

SIGCOMM Preview Session: Datacenter Networking (DCN)

Hakim Weatherspoon
Associate Professor, Cornell University

SIGCOMM
Florianópolis, Brazil
August 22, 2016

Cloud Computing

- The promise of the Cloud
 - A computer utility; a commodity
 - Catalyst for technology economy
 - Revolutionizing for health care, financial systems, scientific research, and society



Google Compute Engine



Windows® Azure™



Cloud Computing

- The promise of the Cloud
 - *ubiquitous, convenient, on-demand network access to a shared pool of configurable computing resources (e.g., networks, servers, storage, applications, and services) that can be rapidly provisioned and released with minimal management effort or service provider interaction.*

NIST Cloud Definition



Google Compute Engine



Windows Azure™



Cloud Computing

- The promise of the Cloud
 - ubiquitous, convenient, *on-demand network access* to a *shared pool* of configurable computing resources (e.g., networks, servers, storage, applications, and services) that can be *rapidly provisioned and released* with minimal management effort or service provider interaction.

NIST Cloud Definition



Google Compute Engine

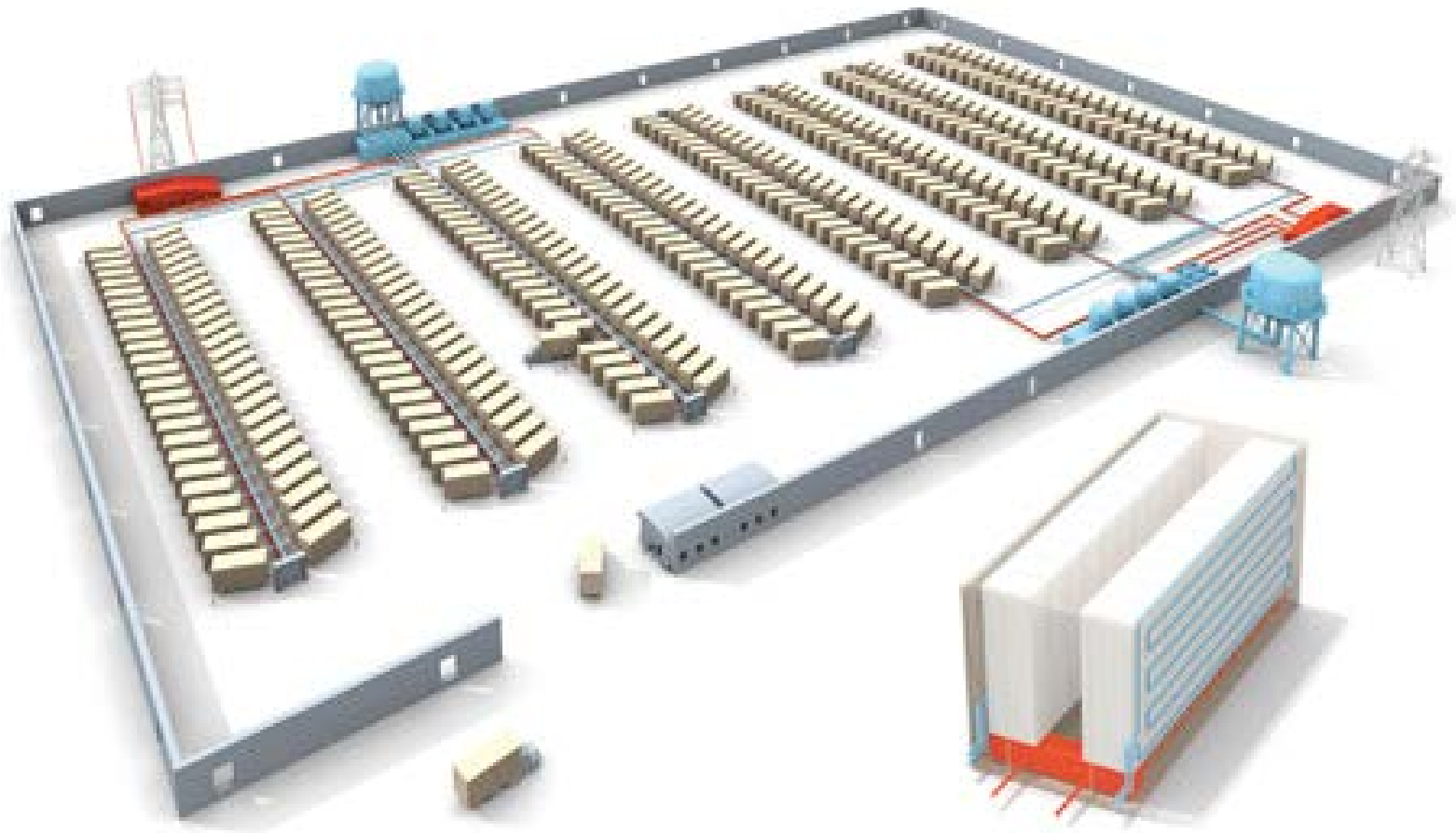


Windows Azure™



Cloud Computing needs Datacenters

- How big is the Cloud?
 - Exabytes: Delivery of petabytes of storage daily



Cloud Computing needs Datacenters

- **How big is the Cloud?**

- Most of the worlds data (and computation) hosted by few companies in datacenters



