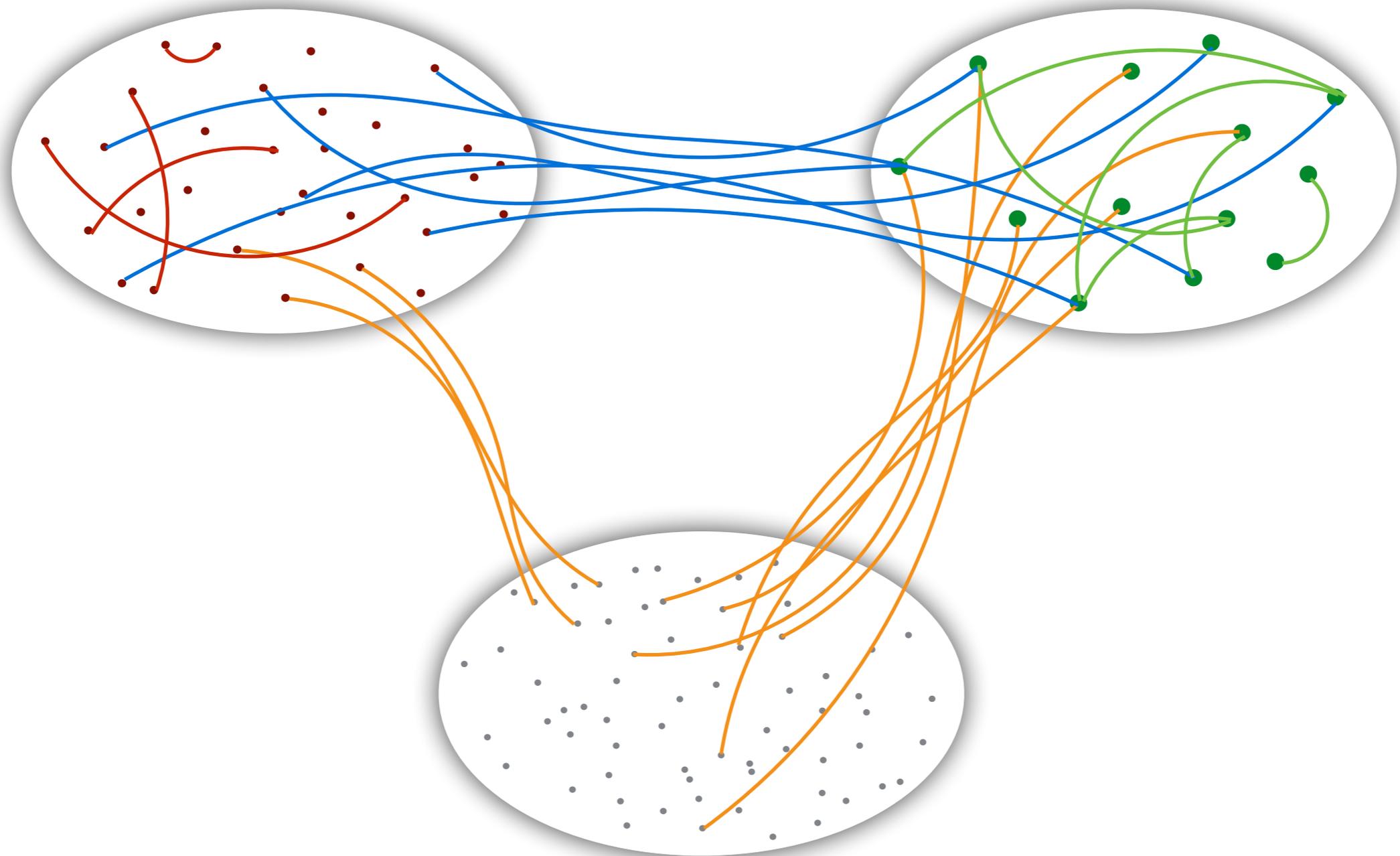
Heterogeneity







Random graphs as a building block



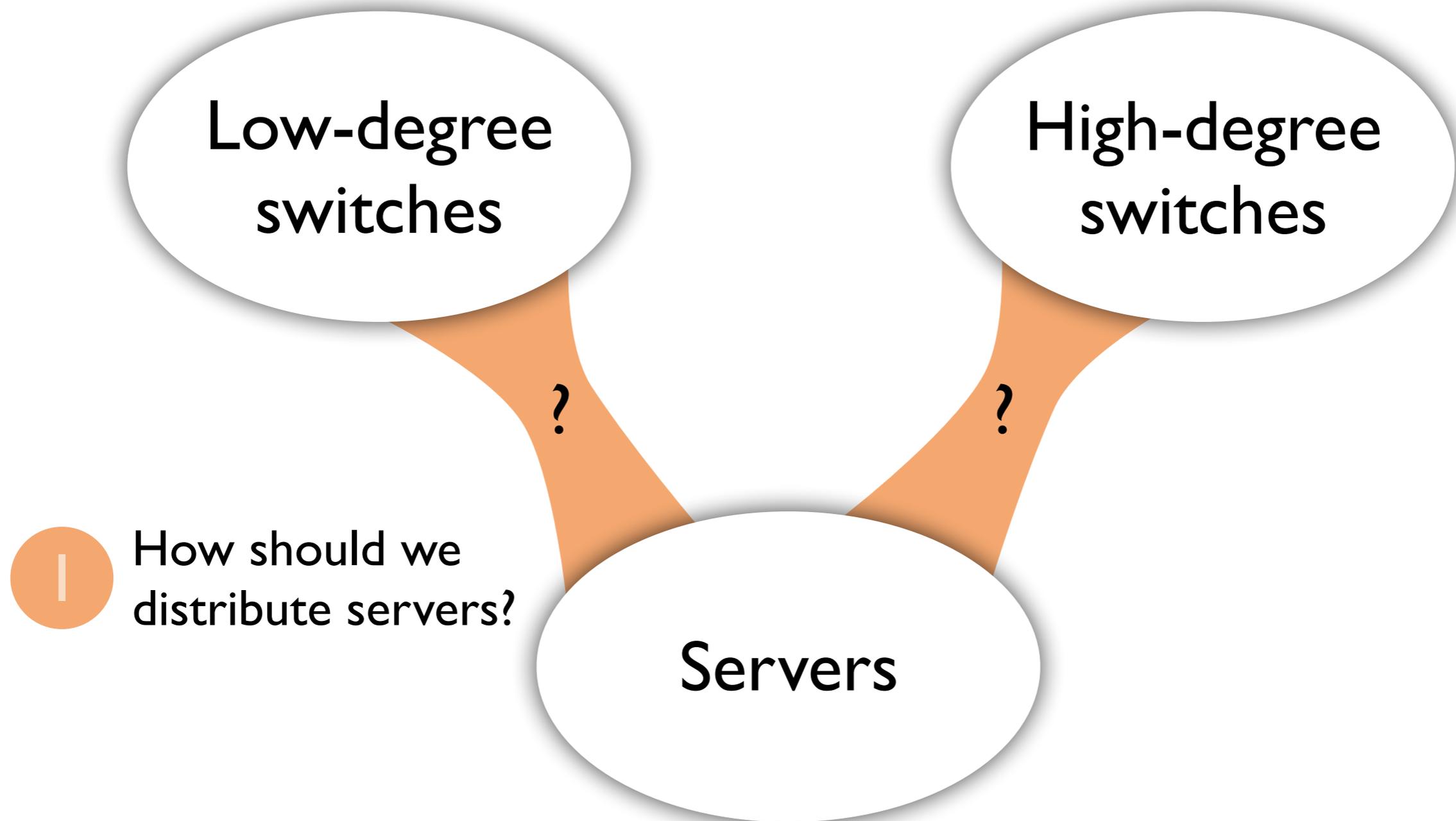
Random graphs as a building block

Low-degree
switches

High-degree
switches

Servers

Random graphs as a building block



Random graphs as a building block

2 How should we interconnect switches?

Low-degree switches

High-degree switches

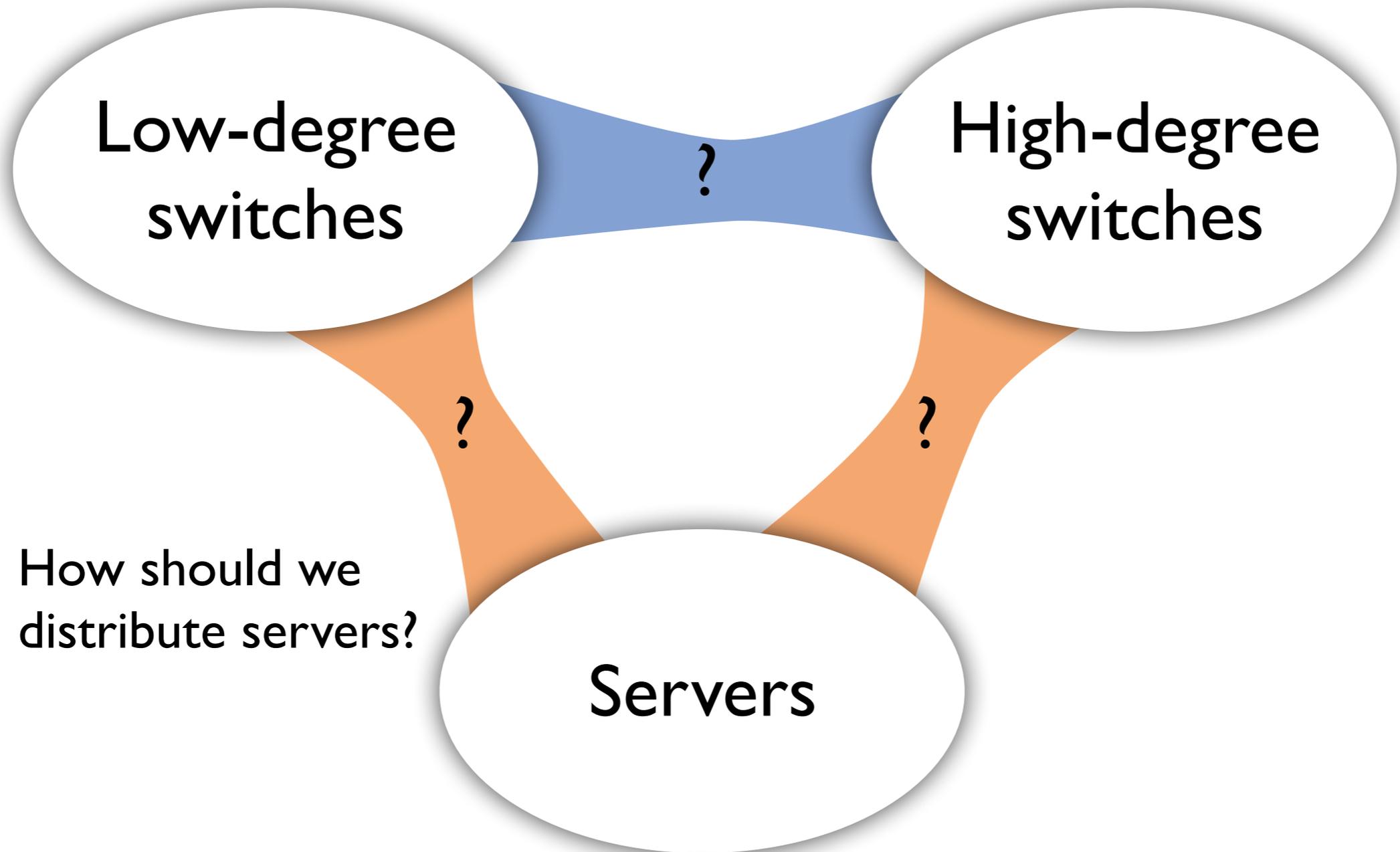
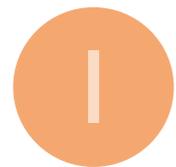
?

?

?

1 How should we distribute servers?

Servers



Random graphs as a building block

2 How should we interconnect switches?

Low-degree switches

High-degree switches

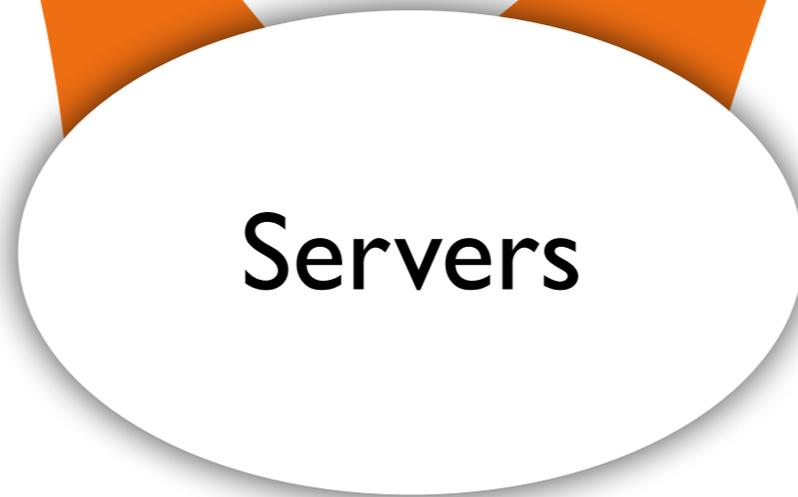
?

?

?

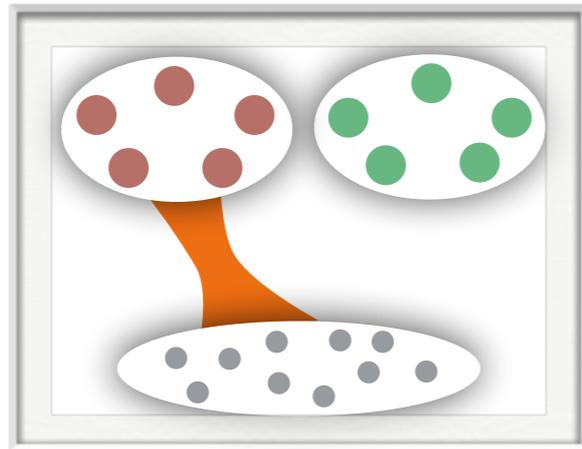
1 How should we distribute servers?

Servers

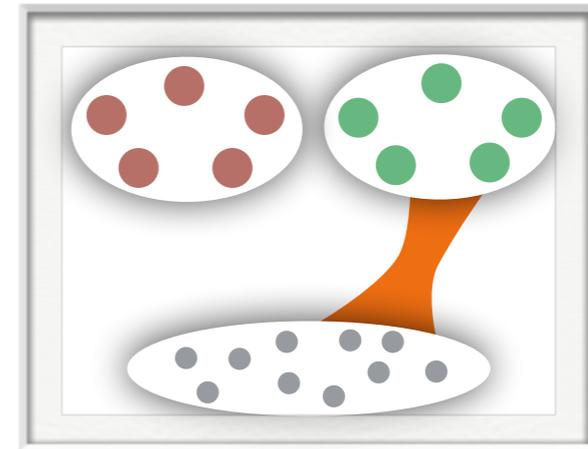
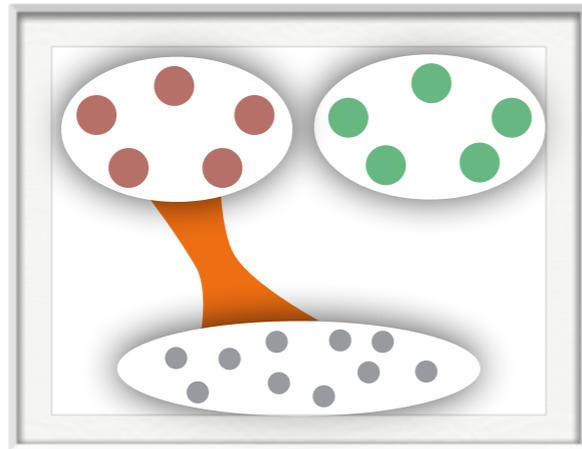


Distributing servers

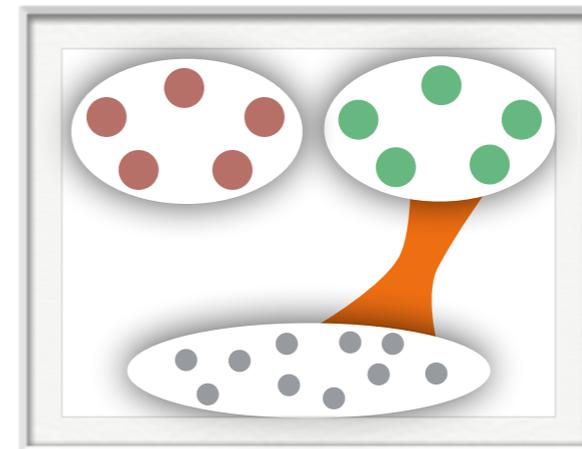
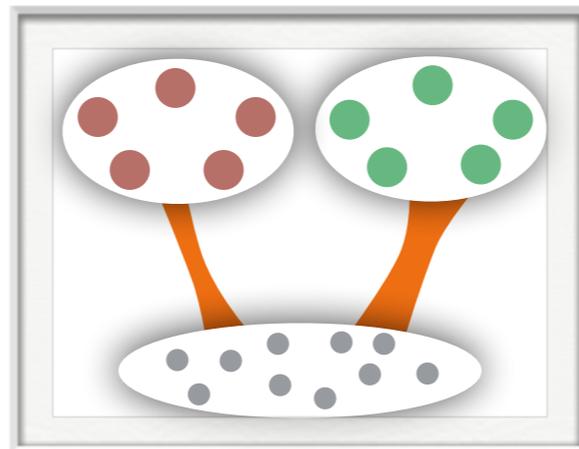
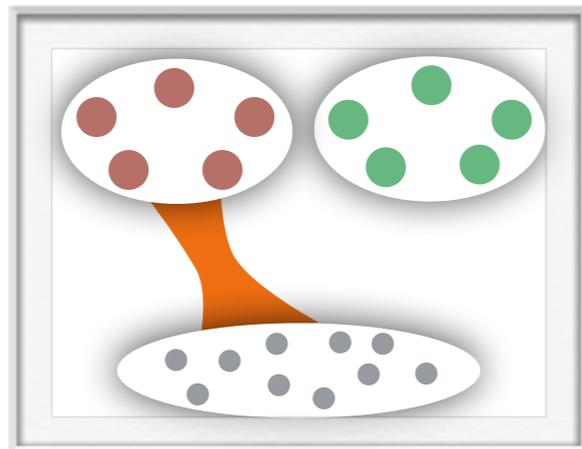
Distributing servers



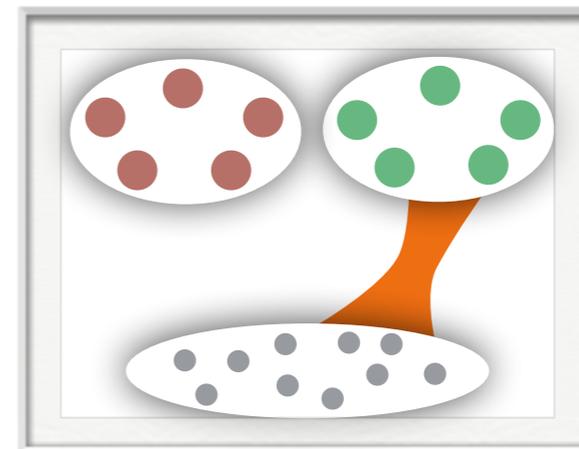
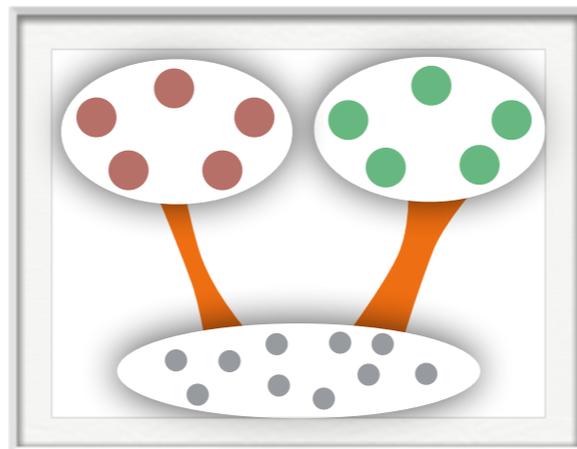
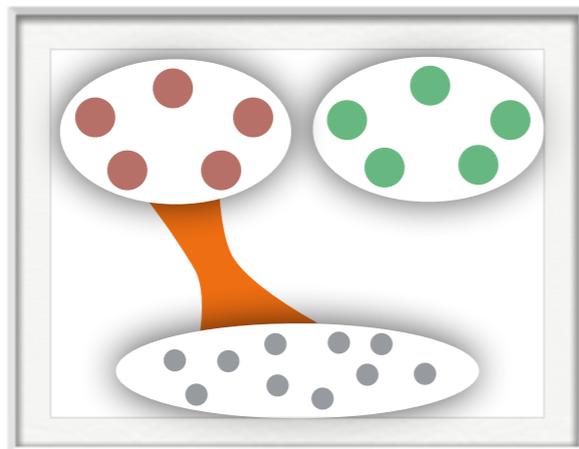
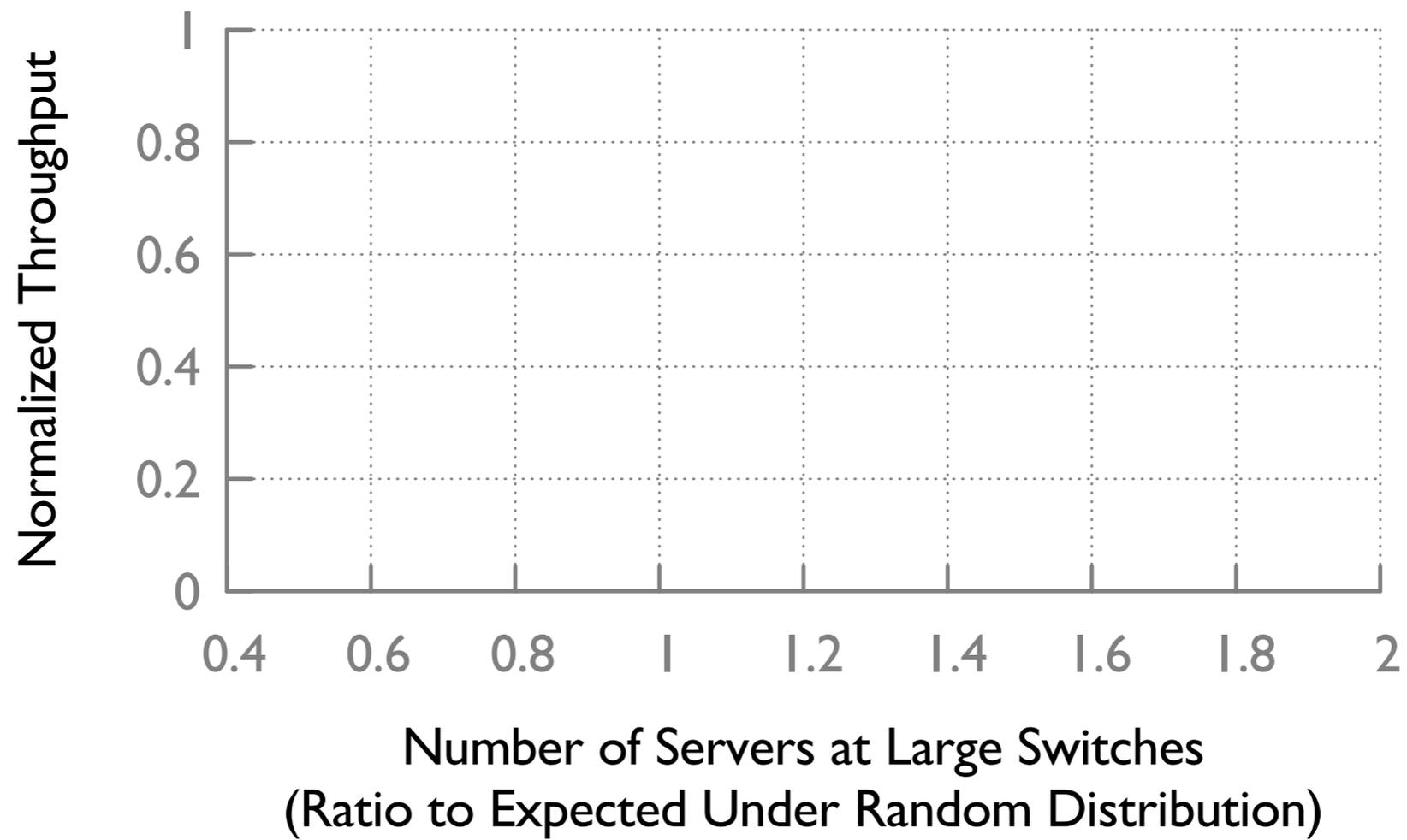
Distributing servers



