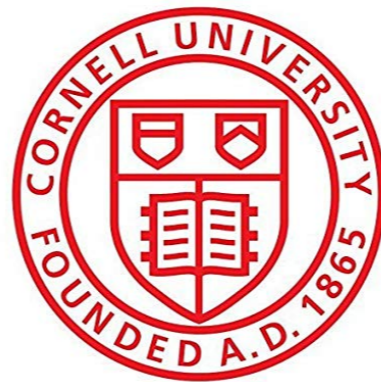


Building High-speed Datacenter Networks in the Post-Moore's Law Era

Vishal Shrivastav
Cornell University

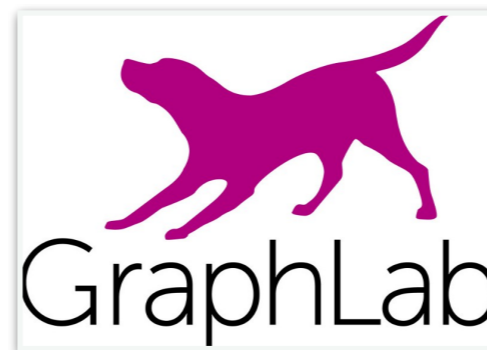
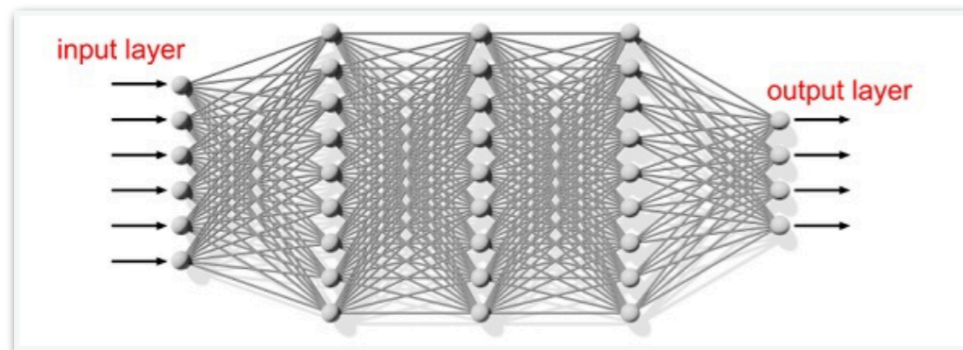


