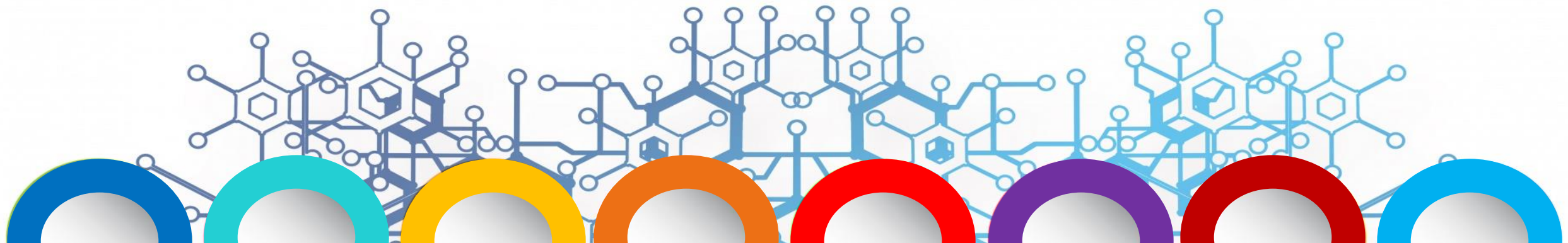
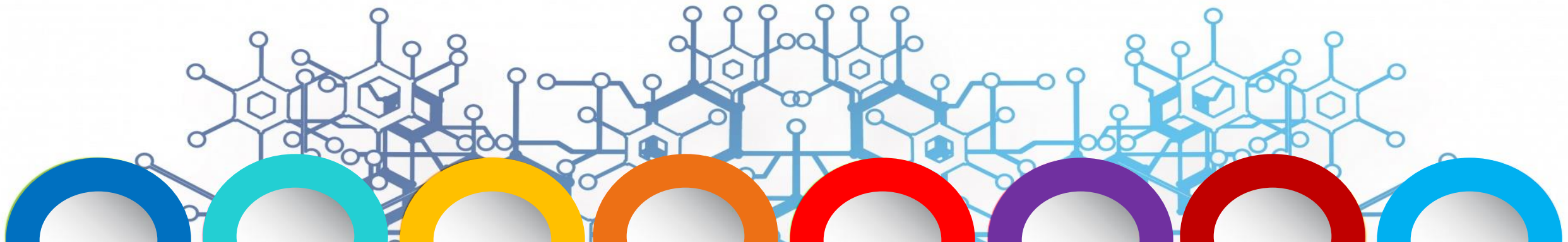


Fundamentals of Distributed Systems & The Cloud:

Network . Compute . Store



Network



Data Center Networks (DCN)

- Tens to hundreds of thousands of hosts, often closely coupled, in close proximity
 - E-business (Amazon)
 - Content-servers (YouTube, Akamai, Apple, Microsoft)
 - Search engines, data mining (Google)



Google Douglas County, Georgia data center



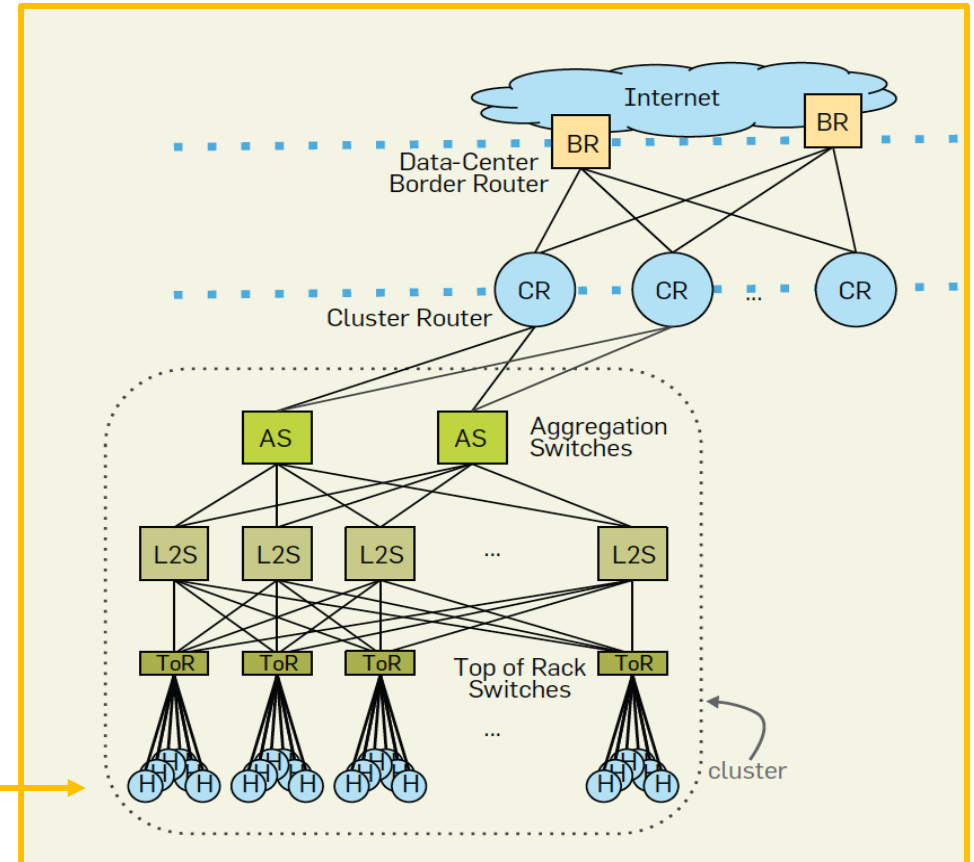
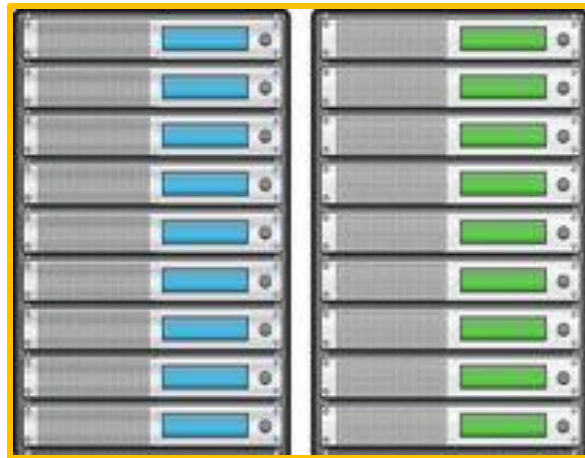
Microsoft, Chicago data center



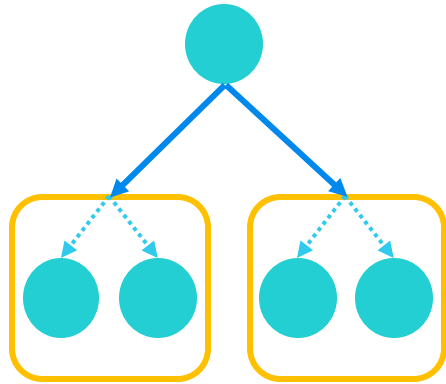
Facebook, Mexico data center

Data Center Hosts

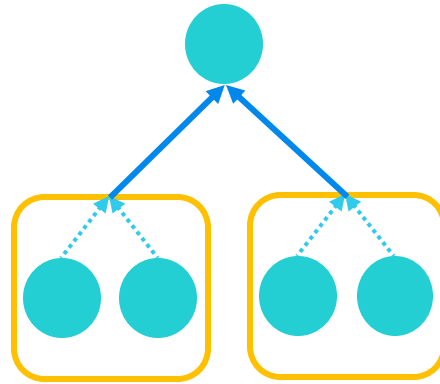
- Commodity Hosts: Blades (Pizza Box)
- Each rack ~20-40 blades
- Top Of the Rack (TOR) switch for each rack
- Each host has a connected interface to the TOR
- Each host has its own DCN IP address