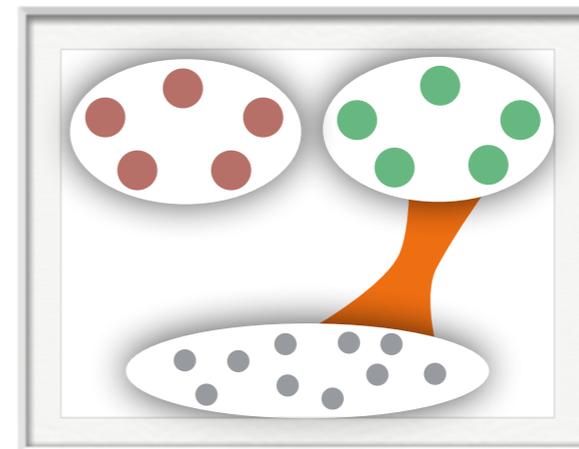
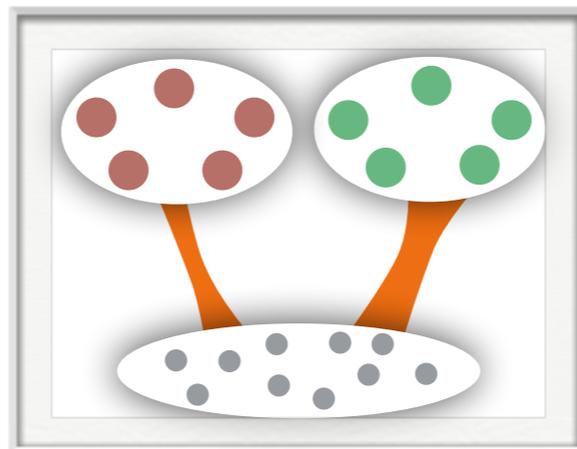
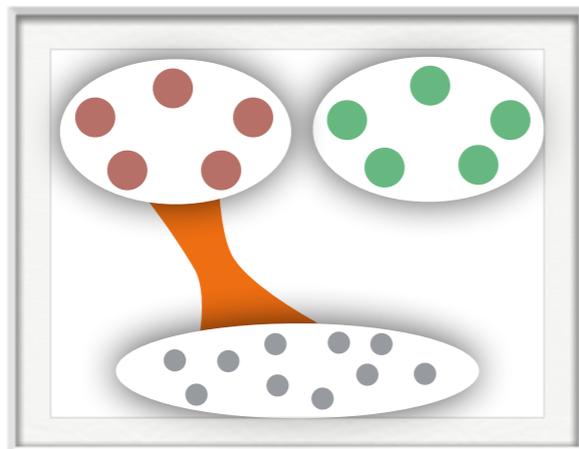
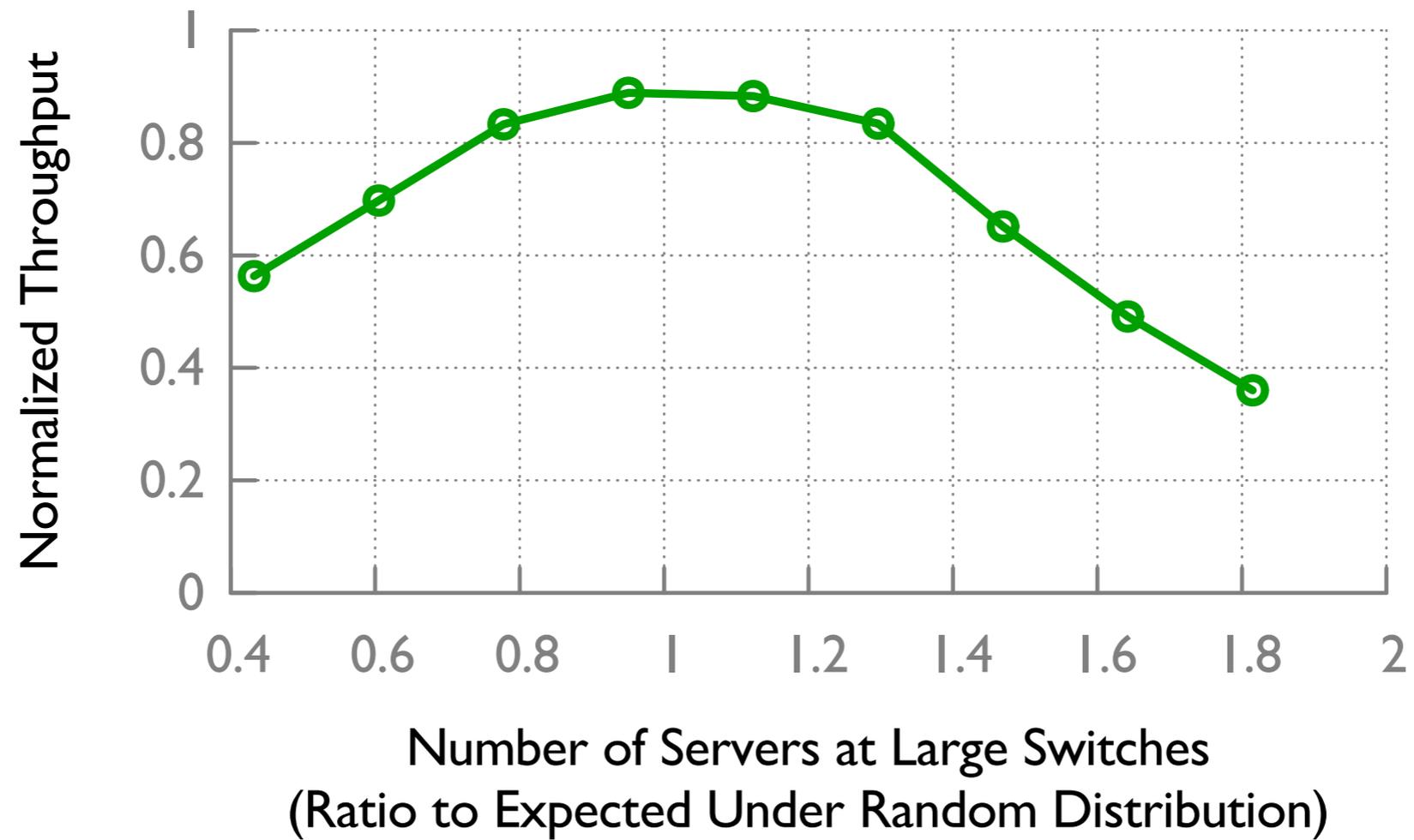
Distributing servers



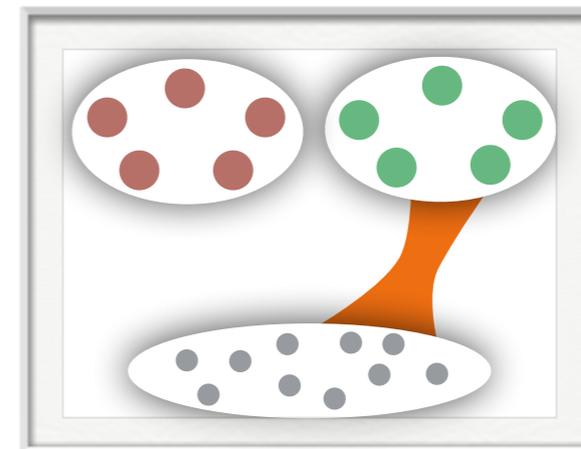
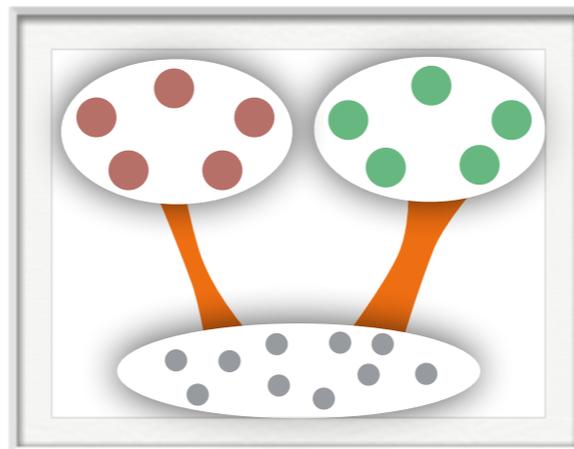
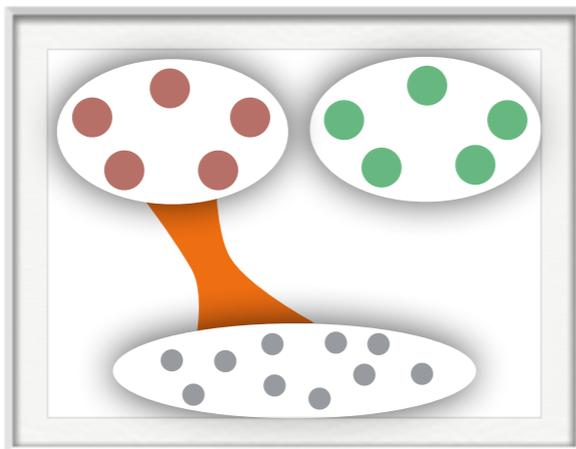
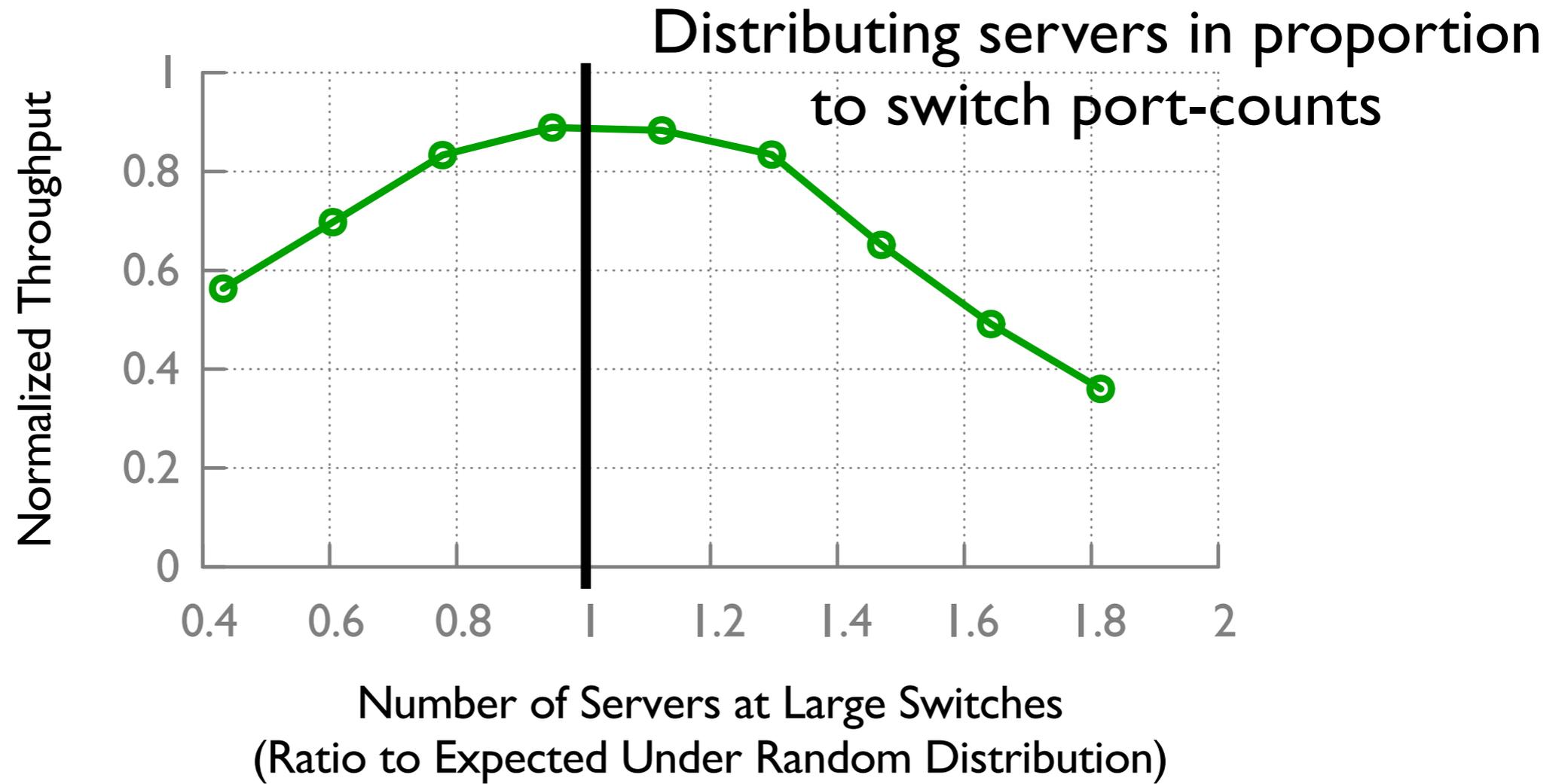
Distributing servers



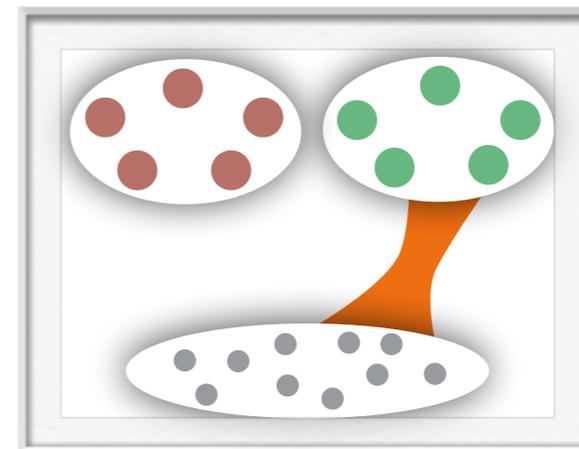
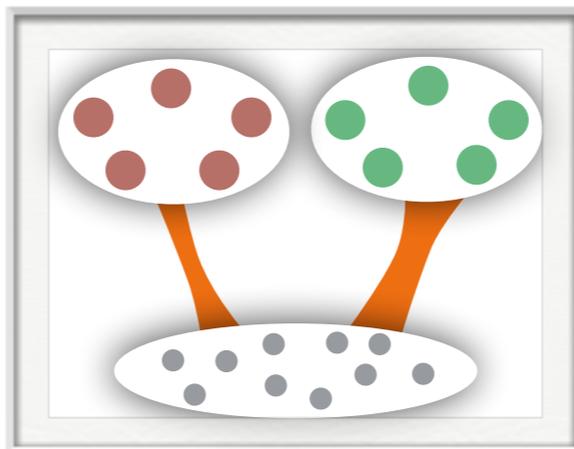
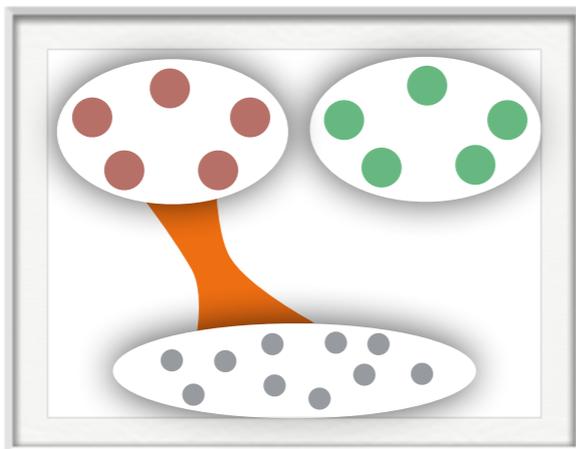
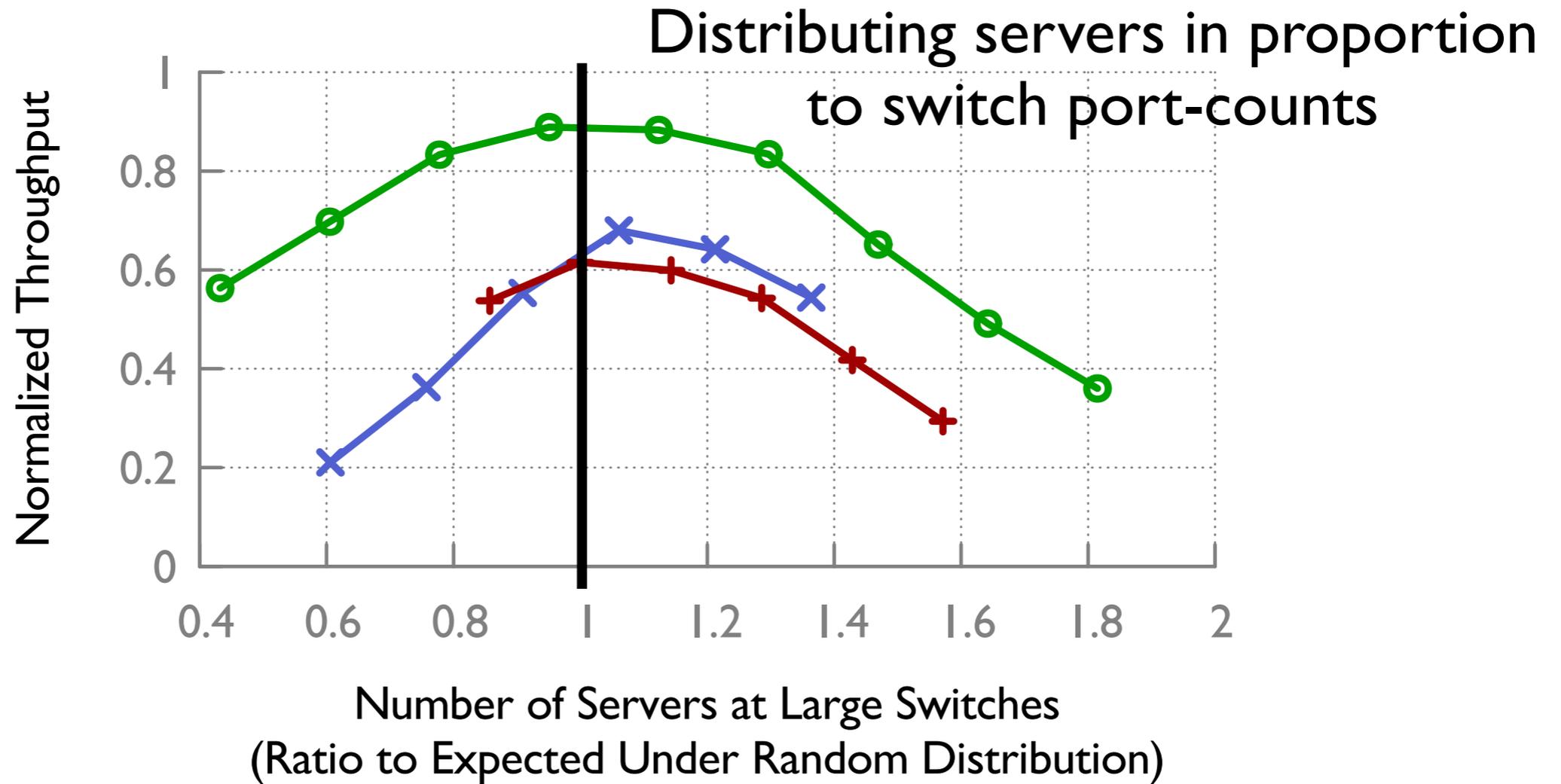
Distributing servers