Datacenters : Workhorses of Modern Internet

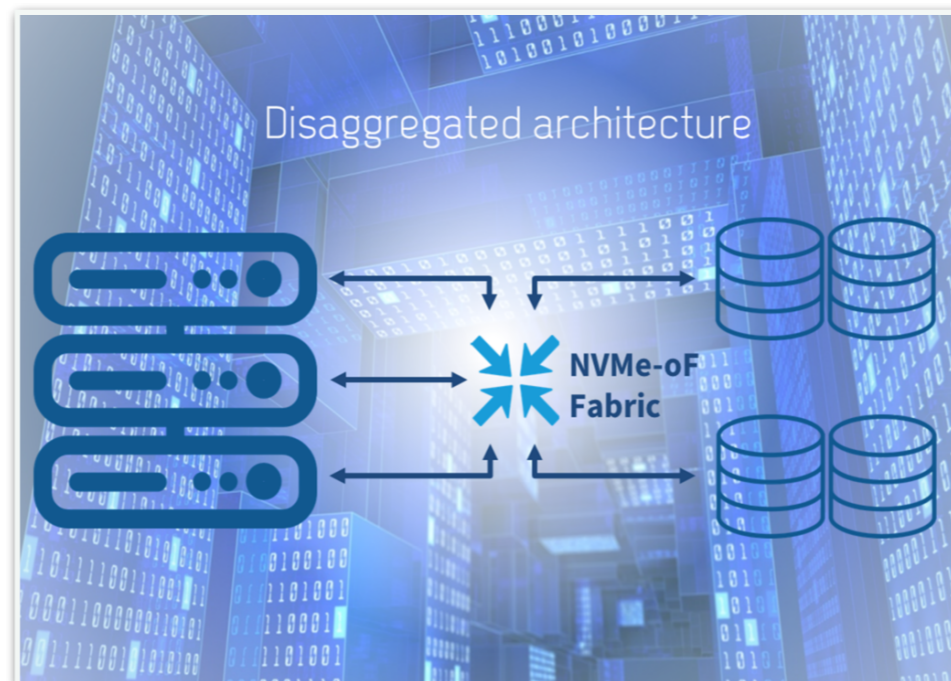


Network Intensive Applications

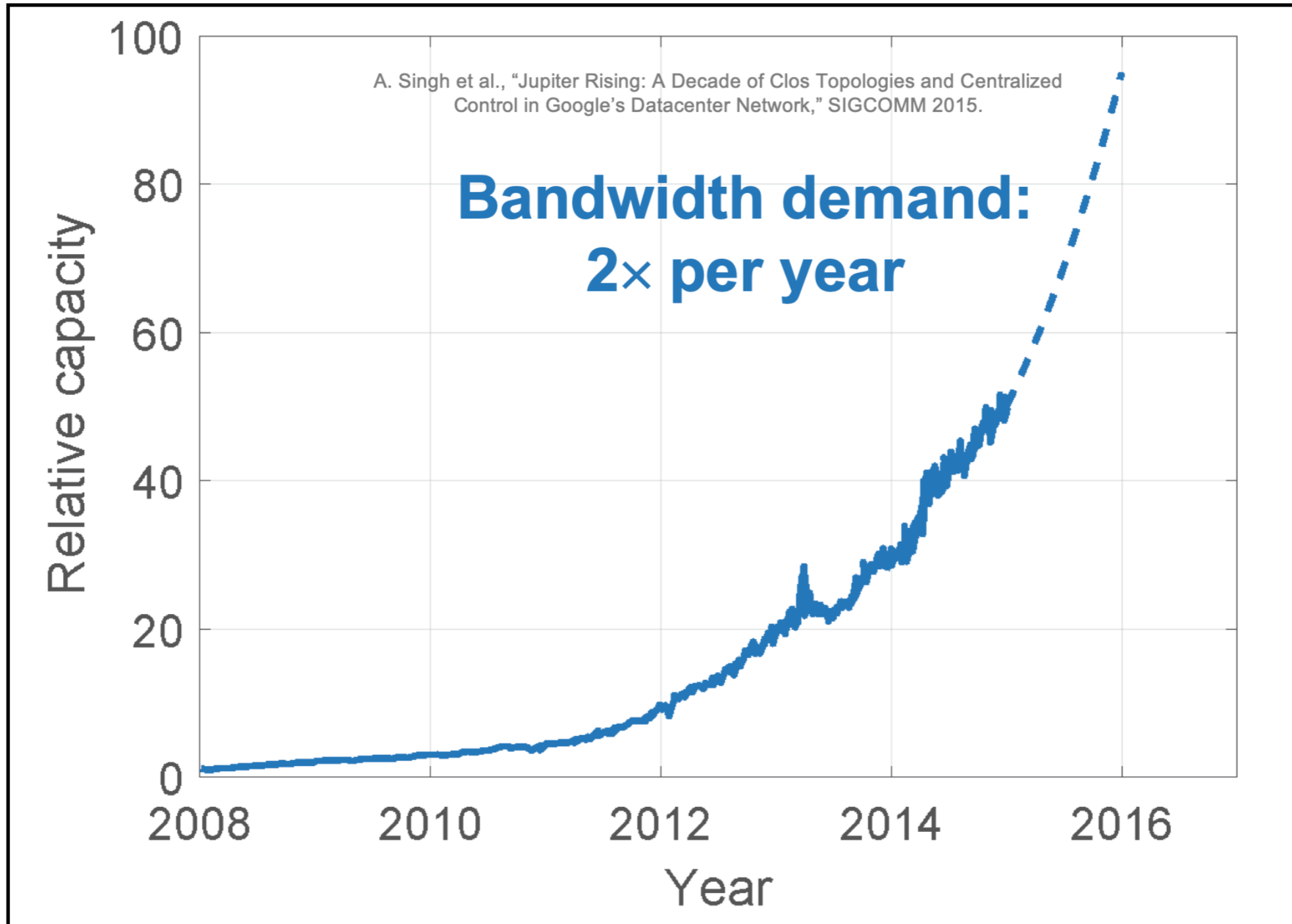
Applications getting more distributed



Resources getting disaggregated



