Advanced Computer Networks

Physical Connectivity at the Rack/Cluster Scale

https://pages.cs.wisc.edu/~mgliu/CS740/F25/index.html

Ming Liu mgliu@cs.wisc.edu

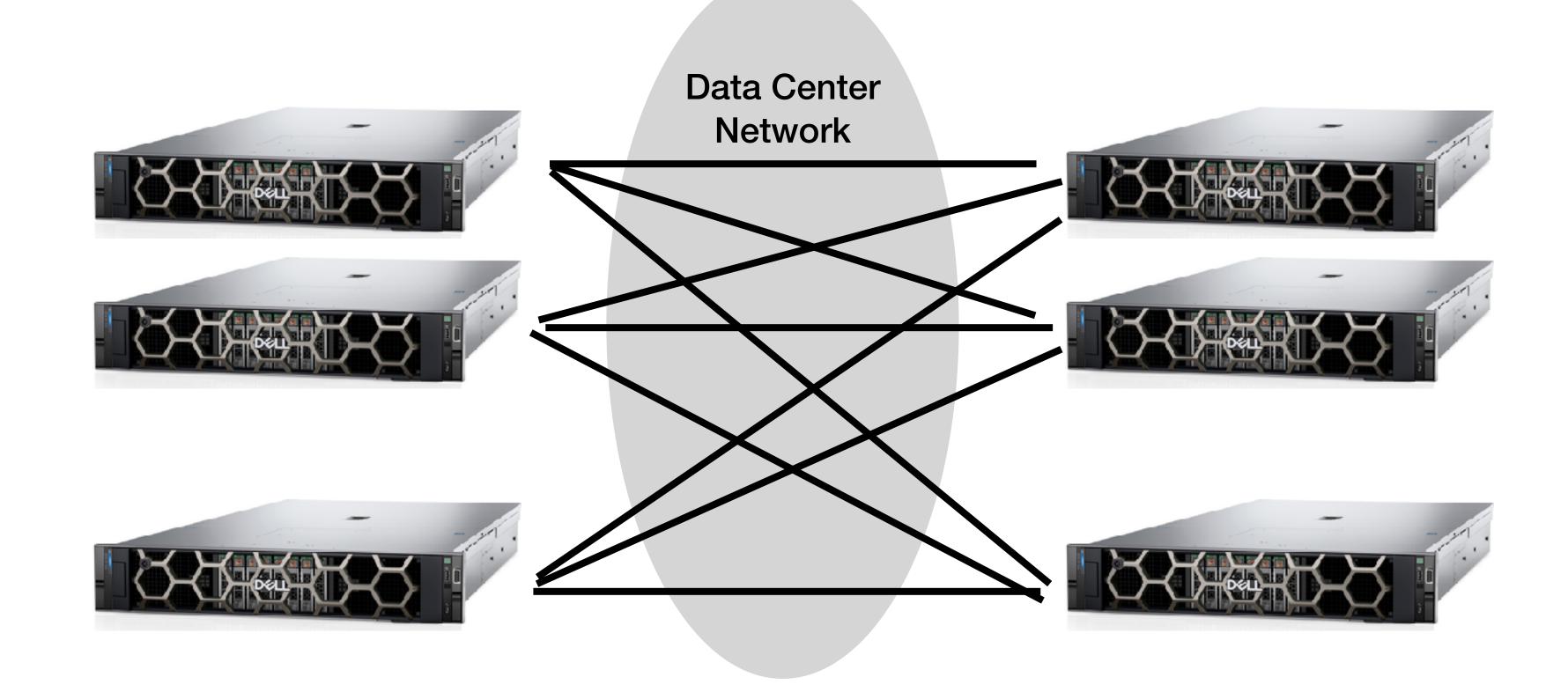
Outline

- Last lecture
 - Course logistics and schedule overview
 - Computer network design history
 - Data center network basics
 - Data center network design requirements
- Today
 - Physical connectivity at the rack/cluster scale
- Announcements
 - Lab1

What are the design requirements for data center networks?

High available server-to-server network connectivity at bandwidth Y among X NIC ports under cost efficiency

High available server-to-server network connectivity at bandwidth Y among X NIC ports under cost efficiency



Today's Focus

High available server-to-server network connectivity at bandwidth Y among X NIC ports under cost efficiency

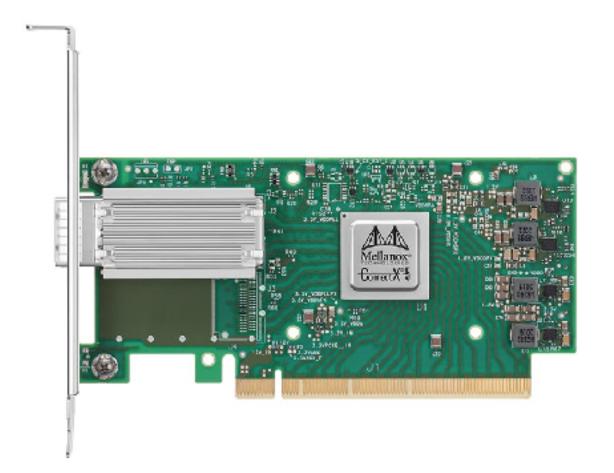
How can we connect two servers?





Physical Connectivity Between Two Servers

- Network Interface Card (NIC)
 - Port bandwidth (1Gbps, 10Gbps, 25Gbps, ...)
 - PCIe lane # and generation







Physical Connectivity Between Two Servers

- Networking Cable
 - Copper (Cat5e, Cat6, Cat6a, Cat7, Cat8) and Fiber (Single/Multi-Mode)
 - Transceiver: a serializer/deserializer(SerDes) converts signals at X GbE
 - Length: reliable data transfer speed, e.g., 1m, 10m, ...



Physical Connectivity Between Two Servers

- We focus on bandwidth in this lecture
 - In practice, server physical location and cost are also important



How can we connect three servers?





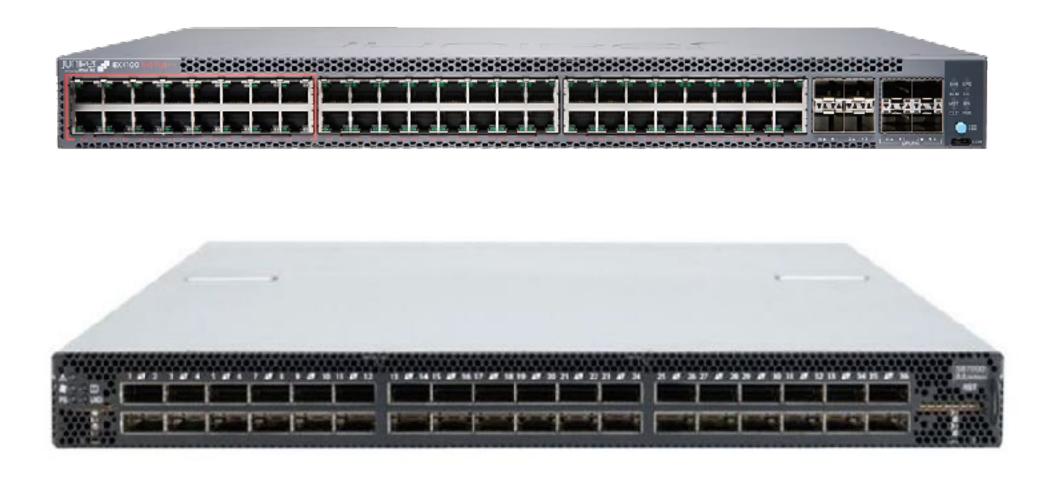


How can we connect three servers?



Physical Connectivity Among Three Servers

- Networking switch
 - A specialized networking gear providing fan-out connectivity
 - Vendors: Broadcom, Cisco, Dell, Arista, Nvidia, Marvell





Physical Connectivity Among Three Servers

- Networking switch
 - A specialized networking gear providing fan-out connectivity
 - Vendors: Broadcom, Cisco, Dell, Arista, Nvidia, Marvell

- Architectural internals
 - Fixed number of ports (K)
 - Switching ASIC for traffic forwarding
 - General-purpose CPU for running switch
 - L2/L3 switching

```
boo(config)#mlag configuration
boo(config-mlag)#show active
mlag configuration
    domain-id Arista
    local-interface Vlan4094
    peer-address 10.10.10.2
    primary-priority 10
    peer-link Port-Channel1000
boo(config-mlag)#show mlag detail | grep State
State : secondary
State changes : 51
boo(config-mlag)#show active
mlag configuration
    domain-id Arista
    local-interface Vlan4094
    peer-address 10.10.10.2
    primary-priority 10
    peer-link Port-Channel1000
boo(config-mlag)#show mlag detail | grep State
State : secondary
State changes : 51
boo(config-mlag)#show mlag detail | grep State
State : disabled
State changes : 53
boo(config-mlag)#show mlag detail | grep State
State : 53
boo(config-mlag)#show mlag detail | grep State
State : 53
boo(config-mlag)#show mlag detail | grep State
State : 53
boo(config-mlag)#show mlag detail | grep State
State : 53
boo(config-mlag)#show mlag detail | grep State
State : 53
boo(config-mlag)#show mlag detail | grep State
State : 53
boo(config-mlag)#show mlag detail | grep State
State : 55
boo(config-mlag)#show mlag detail | grep State
State : 55
boo(config-mlag)#show mlag detail | grep State
State : 55
boo(config-mlag)#show mlag detail | grep State
State : 55
boo(config-mlag)#show mlag detail | grep State
```

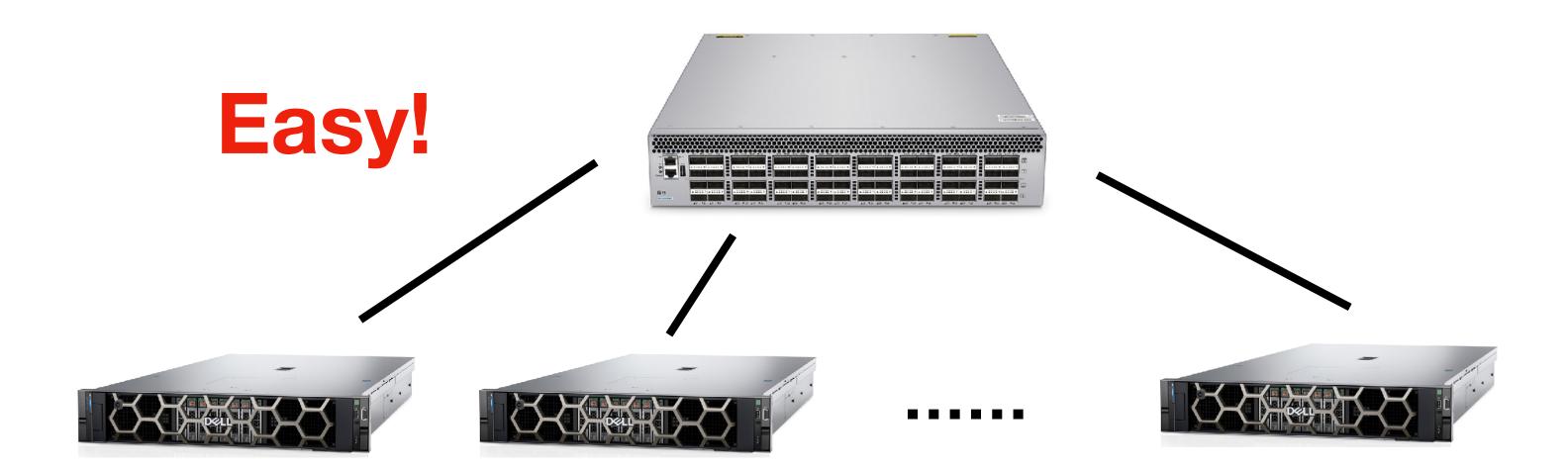
Physical Connectivity Among Three Servers

- Star topology
 - Switch port BW = NIC port BW = Cable BW = Y Gbps



Suppose a switch has K ports, how do we connect K servers?

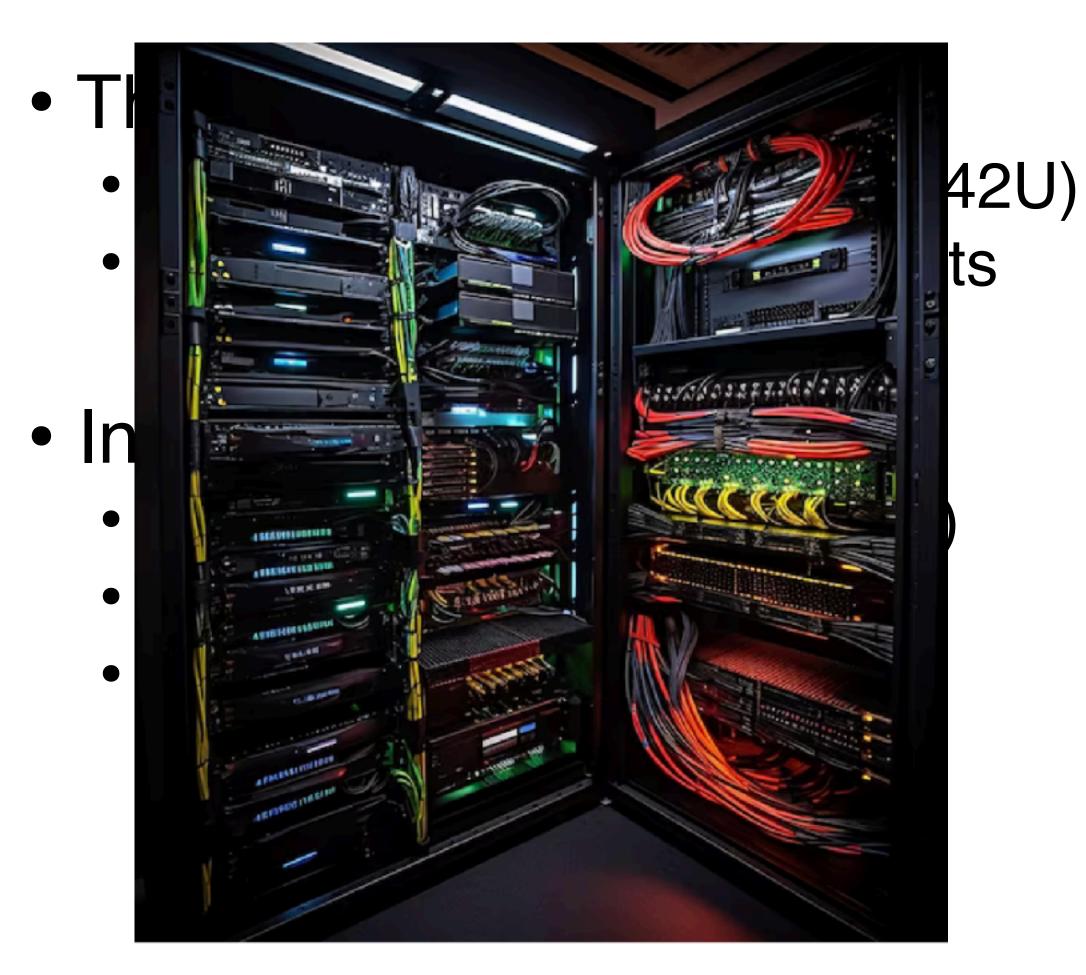
Suppose a switch has K ports, how do we connect K servers?