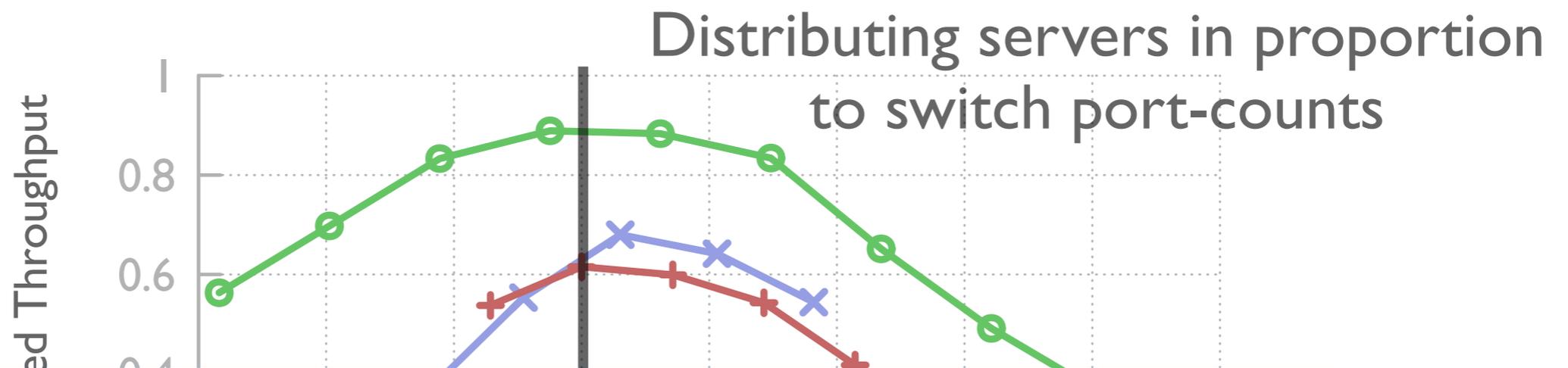
Distributing servers



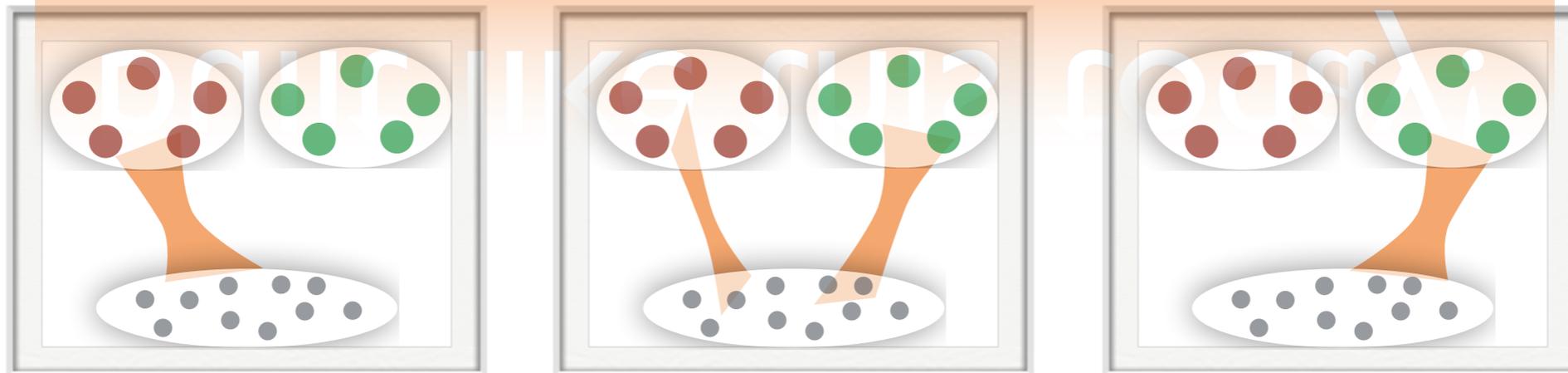
Distributing servers



Distributing servers



Networks aren't built like this today!



Random graphs as a building block

2 How should we interconnect switches?

Low-degree switches

High-degree switches

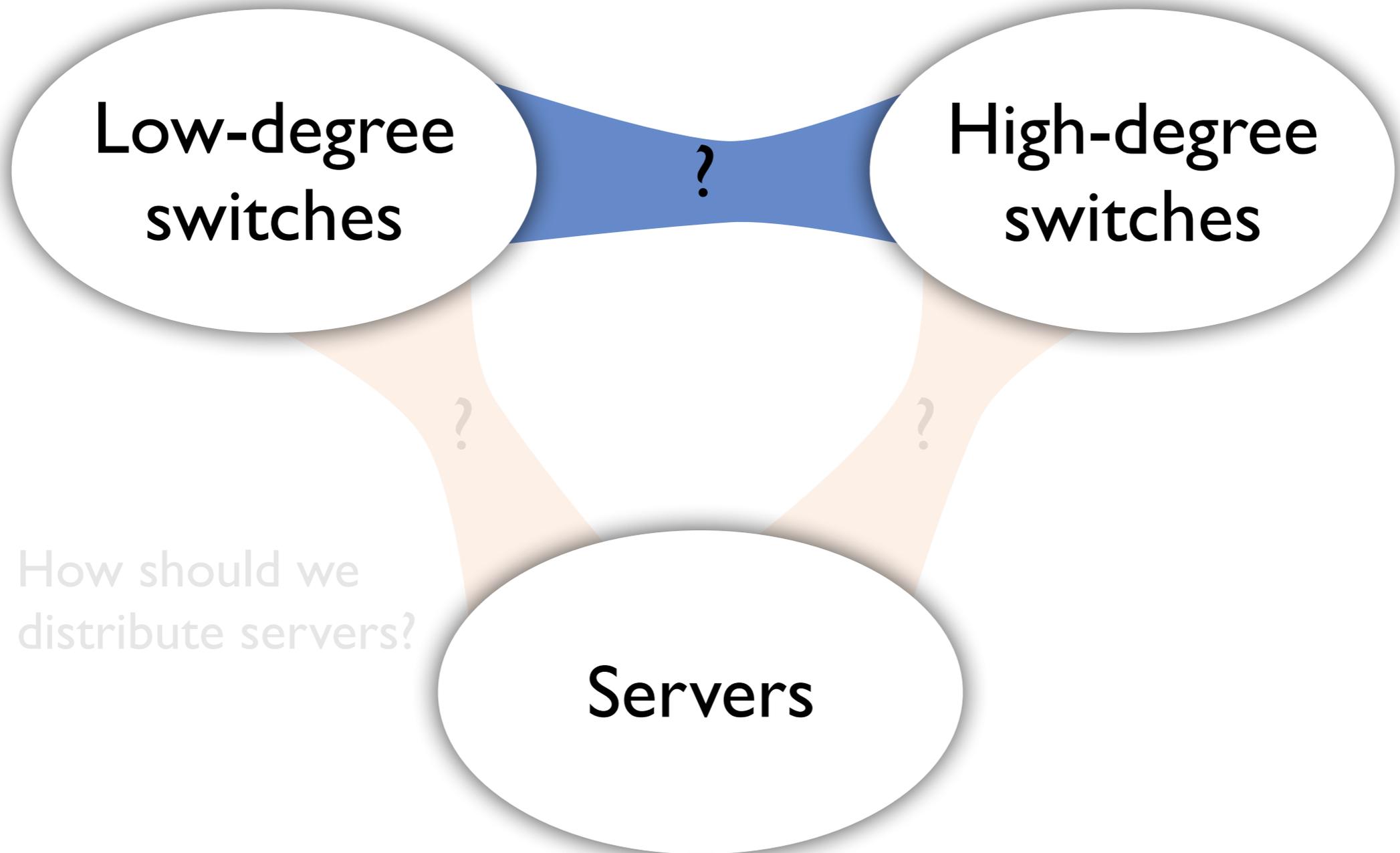
?

?

?

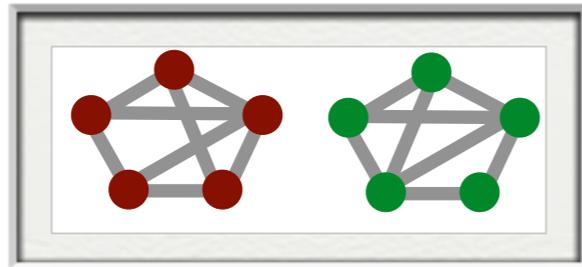
1 How should we distribute servers?

Servers

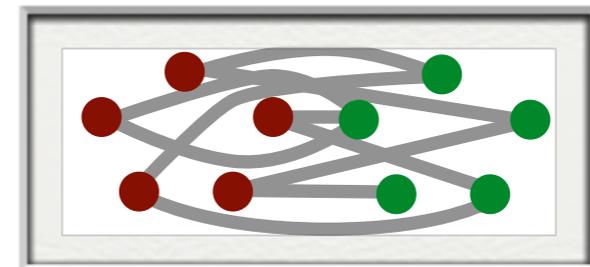
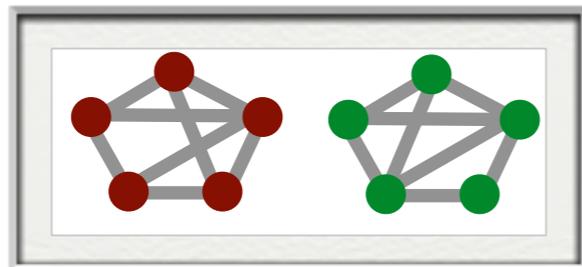


Interconnecting switches

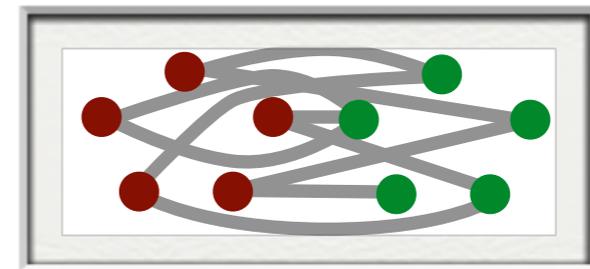
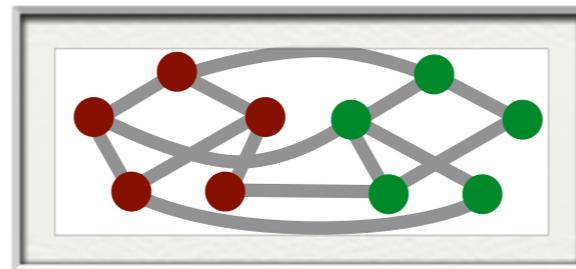
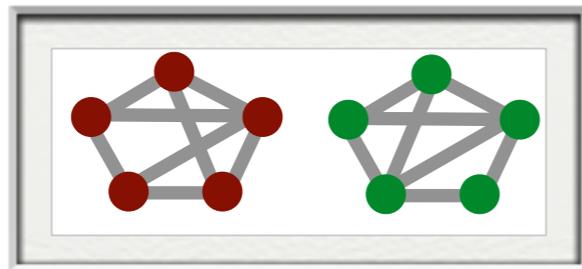
Interconnecting switches



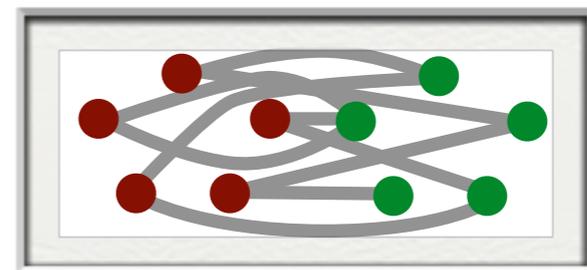
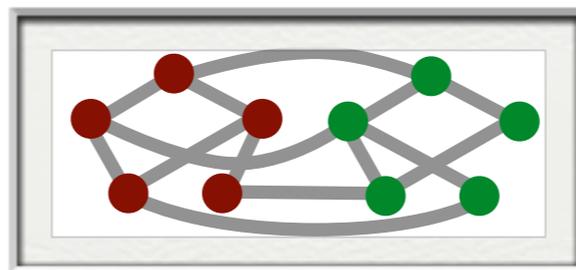
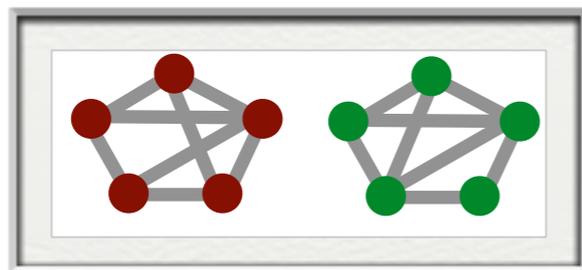
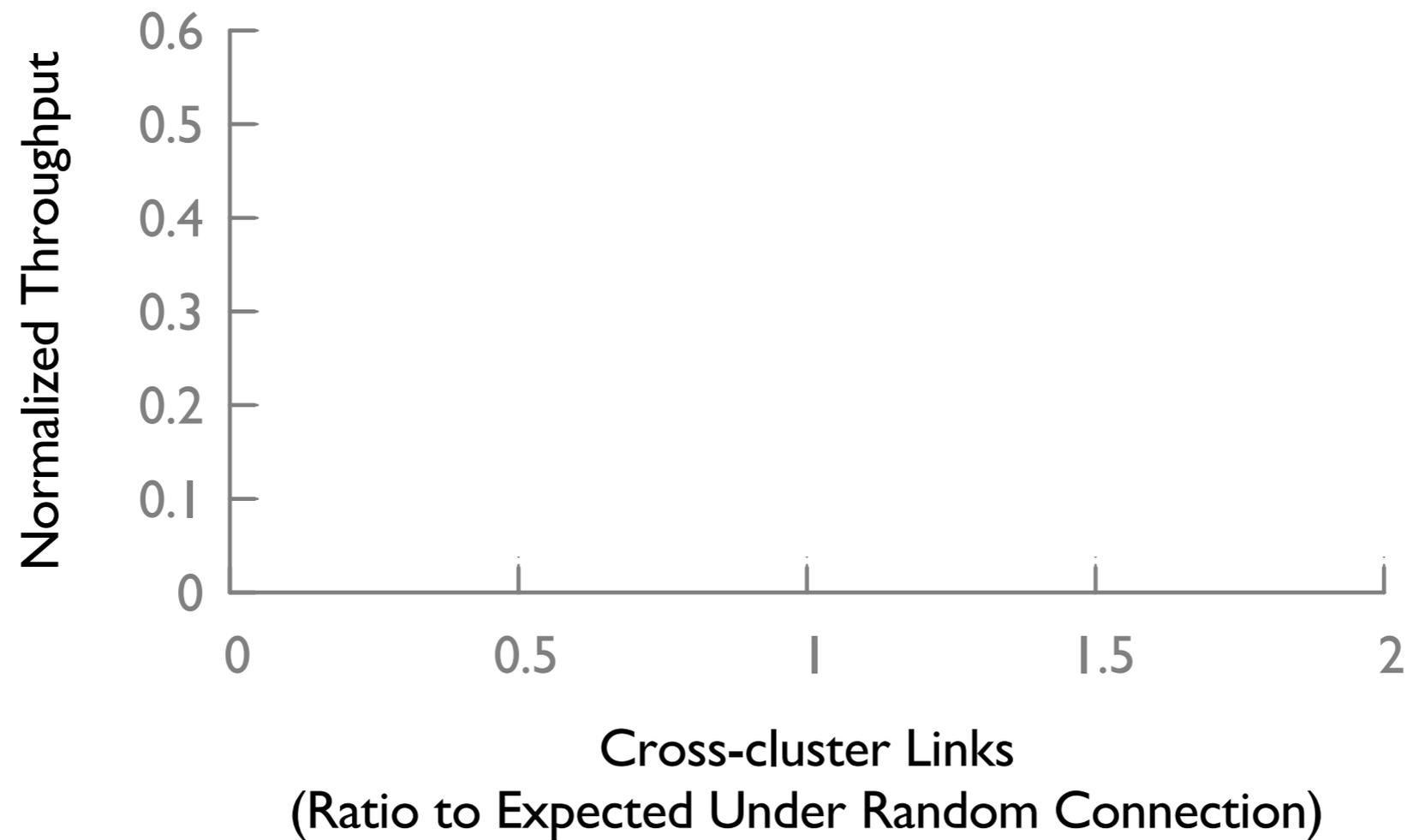
Interconnecting switches



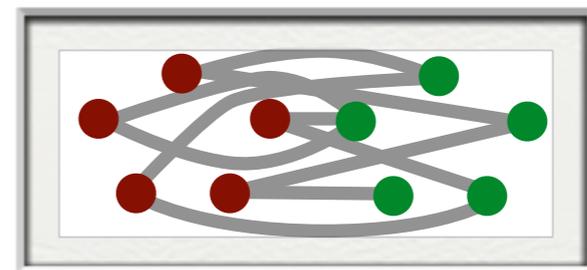
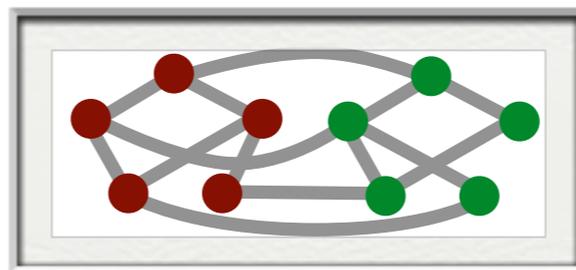
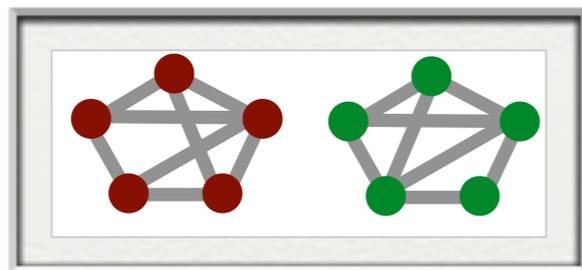
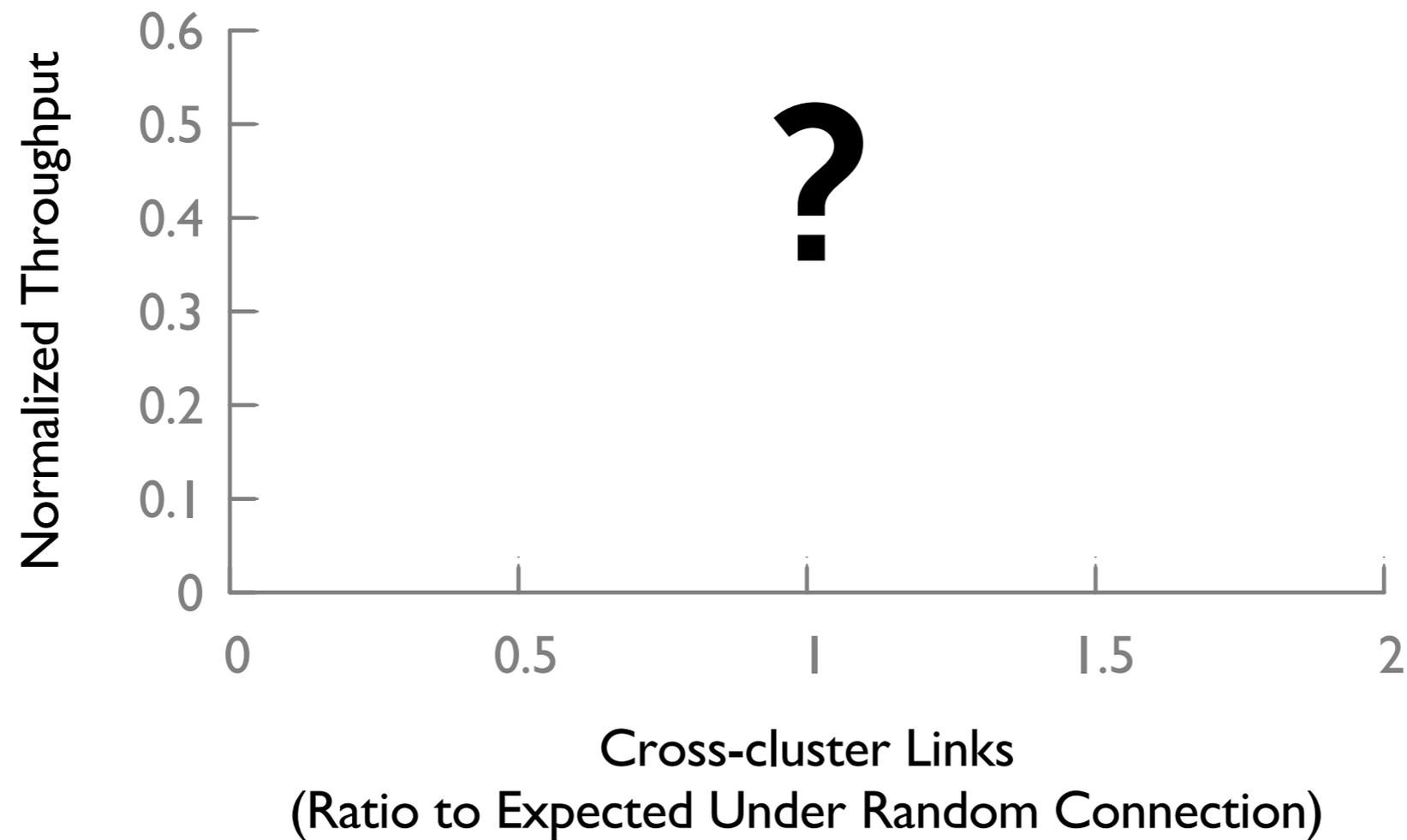
Interconnecting switches