Bandwidth Demand



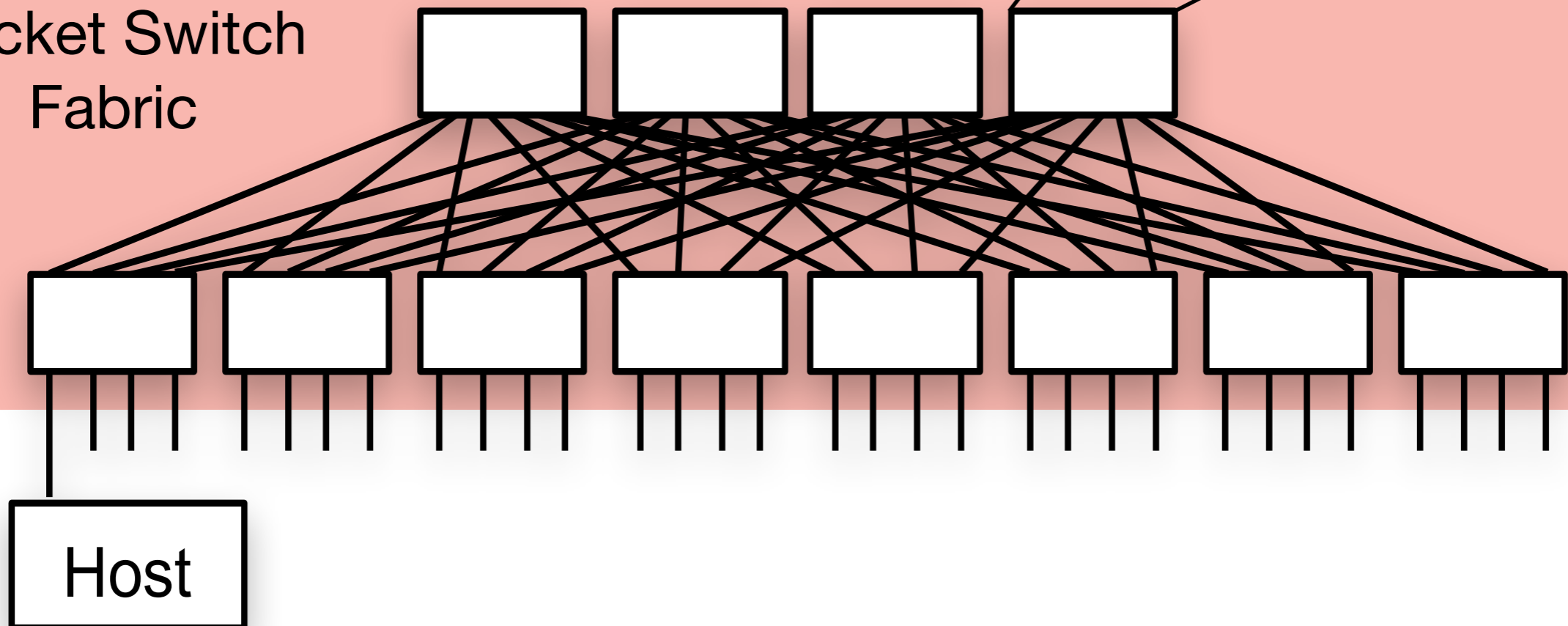
Today's Datacenter Networks

Switching Chip

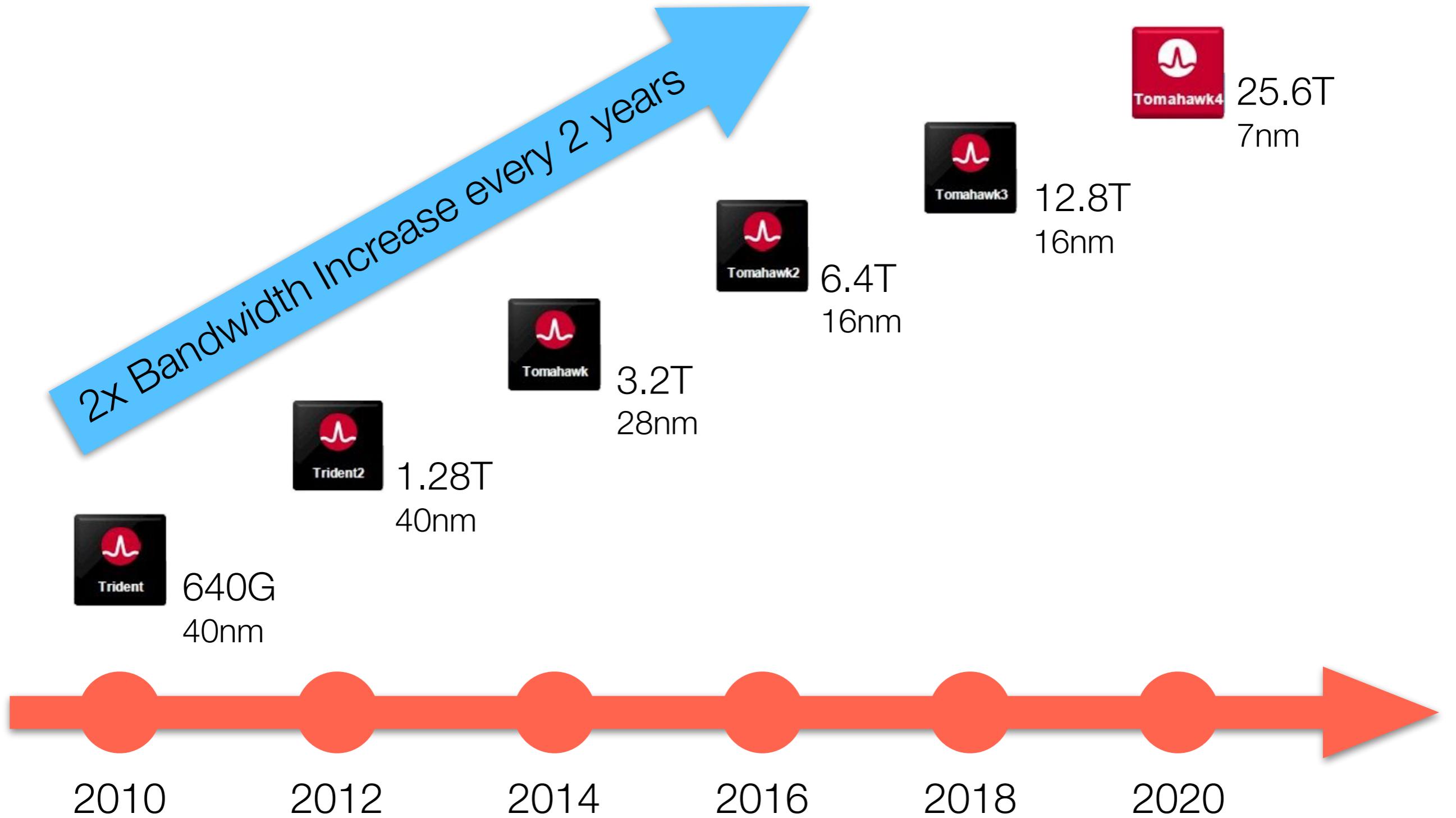
Folded Clos Topology



Packet Switch Fabric