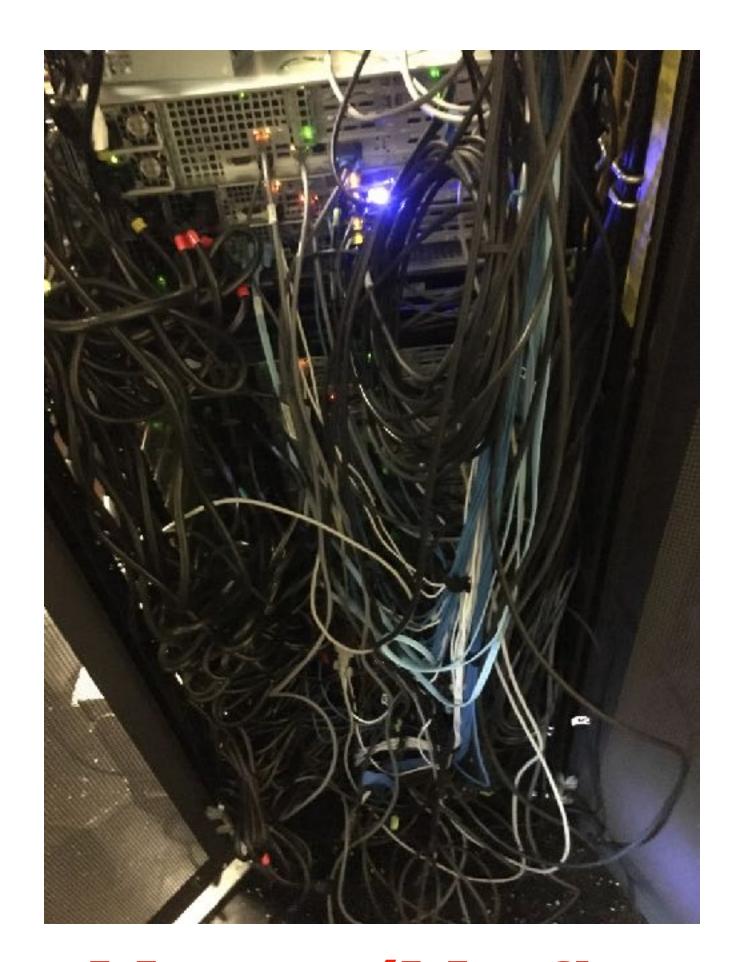
Rack-Scale Network Connectivity

- The size depends on
 - The height of a server rack (42U)
 - The number of switching ports
- Inside a rack
 - PDU (power distribution unit)
 - Servers + Switches
 - Different cables + cable tray

Rack-Scale Network Connectivity



Well-organized

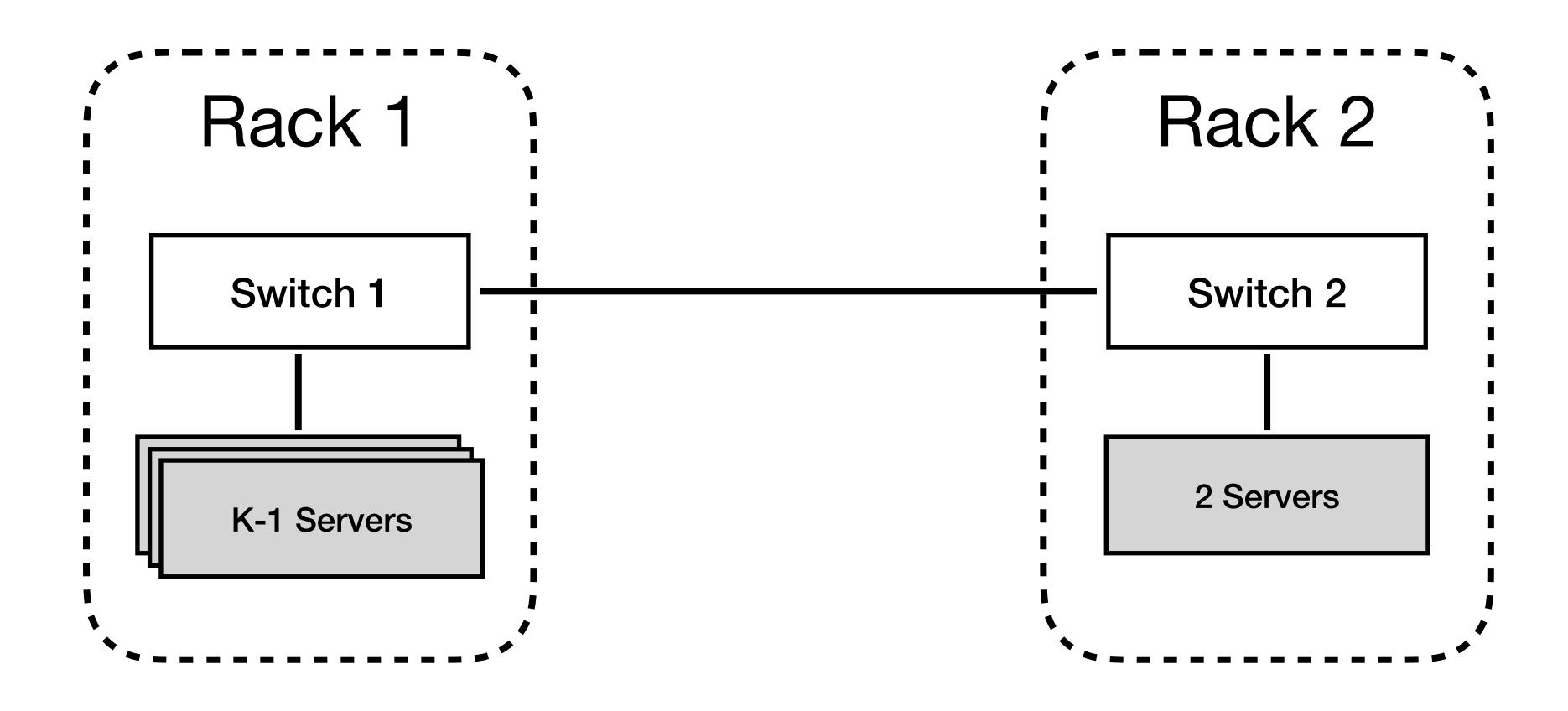


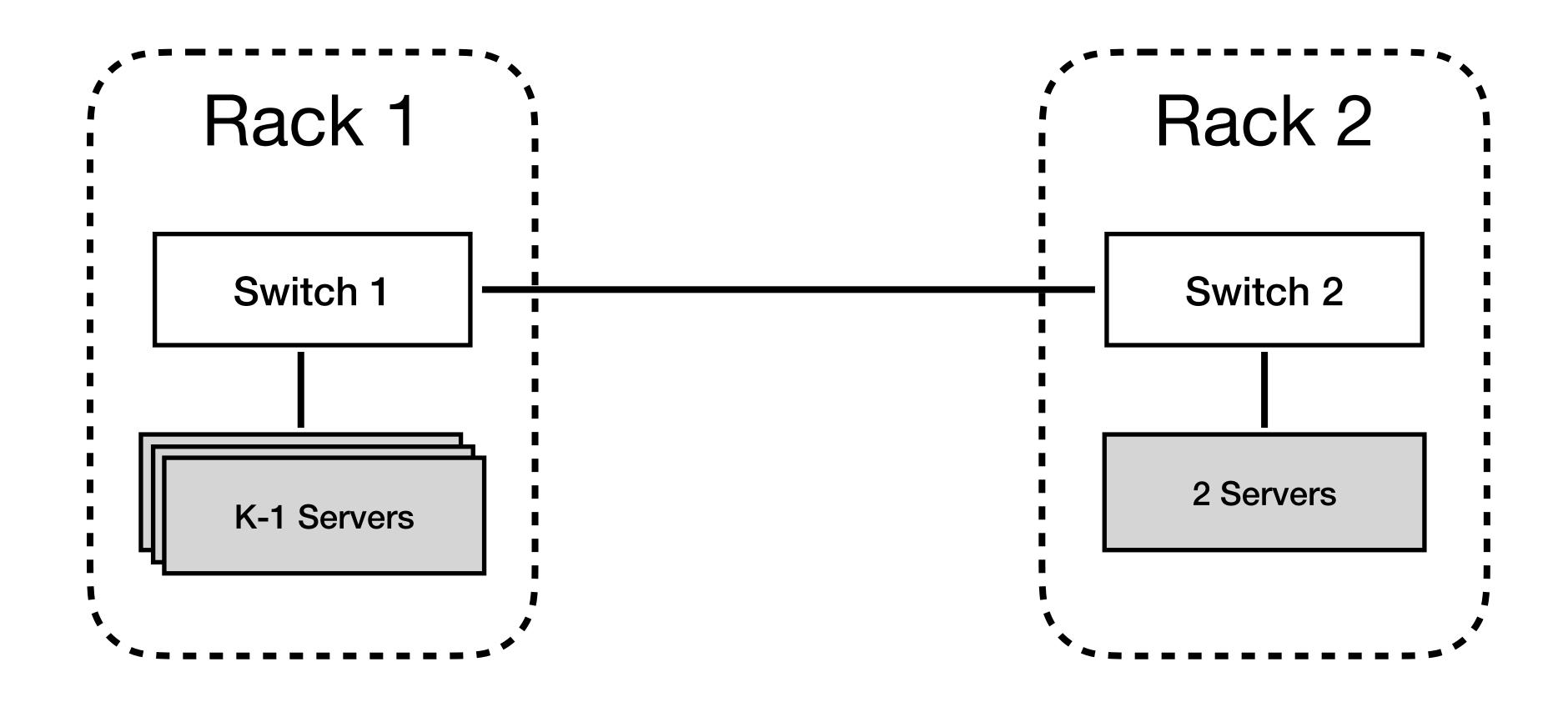
Messy (My first rack experience)

Suppose a switch has ports, how can we connect K+1 servers?

Suppose a switch has ports, how can we connect K+1 servers?

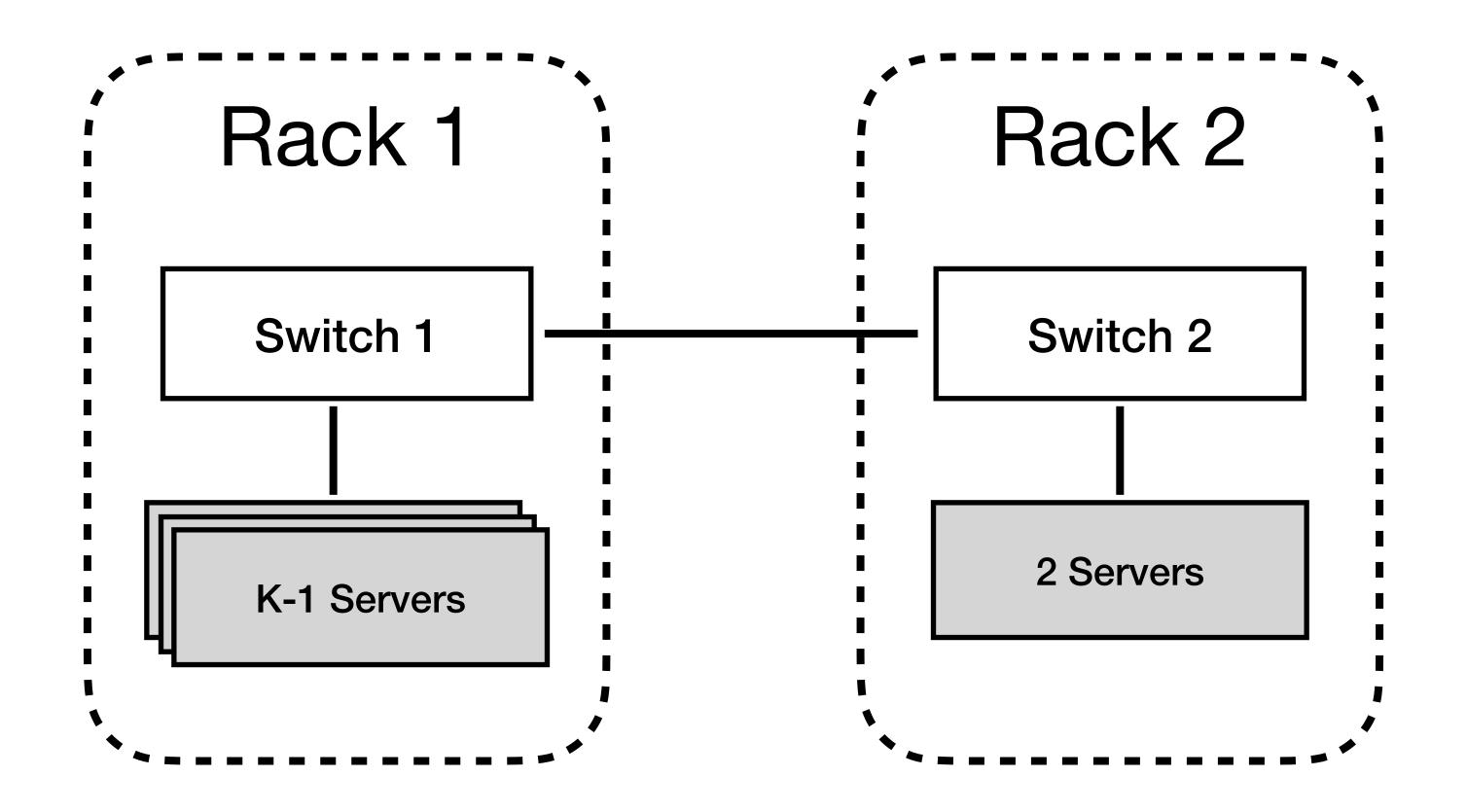
More switches!