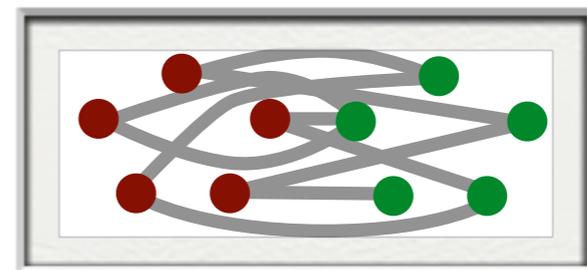
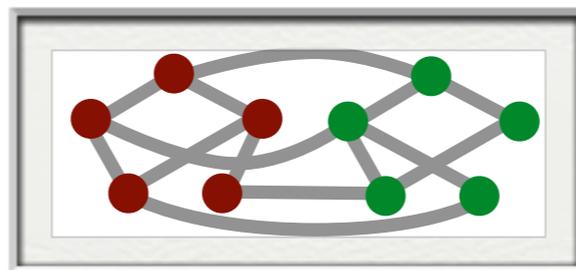
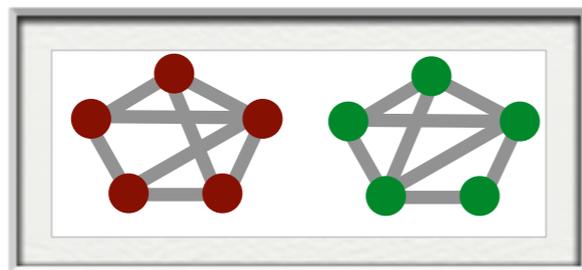
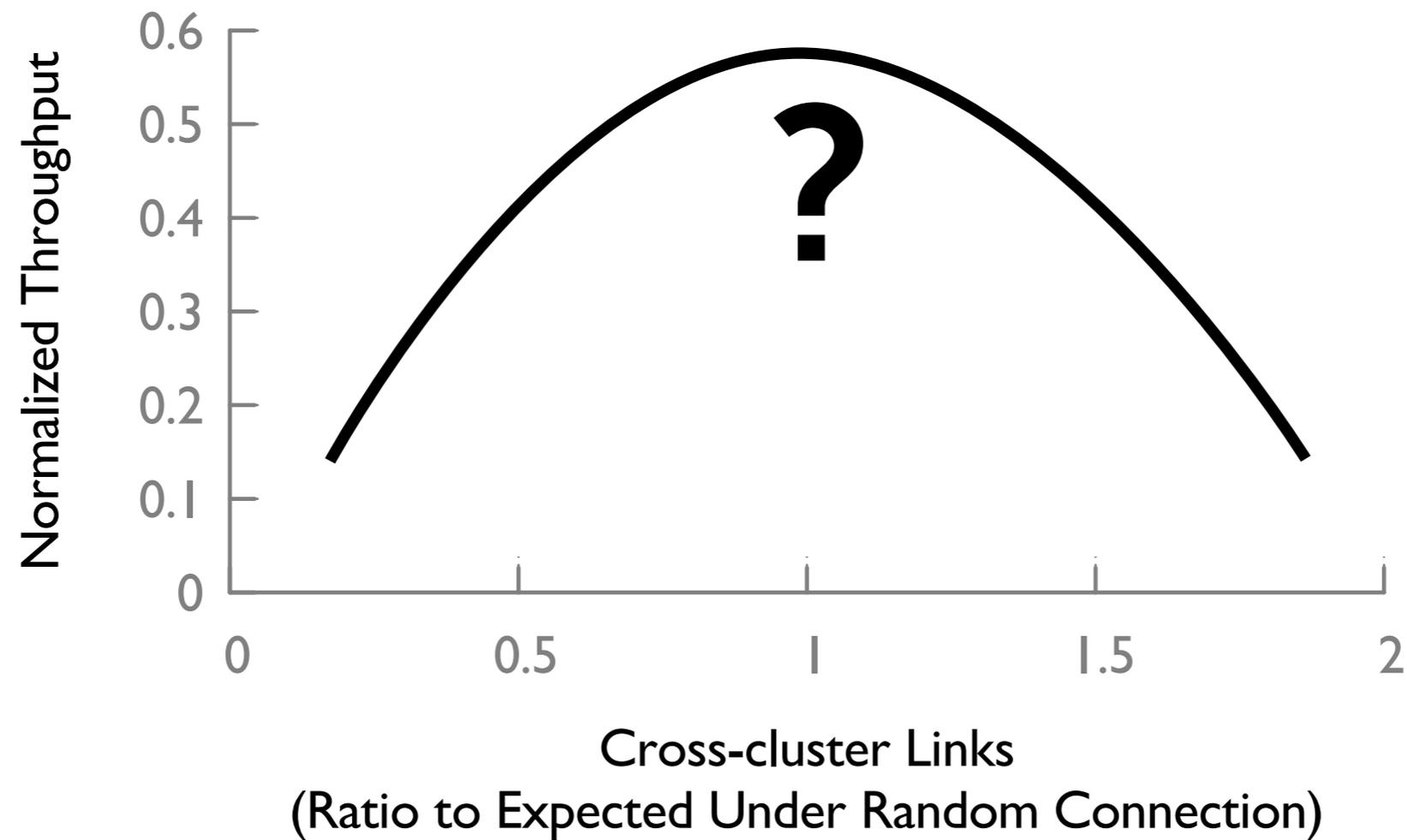
Interconnecting switches



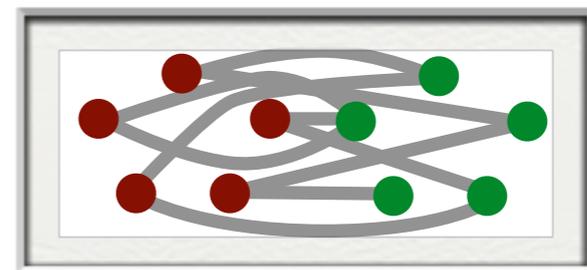
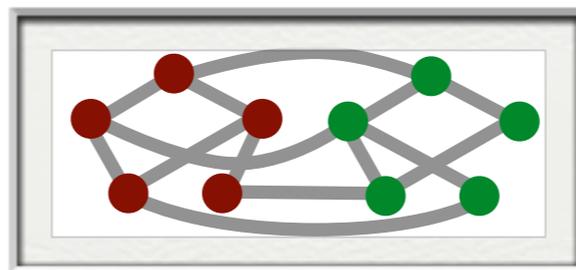
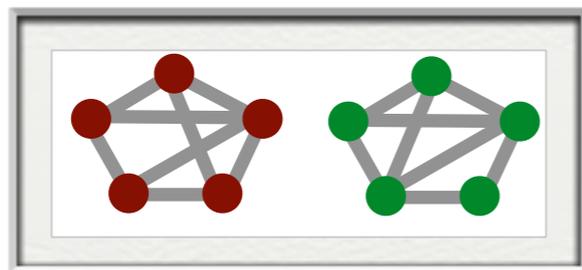
Interconnecting switches



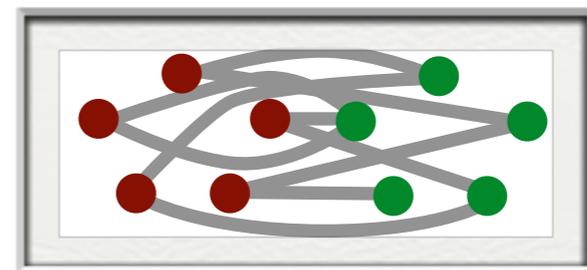
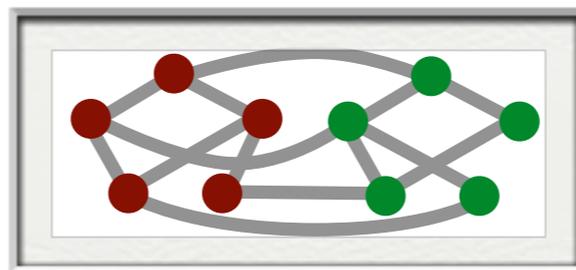
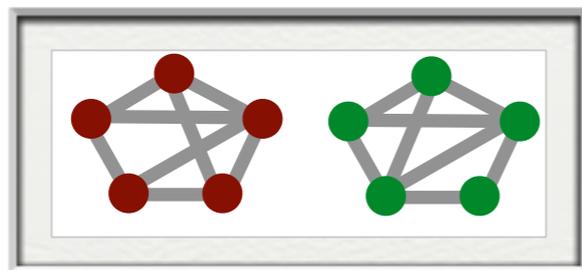
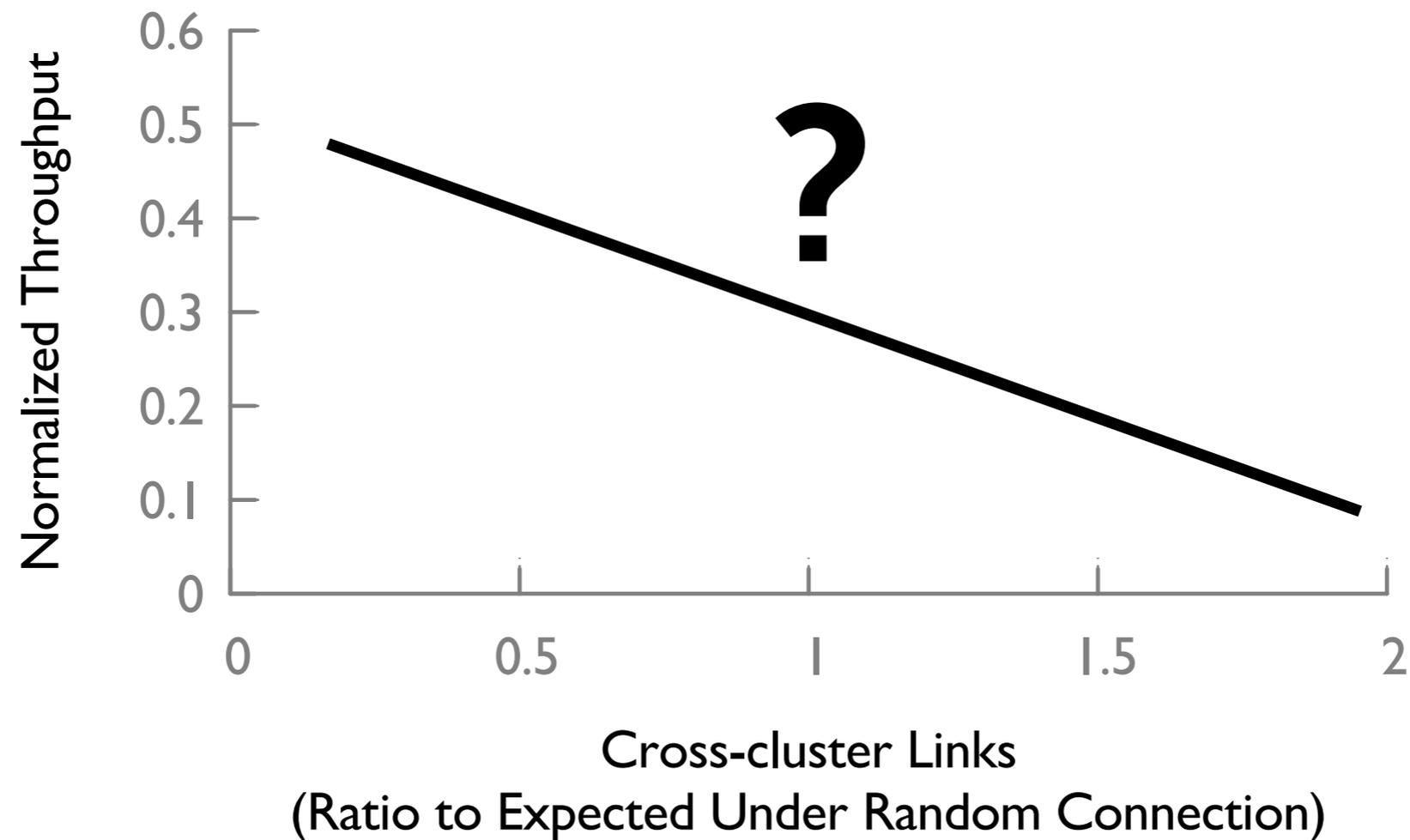
Interconnecting switches



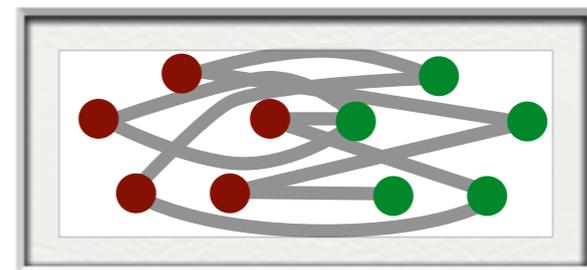
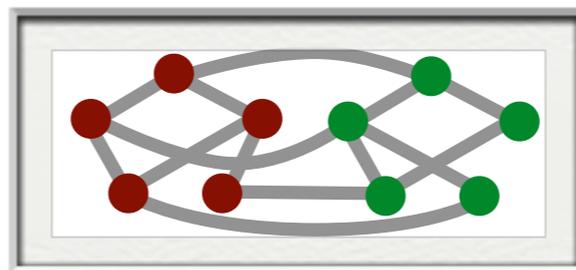
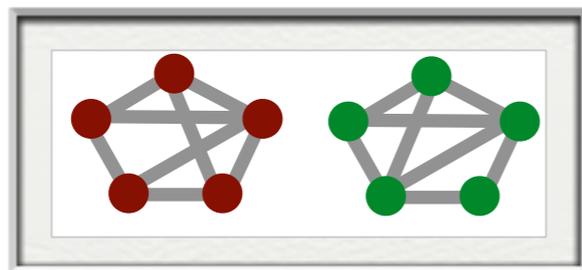
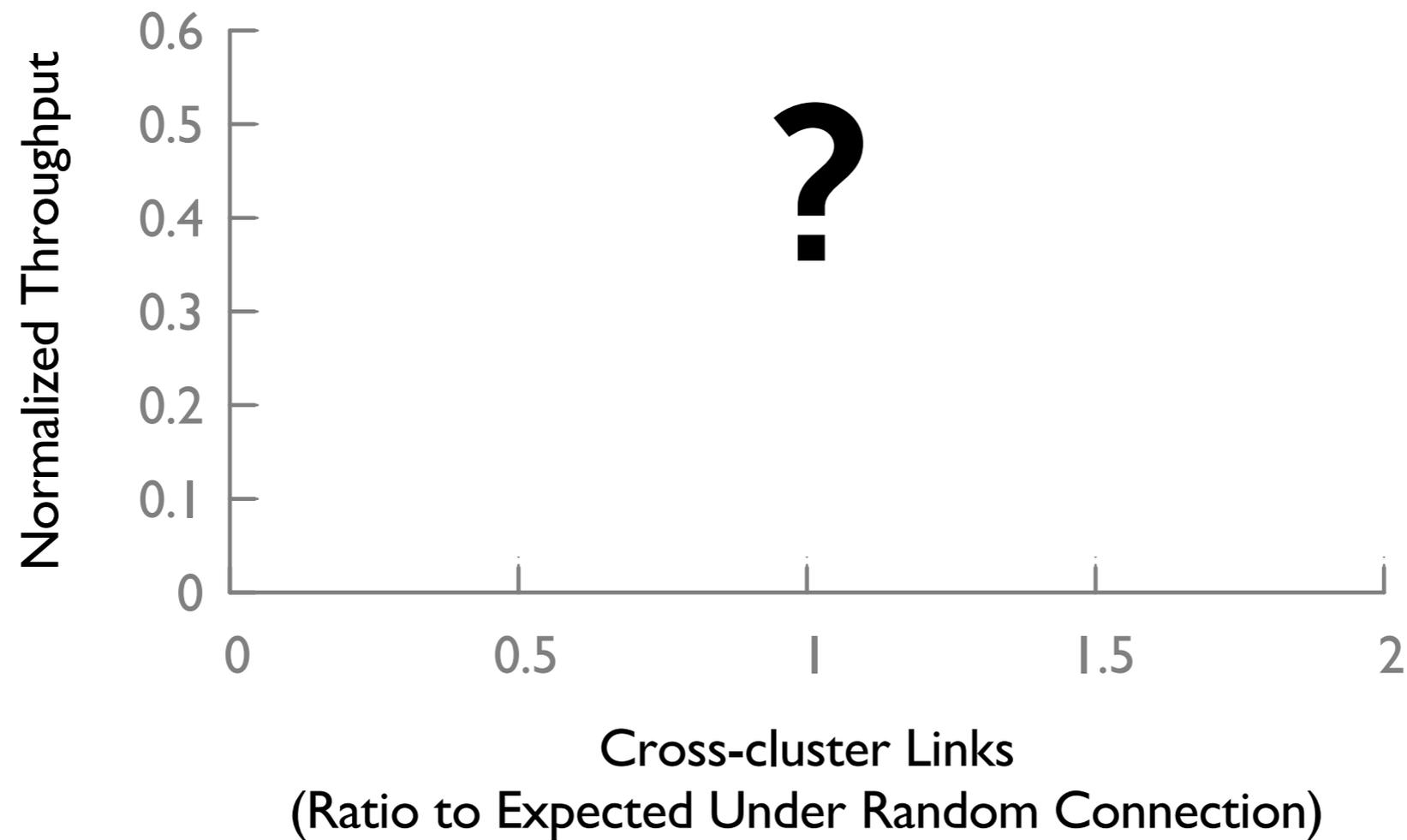
Interconnecting switches



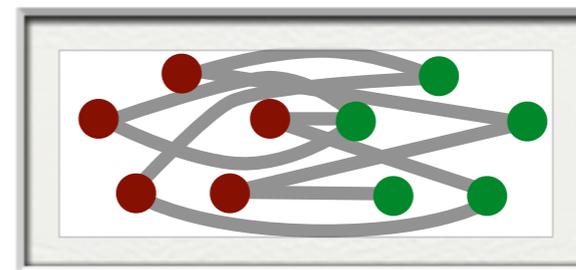
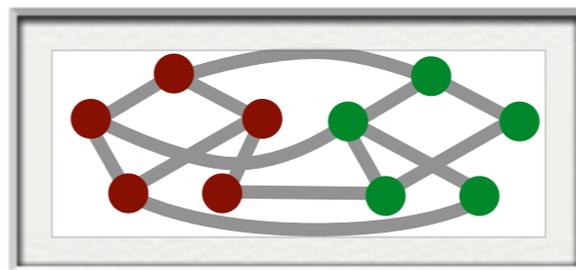
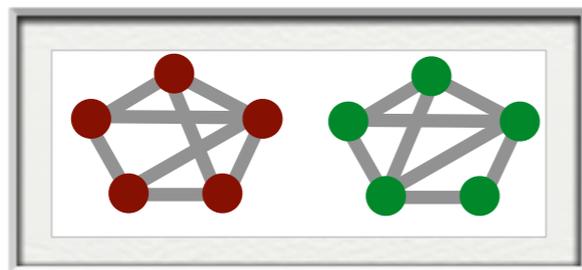
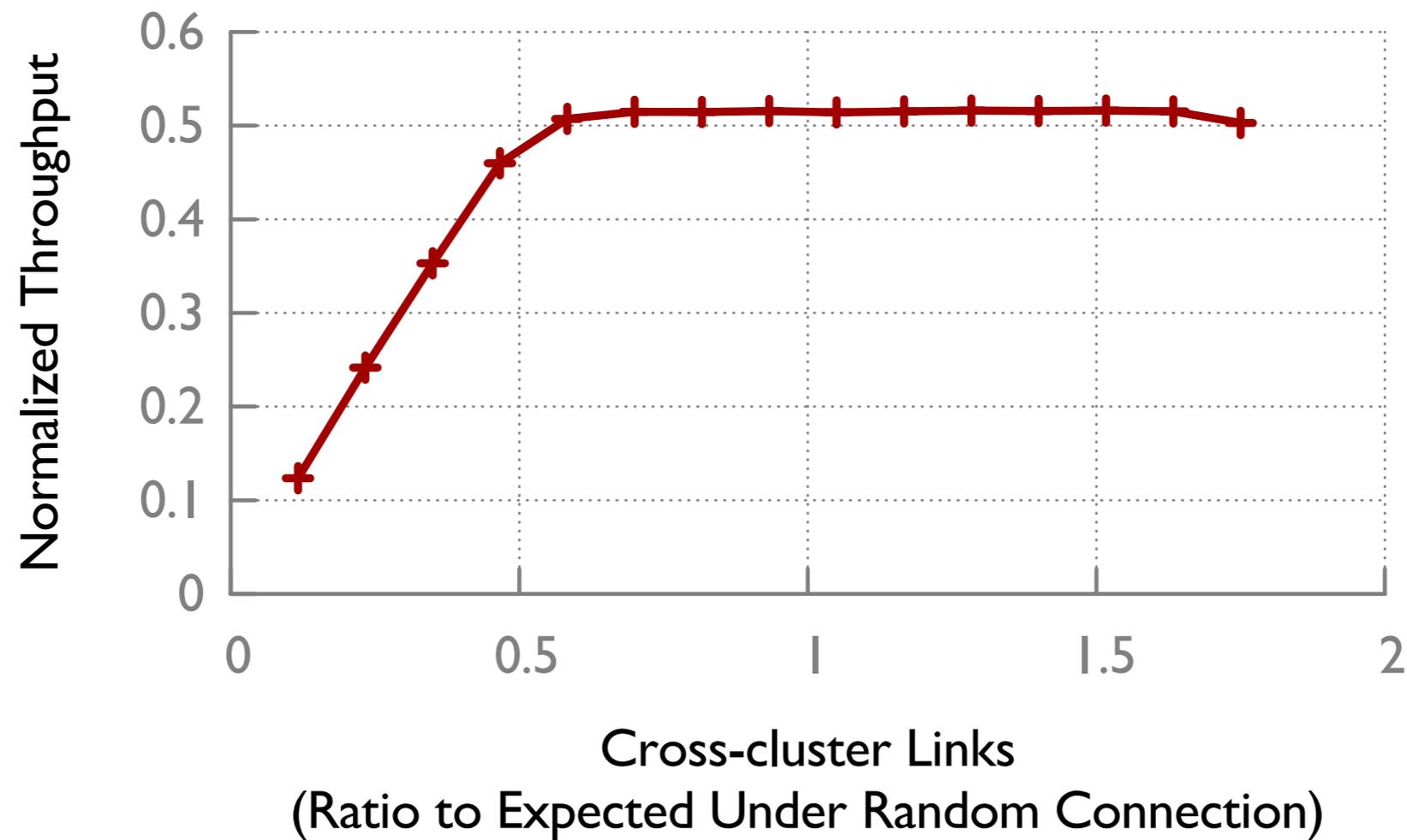
Interconnecting switches