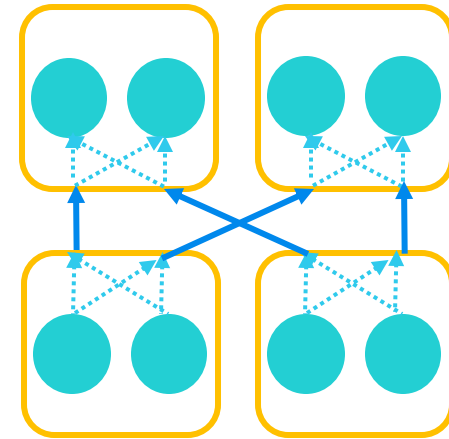
Review: Communication Patterns



Broadcast

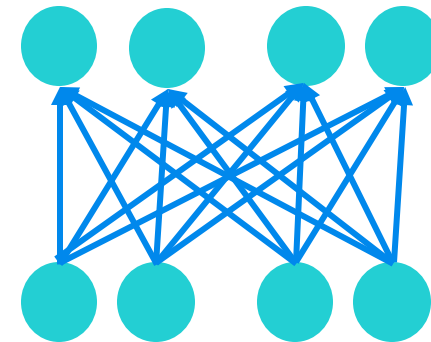
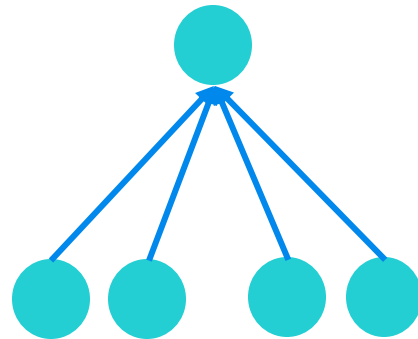
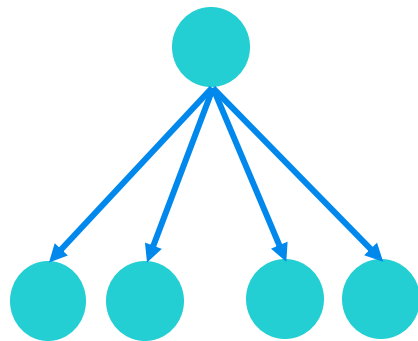


Aggregation



VM-based

Shuffle



Function-based

Communication Patterns

Some Suggested Solutions

- Provide cloud functions with a larger number of cores similar to VM instances
Multiple tasks can combine and share data among them before networking
- Allow developer to explicitly place cloud functions on same VM instance
Offer distributed communication primitives for allocating cloud functions to the same VM instance
- Let applications provide a computation graph
Enables the cloud provider to co-locate the cloud functions to minimize communication overhead

DCN Agility: Any Service, Any Server

- Turn the servers into a single large fungible pool
 - Dynamically expand and contract service footprint as needed
- Benefits
 - Increase service developer productivity
 - Achieve high performance and reliability
 - Lower cost

DCN: Achieving Agility

- Workload management
 - Rapidly installing a service's code on a server
 - Virtual machines, disk images, containers
- Storage Management
 - Server to access persistent data
 - Distributed filesystems (e.g., HDFS, BLOB stores)
- Network
 - Communication among servers, regardless of location in the data center

HDFS: Hadoop Distributed File System

BLOB: Binary Large Object

Reference: Presentation for VL2 Paper. web.mit.edu/6.829/ (2018 offering)

DCN: Routing & Switching

- **Ethernet switching (layer 2)**
 - ✓ Fixed IP addresses and auto-configuration (plug and play)
 - ✓ Seamless mobility, migration, and failover
 - ✗ Broadcast limits scale (ARP)
 - ✗ Spanning Tree Protocol
- **IP routing (layer 3)**
 - ✓ Scalability through hierarchical addressing
 - ✓ Multipath routing through equal-cost multipath
 - ✗ More complex configuration
 - ✗ Can not migrate without changing IP address

Data Center Network

- Load balancer: Application-layer routing

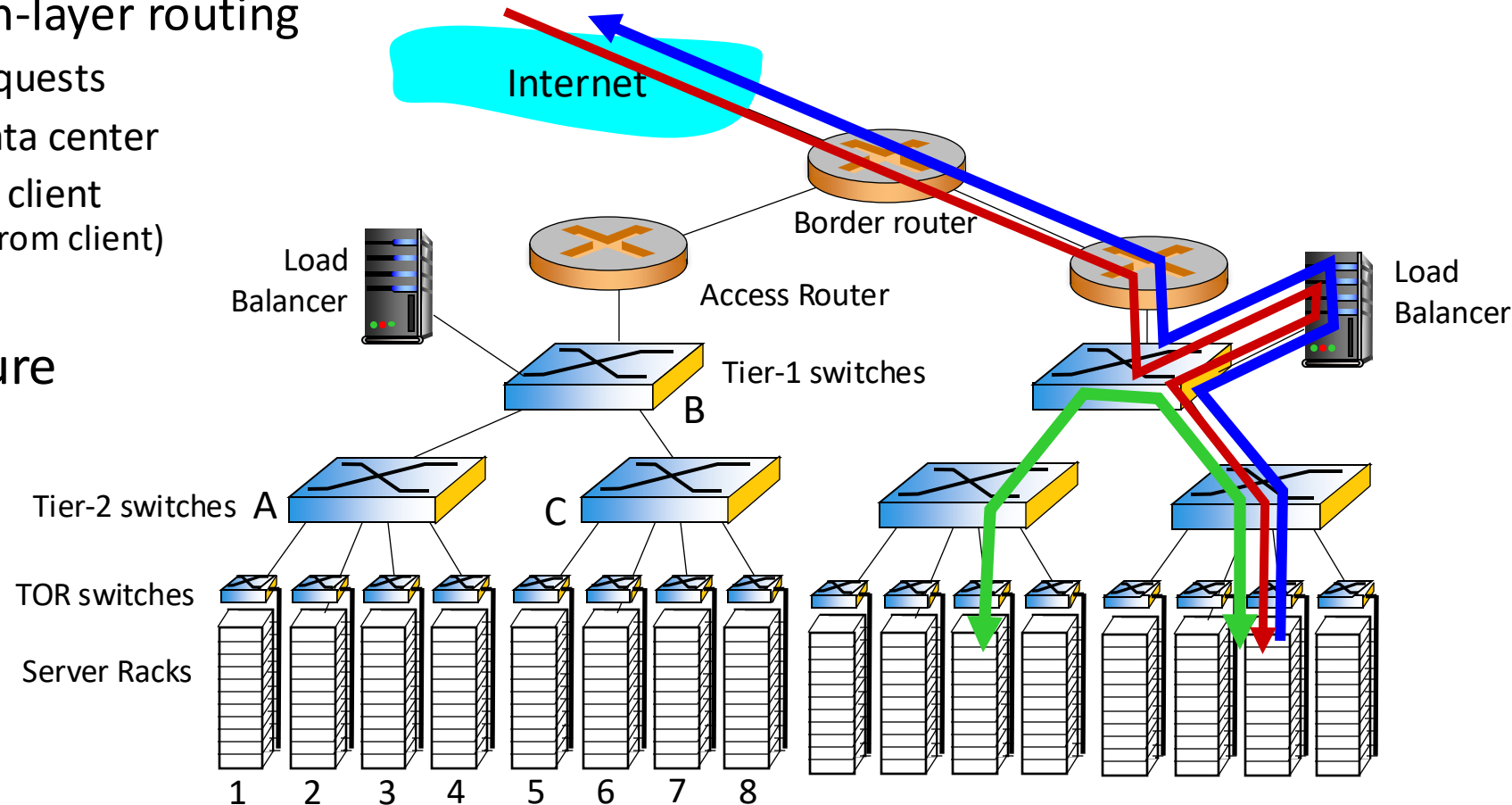
- Receives external client requests
- Directs workload within data center
- Returns results to external client
(Hiding data center internals from client)