Does this work?

- The ingress and egress bandwidth of a switch are unmatched!
 - Egress: the aggregated bandwidth issued to the outside from a switch
 - Ingress: the aggregated bandwidth coming from the outside to the switch

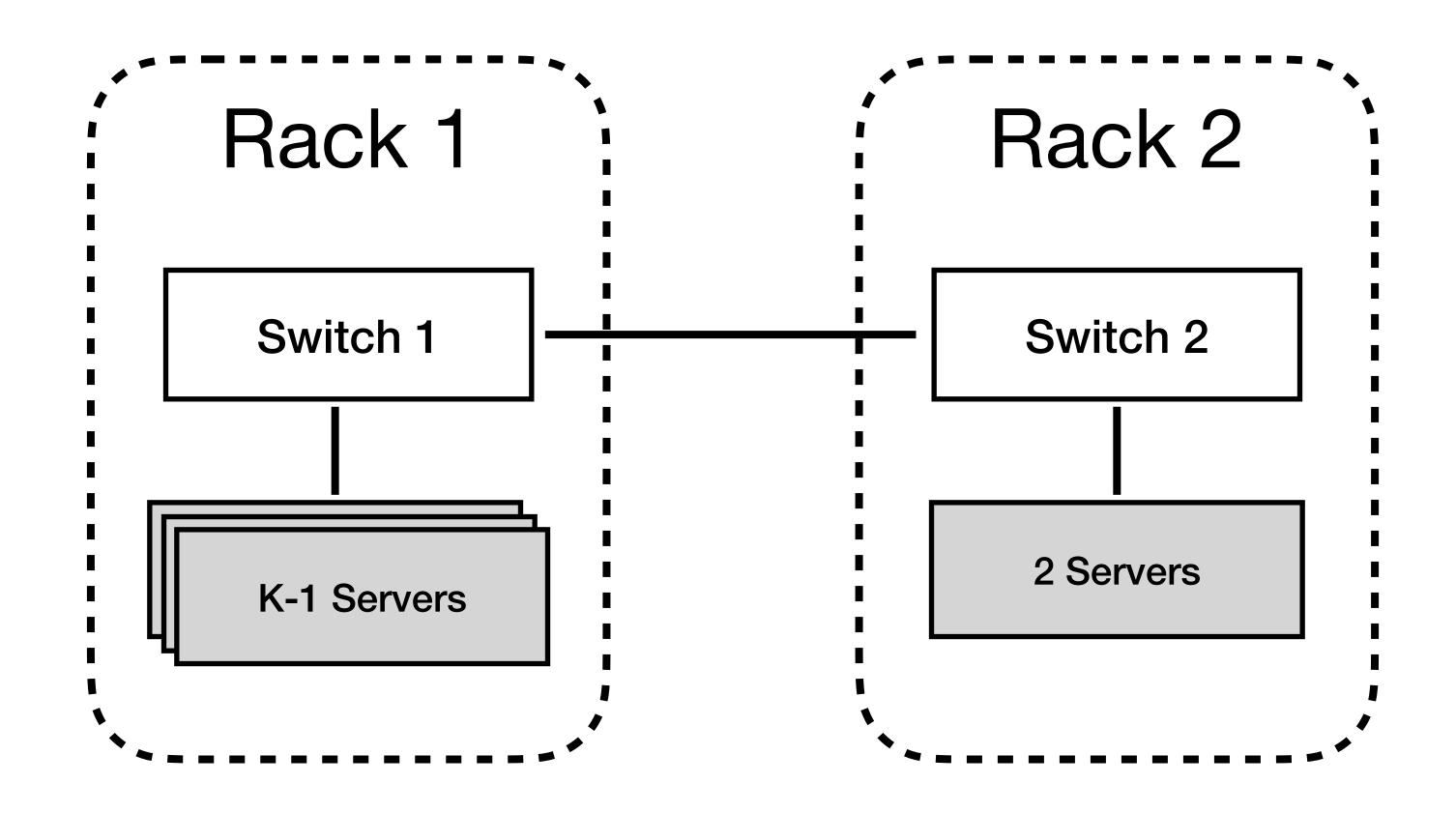


• Switch 1(Rack 1—> Rack 2)

Ingress: (K-1) * Y Gbps

Per-server: 1/(K-1) * Y Gbps

• Egress: 1 * Y Gbps

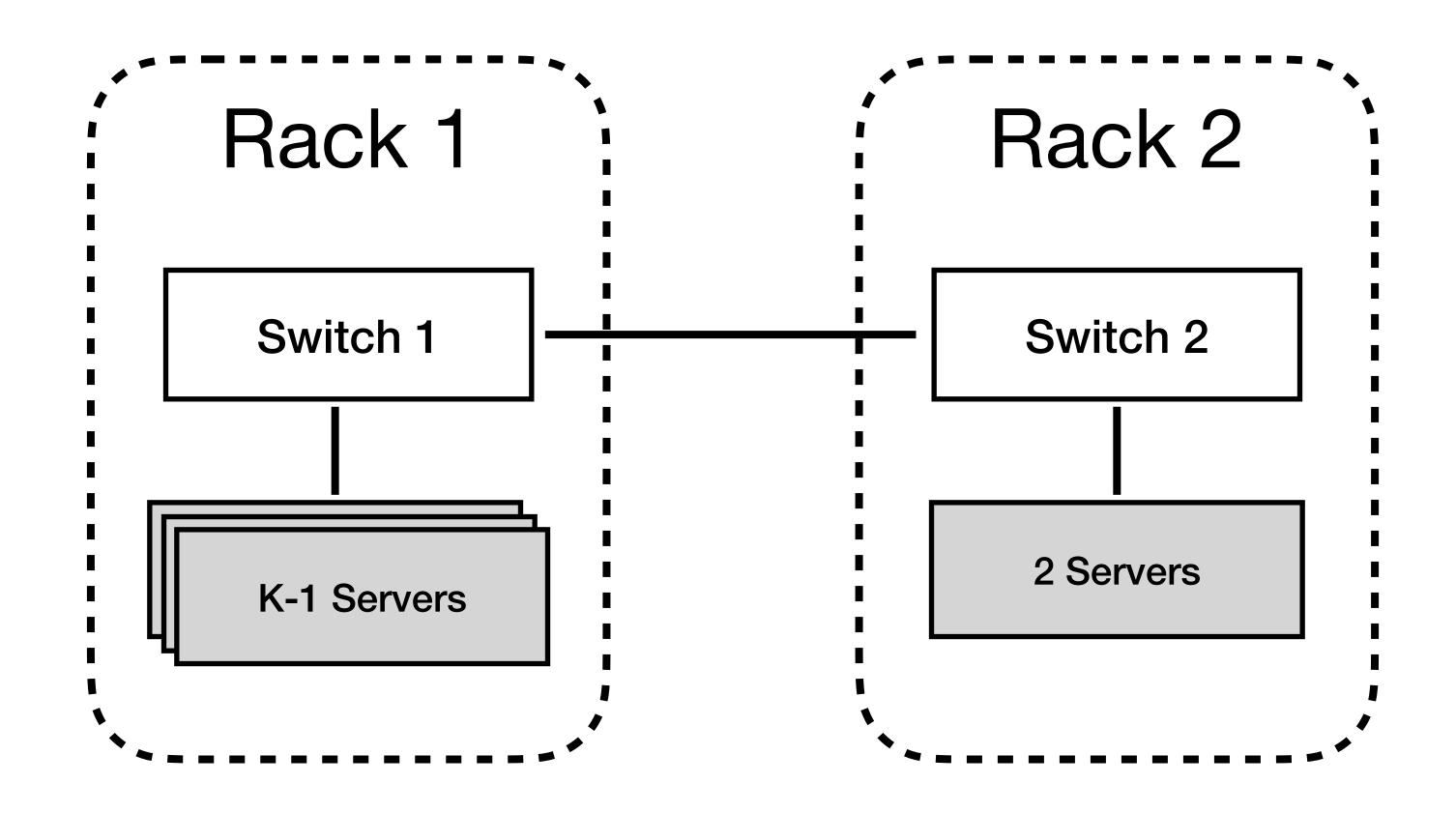


• Switch 1(Rack 2—> Rack 1)

