## **Datacenters**

### CS 168 – Spring 2024

### Datacenters

- How does one design datacenter networks?
- Consider new assumptions (i.e., different from what we learned in Internet).
  - Single administrative control over the network topology, traffic, and to some degree end hosts
  - Much more **homogenous**
  - Strong emphasis on **performance**
  - Few backwards compatibility requirements, clean-slate solutions welcome!
  - 0 ...

# We will focus on DC topology in the rest of the discussion

## **Bisection Bandwidth**

• We want a network with high bisection bandwidth:

- Pick the number of links we must cut in order to partition a network into two halves
- Bisection bandwidth is the sum of those bandwidths.



## **Bisection Bandwidth**

- Full bisection bandwidth: Nodes in one partition can communicate simultaneously with nodes in the other partition at full rate.
  - Given N nodes, each with access link capacity R, bisection bandwidth = N/2 x R
- Oversubscription, informally, how far from the full bisection bandwidth we are
  - Formally: ratio of worst-case achievable bandwidth to full bisection bandwidth.

## **Bisection Bandwidth**



Bisection Bandwidth: 200G Full Bisection Bandwidth: (8/2)\*100G = 400G

Oversubscription: 200/400 = 2x

## **Big switch abstraction**

- We want an abstraction of big switch;
- Naive solutions:
  - Server-to-server full-mesh;
  - A physical big switch?
- Either impractical or too costly (\$\$\$); 10k hosts \* 10k hosts = ~100M links
- Can we do better?



# Design 1: Fat tree (scale-up)

- Borrowed straight from the High Performance Computing (i.e., supercomputers) community
- Use of big, non-commodity switches



- Problem? Scales badly in terms of cost
  ..only a few switch vendors can make big switches
- Scales badly in terms of fault-tolerance

cma	itch	



big switch

## **Design 2: Clos (scale-out)**

- Replace nodes in the fat tree with groups of cheap commodity switches
  - cheaper
  - high redundancy (bisection width)
- Allows oversubscription ratio of 1



## A Closer look at Clos

- The concept of redundancy and the clos topology itself have many variations!
- We focus on the 3-tiered clos topology
- Homogenous switches (k ports) and links
- Each non-core switch has k/2 ports pointing north and k/2 pointing south



## **Caveats about clos**

- Oversubscription ratio of 1 only with optimal load balancing
- Non-trivial to incrementally build and/or expand the network
  - e.g., port count k is fixed



- ECMP Equal Cost Multi-Path
  - Goal: use multiple paths in network topology that are <u>equal</u> <u>cost</u>
  - Idea: Load-balance packets across different forwarding paths
- ECMP Hash Function:
  - f(src\_ip, dst\_ip, proto, src\_port, dst\_port)
  - "Per-flow" load balancing

# (Over/Under)lay

- With VMs, many hosts spin up frequently
- Addressing could become a mess!
  - And won't scale!
- Thus, we build **overlay** networks



## **How? Encapsulation**

• Encapsulation: put another header on the packet



• Decapsulation: remove extra headers that were added for encapsulation



## Worksheet - Q1

- 1 Datacenter True/False
  - (1) Achieving low latency for small jobs (mice) is critical for datacenter networks.
  - (2) Clos topology enables the use of commodity hardware in datacenter networks.
  - (3) ECMP leads to reordering of packets in the same flow.
  - (4) Over-subscription is the ratio of worst-case achievable aggregate bandwidth between end-hosts to total bisection bandwidth.

## Worksheet - Q1

- 1 Datacenter True/False
  - (1) Achieving low latency for small jobs (mice) is critical for datacenter networks.

Solution: True: It is necessary to prevent large slow flows from starving small latency sensitive tasks.

(2) Clos topology enables the use of commodity hardware in datacenter networks.

#### Solution: True

(3) ECMP leads to reordering of packets in the same flow.

#### Solution: False. Packets from the same flow take the same path.

(4) Over-subscription is the ratio of worst-case achievable aggregate bandwidth between end-hosts to total bisection bandwidth.

### Solution: True

2 Clos-based Topology



Given the datacenter network above:

- (1) How many paths from an arbitrary server upward to an arbitrary root switch? Consider only minimallength paths.
- (2) How many paths can be set up between M1 to M3 in this topology? M1 to M5? Consider only minimal-length paths.
- (3) How many of the M1 to M5 paths go through switch S2?

### 2 Clos-based Topology



Given the datacenter network above:

(1) How many paths from an arbitrary server upward to an arbitrary root switch? Consider only minimallength paths.

#### Solution: 1.

(2) How many paths can be set up between M1 to M3 in this topology? M1 to M5? Consider only minimal-length paths.

#### Solution: 2;4.

(3) How many of the M1 to M5 paths go through switch S2?

#### Solution: 1.

### 3 Bisection Bandwidth and ECMP

Consider the same network topology in Q2. Assume all links have capacity 10 Gbps and ECMP is used for multipath routing.

- (1) What is the bisection bandwidth of this topology?
- (2) Assume ECMP makes the best possible hashing/load balancing. If switch S1 and S3 fail, what is the available aggregate bandwidth from each four-server Pod to any other such Pod (a Pod contains four servers closest to each other)?
- (3) Assume M1, M2 are sending flows to M5; and M3, M4 are sending flows to M6. TCP makes optimal use of the available bandwidth and ECMP makes the worst possible hashing. What is the approximate average data rate server M1 can send to server M5?

3 Bisection Bandwidth and ECMP

Consider the same network topology in Q2. Assume all links have capacity 10 Gbps and ECMP is used for multipath routing.

(1) What is the bisection bandwidth of this topology?

### Solution: 80 Gbps.

(2) Assume ECMP makes the best possible hashing/load balancing. If switch S1 and S3 fail, what is the available aggregate bandwidth from each four-server Pod to any other such Pod (a Pod contains four servers closest to each other)?

### Solution: 20 Gbps.

(3) Assume M1, M2 are sending flows to M5; and M3, M4 are sending flows to M6. TCP makes optimal use of the available bandwidth and ECMP makes the worst possible hashing. What is the approximate average data rate server M1 can send to server M5?

Solution: 2.5 Gbps.



- (1) If routing a packet from VM 1 to VM 6, what is the expected path?
- (2) What encapsulation and decapsulation should the servers on the path from VM 1 to VM 6 perform? Fill in the following table.

Router/	Server	Adds/Deletes Overlay	<b>Overlay Address</b>	<b>Underlay Address</b>

(3) How many destinations does **router 1** ave to keep track of without encapsulation? How about with encapsulation? Consider how this isolation impacts VM scalability in a data center!

(1) If routing a packet from VM 1 to VM 6, what is the expected path?

### Solution: S1-S5-S4

(2) What encapsulation and decapsulation should the servers on the path from VM 1 to VM 6 perform? Fill in the following table.

Server	Adds/Deletes Overlay	<b>Overlay Address</b>	Underlay Address

Solution:	Server	<b>Adds/Deletes Overlay</b>	<b>Overlay Address</b>	<b>Underlay Address</b>
	Server 1	Adds	4.4.4/32	192.0.2.3
	Server 5	Neither	4.4.4/32	192.0.2.3
	Server 4	Deletes	4.4.4/32	192.0.2.3

(3) How many destinations does server 5 have to keep track of without encapsulation? How about with encapsulation? Consider how this isolation impacts VM scalability in a data center!

Solution: Without = all of them = 11, With = just the underlay / physical addresses = 5