Cloud Computing needs Datacenters

- How big is the Cloud?

- Most of the worlds data (and computation) hosted by few companies in datacenters



Inside of a Datacenter

- 10s to 100s of thousands of servers
- Exabytes (1000s of petabytes) of storage
- Infrastructure-as-a-Service (IaaS)
 - Amazon EC2, Google Compute Engine, Microsoft Azure
- Single “application” spread across many thousands of servers (e.g. Amazon.com)
 - Application components such as caches, web servers, data bases, distributed file servers,...
 - Each component is “scaled” to meet the needs of millions (or billions) of users

Why Study DCN

- Scale
 - Google: 0 to 1B users in ~15 years
 - Facebook: 0 to 1B users in ~10 Years
 - *Must operate at the scale of $O(1M+)$ users*
- Cost:
 - To build: Google (\$3B/year), MSFT (\$15B/total)
 - To operate: 1--2% of global energy consumption*
 - *Must deliver apps using efficient HW/SW footprint*

* LBNL, 2013

What defines a datacenter network (DCN)

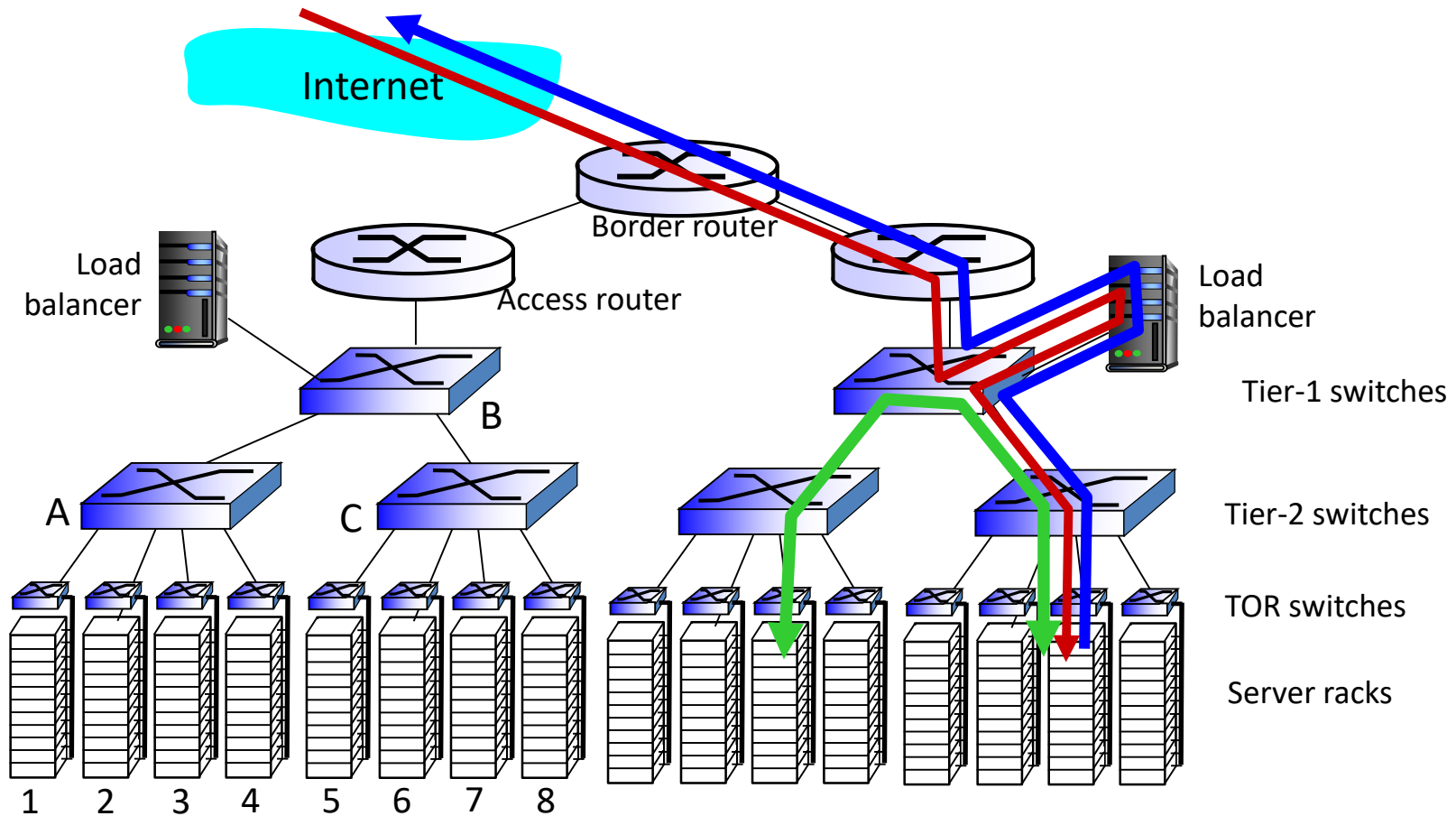
The Internet	Datacenter Network (DCN)
Many autonomous systems (ASes)	One administrative domain
Distributed Control/routing	Centralized Control and route selection
Single shortest path routing	Many paths from source to destination
Hard to measure	Easy to measure
Standardized Transport (TCP and UDP)	Many transports (DCTCP, qFabric,...)
Innovation requires consensus (IETF)	Single company can innovate
“Network of networks”	“Backplane of giant supercomputer”