Ingress: 1 * Y Gbps

Per-server: 1/(K-1) * Y Gbps

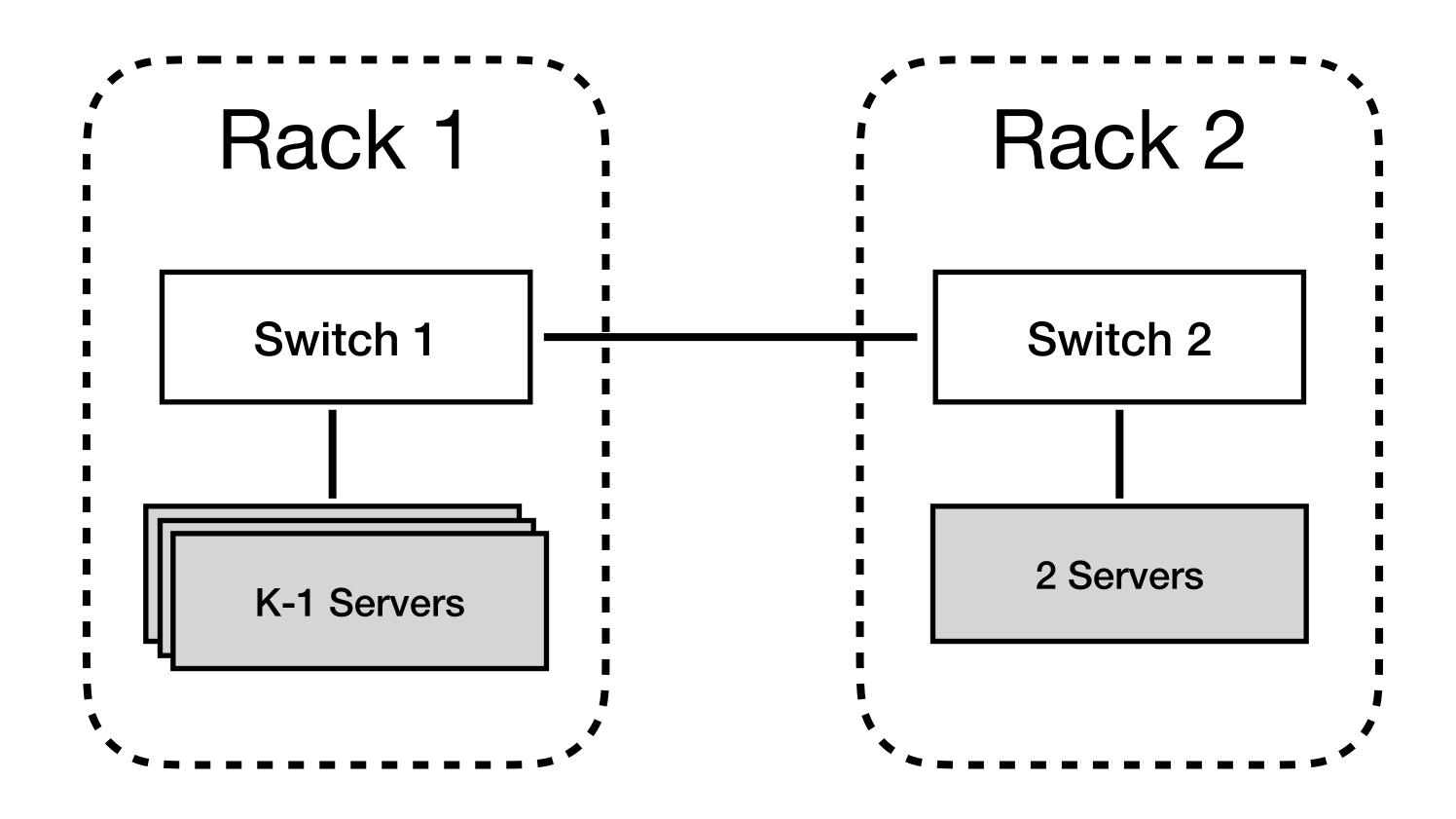
Egress: (K-1) * Y Gbps



- Switch 2(Rack 1 -> Rack 2)
 - Ingress: 1 * Y Gbps

• Egress: 2 * Y Gbps

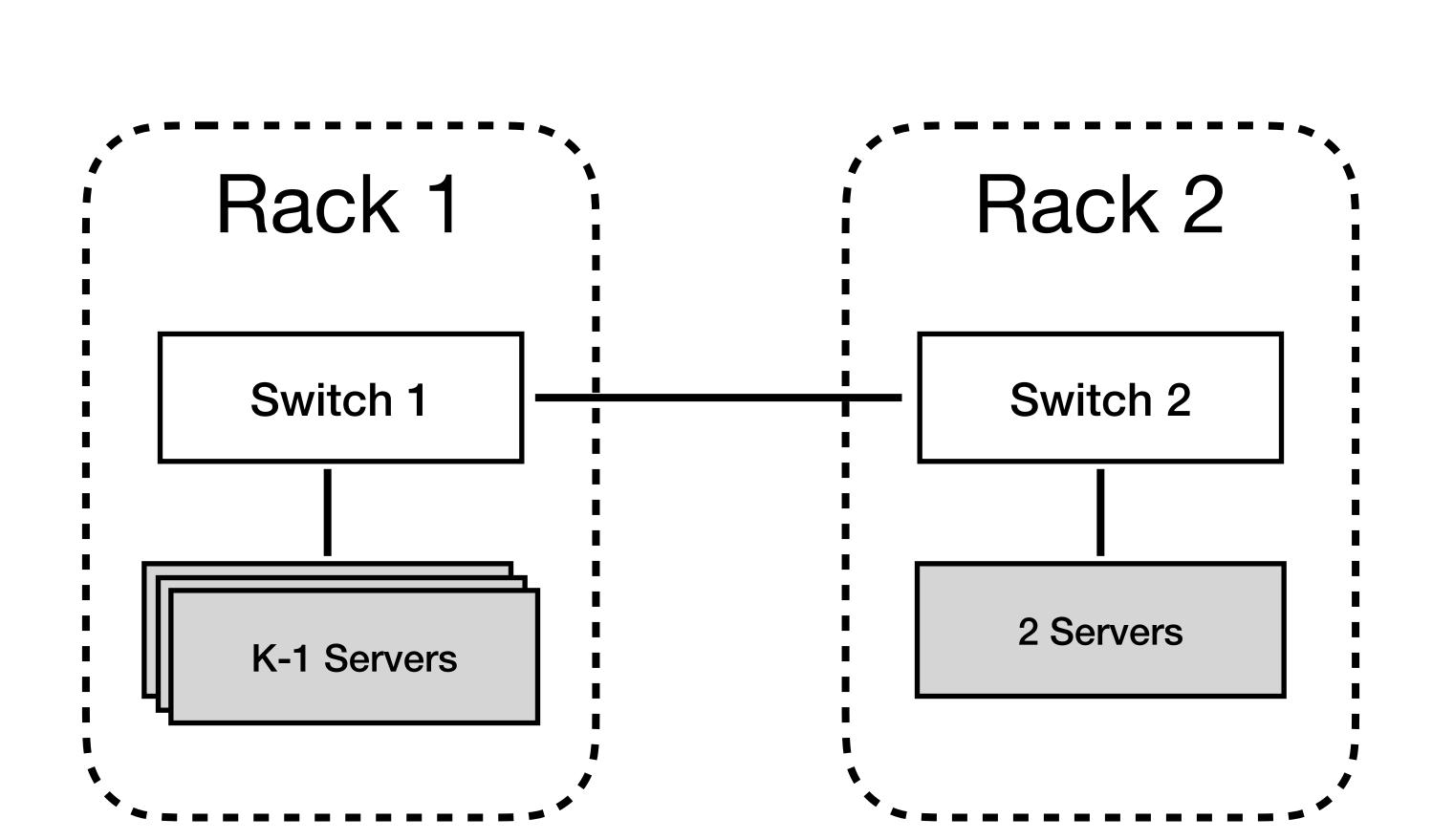
Per-server: 1/2 * Y Gbps

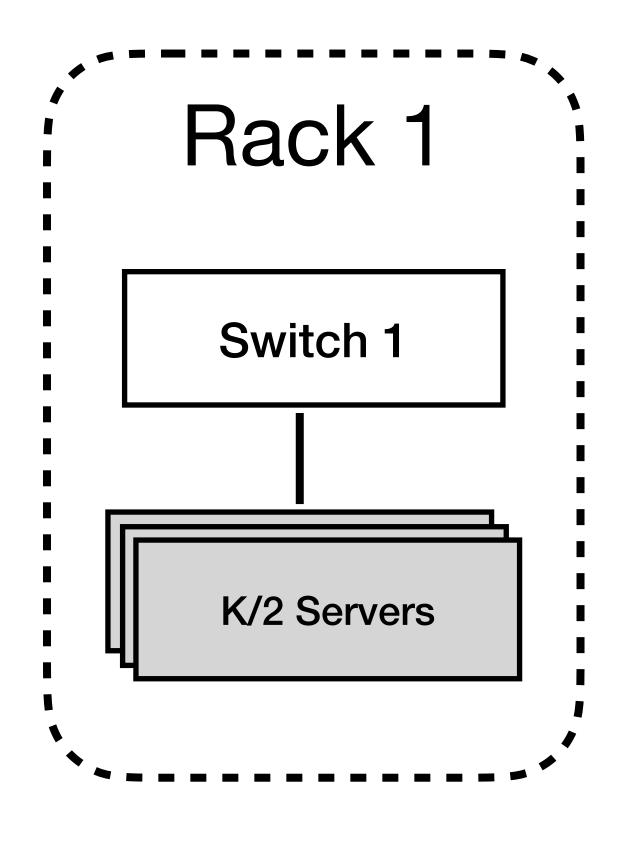


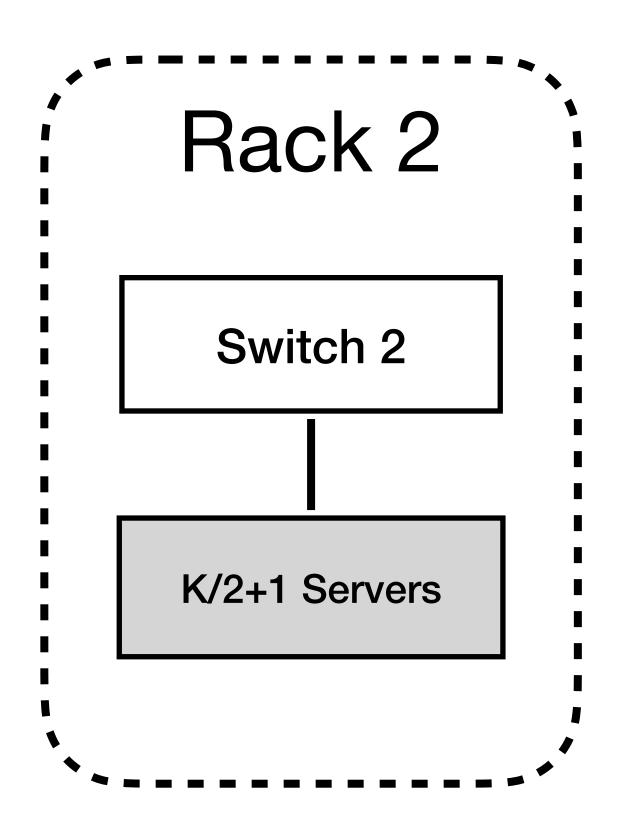
Per-server: 1/2 * Y Gbps

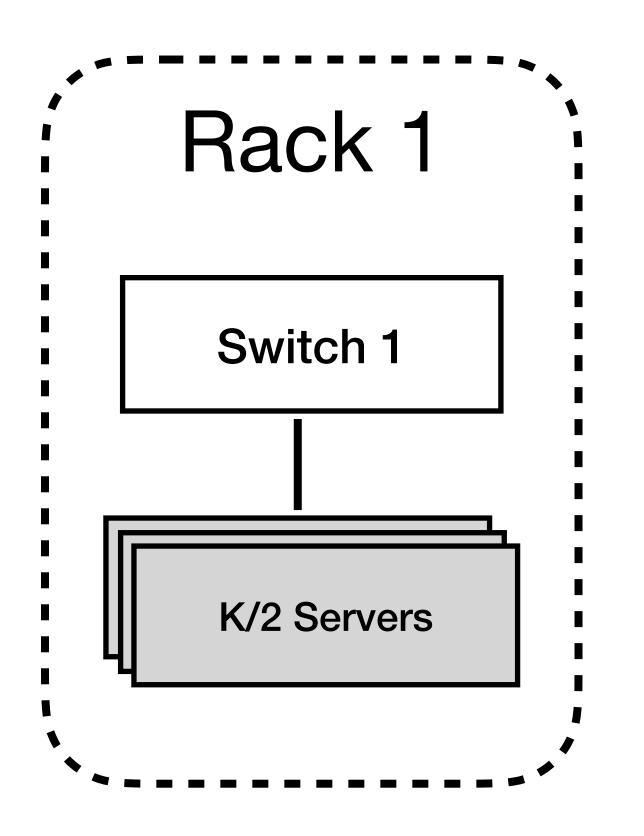
- Switch 2(Rack 2 -> Rack 1)
 - Ingress: 2 * Y Gbps

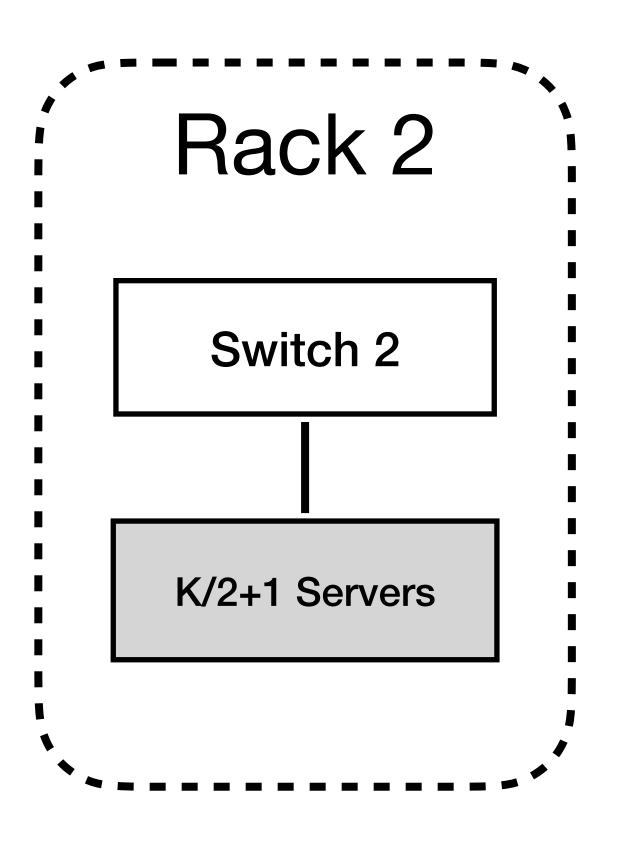
• Egress: 1 * Y Gbps



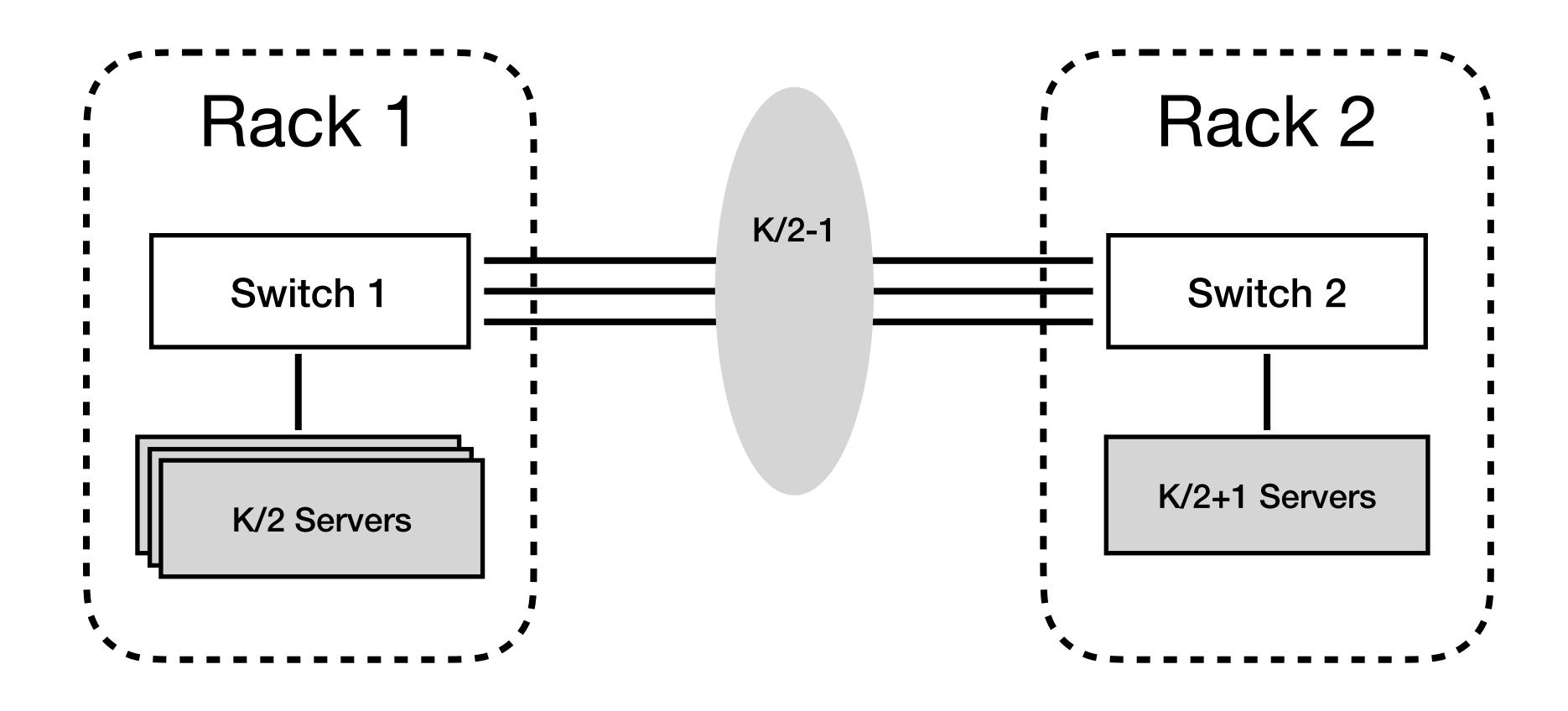








How to connect switch 1 and switch 2?