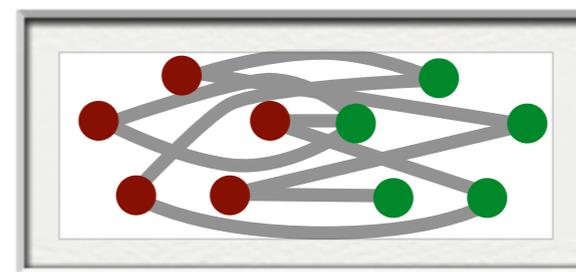
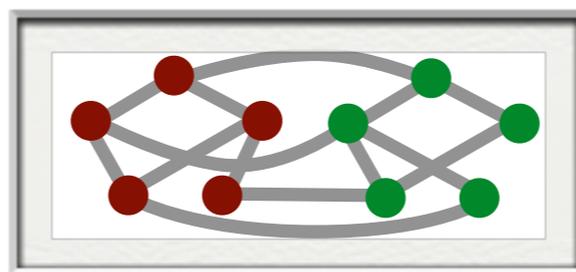
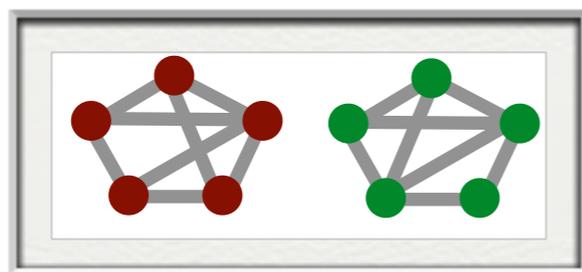
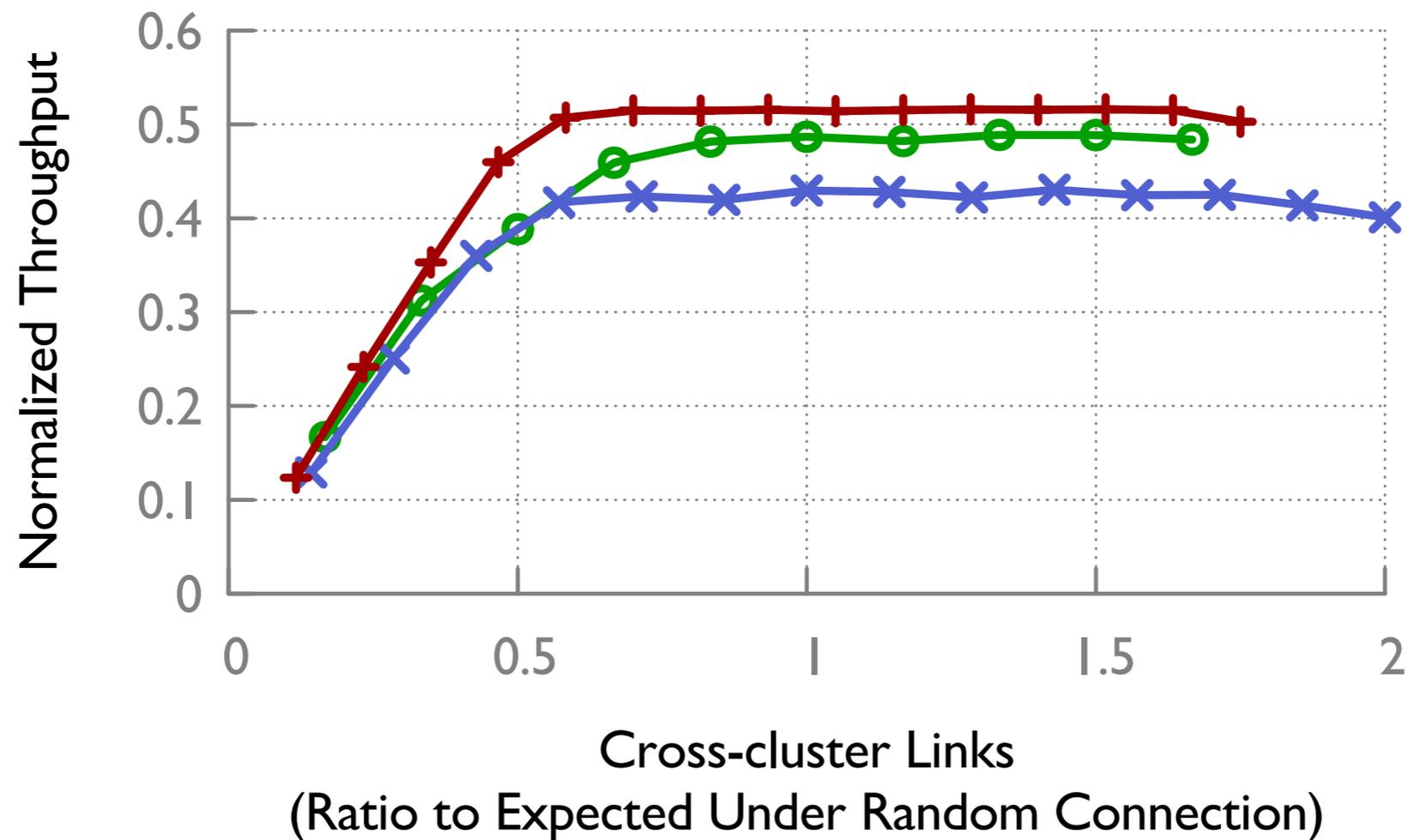
Interconnecting switches



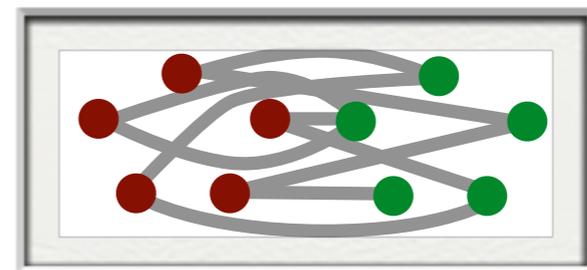
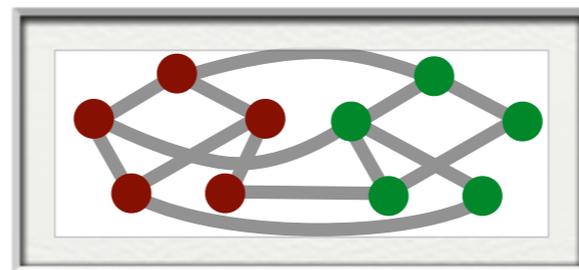
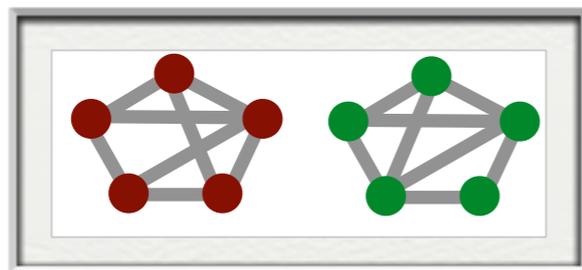
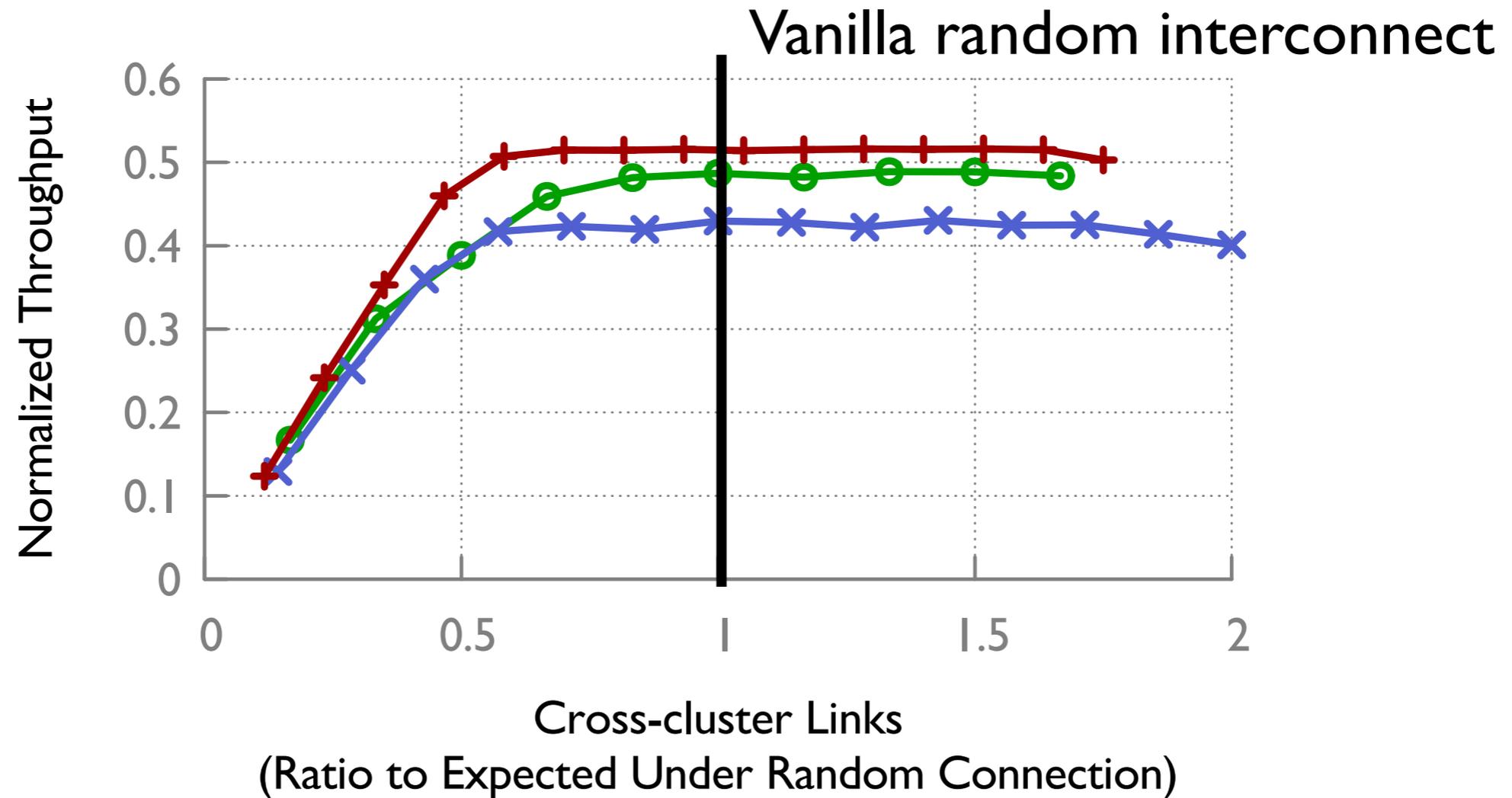
Interconnecting switches



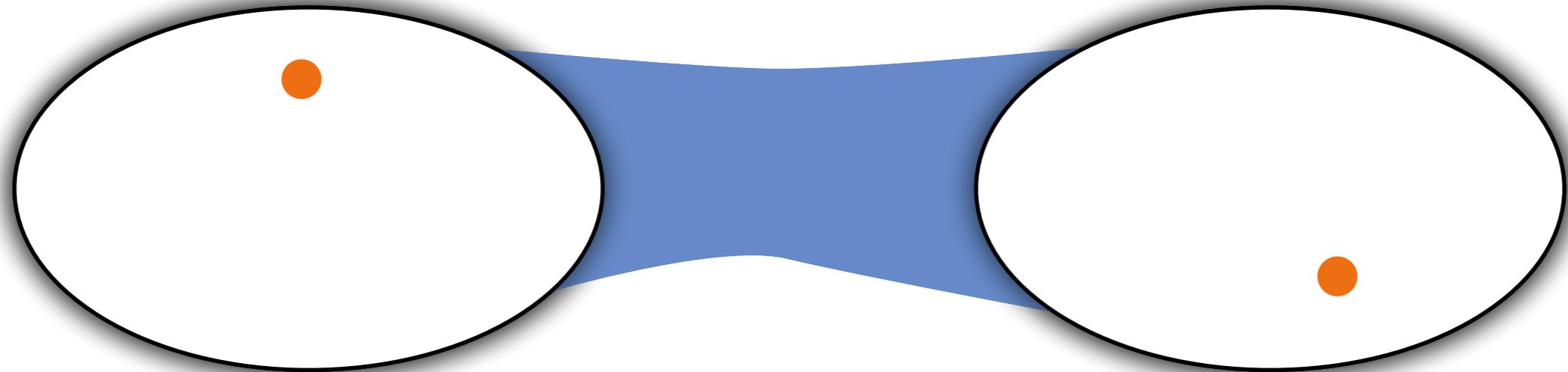
Interconnecting switches



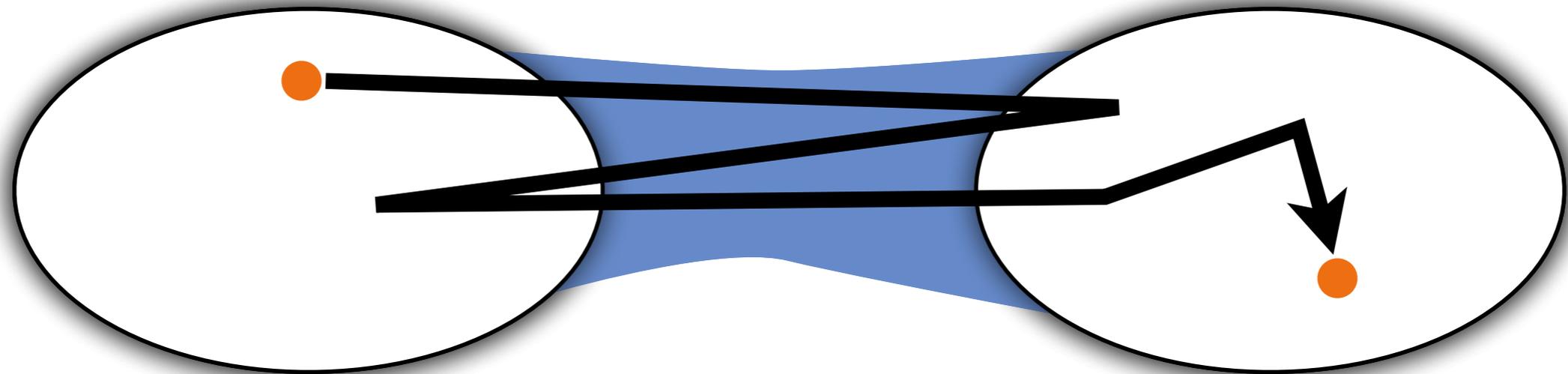
Interconnecting switches



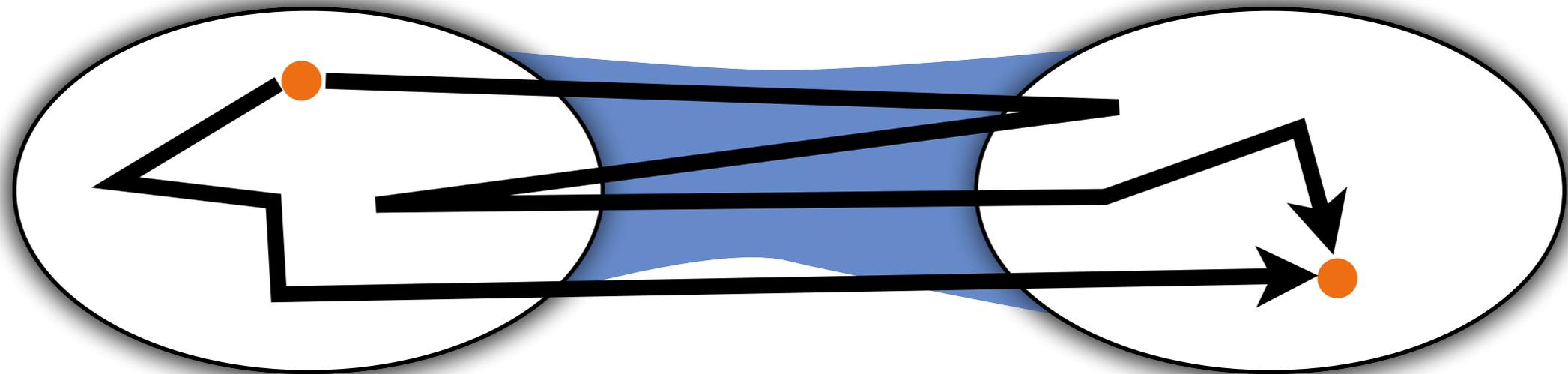
Intuition



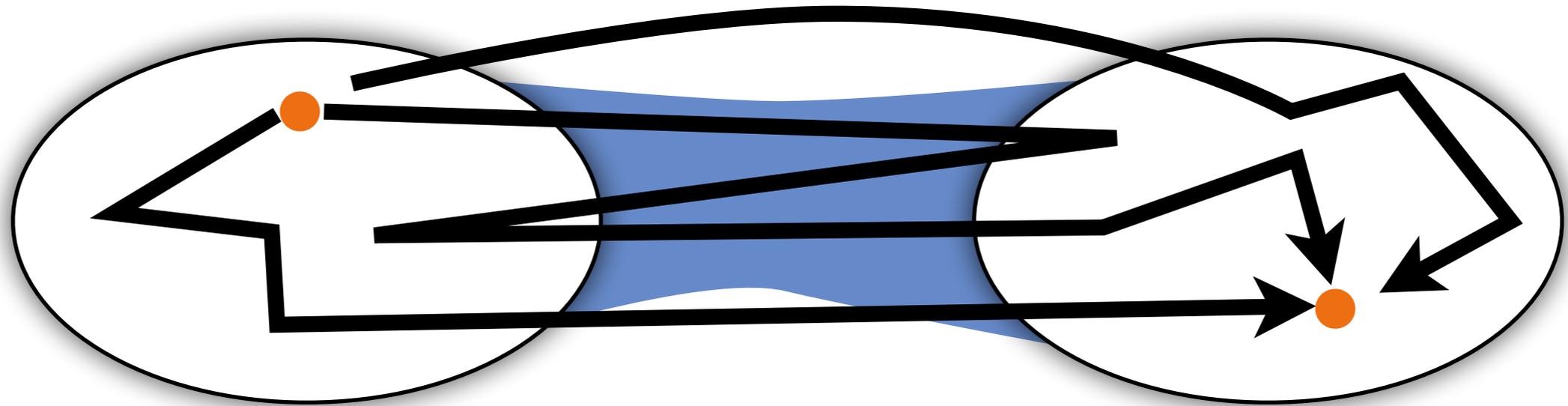
Intuition



Intuition

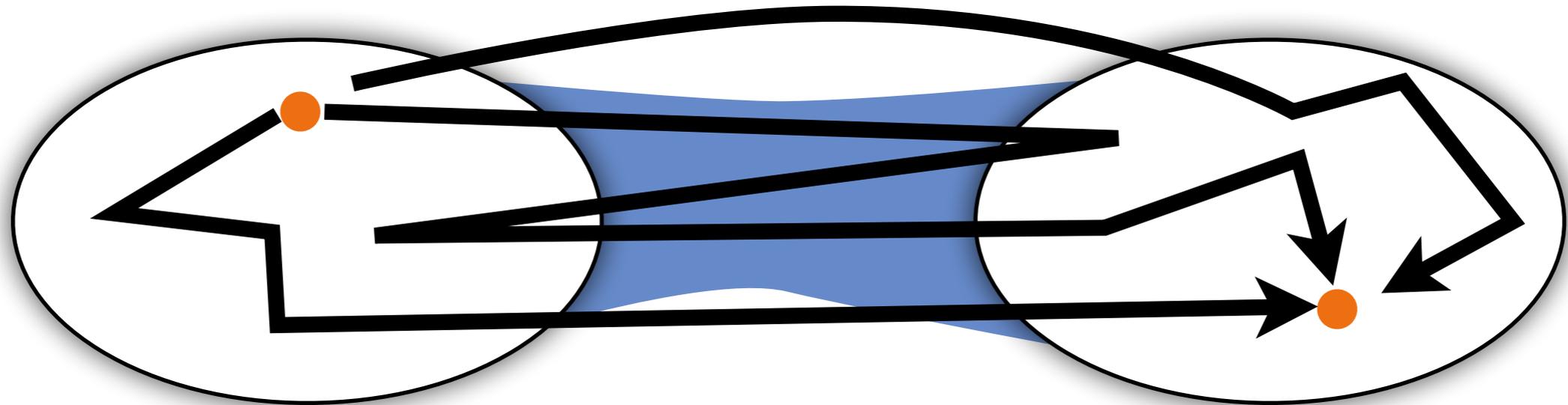


Intuition



Still need one crossing!