- Conventional Architecture

- **Good:** Scalable

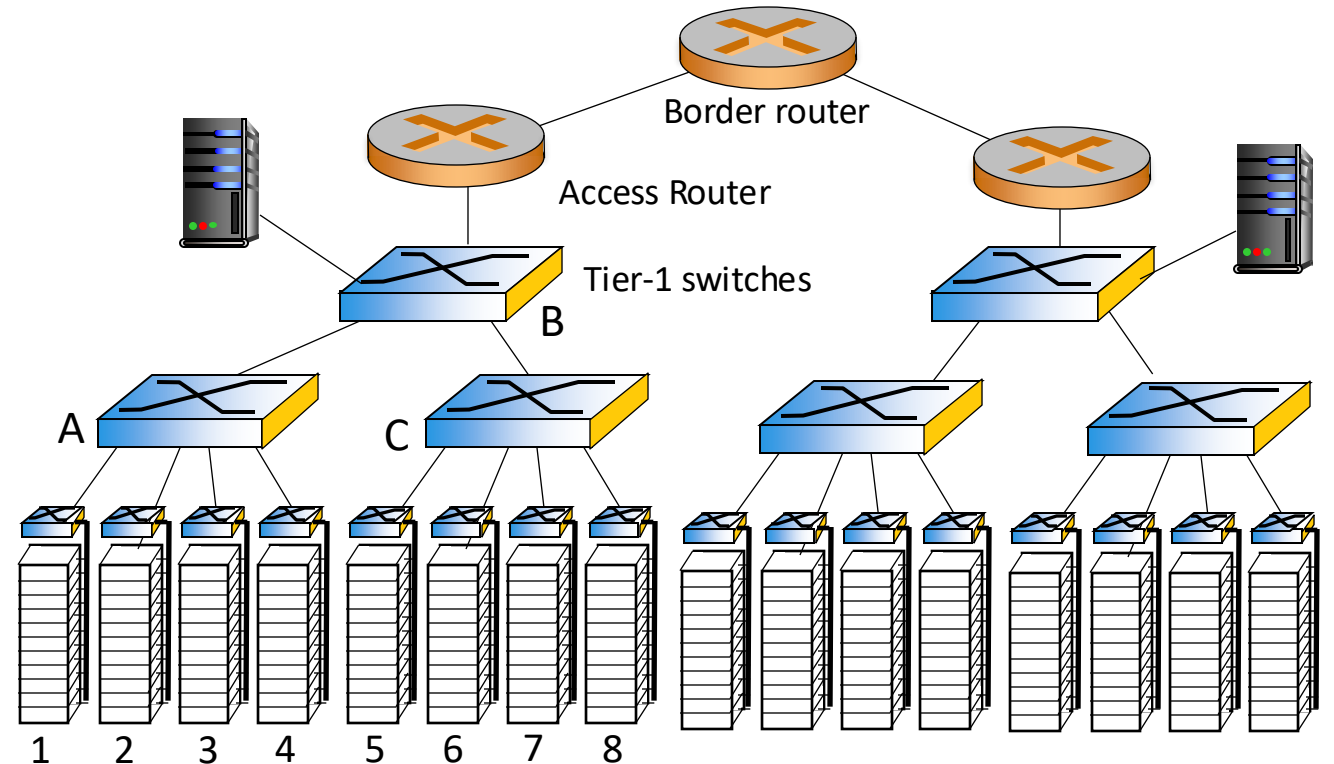
- **Not So Good**

Limited host-to-host capacity
No Service Agility



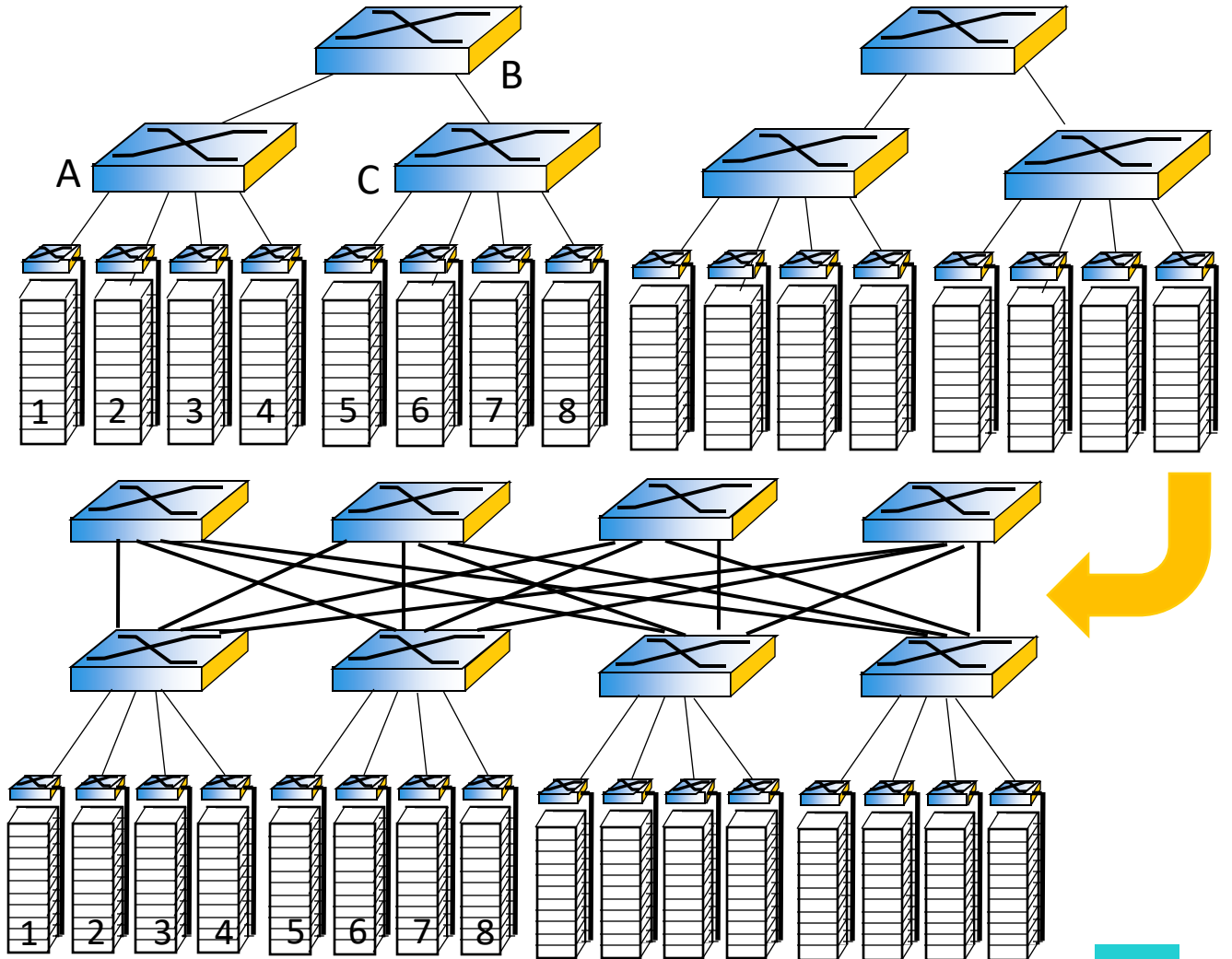
Example

- Host-TOR: 1Gbps
- Between switches: 10Gbps
- 40 Flows
 - 10 Hosts in Rack 1 - 10 Hosts in Rack 5
 - 10 Hosts in Rack 2 - 10 hosts in Rack 6
 - 10 Hosts in Rack 3 - 10 hosts in Rack 7
 - 10 Hosts in Rack 4 - 10 hosts in Rack 8
- Performance?



Example

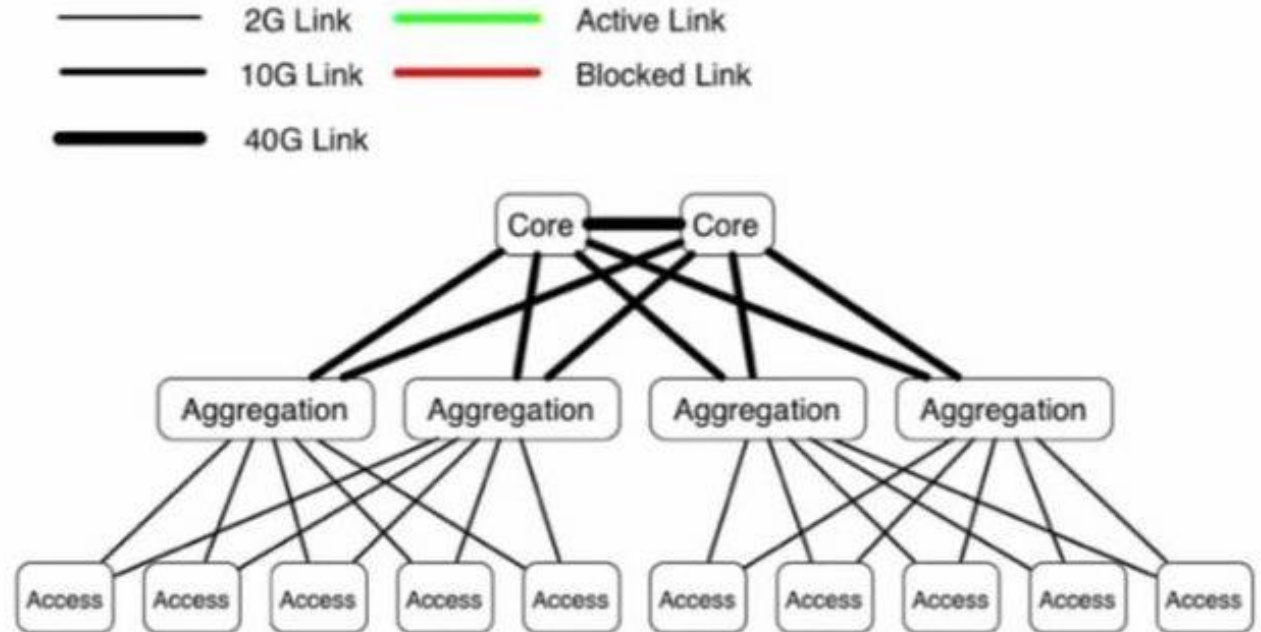