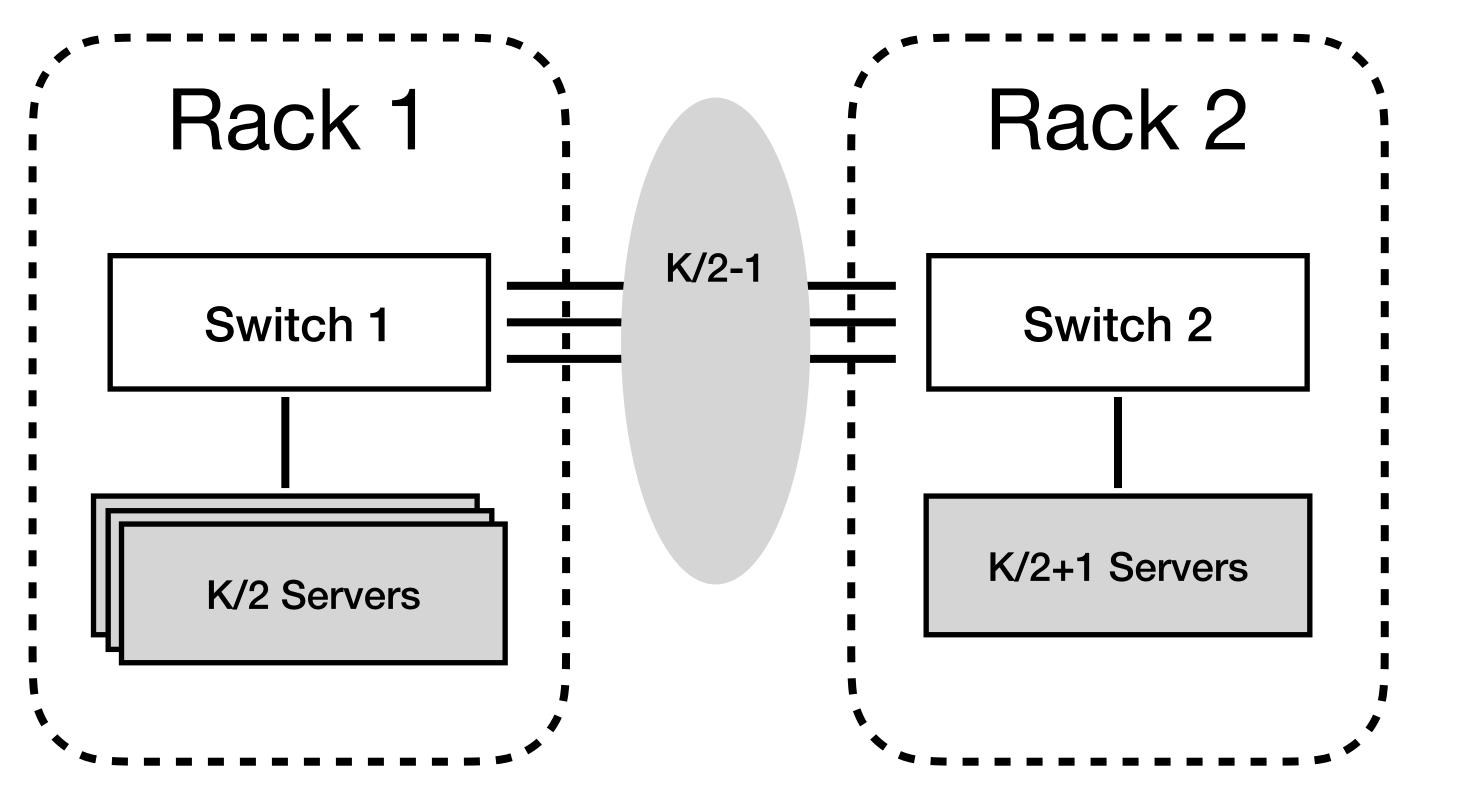


• Switch 1(Rack 1—> Rack 2)

Ingress: K/2 * Y Gbps

Per-server: (K-2)/K * Y Gbps

• Egress: (K/2 - 1) * Y Gbps

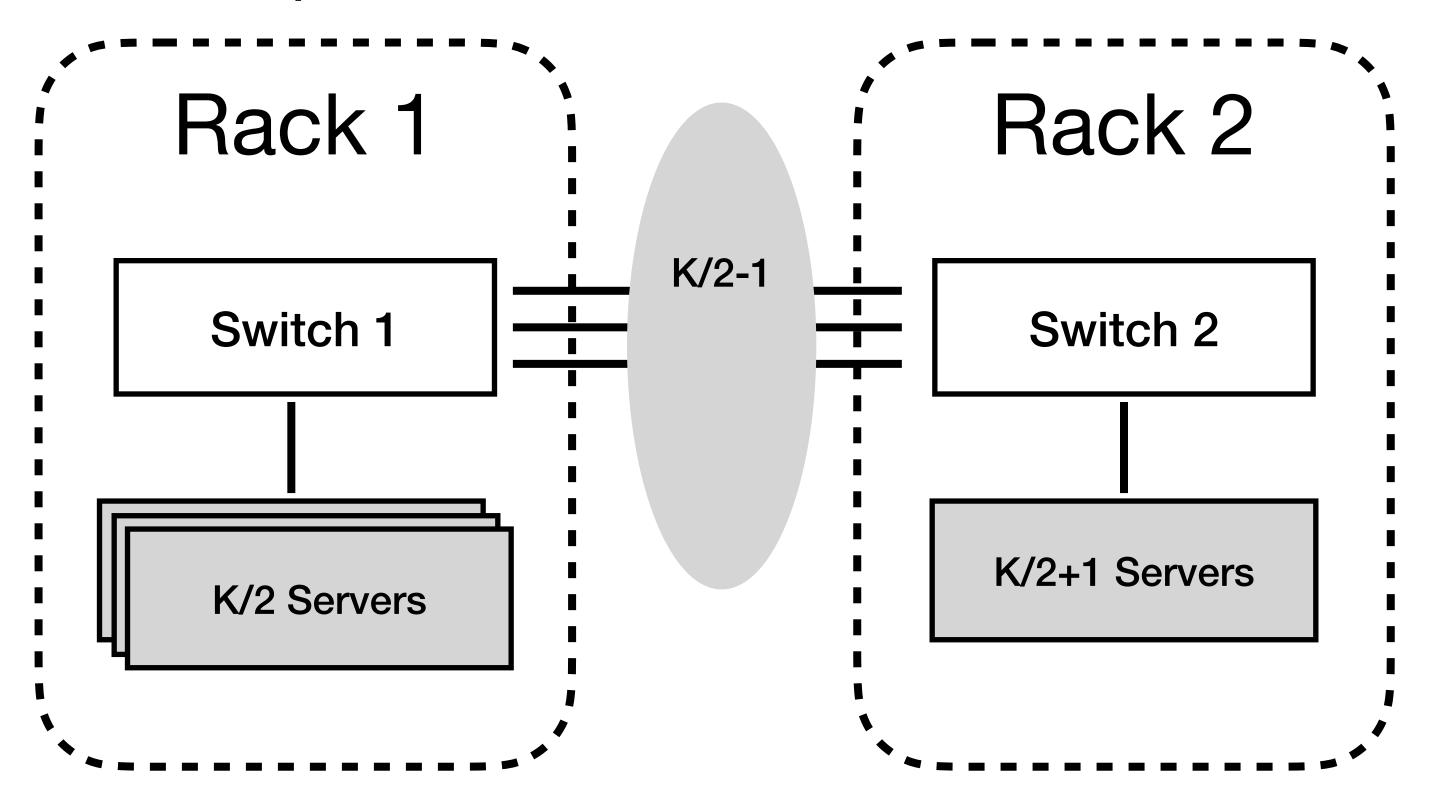


• Switch 1(Rack 2—> Rack 1)

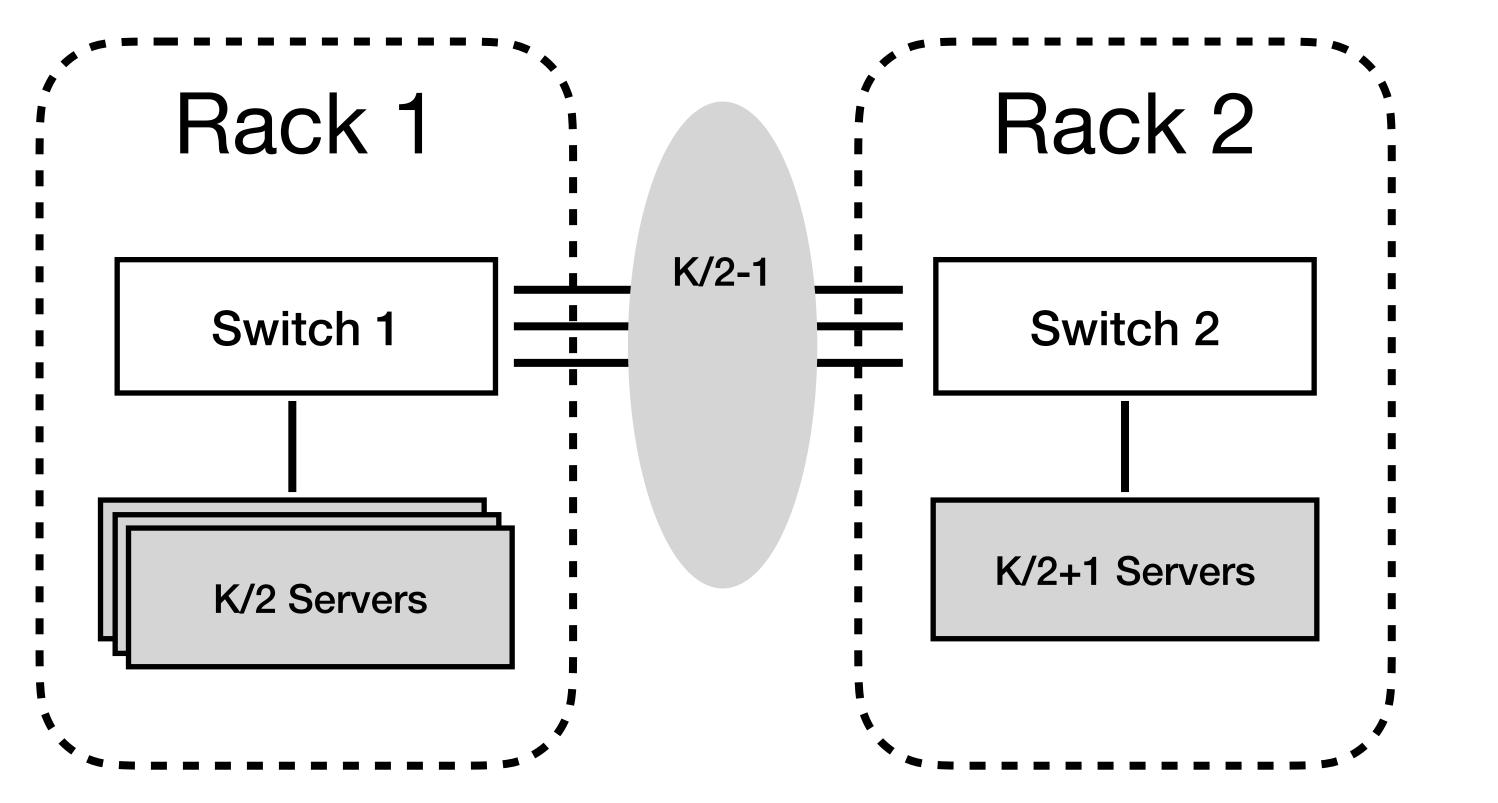
Ingress: (K/2-1) * Y Gbps

Per-server: (K-2)/K * Y Gbps

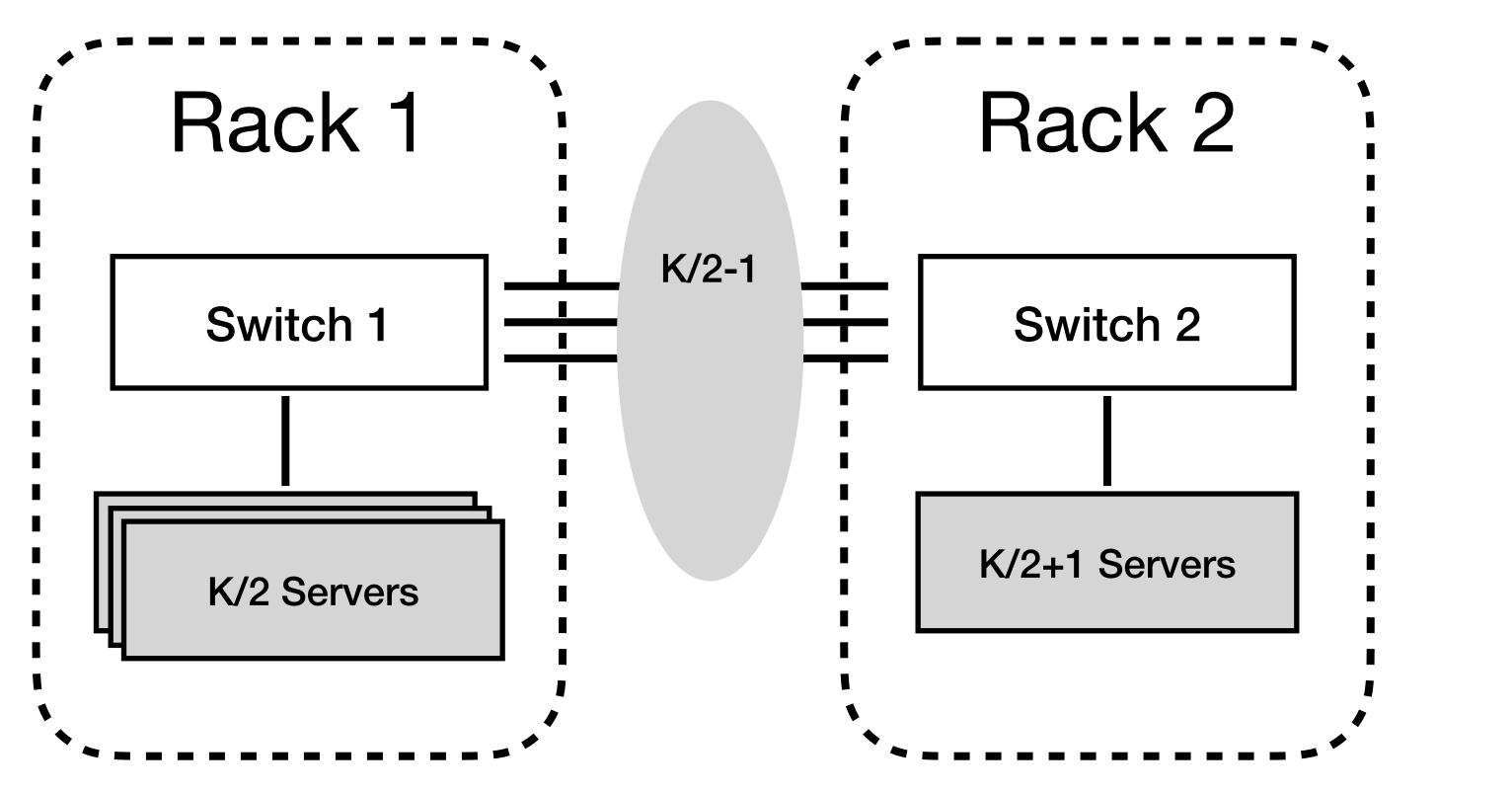
Egress: K/2 * Y Gbps



- Switch 2(Rack 1 -> Rack 2)
 - Ingress: (K/2 1) * Y Gbps Per-server: (K-2)/(K+2) * Y Gbps
 - Egress: (K/2 + 1) * Y Gbps