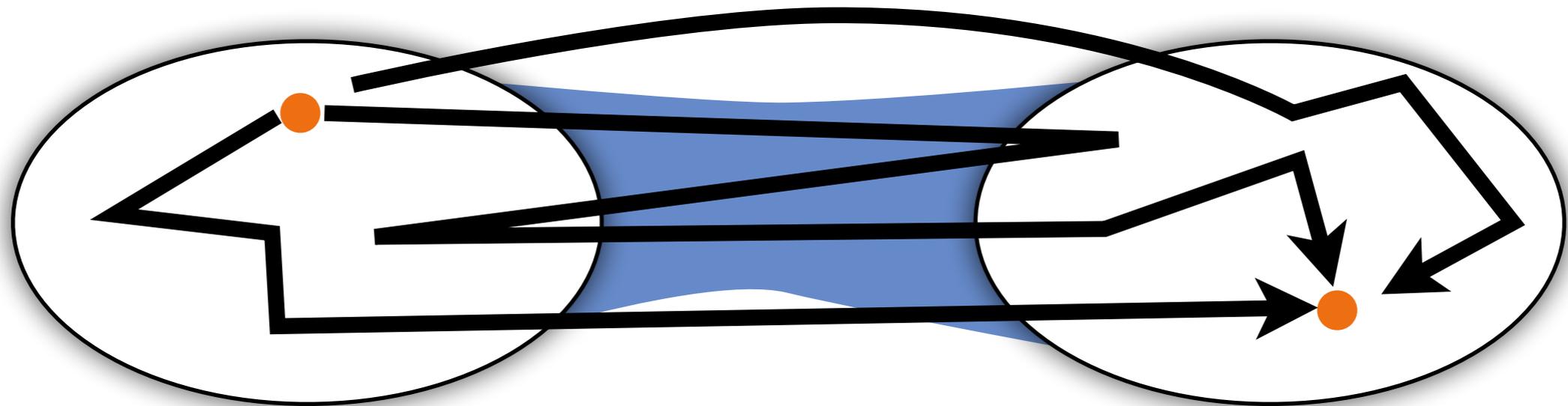
Intuition



Still need one crossing!

$$\Theta \left(\frac{1}{APL} \right)$$

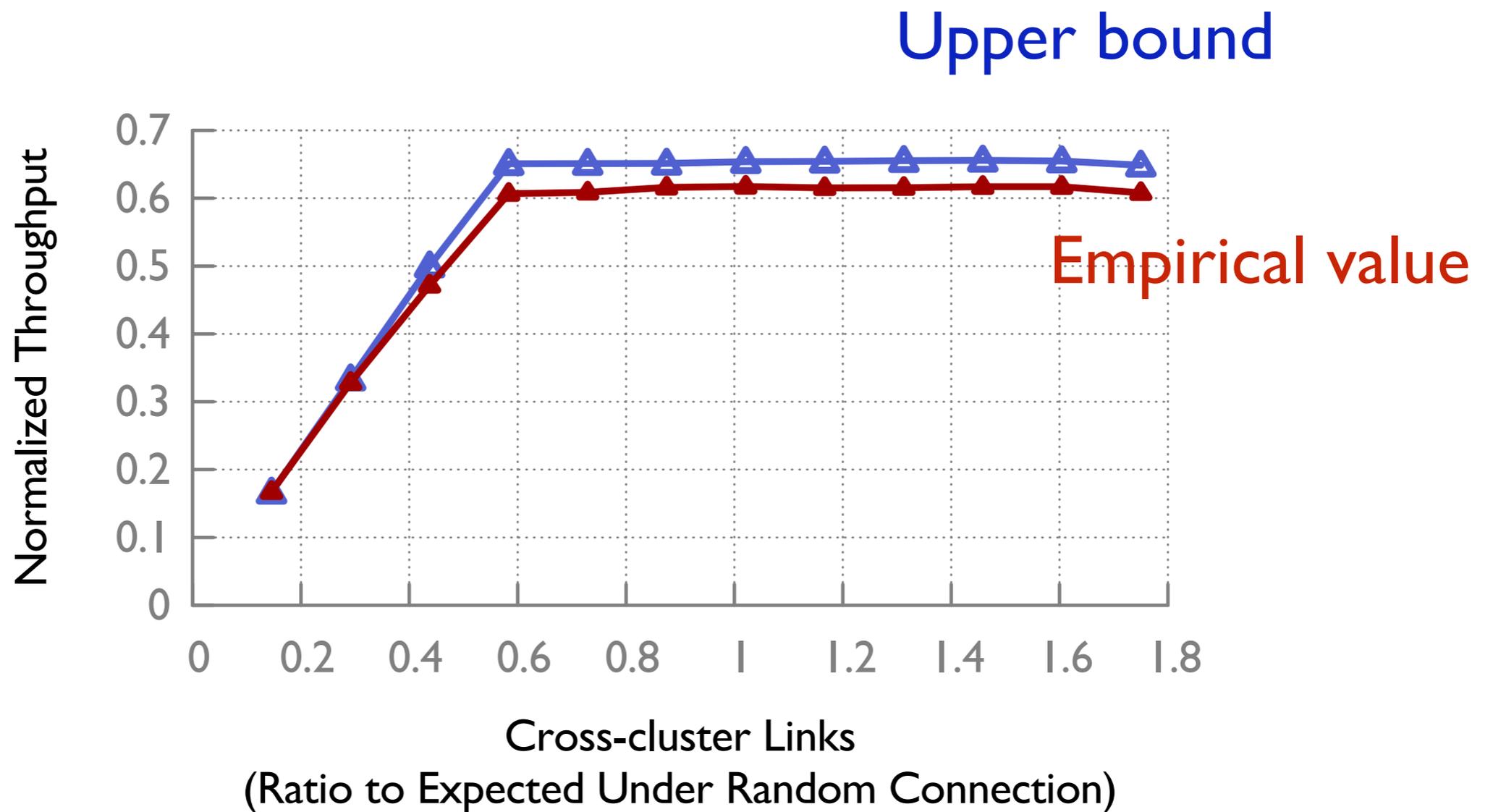
Intuition



Still need one crossing!

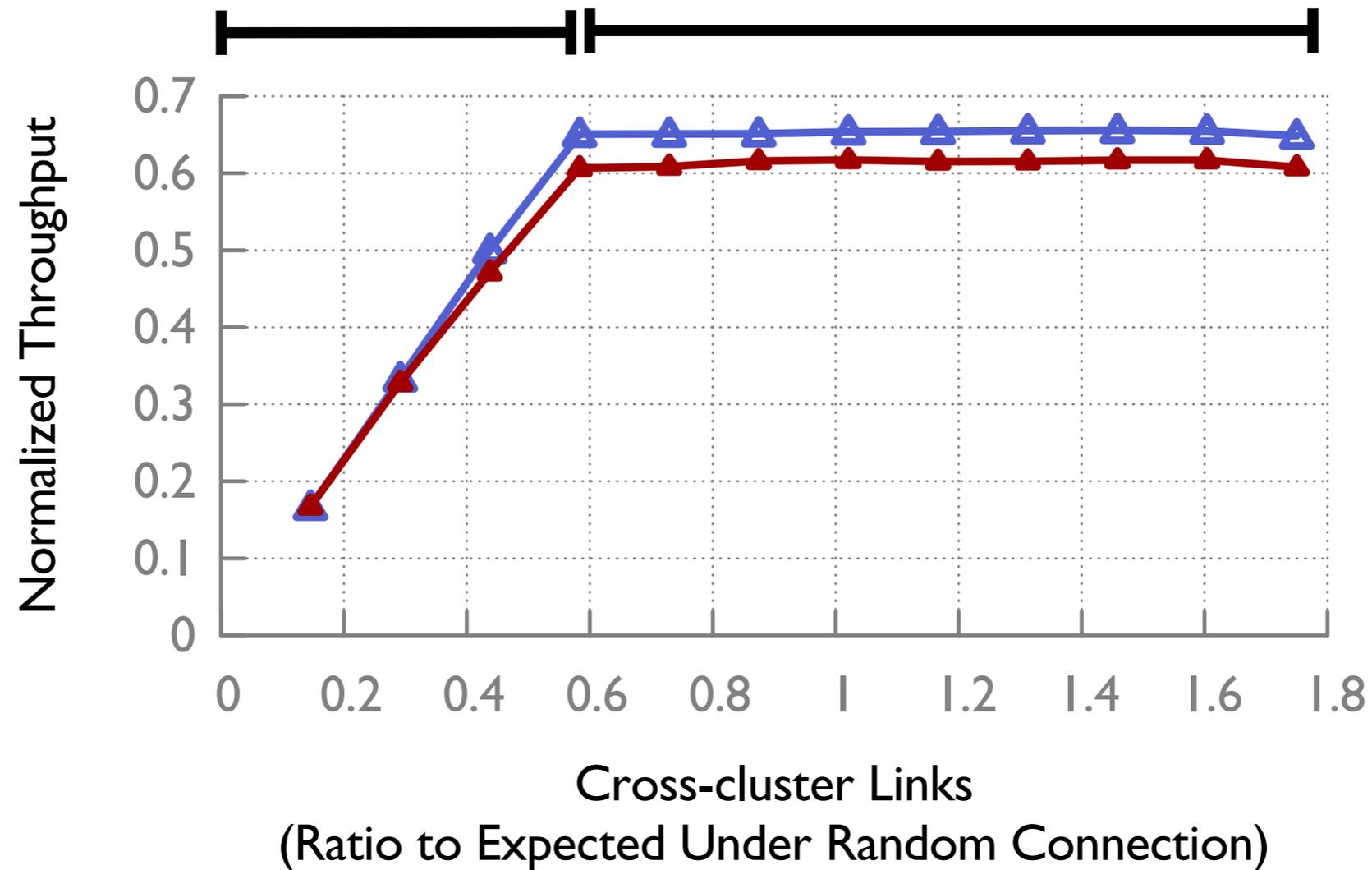
Throughput should drop when less than $\Theta\left(\frac{1}{APL}\right)$ of total capacity crosses the cut!

Explaining throughput

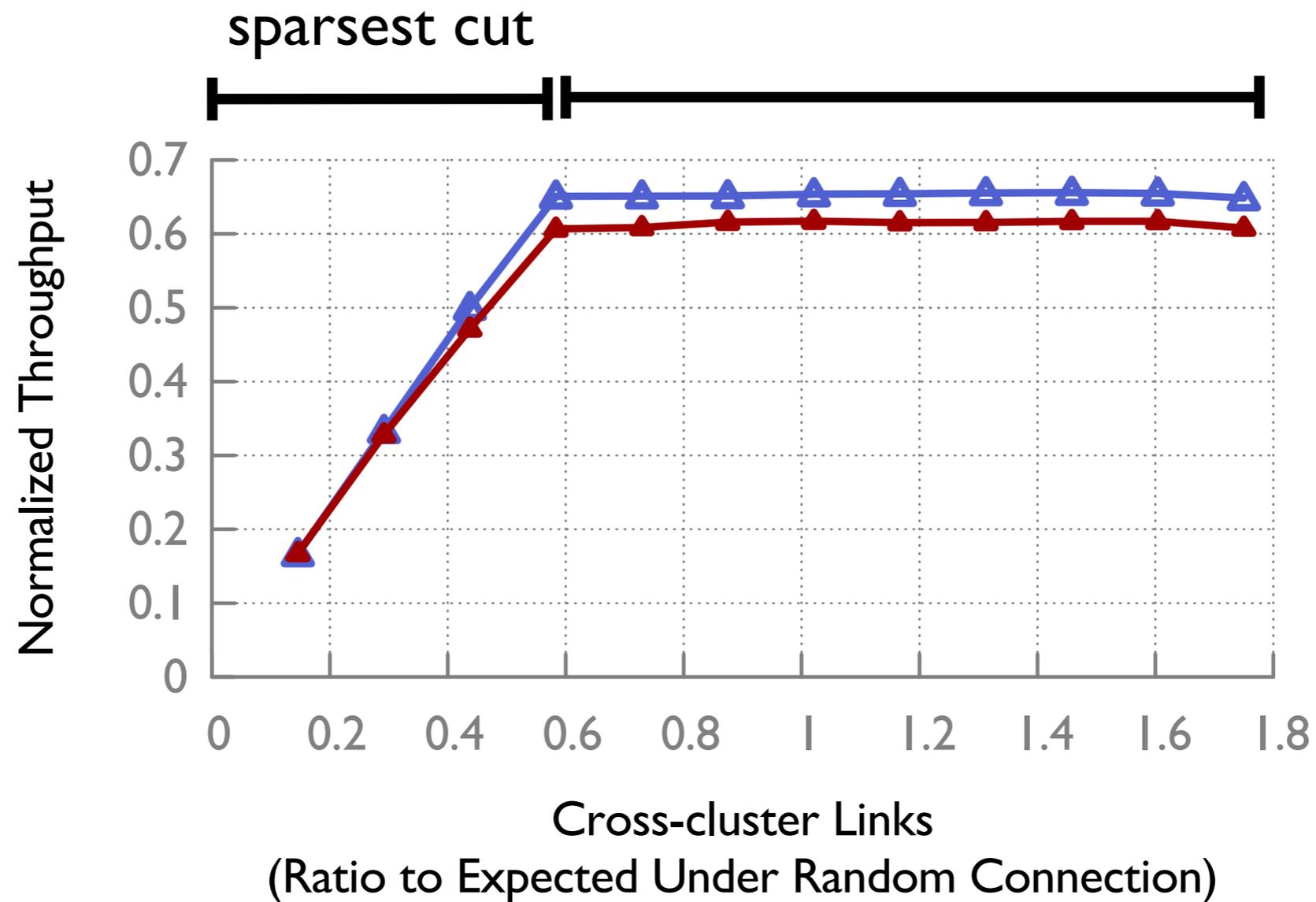


And constant-factor matching lower bounds in special case

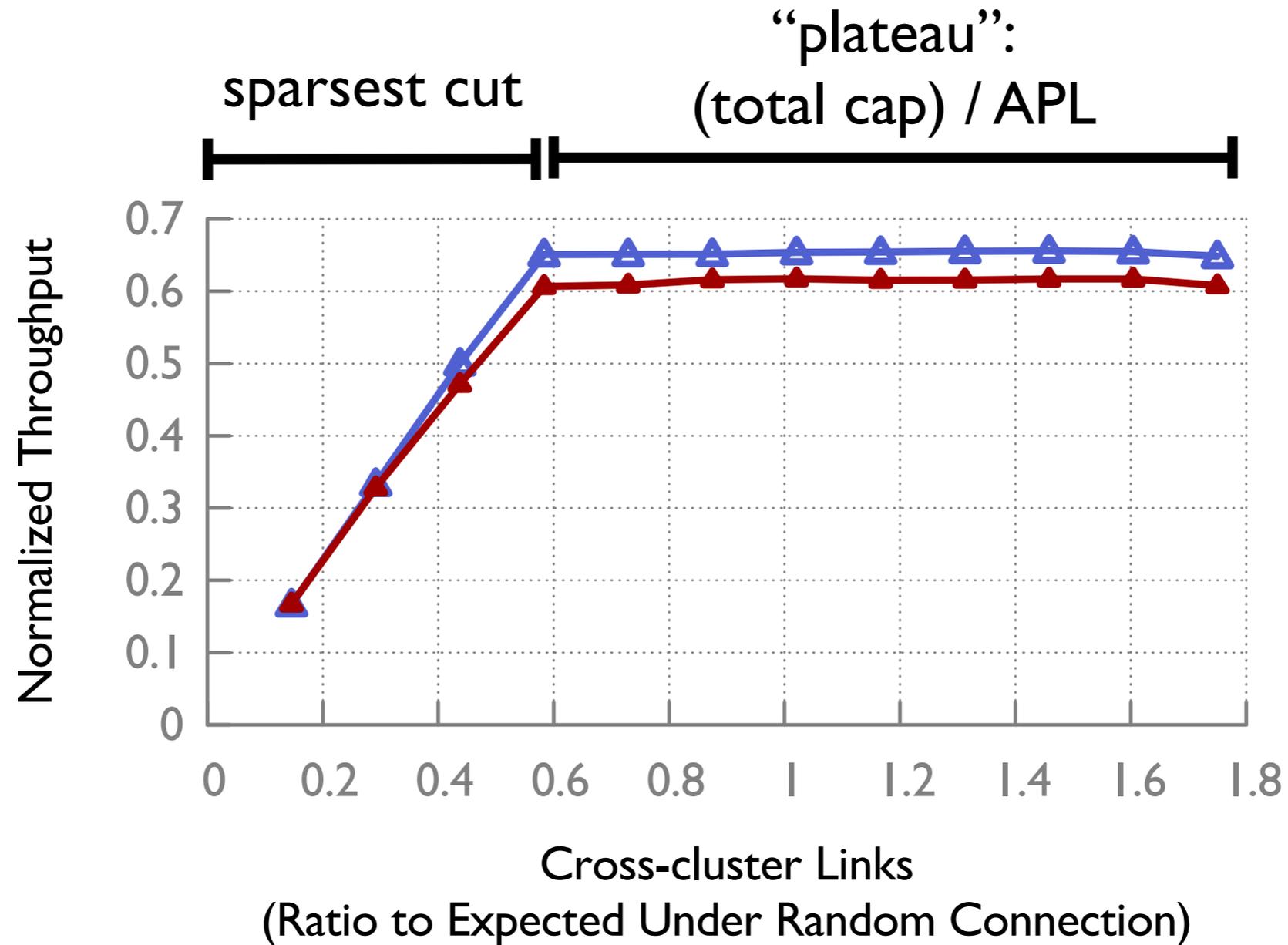
Two regimes of throughput



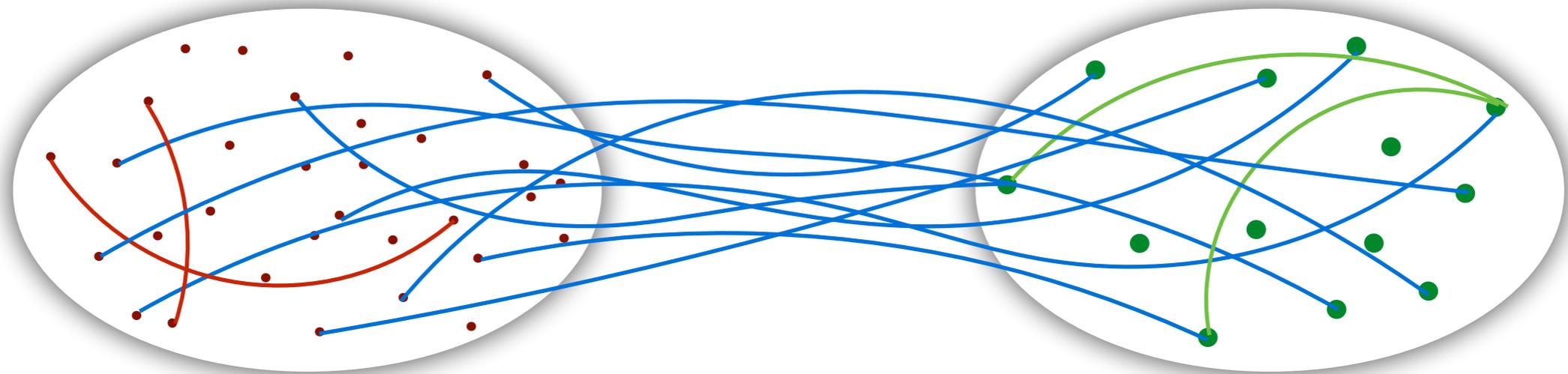
Two regimes of throughput



Two regimes of throughput

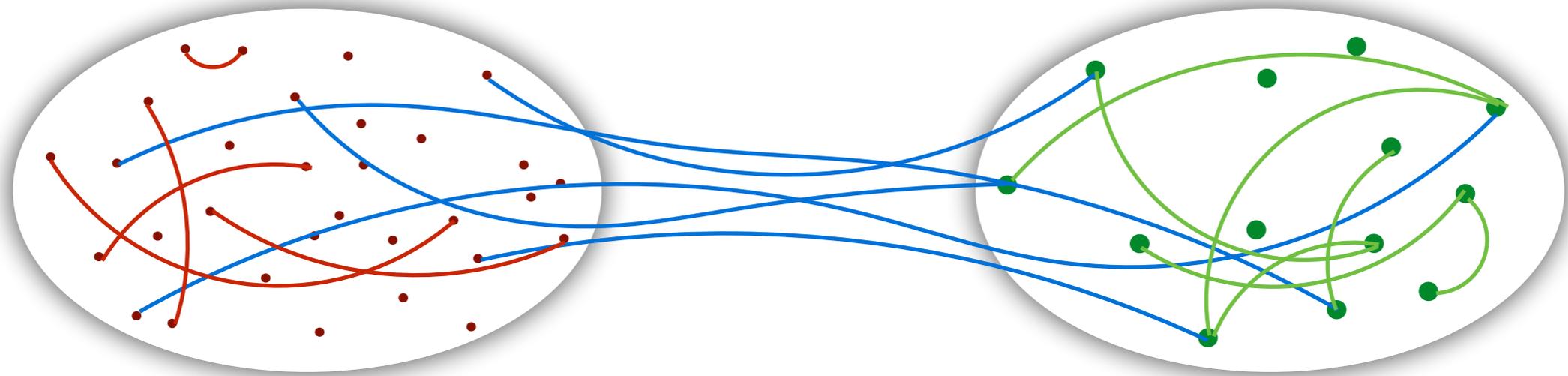


Implications



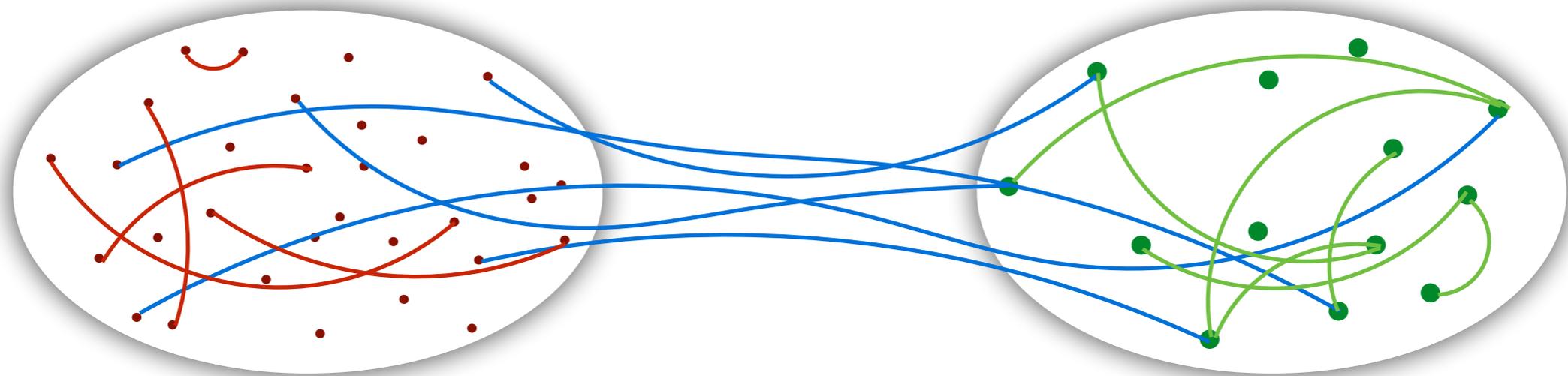
A wide range of connectivity options

Implications



A wide range of connectivity options

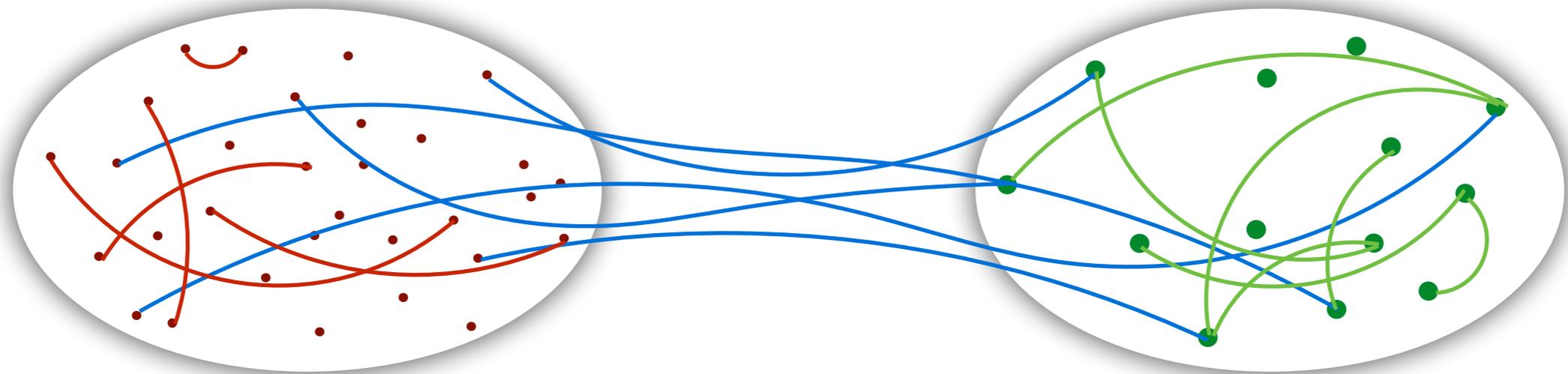
Implications



A wide range of connectivity options

Bisection bandwidth \neq throughput

Implications



A wide range of connectivity options

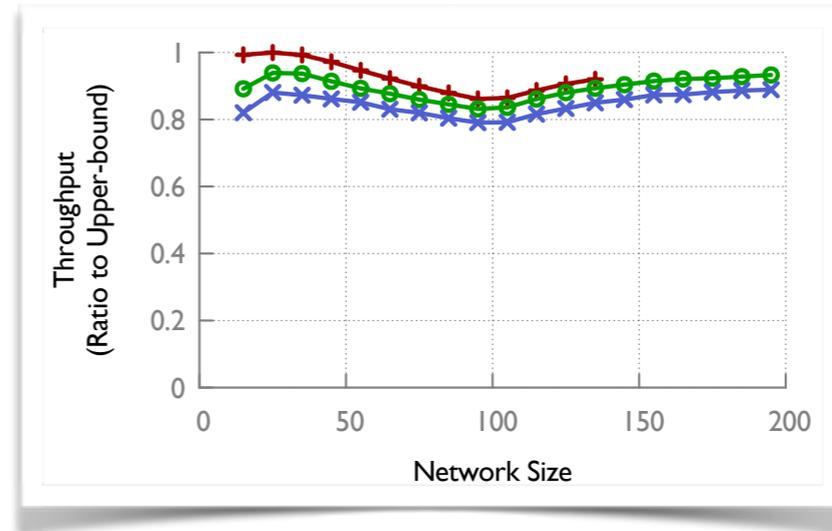
Bisection bandwidth \neq throughput

Greater freedom in cabling

Quick recap!