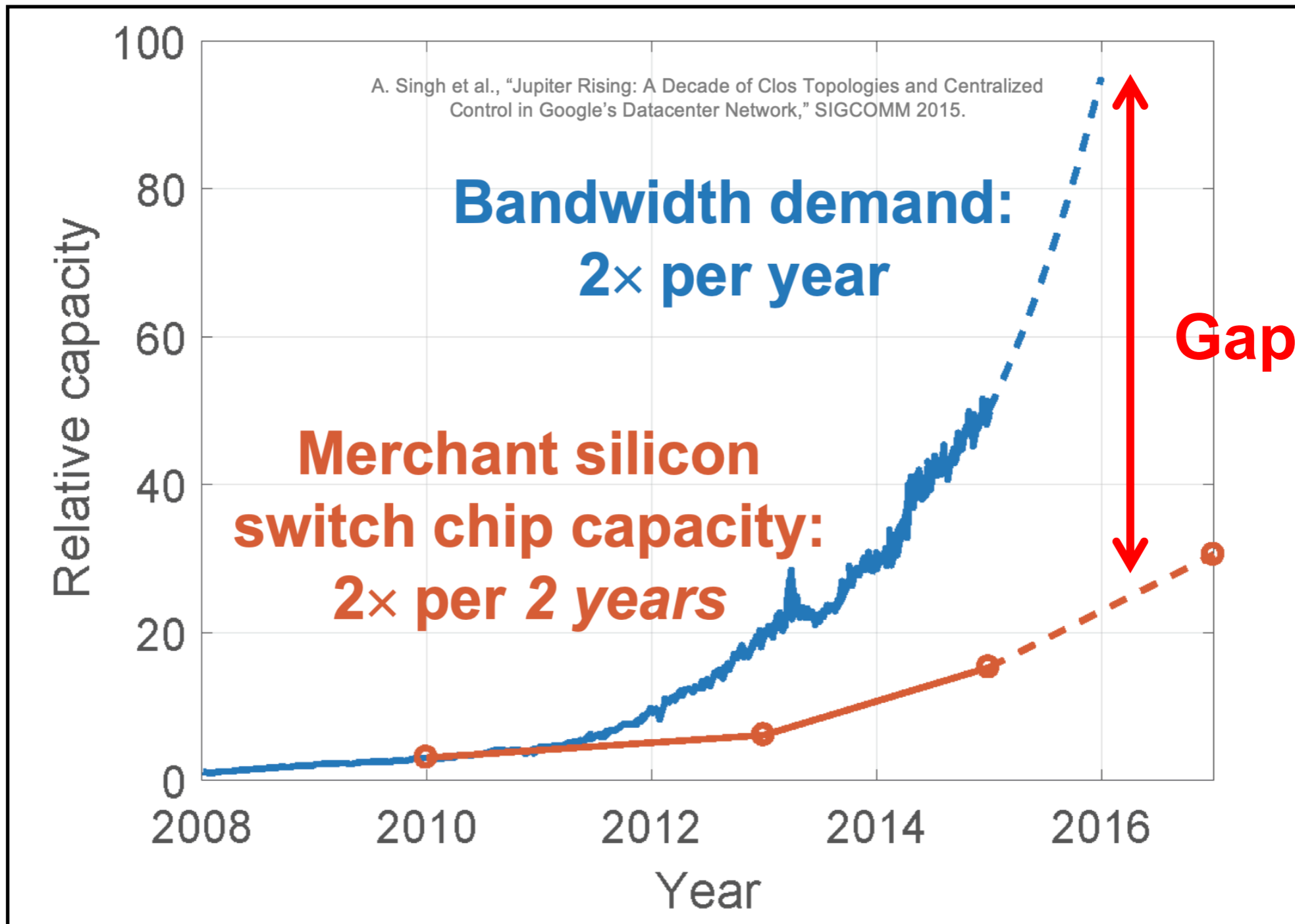


Switching Capacity



[Source: Broadcom]

Growing Gap



Scale-Out to the Rescue?



Fully Provisioned Clos Network with 8,192 hosts

Topology	# Hops	# Switch Chips	# Transceivers	# Wires
2 tiers	4	192	24K	12K