- Switch 2(Rack 2—> Rack 1)
 - Ingress: (K/2 + 1) * Y Gbps Per-server: (K-2)/(K+2) * Y Gbps
 - Egress: (K/2 1) * Y Gbps



Switch 2(Rack 2—> Rack 1)
 Ingress: (K/2 + 1) * Y Gbps
 Egress: (K/2 - 1) * Y Gbps
 Rack 1

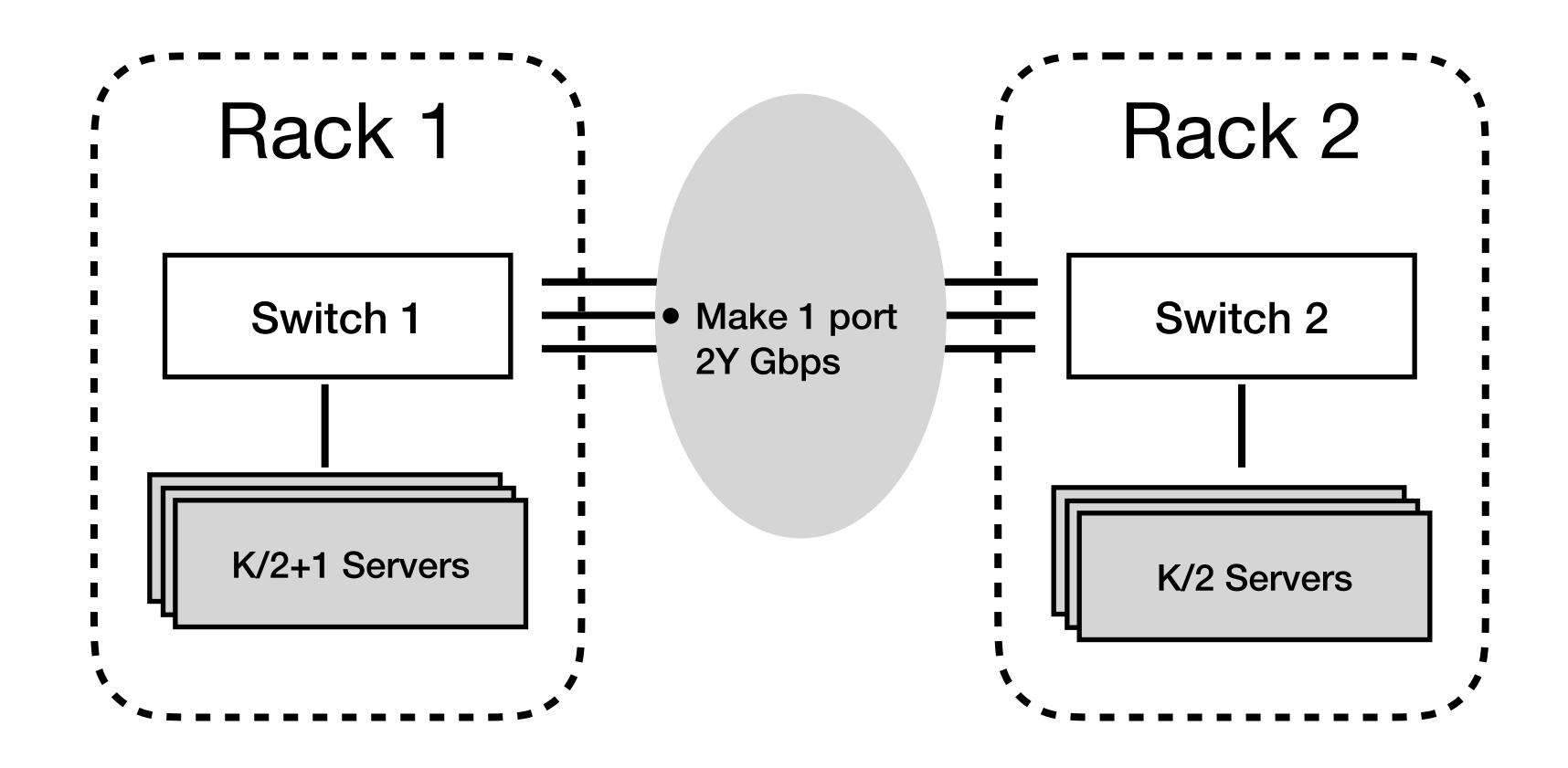
Rack 2

 Rack 2

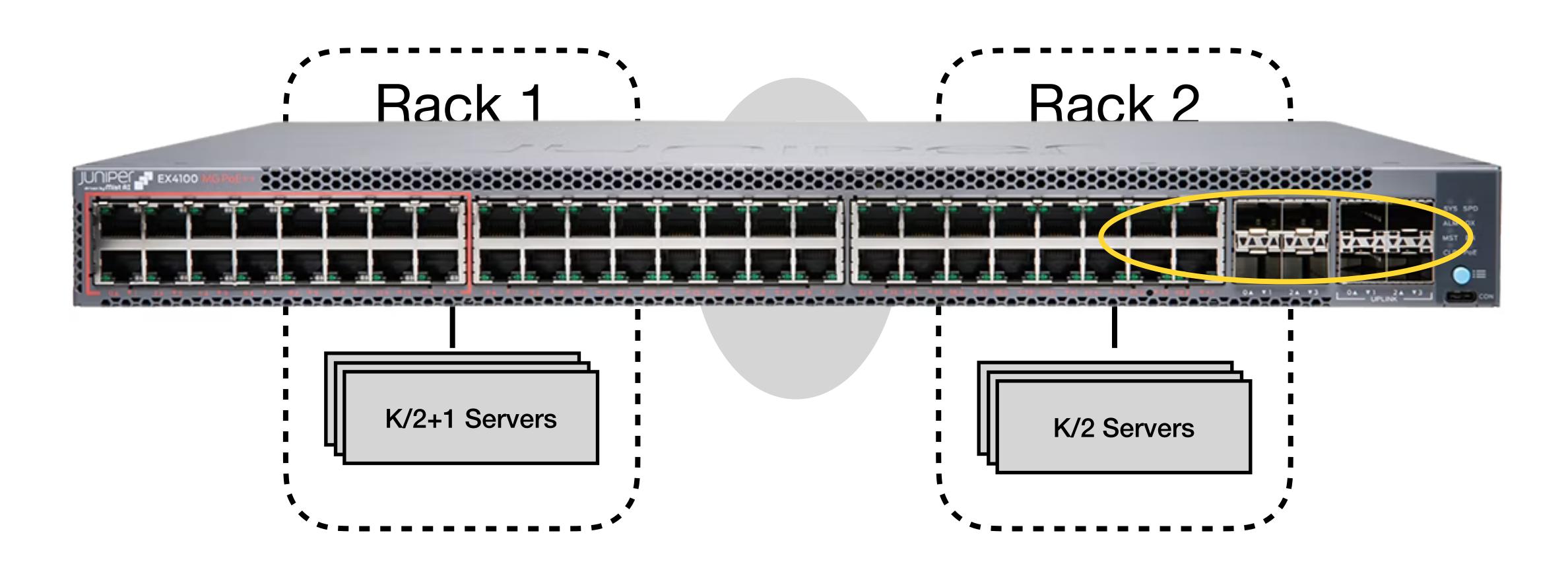
But server bandwidth is still not fully used!

Key: Match ingress and egress bandwidth at each switching point!

- #1: scale-out strategy
 - Enhance the switch
 - Slim (slow) port + Fat (fast) port



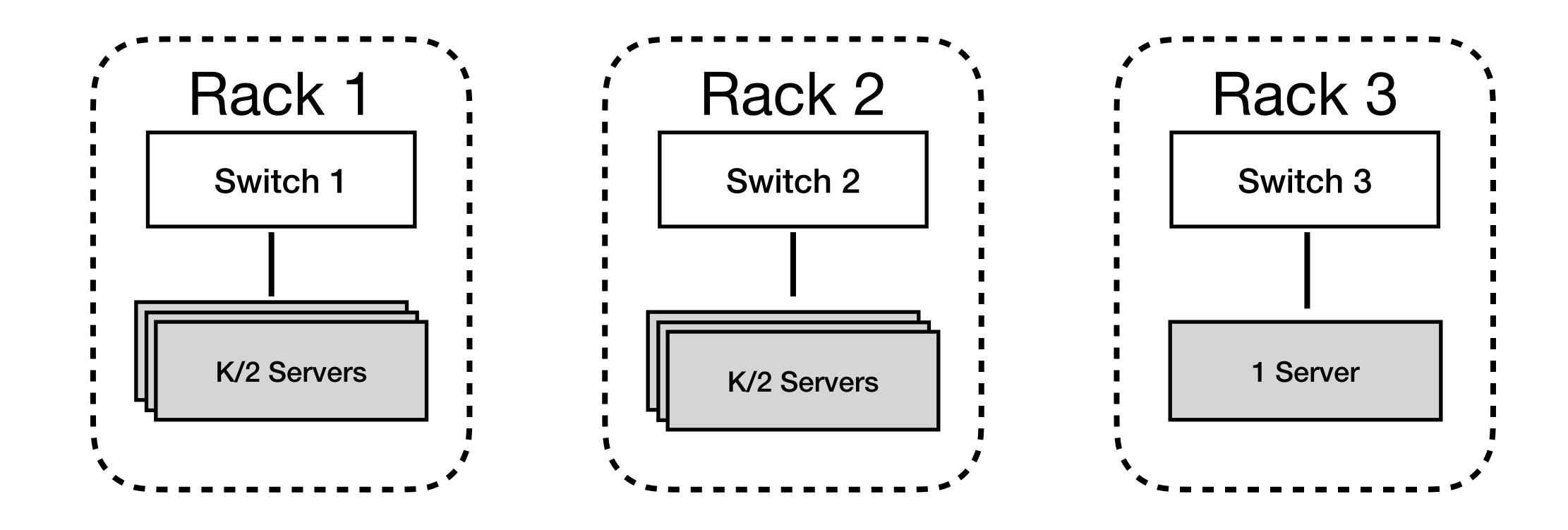
- #1: scale-out strategy
 - Enhance the switch
 - Slim (slow) port + Fat (fast) port



- #1: scale-out strategy
 - Enhance the switch
 - Slim (slow) port + Fat (fast) port



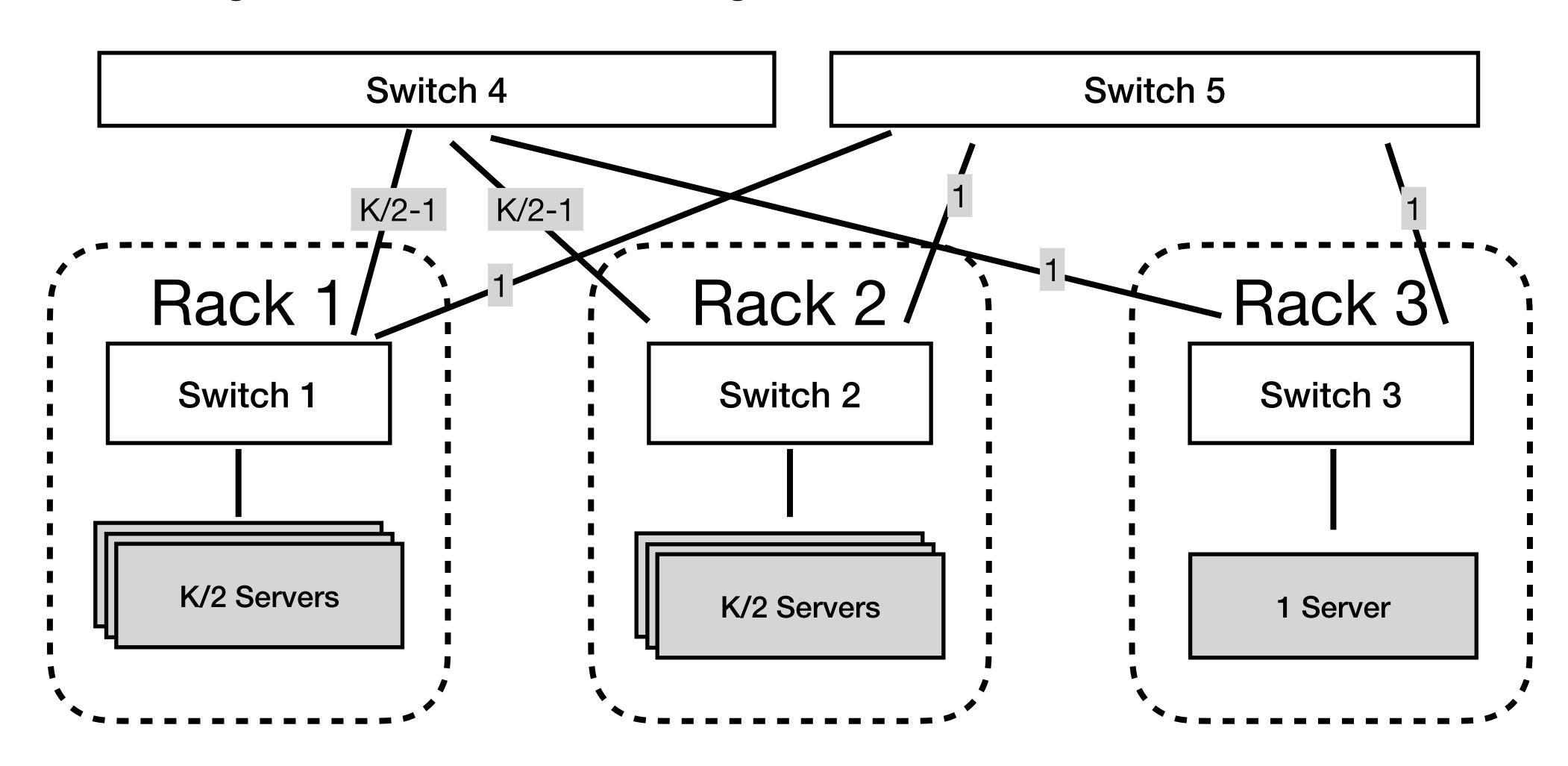
- #2: scale-up strategy
 - Adding more intermediate stages