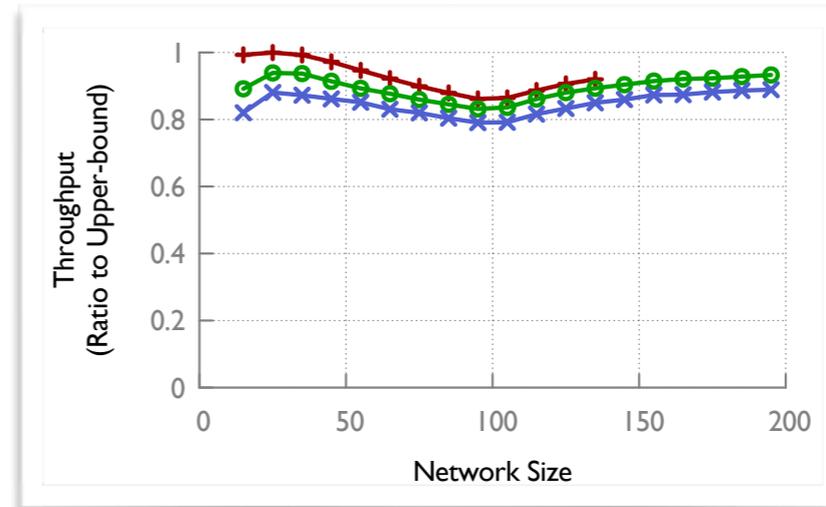


How close can we get to optimal network capacity?

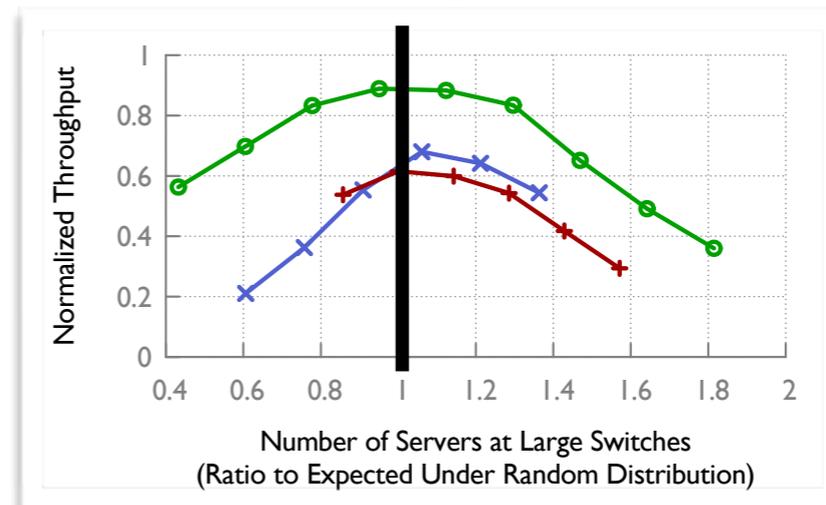




How close can we get to optimal network capacity?

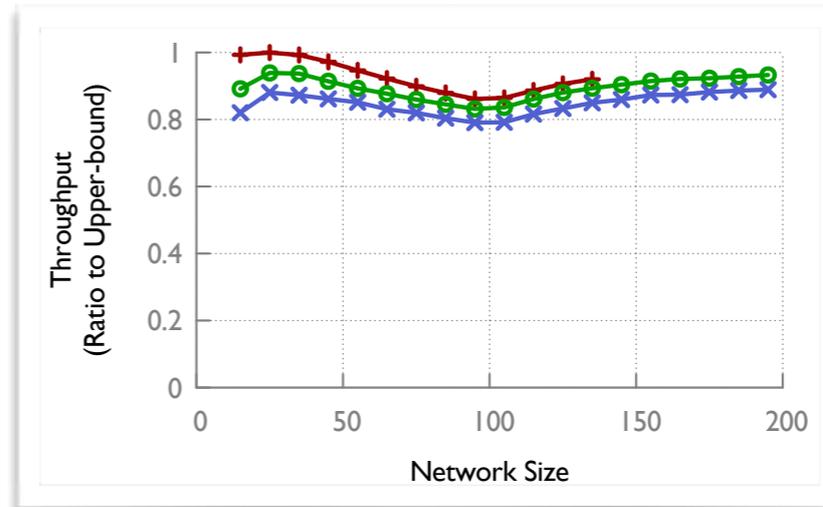


How should we distribute servers?



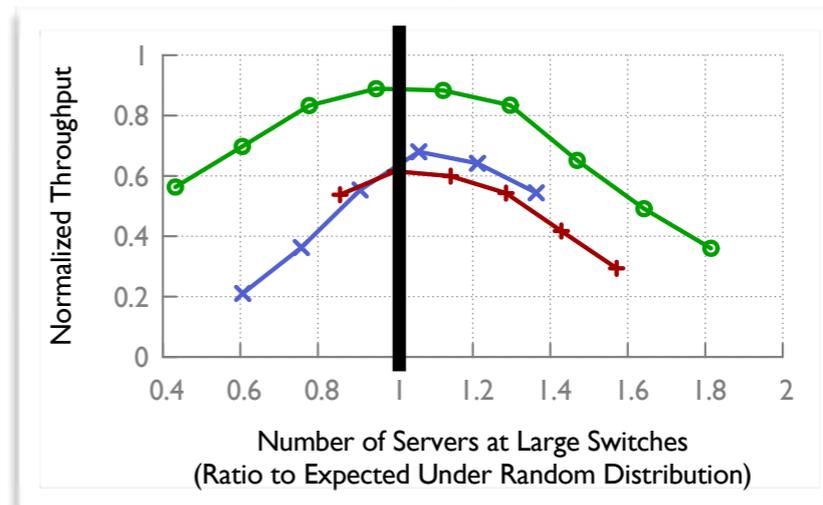
0

How close can we get to optimal network capacity?



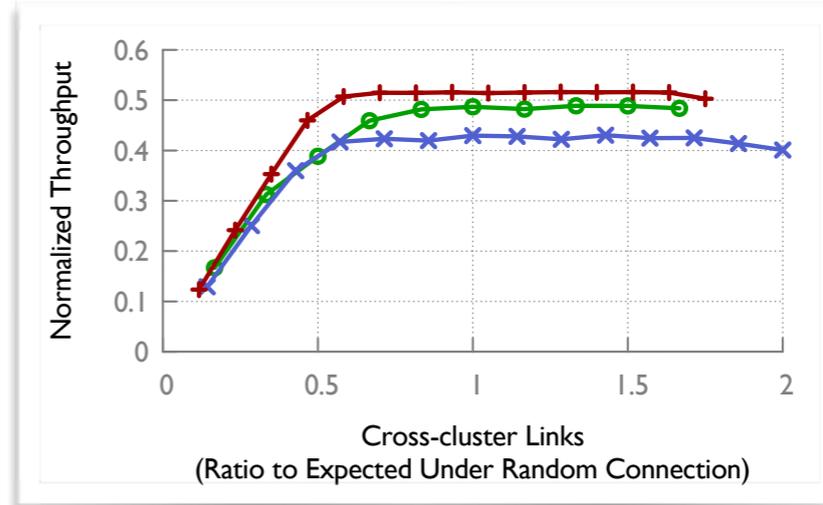
1

How should we distribute servers?



2

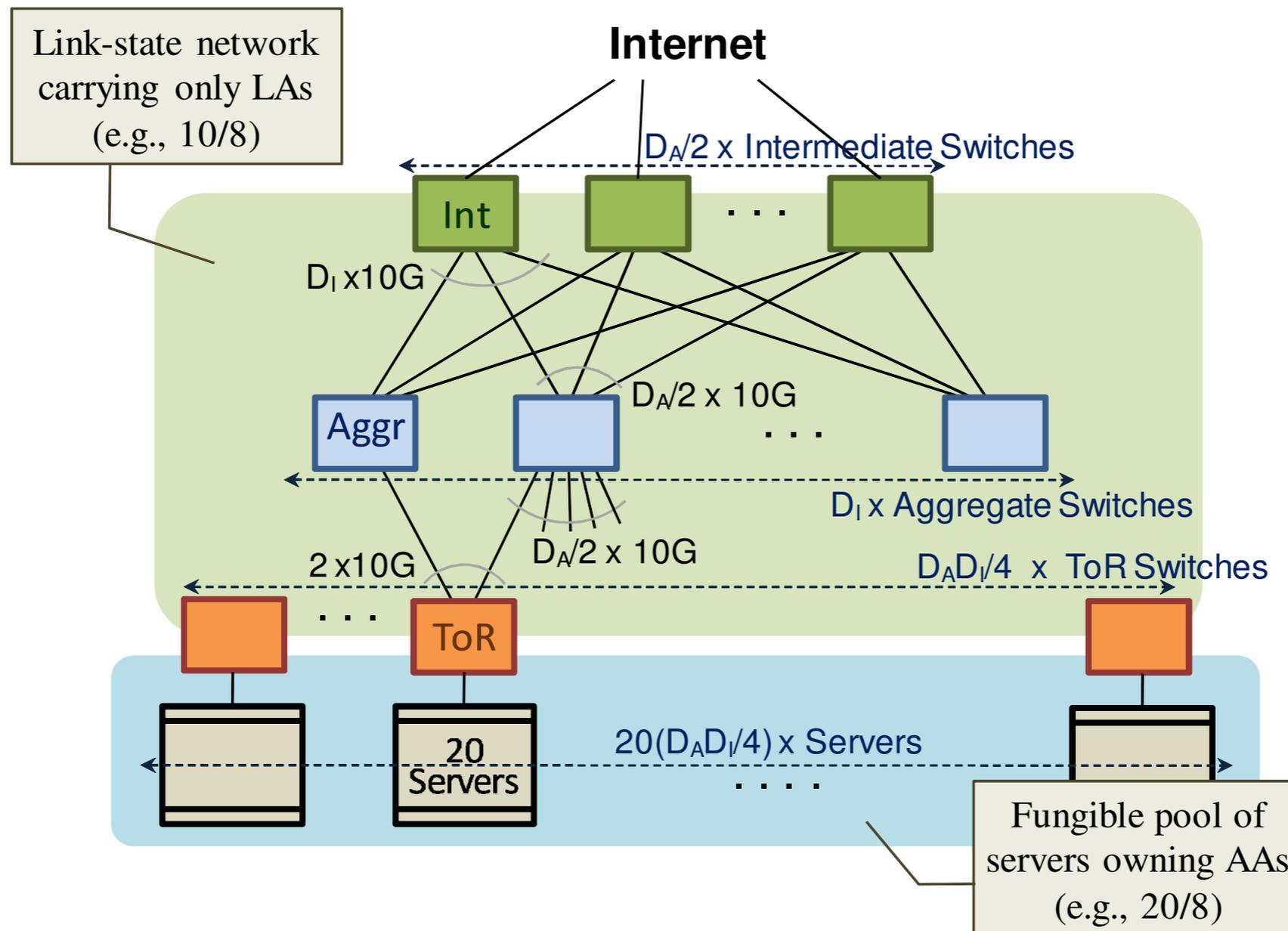
How should we interconnect switches?



Improving a REAL heterogeneous topology

The VL2 topology

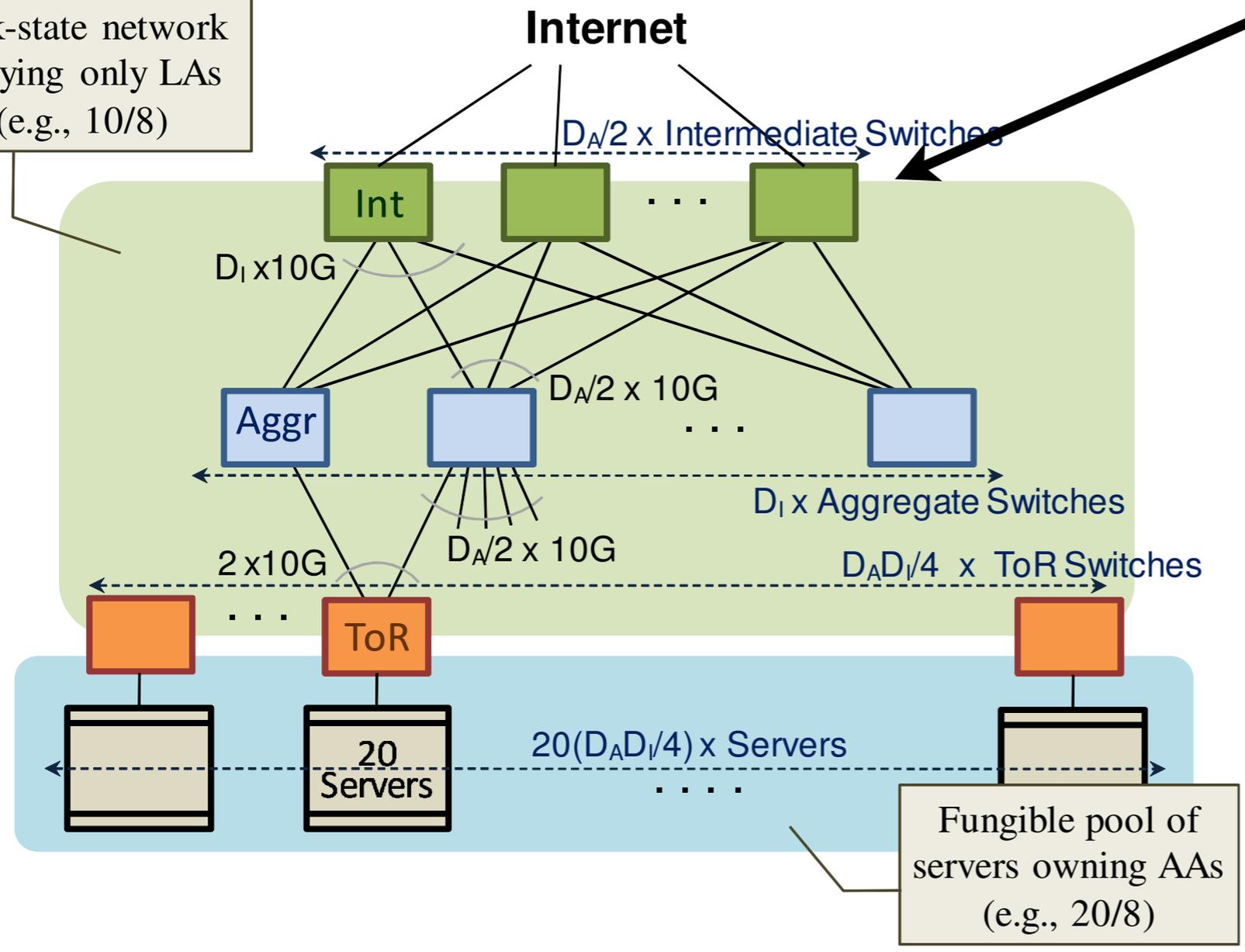
[Greenburg, Hamilton, Jain, Kandula, Kim, Lahiri, Maltz, Patel, Sengupta, SIGCOMM'09]



The VL2 topology

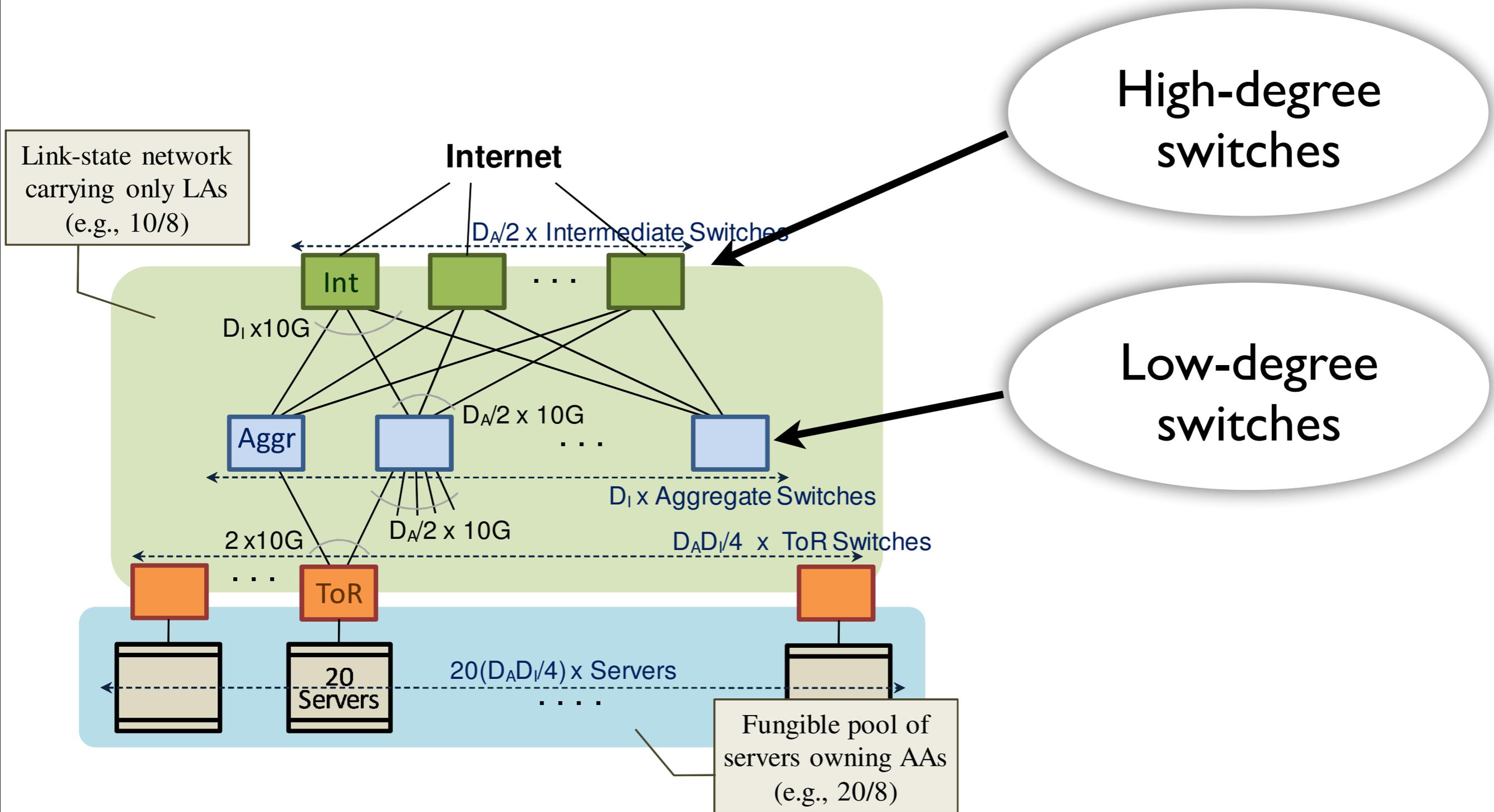
High-degree switches

Link-state network carrying only LAs (e.g., 10/8)