3 tiers	6	1,280	40K	20K

Goal

A **datacenter switching fabric** that provides **unlimited bandwidth scaling** at **low power, low cost, high performance**

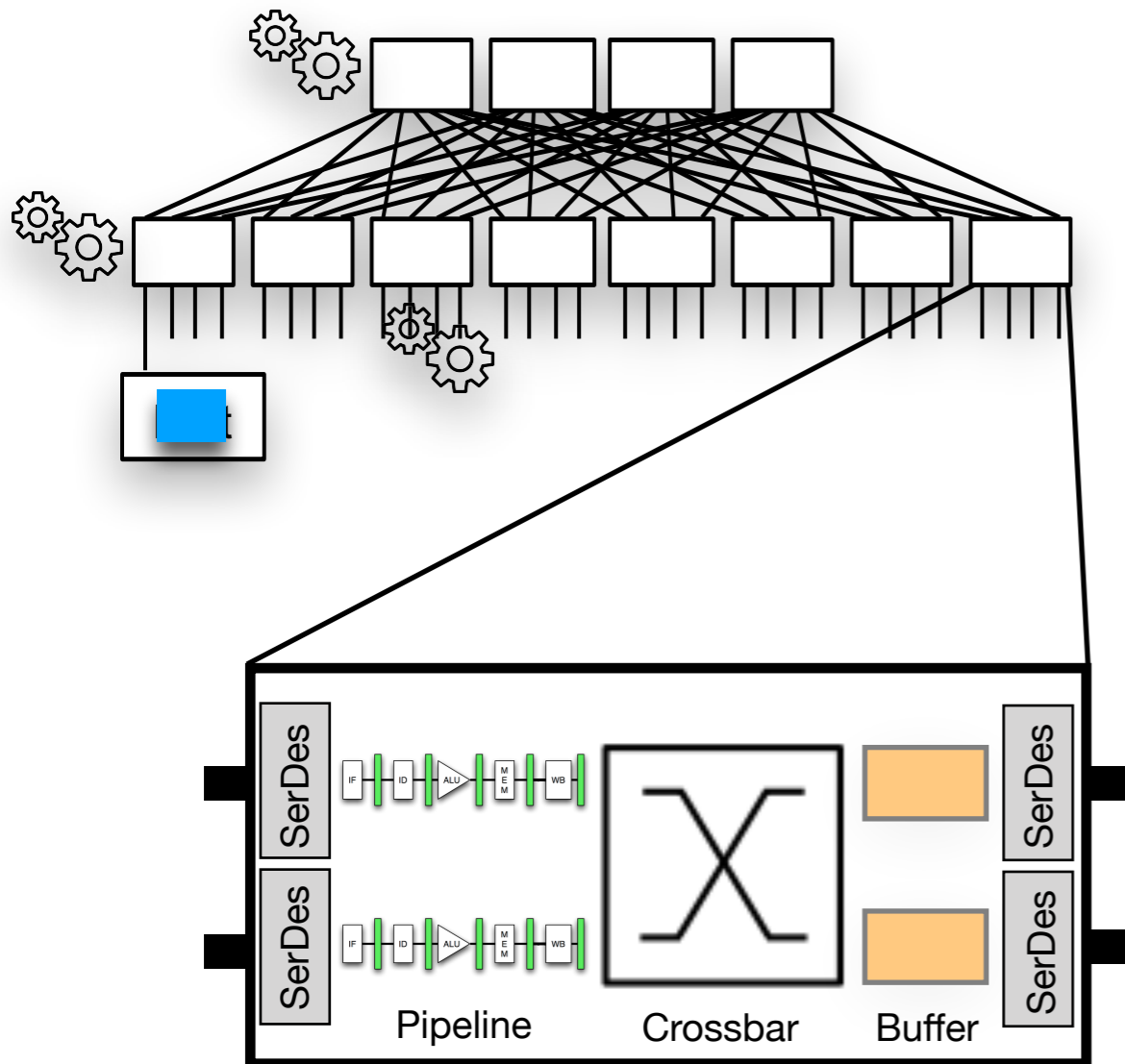


Circuit Switching?