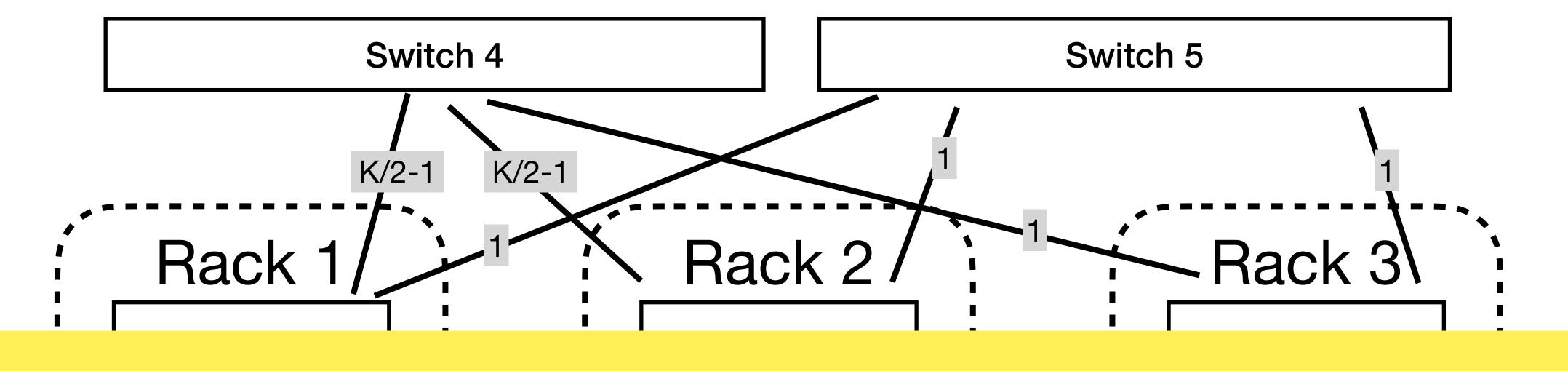
- #2: scale-up strategy
 - Adding more intermediate stages

Switch 5 Switch 4 Rack 3 Rack 1 Rack 2 Switch 2 Switch 3 Switch 1 K/2 Servers K/2 Servers 1 Server

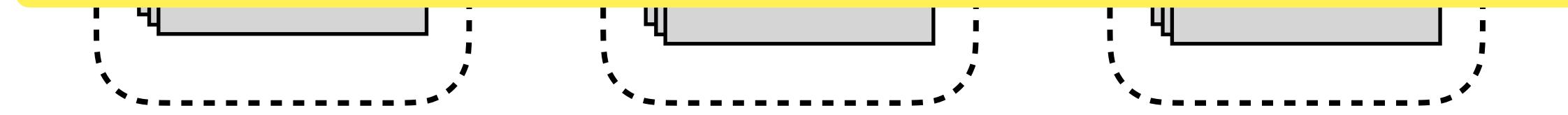
- #2: scale-up strategy
 - Adding more intermediate stages



- #2: scale-up strategy
 - Adding more intermediate stages



Adding more communication paths!



How can we connect X servers?

A Multistage Switching Network

- Clos networks, originally proposed in the telecommunications
 - Invented by Edson Erwin in 1938 and formalized by Charles Clos in 1952
- Fat-Tree topology
 - First proposed for parallel supercomputers

Fat-Trees: Universal Networks for Hardware-Efficient Supercomputing

CHARLES E. LEISERSON, MIMBER, IEEE

Abstract — This paper presents a new class of universal routing further from the leaves. In physical structure, a fat-tree networks called fur-trees, which might be used to interconnect the resembles, and is based on, the tree of meshes graph due to processors of a general-purpose parallel supercomputer. A fai-tree routing network is parameterized not only in the number of pendently from number of processors, substantial hardware can
be saved over, for example, hypercube-based networks, for such
parallel processing applications as finite-element analysis, but
without reserting to a special-purpose architecture.

Of some processing applications are finite-element analysis, but
without reserting to a special-purpose architecture.

Most retworks that have been repossed for parallel pro-

M CST routing networks for parallel processing super-computers have been analyzed in terms of perhow long it takes to route permutations, and cost is measured by the number of switching components and wires. This any planar interconnection strategy requires only O(s) volpaper presents a new routing network called fat-trees, but une. Thus, a natural implementation of a parallel finiteanalyzes it in a somewhat different model. Specifically, we element algorithm would waste much of the communication use a three-dimensional VLSI model in which pin boundednyss has a direct analog as the bandwidth limitation imposed by the surface of a closec three-dimensional region. Performance is measured by how long it takes to route an arbitruy set of messages, and vost is measured as the voume of that for a given physical volume of hardware, no network is a physical implementation of the network. We prove a wei-

Unlike a computer scientist's traditional notion of a tree, fat-trees are more like real trees in that they get thicker

Manuscript received February 1, 1985; revised May 30, 1985. This work wis supported in part by the Defense Advanced Sesearch Project Agency
under Couract N00014-80C-0622. A preliminary version of this paper was
prosented at the IEEE 1998 International Conference on Parallel Processing.

St. Charles, IL., Aug. 1982.

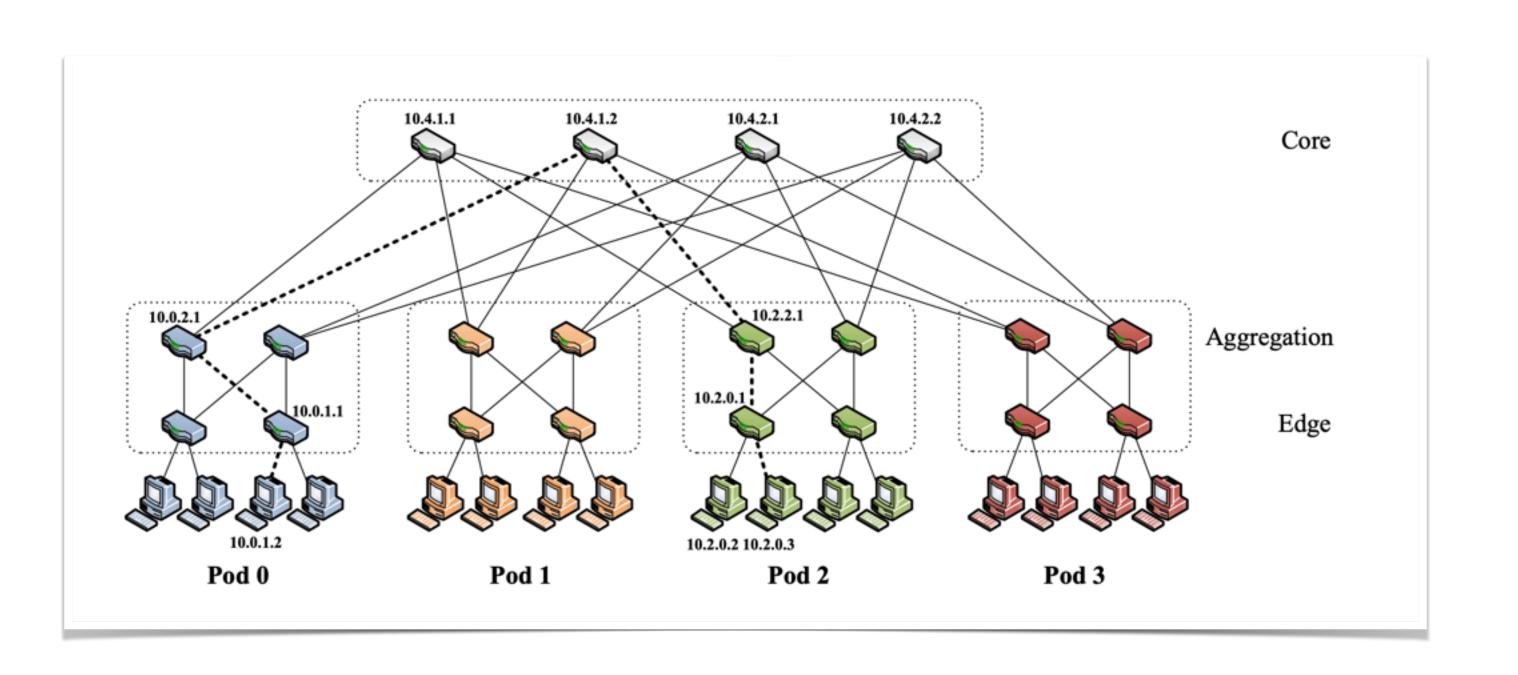
The author is with the Laboratory for Computer Science, Massachusetts assistant of Technology, Cambridge, MA UTLIFF.

resembles, and is based on, the tree of meshes graph due to Leighton [12], [14]. The processors of a fat-tree are located processors, but also in the amount of simultaneous communication it can support. Since communication can be scaled inde- are switches. Going up the fat-tree, the number of wires

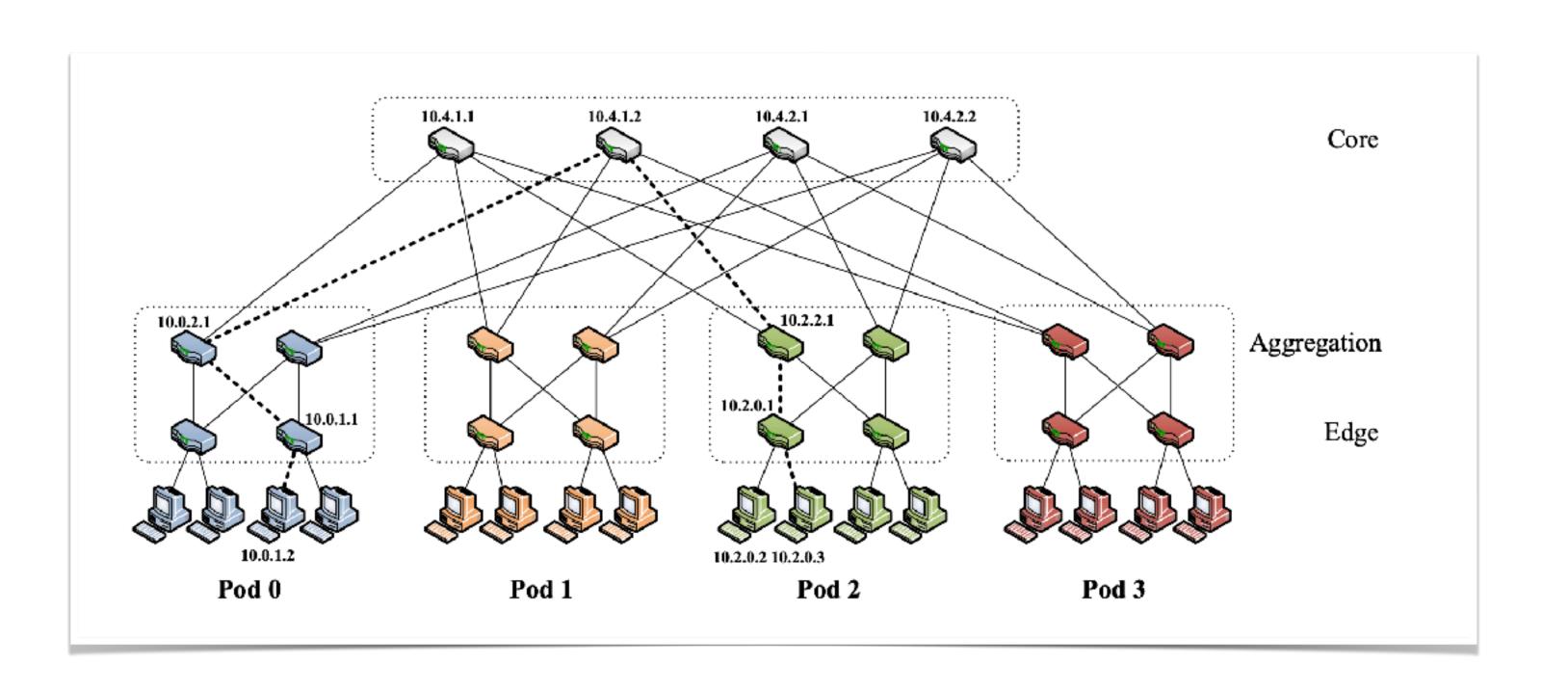
Most retworks that have been proposed for parallel propoof that a fai-tree of a given size is nearly the best routing cessing are based on the Boolean hypercube, but these net-network of that size. This universality theorem is proved using a works suffer from wirability and packaging problems and require nearly order x32 physical volume to interconnect x processors. In his influential paper on "ultracomputers" [27] Schwartz demonstrates that many problems can be solved the preceding section would appear to be the very large num ber of intercabinet wires which it implies." Schwartz ther goes on to consider a "layered" architecture, which seem easier to build, but which may not have all the nice properties of the original architecture.