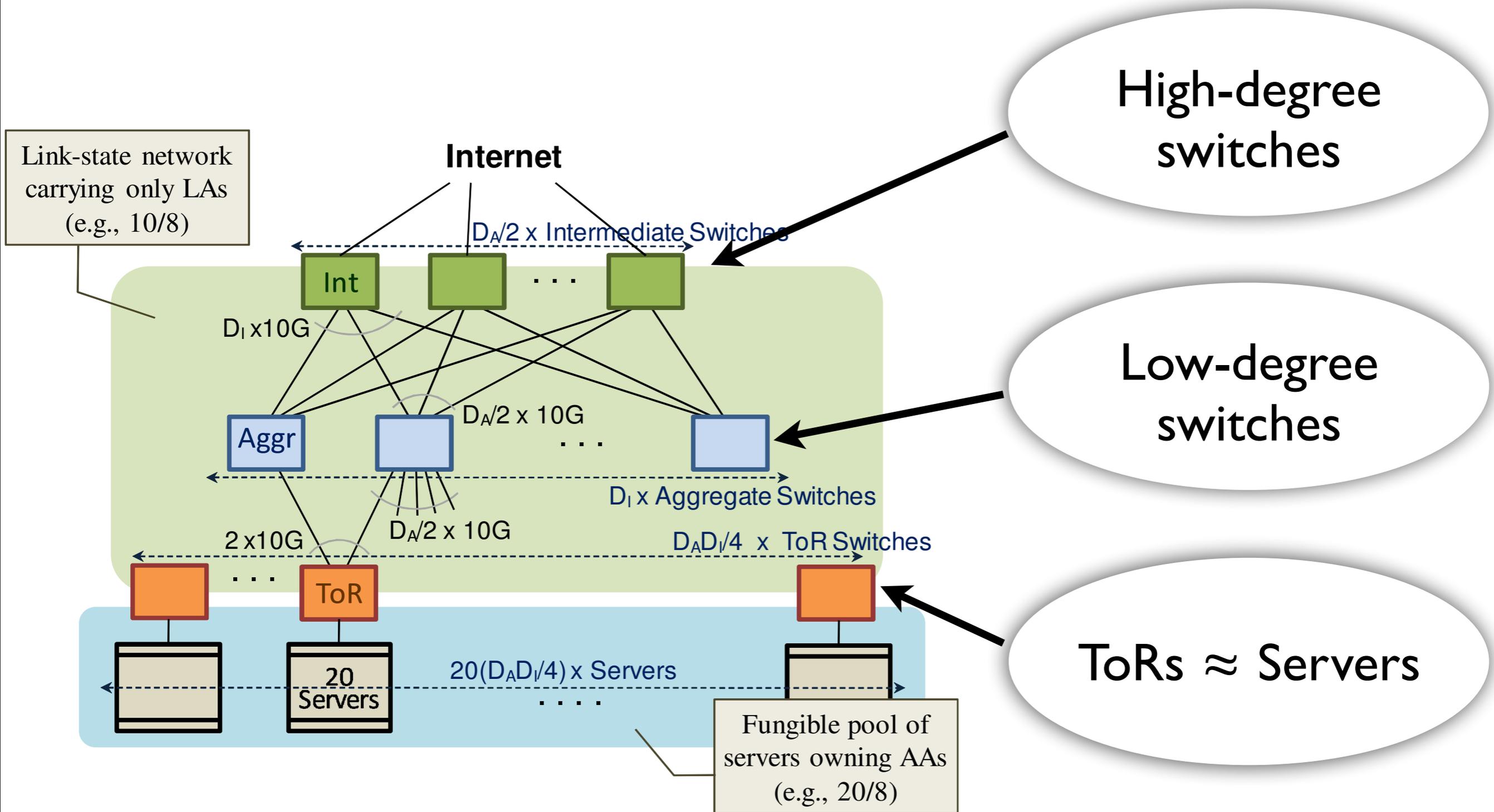


Fungible pool of servers owning AAs (e.g., 20/8)

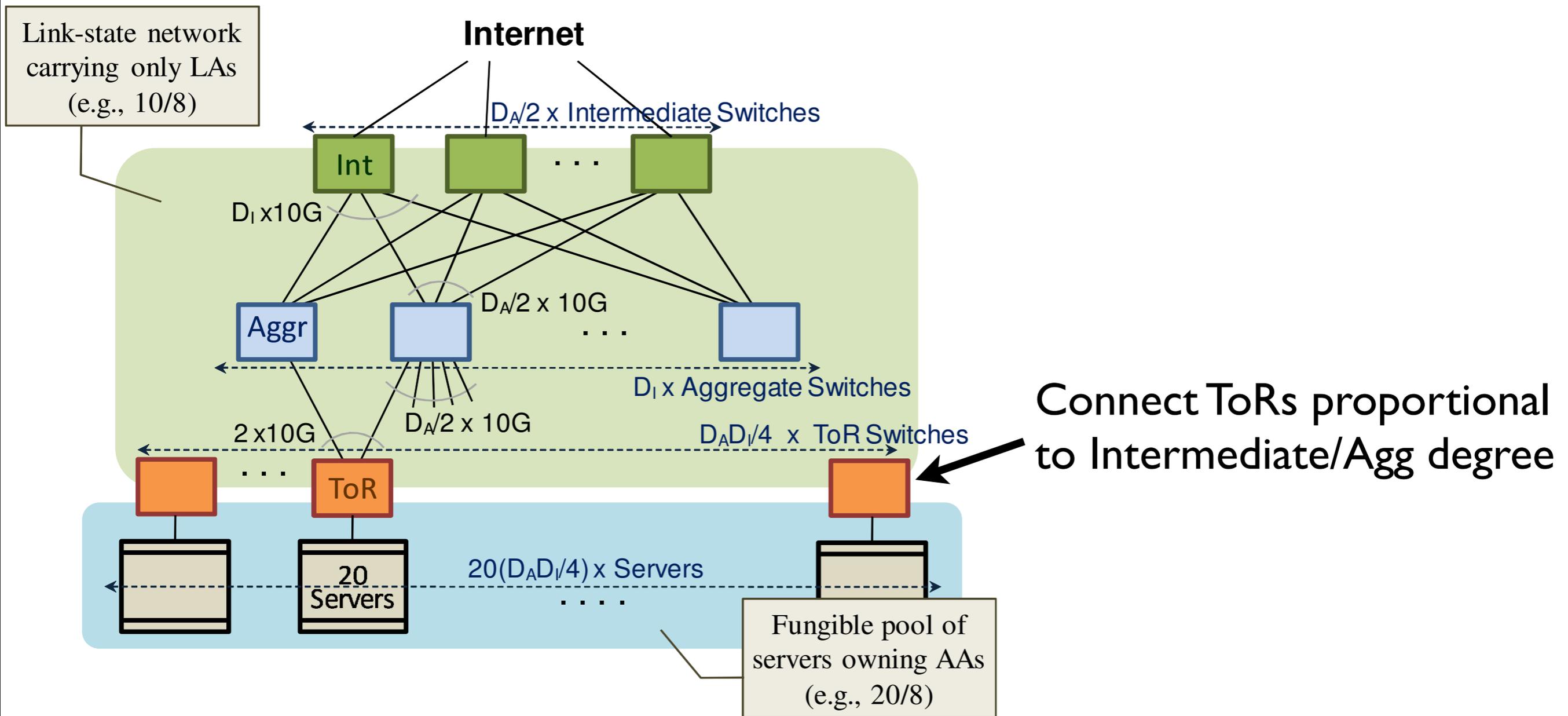
The VL2 topology



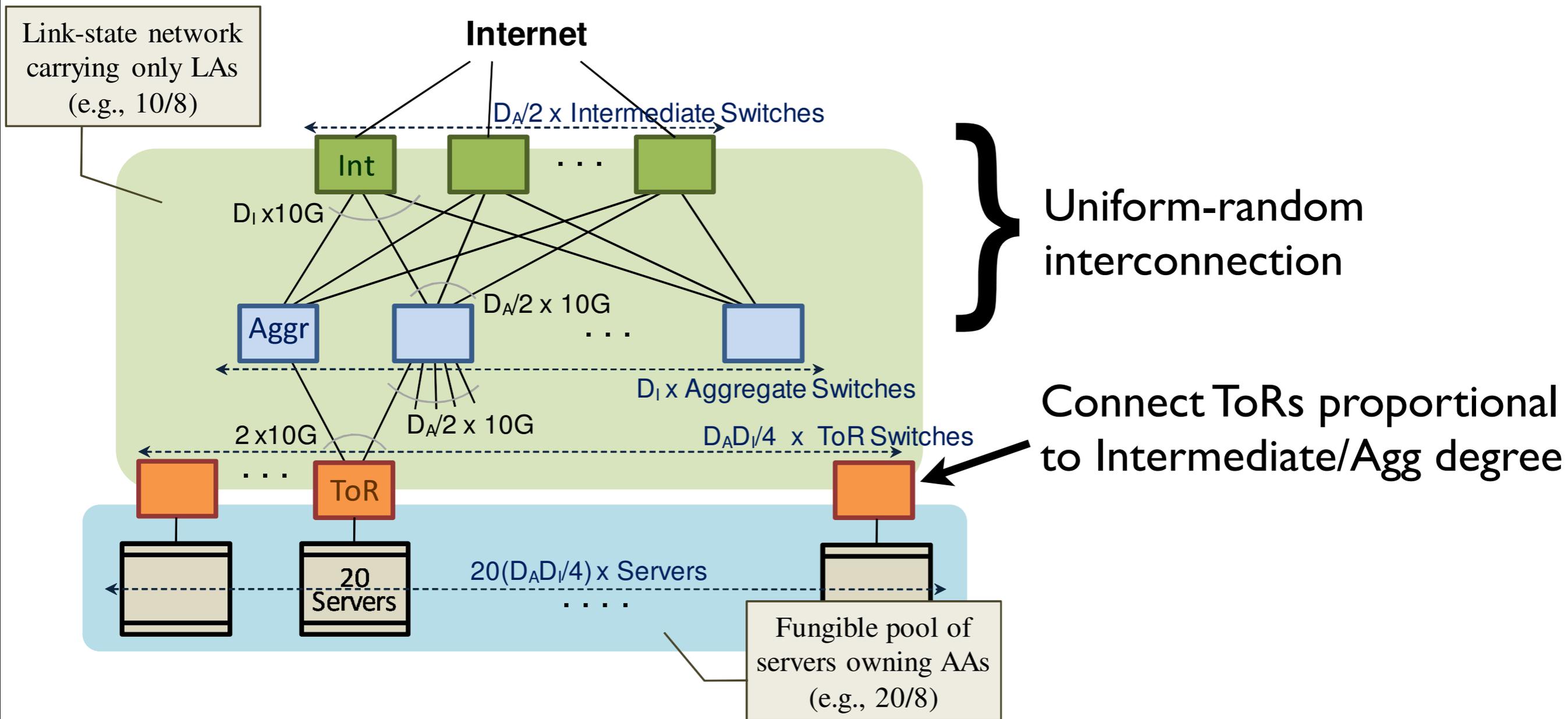
The VL2 topology



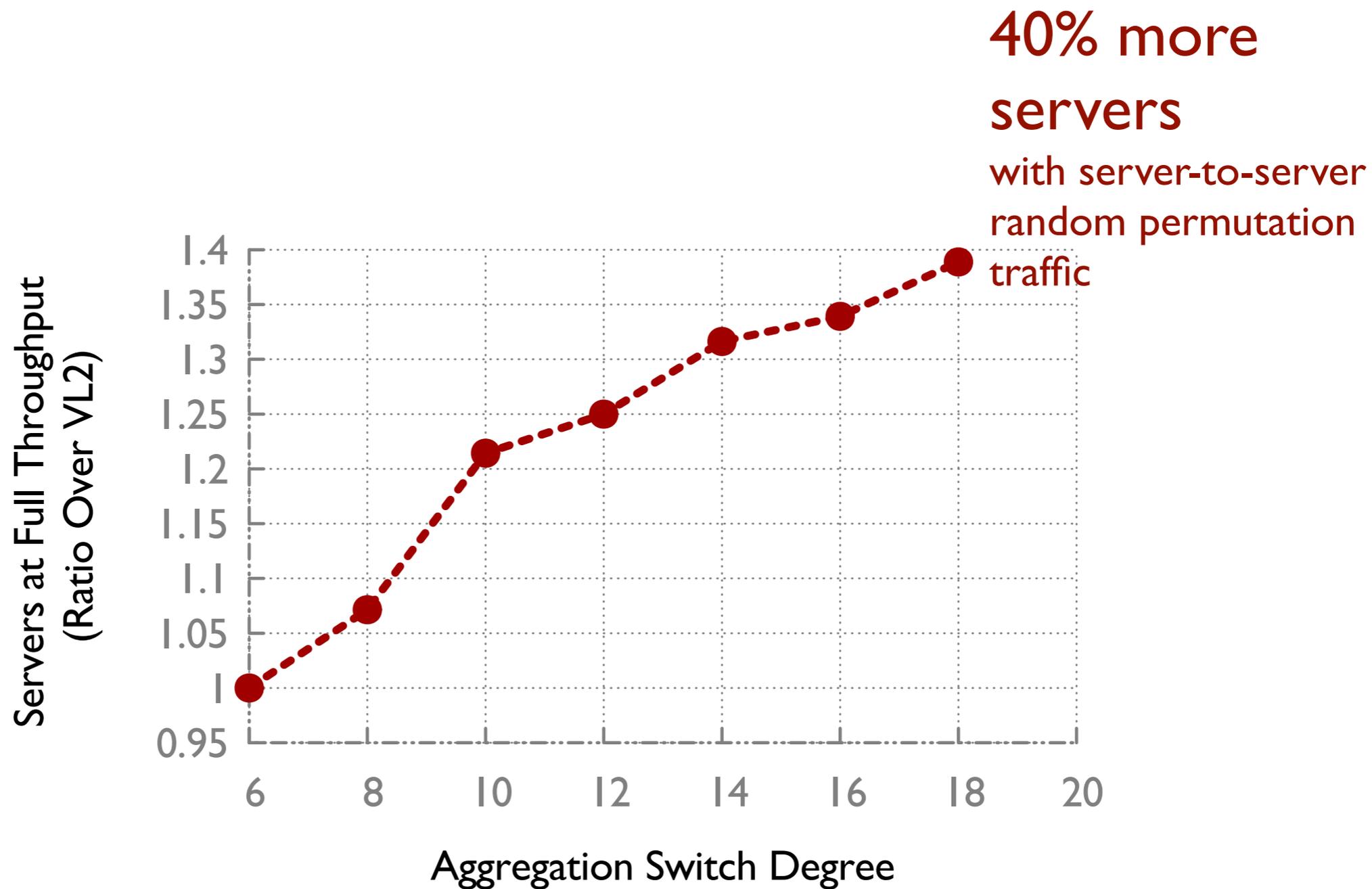
Rewiring VL2



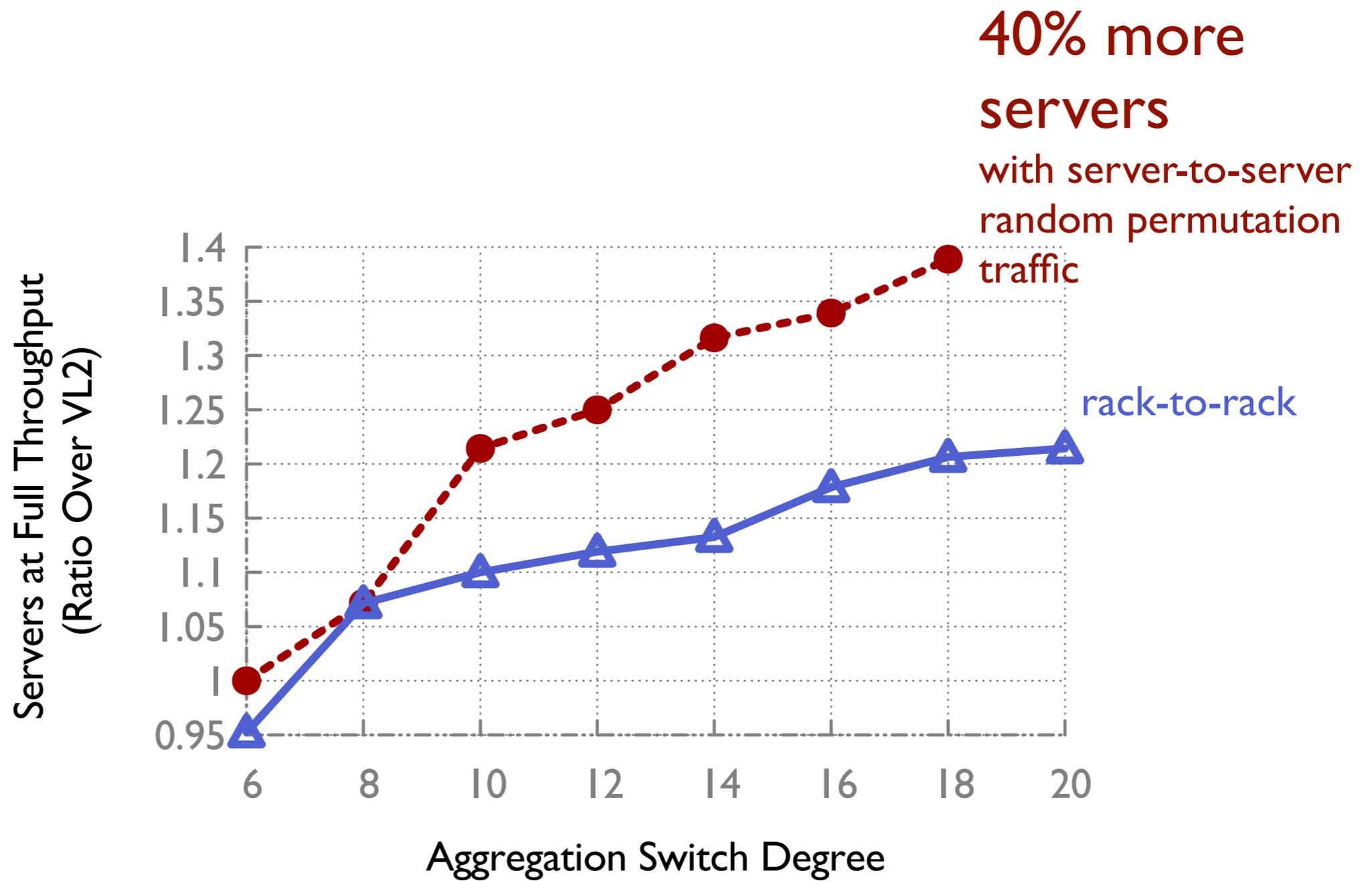
Rewiring VL2



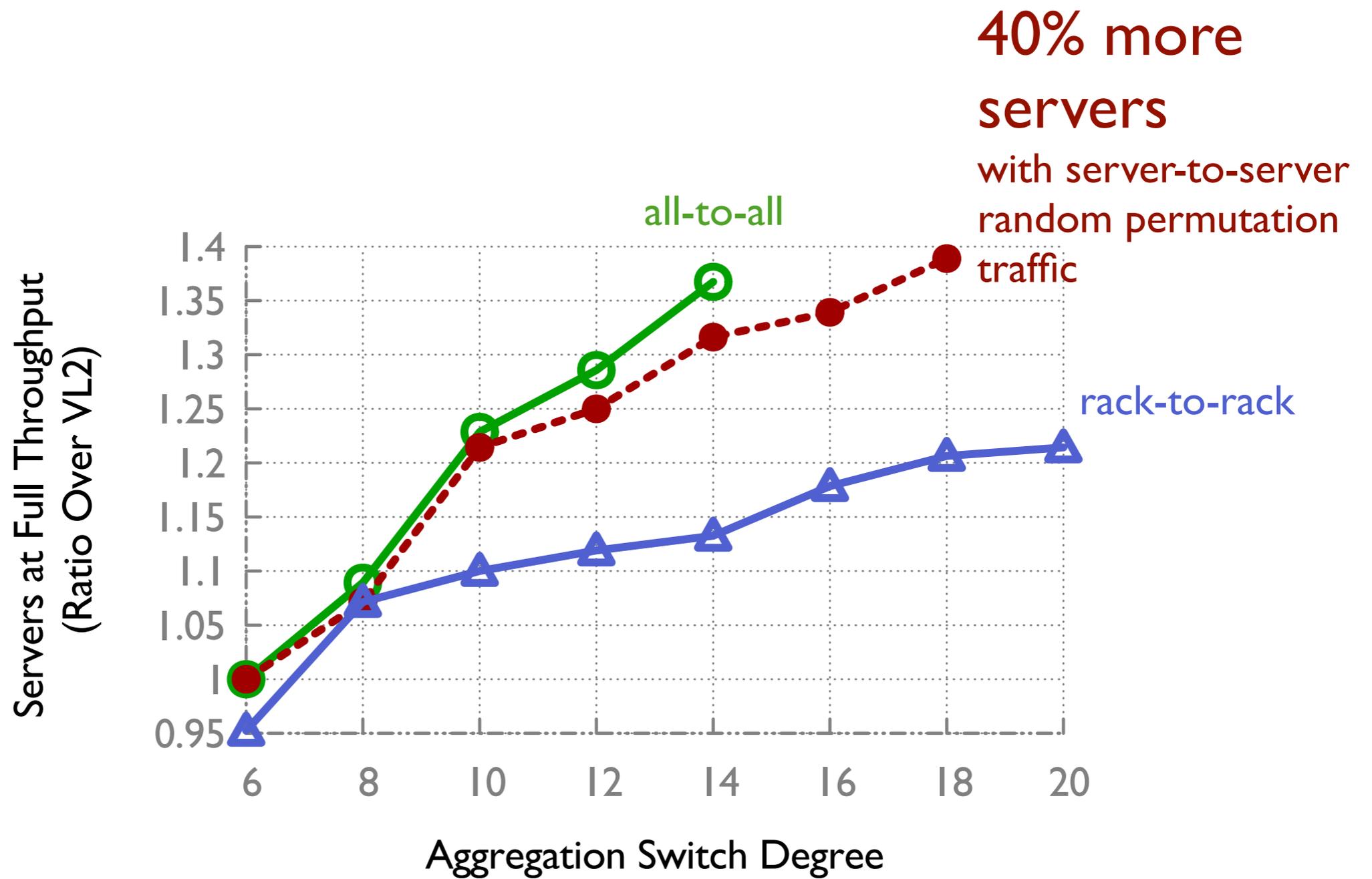
Rewiring VL2



Rewiring VL2



Rewiring VL2



How do we design
throughput optimal
network topologies?

<https://github.com/ankitsingla/topobench>