Circuit Switching: Opportunities & Challenges

Packet Switching

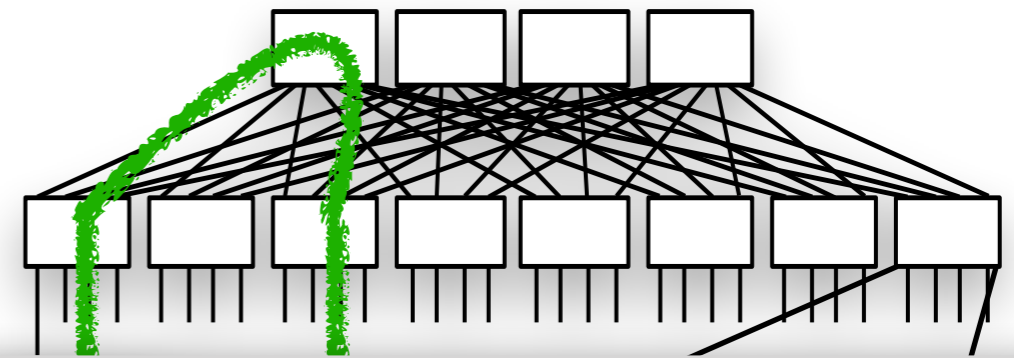
(Per packet, Per hop Decisions)



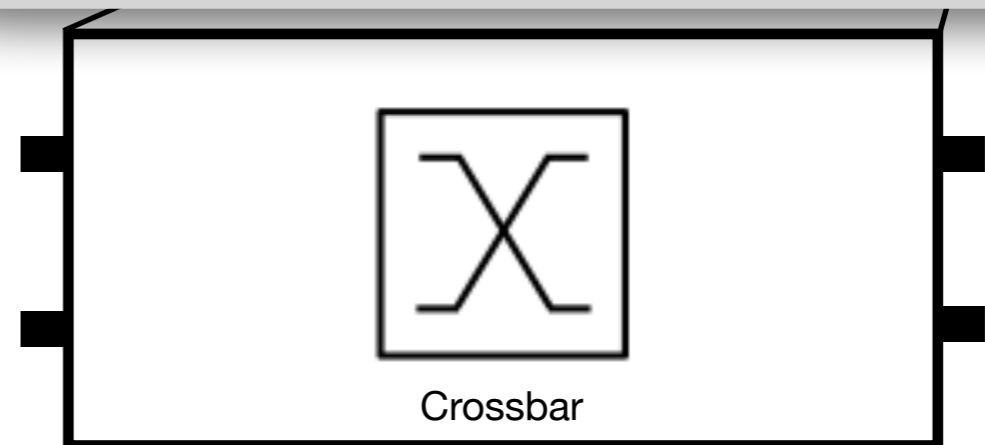
Packet Switch

VS.

Circuit Switching



Lower Power
Lower \$\$
Unlimited Bandwidth Scaling (w/ optics)



Circuit Switch

Circuit Switching: Opportunities & Challenges



Circuit Configuration

Data Plane

setting up physical circuits

milliseconds

nanoseconds

Control Plane

circuit scheduling

Traffic demand



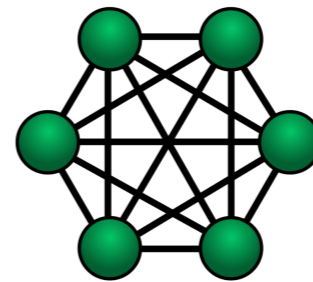
Centralized

Circuits

???

ns-scale circuit scheduling with high performance?

Shoal



“Shoal: A Network Architecture for Disaggregated Racks”

Vishal Shrivastav, Asaf Valadarsky, Hitesh Ballani, Paolo Costa, Ki Suh Lee,
Han Wang, Rachit Agarwal, Hakim Weatherspoon