- Host-TOR: 1Gbps
- Between switches: 10Gbps
- 40 Flows
 - 10 Hosts in Rack 1 - 10 Hosts in Rack 5
 - 10 Hosts in Rack 2 - 10 hosts in Rack 6
 - 10 Hosts in Rack 3 - 10 hosts in Rack 7
 - 10 Hosts in Rack 4 - 10 hosts in Rack 8
- Performance?



DCN Architectures: Fat Tree

- Every layer above a layer has to support sum of the capacity of the lower layer

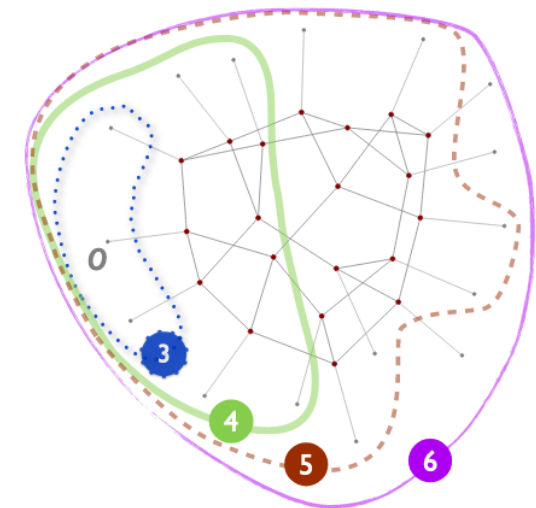
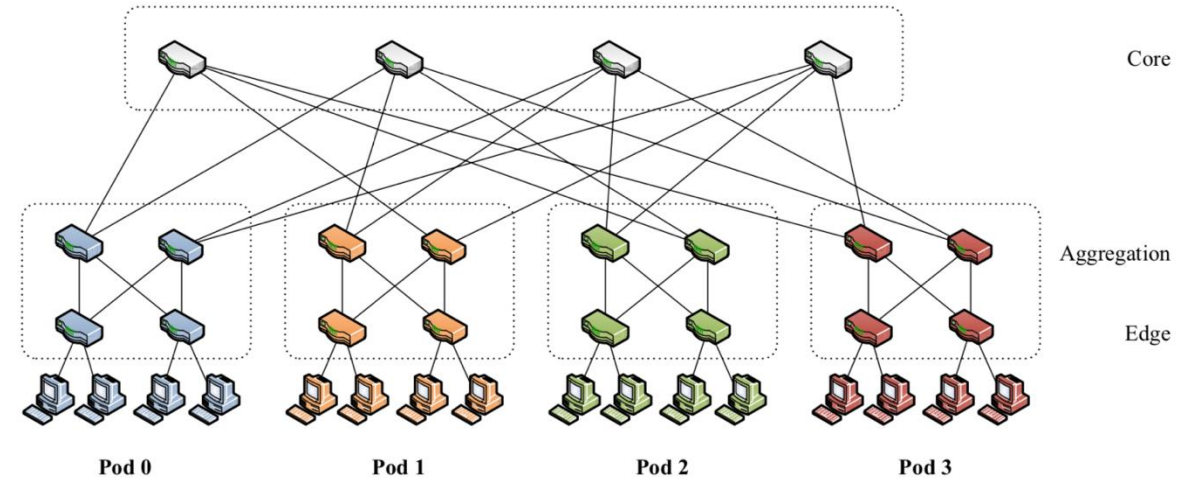
- Scale up!
 - Expensive for higher layers