DCN Research “cheat sheet”

- How would you design a network to support 1M endpoints?
- If you could...
 - Control all the endpoints and the network?
 - Violate layering, end-to-end principle?
 - Build custom hardware?
 - Assume common OS, Dataplane functions?
- Top-to-bottom rethinking of the network

DCN Topologies: Background

DCN Topologies: Tree-based



DCN Topologies: Tree-based

Issues with Traditional Data Center Topology

◎ *Oversubscription:*

- Ratio of the worst-case achievable aggregate bandwidth among the end hosts to the total bisection bandwidth of a particular communication topology
- Lower the total cost of the design
- Typical designs: factor of 2:5:1 (400 Mbps) to 8:1 (125 Mbps)

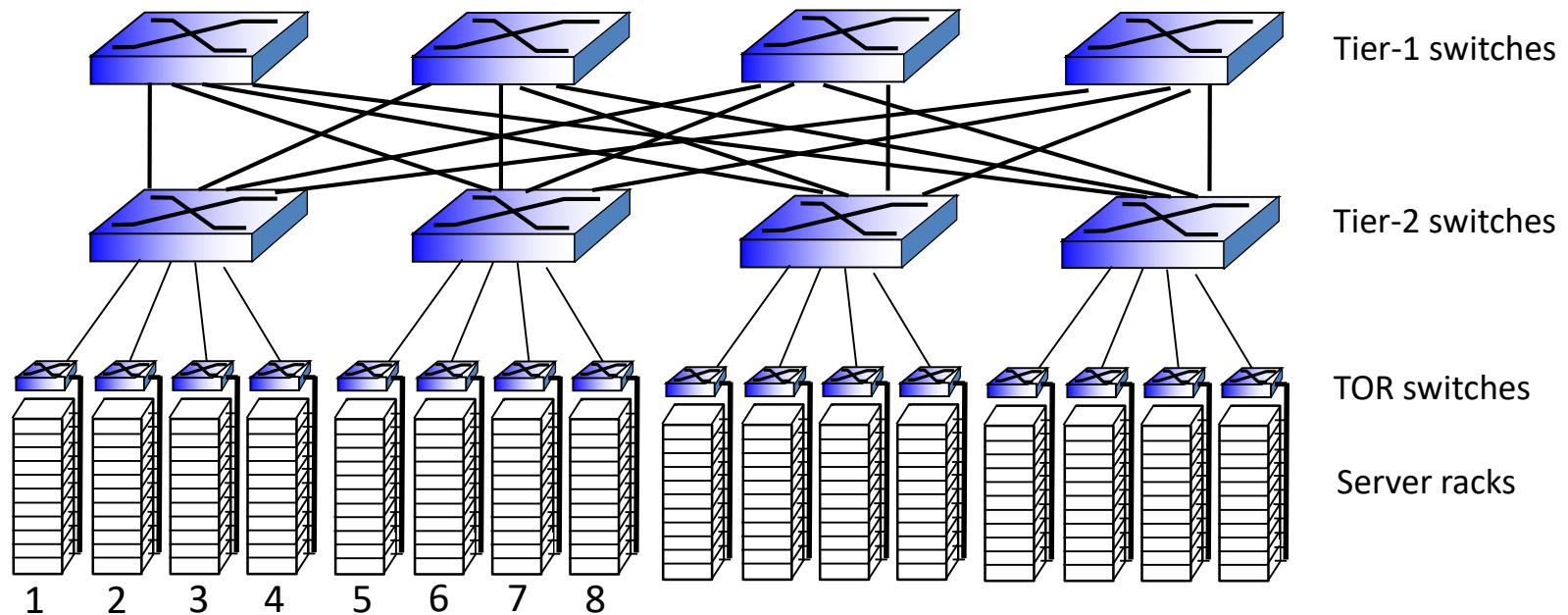
◎ *Cost:*

- Edge: \$7,000 for each 48-port GigE switch
- Aggregation and core: \$700,000 for 128-port 10GigE switches
- Cabling costs are not considered!

DCN Topologies: Folded-Clos multi-rooted tree

“FatTree” overcomes limitations

- ❖ rich interconnection among switches, racks:
 - increased throughput between racks (multiple routing paths possible)
 - increased reliability via redundancy





Paper Previews: Datacenters 1

- Session, Wednesday, 10:40am to 12:20pm
- Globally Synchronized Time via datacenter networks
 - Ki Suh Lee, Han Wang, Vishal Shrivastav, Hakim Weatherspoon
 - Datacenter Time Protocol, DTP, provides tight and bounded precision for synchronizing time throughout an entire datacenter. The paper shows that DTP uses the physical layer to bound precision by 25.6 nanoseconds for directly connected nodes, 153.6 nanoseconds for a datacenter with six hops. I.e. no two clocks differ by more than 153.6ns throughout the entire datacenter.



Paper Previews: Datacenters 1

- Session, Wednesday, 10:40am to 12:20pm
- Robotron: Top-down network management at Facebook
 - Yu-Wei Eric Sung, Xiaozheng Tie, Startsky H.Y. Wong, Hongyi Zeng
 - Network management via separating intent (expressed by Engineers) from implementation (translated by the system), which making the system more robust. Further, Robotron monitors operational state to ensure conformance to desired state.



Paper Previews: Datacenters 1

- Session, Wednesday, 10:40am to 12:20pm
- RDMA over Commodity Ethernet at Scale
 - Chuanxiong Guo, Haitao Wu, Zhong Deng, Gaurav Soni, Jianxi Ye, Jitu Padhye, Marina Lipshteyn