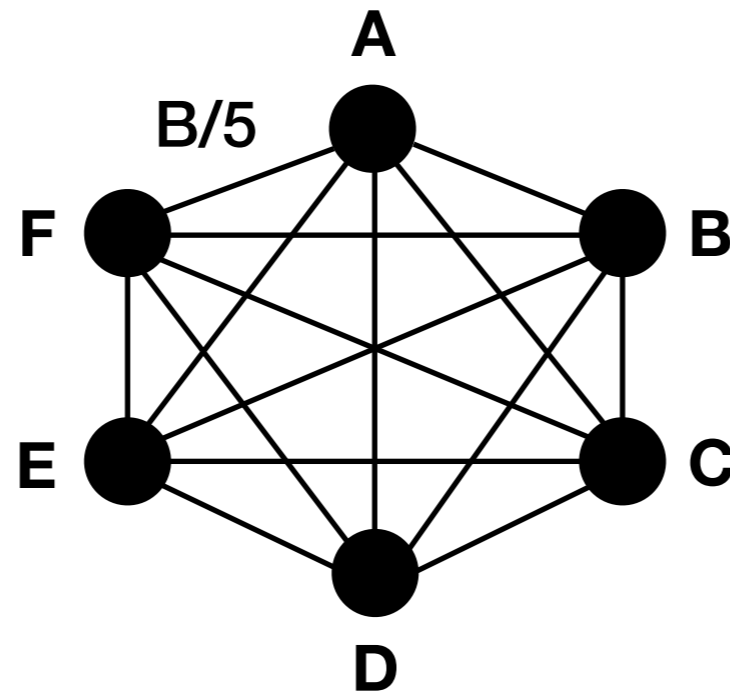
In USENIX NSDI 2019

Shoal is the first **network stack** and **switching fabric** design for a **fast circuit-switched** network

1. **Physical Layer** : Fast circuit scheduling mechanism
2. **Routing** : Bounded worst-case throughput
3. **Congestion Control** : Bounded worst-case queuing

Achieves comparable or better performance than several recent packet-switched network designs

Physical Layer



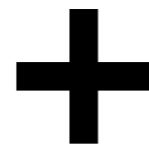
Full Mesh
Virtual Topology

N-1 timeslots
(epoch)

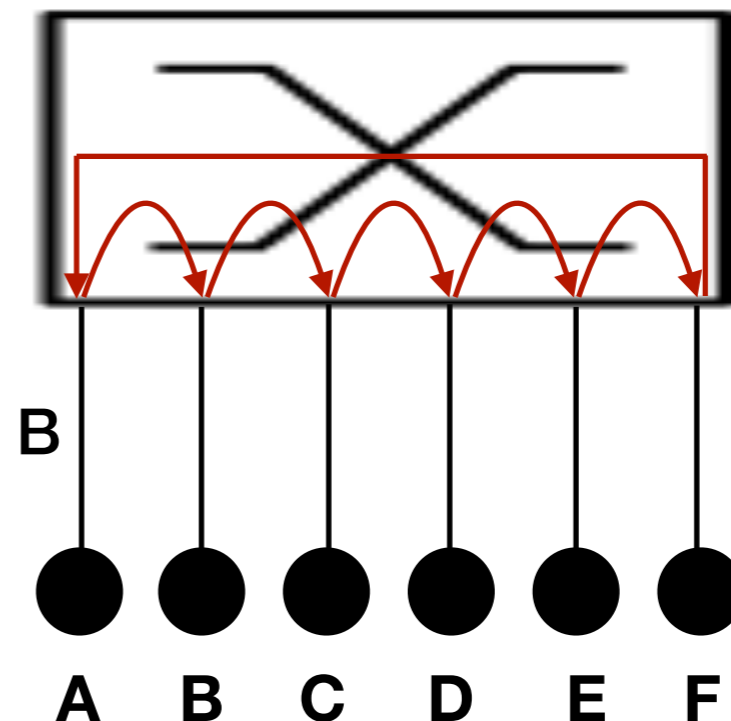
	1	2	3	4	5
A	B	C	D	E	F
B	C	D	E	F	A
C	D	E	F	A	B
D	E	F	A	B	C
E	F	A	B	C	D
F	A	B	C	D	E