- Why Fat-Tree?
 - Identical bandwidth at bisections
 - Same aggregated bandwidth at each layer
 - All devices **can** transmit at line speed if packets are distributed properly along available paths

Other DCN Architecture Proposals

- **MDC (Modular Data Center) Network**
 - Container-internal network
 - Core network to connect containers
- **BCube (Container-based)**
 - A standard 12-meter shipping container
 - A few thousand hosts
 - Graceful performance degradation over time
 - Replaceable
- **PortLand: Plug and play large scale data center networks**
- **JellyFish**
 - Network interconnect
 - Degree-bounded random graph topology among TOR switches
 - Different degrees of oversubscription



Reference: PortLand: A Scalable Fault-Tolerant Layer 2 Data Center Network Fabric (SIGCOMM 2009)

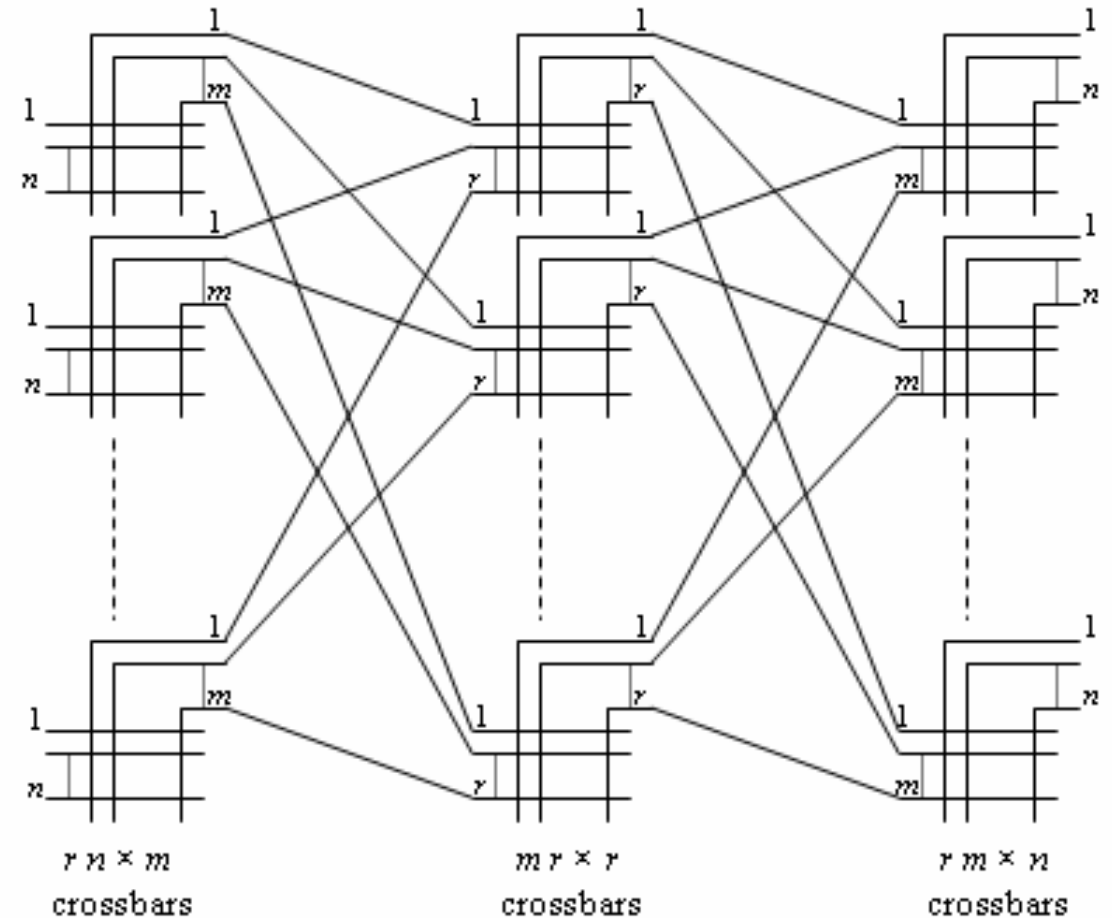
Reference: Jellyfish: Networking Data Centers Randomly (NSDI 2012)

Reference: Computer Networking : A Top-Down Approach. James F. Kurose, Keith W. Ross, 7th Edition, Pearson, 2017

Reference & Figure: BCube: A High Performance, Server-centric Network Architecture for Modular Data Centers (SIGCOMM 2009)

Clos Topology

- Historical: Telephone Circuit Switching
- Parameters
 - Stages
 - Inputs
 - Outputs
- Non-blocking properties