 - Challenges and approaches to using RDMA¹ (remote direct memory access) over commodity Ethernet (RoCEv2). Paper shows that RoCEv2 scale and issues can be addressed and that RDMA can replace TCP in the datacenter.

¹RDMA supports zero-copy networking by enabling the network adapter to transfer data directly to or from application memory, eliminating the need to copy data between application memory and the data buffers in the operating system. Such transfers require no work to be done by CPUs, caches, or context switches, and transfers continue in parallel with other system operations. When an application performs an RDMA Read or Write request, the application data is delivered directly to the network, reducing latency and enabling fast message transfer. However, this strategy presents several problems related to the fact that the target node is not notified of the completion of the request (1-sided communications). -- Wikipedia



Paper Previews: Datacenters 1

- Session, Wednesday, 10:40am to 12:20pm
- ProjecToR: Agile Reconfigurable Data Center Interconnect
 - Monia Ghobadi, Ratul Mahajan, Amar Phanishayee, Nikhil Devanur, Janardhan Kulkarni, Gireeja Ranade, Pierre-Alexandre Blanche, Houman Rastegarfar, Madeleine Glick, Daniel Kilper
 - Explores use of free-space optics¹ for building datacenter interconnects using digital micromirror devices (DMD) and mirror assembly combination as a transmitter and photodetector on top of the rack as a receiver.

¹Free-space optical communication (FSO) is an optical communication technology that uses light propagating in free space to wirelessly transmit data for telecommunications or computer networking. "Free space" means air, outer space, vacuum, or something similar. This contrasts with using solids such as optical fiber cable or an optical transmission line. The technology is useful where the physical connections are impractical due to high costs or other considerations.--
Wikipedia

Final thoughts

- DCN Is an exciting, fun research area
- While many papers are from Microsoft, Google, Facebook, ...
 - YOU have the ability to have enormous impact
 - Many Projects are open--source
 - E.g., <http://sonic.cs.cornell.edu>
- Rethink the entire network stack!
 - Hardware, software, protocols, OS, NIC, ...