End-hosts

static, pre-defined
schedule

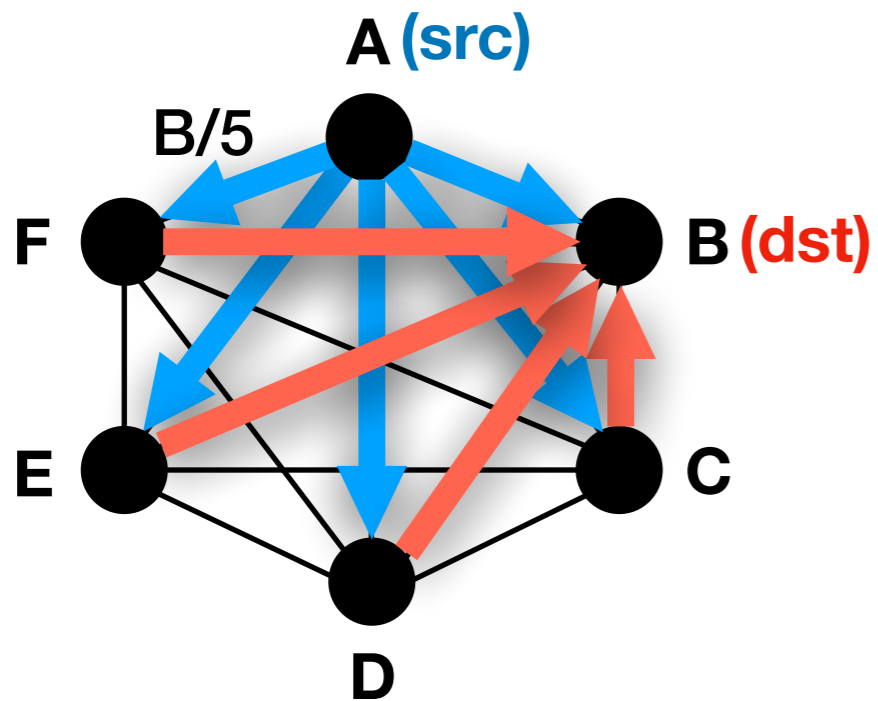


Synchronous
System

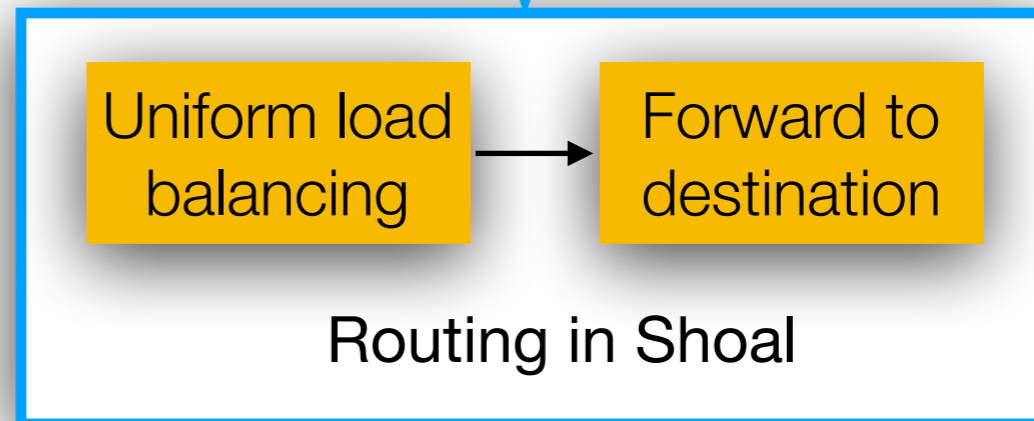


Physical
Topology

Routing



	Best-case throughput	Worst-case throughput	
Direct Routing	1	$1/(N-1)$	✗
Valiant LB	1	0.5	✓



Practical Challenges

Shoal is a synchronous system

Requires ns-precision network-wide synchronization

“Globally Synchronized Time via Datacenter Networks”

Vishal Shrivastav, Ki Suh Lee, Han Wang, Hakim Weatherspoon

In Journal of Transactions on Networking (ToN) 2019

*** conference version in ACM SIGCOMM 2016**

Circuit set-up changes at ns-scale

Requires fast, ns-precision scheduling at end-hosts

“Fast, Scalable, and Programmable Packet Scheduler in Hardware”

Vishal Shrivastav

In ACM SIGCOMM 2019

Conclusion

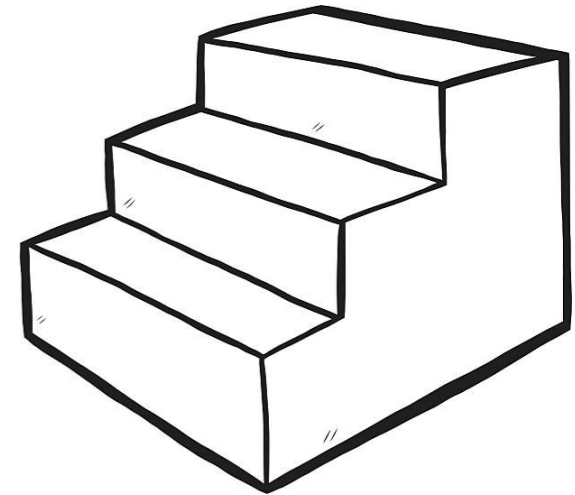
Unlimited
Bandwidth Scaling
at
Low Power, Low Cost, High Performance



**Fast, High Performance
Circuit Scheduling**



Shoal



Circuit Switching

Thank you!

Email: vishal@cs.cornell.edu

Webpage: <http://www.cs.cornell.edu/~vishal/>