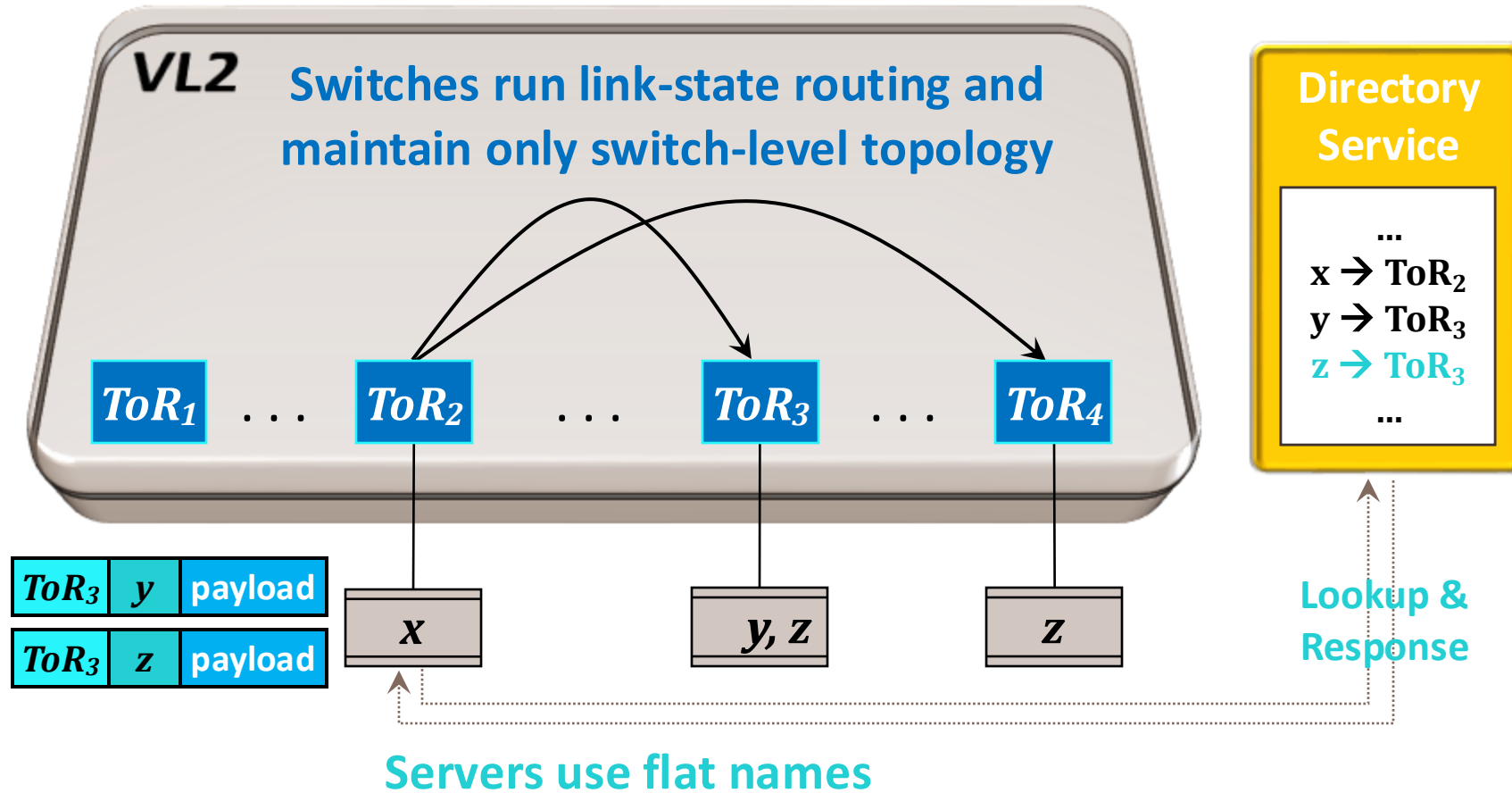
VL2: Solutions

Objective	Approach	Solution
1. Layer-2 semantics	Employ flat addressing	Name-location separation & resolution service
2. Uniform high capacity between servers	Guarantee bandwidth for hose-model traffic	Flow-based random traffic indirection (Valiant LB)
3. Performance Isolation	Enforce hose model using existing mechanisms only	TCP

TCP: Transmission Control Protocol, LB: Load Balancing

Reference: web.mit.edu/6.829/www/currentsemester/materials/datacenter-networking.pptx

VL2: Addressing & Routing

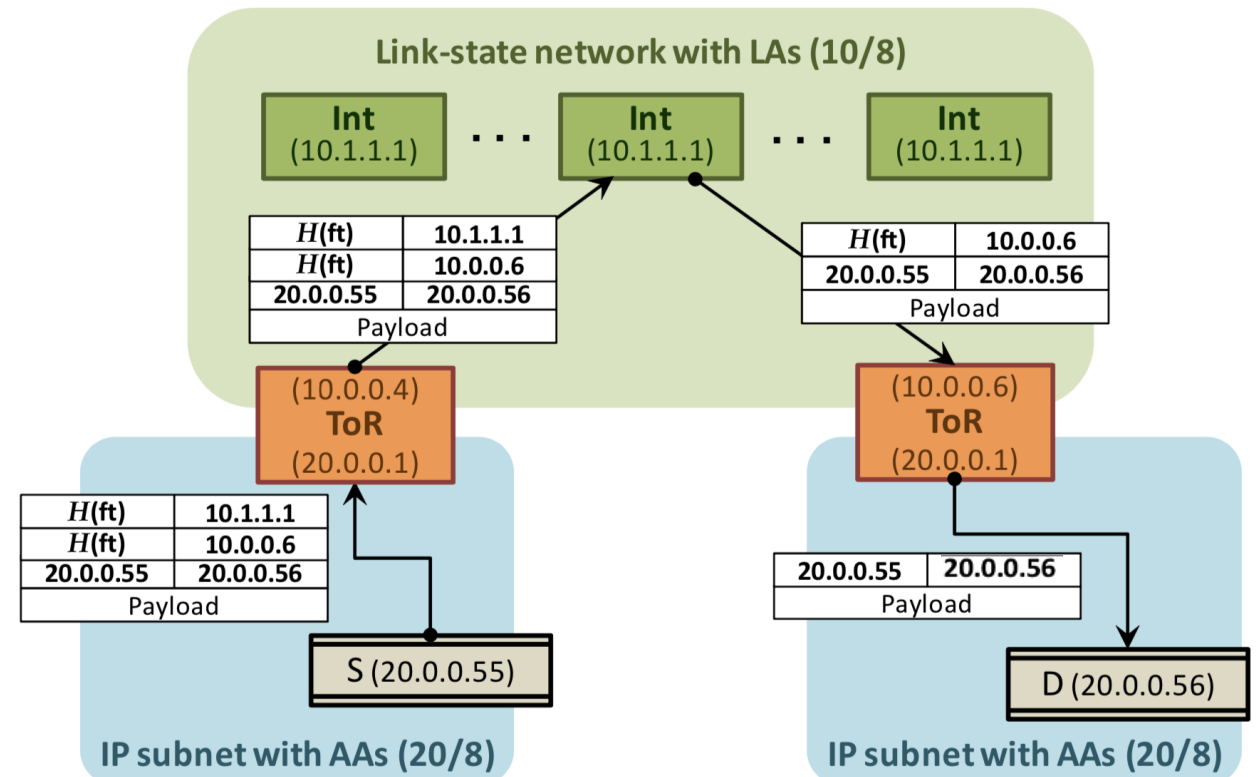


Reference: VL2: A Scalable and Flexible Data Center Network (SIGCOMM 2009)

Figure: web.mit.edu/6.829/www/currentsemester/materials/datacenter-networking.pptx

VL2: Summary

- Network built from low-cost switch ASICs arranged into a **Clos topology**
- **Valiant Load Balancing (VLB)**: Spread traffic uniformly across network paths without central coordination or traffic engineering (sending server picks random path for each flow)
- **Flat addressing**: Allow service instances to be placed anywhere in the network
- **End-system based** address resolution to scale to large server pools without introducing complexity to the network control plane