On the other hand, there are many applications that do not require the full communication potential of a hypercube based network. For example, many finite-element problem formance and cost. Performance is typically measured by are planar, and planar graphs have a bisection width of size

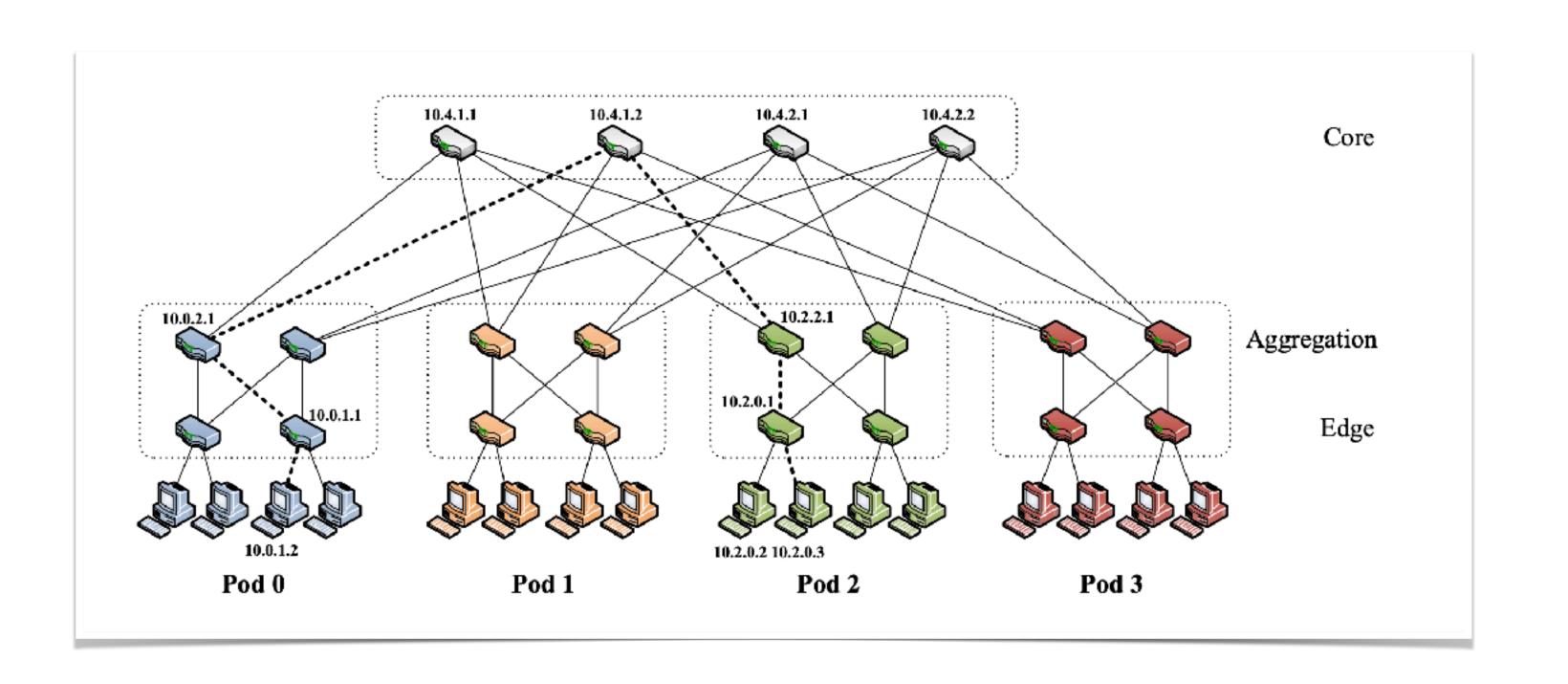
> Fat-trees are a family of general-purpose interconnection implementation. Section III shows how communication on a fat-tree can be scheduled off-line in a near-optimal fashion Section IV defines the class of universal fat-trees and in vestigates their hardware cost in a three-dimensional VLSI model. Section V contains several combinatorial theorem concerning the recursive decomposition of an arbitrary routsal routing networks. Finally, Section VII offers some remarks about the practicality of fat-trees.



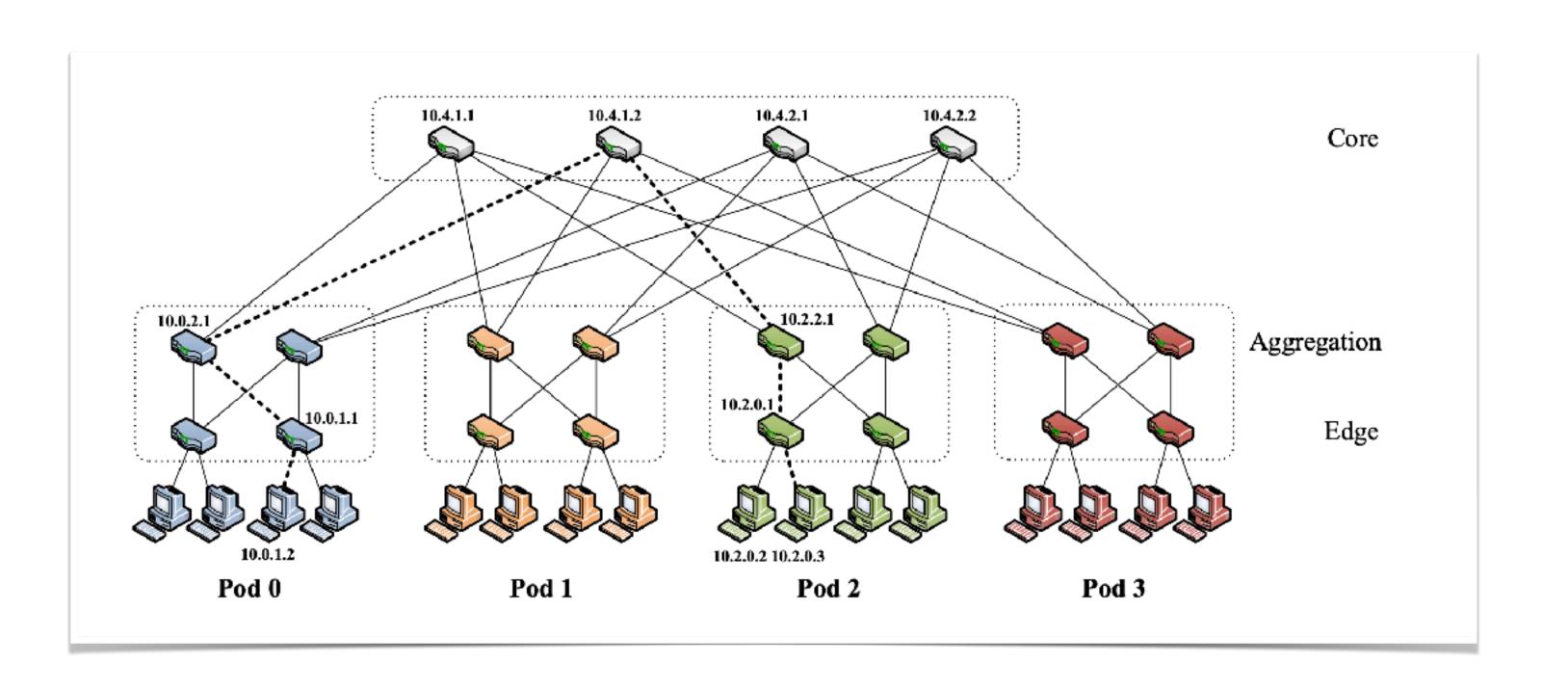
• K-ary fat tree: three-layer topology (edge, aggregation, and core)



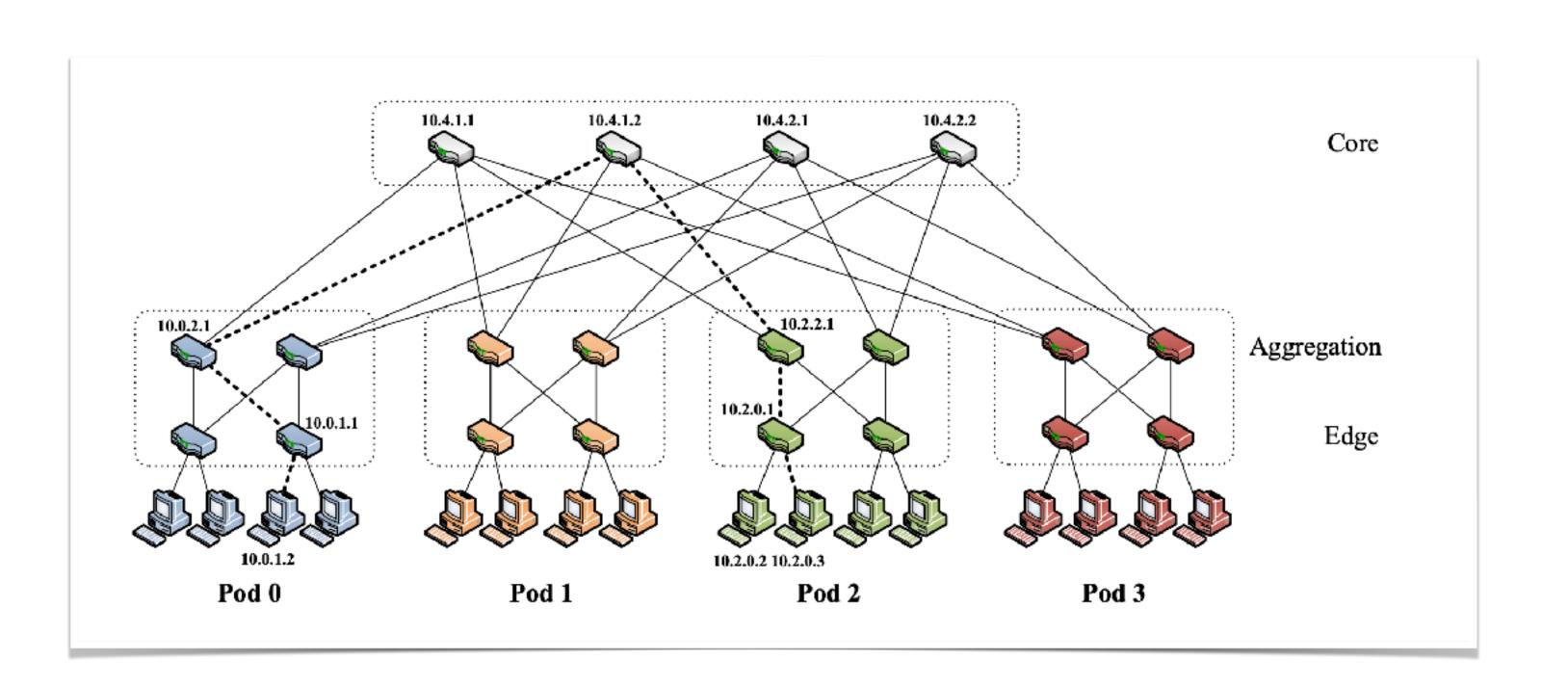
- K-ary fat tree: three-layer topology (edge, aggregation, and core)
 - Each edge switch connects to K/2 servers and K/2 aggregation switches



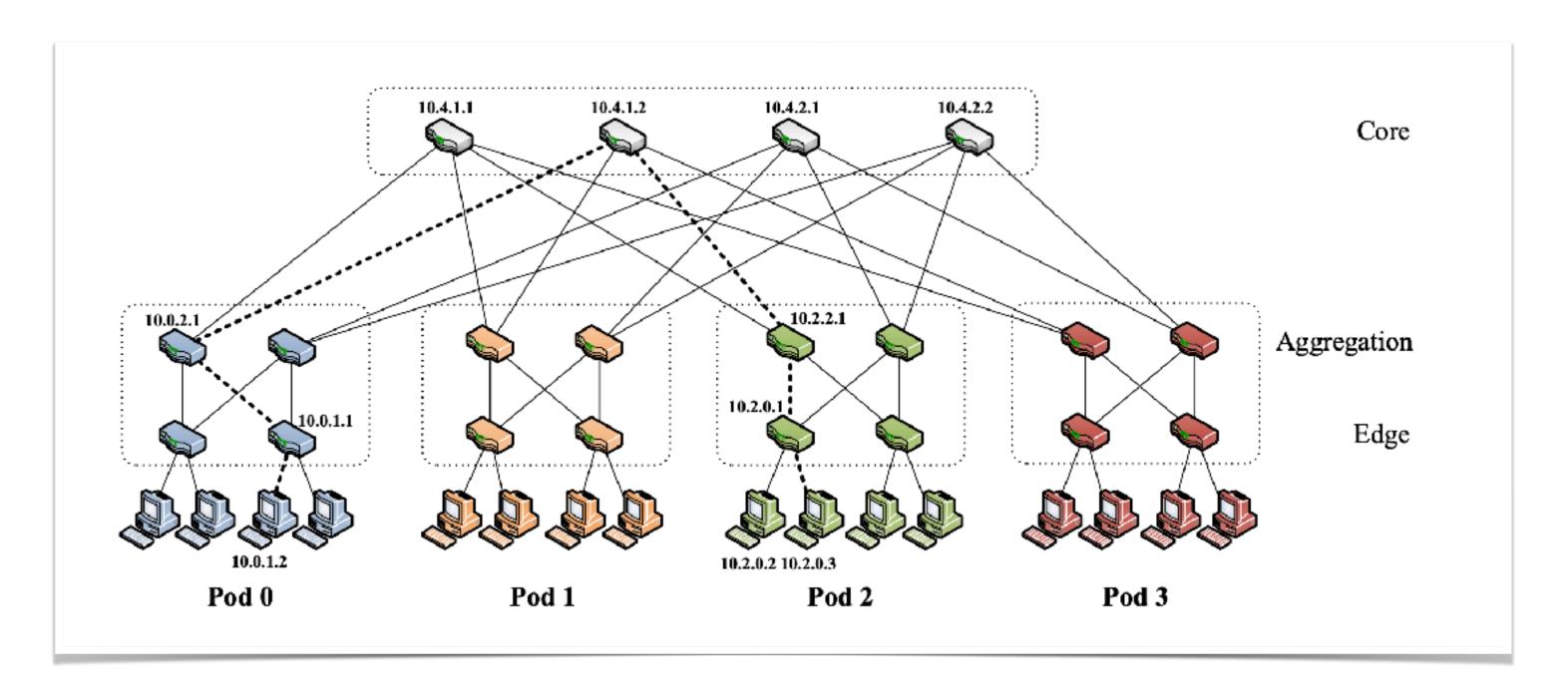
- K-ary fat tree: three-layer topology (edge, aggregation, and core)
 - Each edge switch connects to K/2 servers and K/2 aggregation switches
 - Each aggregation switch connects to K/2 edge and K/2 core switches



- K-ary fat tree: three-layer topology (edge, aggregation, and core)
 - Each edge switch connects to K/2 servers and K/2 aggregation switches
 - Each aggregation switch connects to K/2 edge and K/2 core switches
 - (K/2)^2 cores switches



- K-ary fat tree: three-layer topology (edge, aggregation, and core)
 - Each edge switch connects to K/2 servers and K/2 aggregation switches
 - Each aggregation switch connects to K/2 edge and K/2 core switches
 - (K/2)² cores switches
 - Support K³/4 servers



A Generic Workflow to Construct Fat-Tree Networks

- Step 1: Determine the networking configuration
 - E.g., bandwidth of server NIC and switch port, switching port #
- Step 2: Add intermediate switching stages to match the BW
 - Ingress BW == Egress BW at any switching point (Bandwidth rule)
- Step 3: Apply the scale-out strategy to merge connections
 - Use the Slim and Fat port ratio to decide (Scale-out rule)
- Step 4: Apply the scale-up strategy to add communication paths
 - Added path # relates to switching hops # in the next stage (Scale-up rule)

Summary

- Today
 - Physical connectivity at the rack/cluster scale

- Next lecture
 - Physical connectivity for inter-data centers