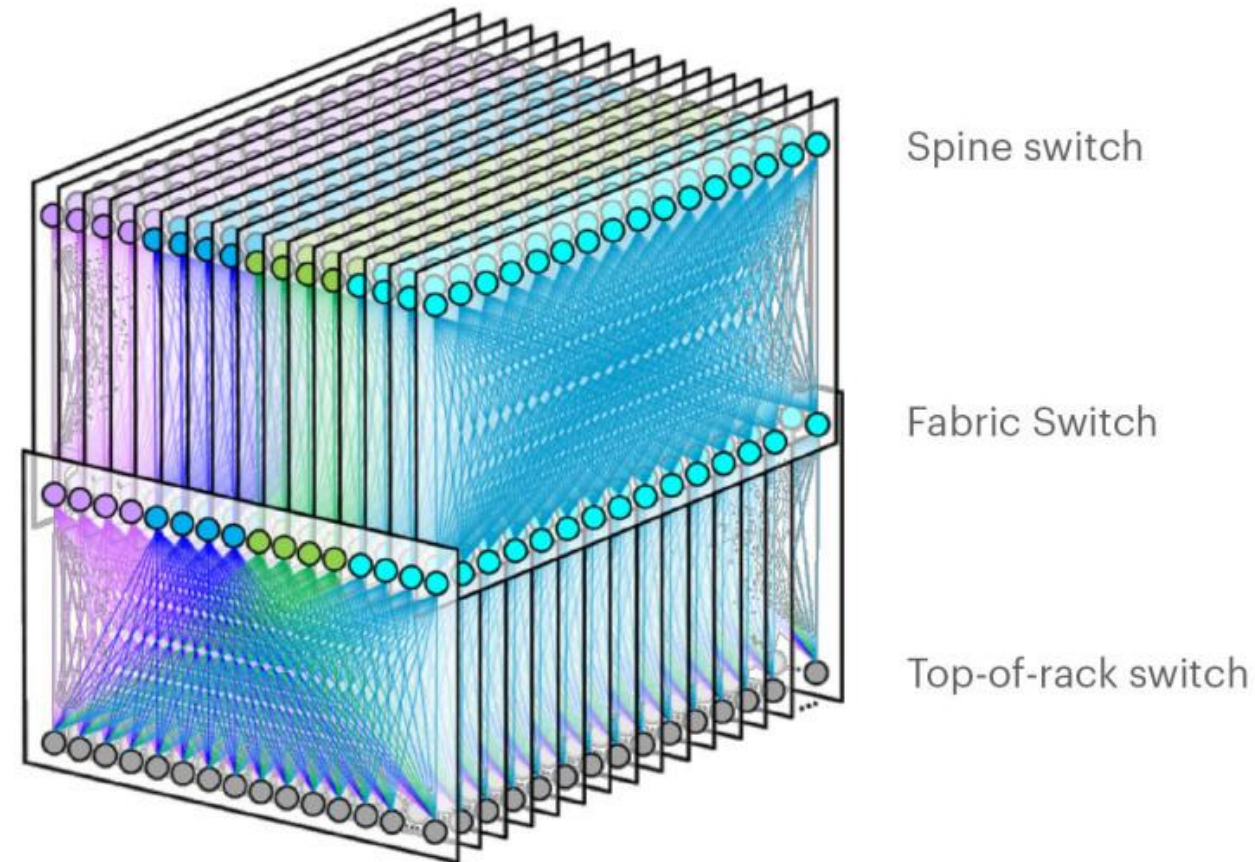


VL2: Summary

- Benefits
 - Uniform High Capacity
 - Max rate of server-to-server traffic flow limited only by available capacity of the NICs of servers involved
 - Assigning servers to a service is independent of network topology
 - Performance Isolation
 - Traffic of one service is not affected by the traffic of any other service (just as if each service was connected by a separate physical switch)
 - Layer 2 Semantics
 - Just as if servers on a LAN: IP address can connect to any port of Ethernet switch due to flat addressing
 - Virtual machines able to migrate to any server while keeping the same IP address

Data Center Network Elements

Facebook F16 data center network topology:



<https://engineering.fb.com/data-center-engineering/f16-minipack/> (posted 3/2019)

Jupiter (Google)

- Software Defined Networking
 - Logically centralized and hierarchical control plane
- Clos Topology
 - Non-blocking multistage switching topology
- Merchant Switch Silicon
 - Cost-effective, commodity general-purpose Ethernet switching

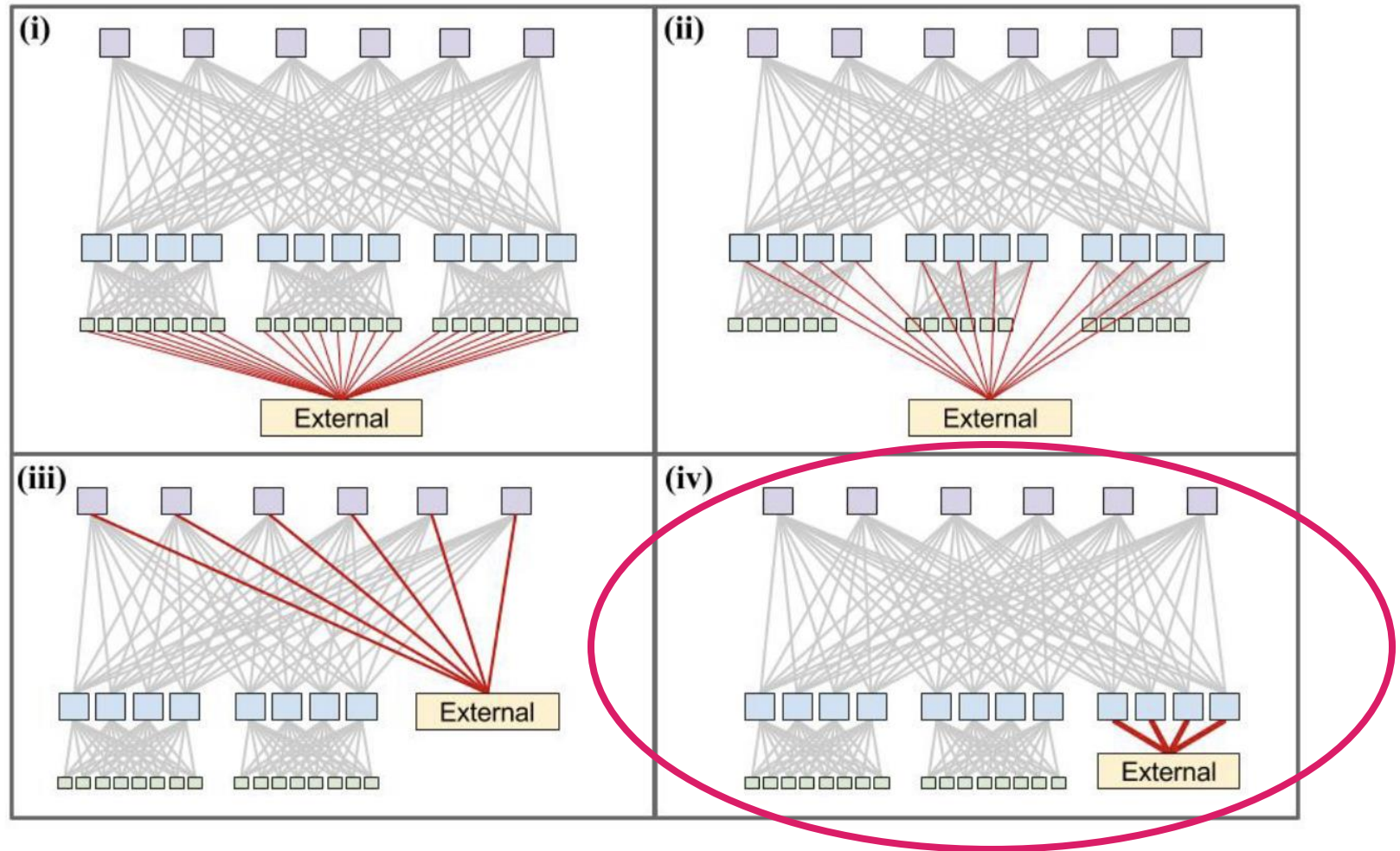
Google Cloud Blog (1): <https://cloudplatform.googleblog.com/2015/06/A-Look-Inside-Googles-Data-Center-Networks.html> (2015)

Google Cloud Blog (2): <https://cloud.google.com/blog/topics/systems/the-evolution-of-googles-jupiter-data-center-network> (2022)

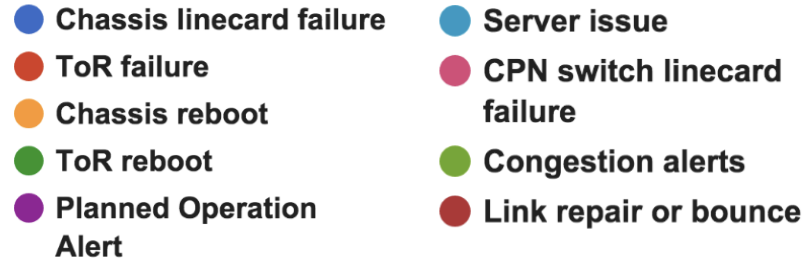
Jupiter

External Connectivity

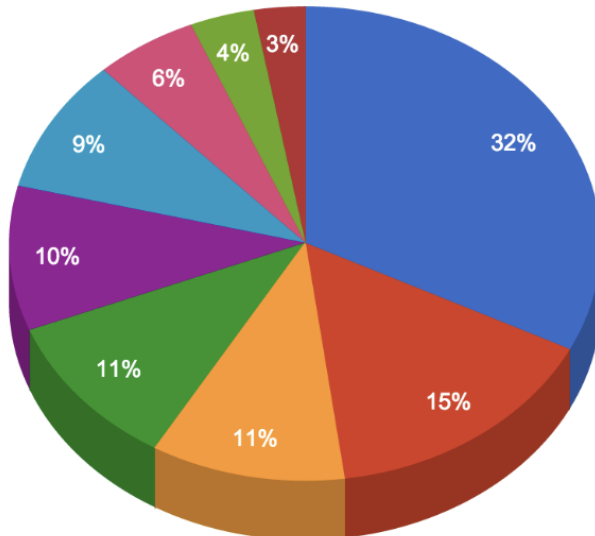
Can it handle burst
external bandwidth?



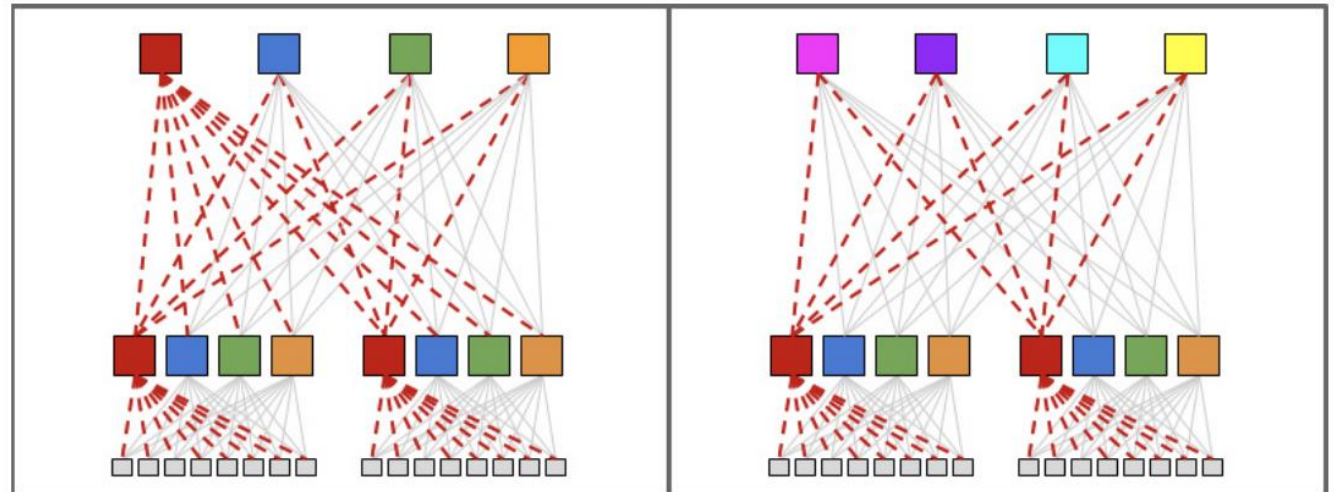
Jupiter



Alerts in a Cluster over 9 months

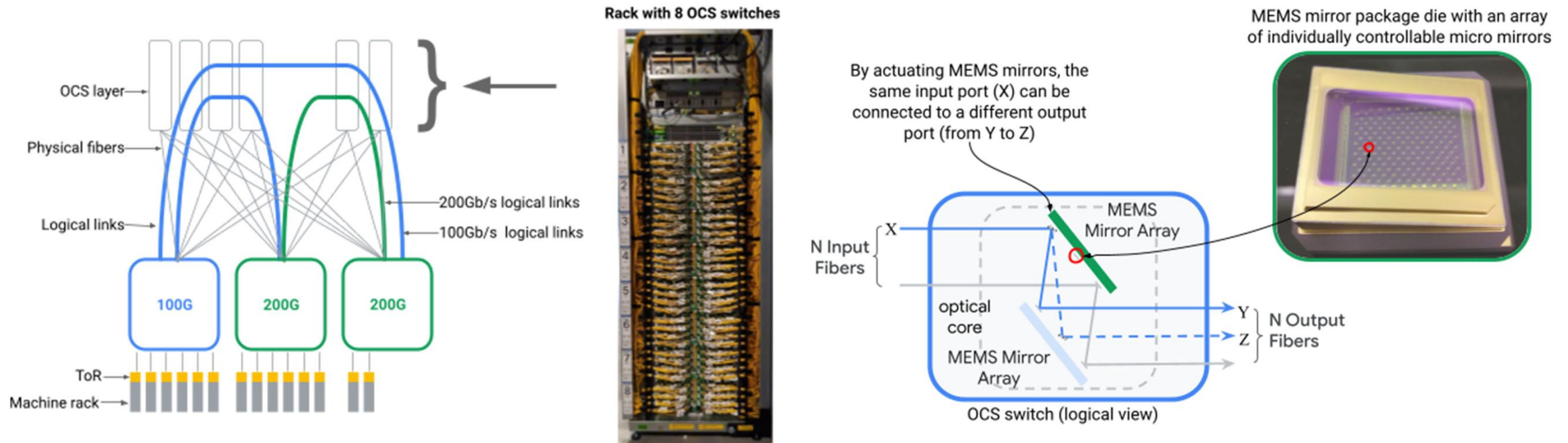


Fabric Software Redundancy Upgrades



Jupiter Evolving

Direct-connect topology & Optical Circuit Switching (OCS)



5x Higher Speed, 30% Capex Reduction, 41% Power Reduction

Jupiter Evolving

- Optical Circuit Switches (OCS)
- Direct **mesh-based** network topologies for higher performance, lower latency, lower cost, and lower power consumption
- Real-time topology and traffic engineering to simultaneously adapt network connectivity and pathing to match application priority and communication patterns
- Hitless network upgrades with localized add and removal of capacity

Cloud Resources

- Cloud Building Blocks
 - Compute
 - Store
 - Network
- Resources
 - CPU (Example: 4 Cores)
 - Memory (Example: 16GB RAM)
 - Storage (Example: 1TB HDD)
 - **Network (Example: 20 Gbps BW)**
 - I/O (Example: 1TB SATA)

Network Resources in The Cloud

- Virtual Private Cloud (Cloud VPC)
 - Connectivity for VM instances
 - Built-in Internal TCP & UDP Load Balancing
 - Proxy systems for Internal HTTP(S) Load Balancing
 - Distributes traffic
- Virtual Private Network (Cloud VPN)
 - Connect VPC network to on-premises or another cloud by a secure virtual private network
- Cloud NAT
 - Software-defined **network address translation** support
- Cloud Router
 - Border Gateway Protocol (BGP) exchange of routes between Virtual Private Cloud (VPC) and peer network
- Cloud Interconnect
 - Connect VPC network to an on-premises network by high-speed physical connection

Acknowledgement

The list of resources used in preparation of this slide set are provided on:

<https://canvas.sfu.ca/courses/88212/pages/references>

Pictures and quoted resources are mentioned in each use.

