



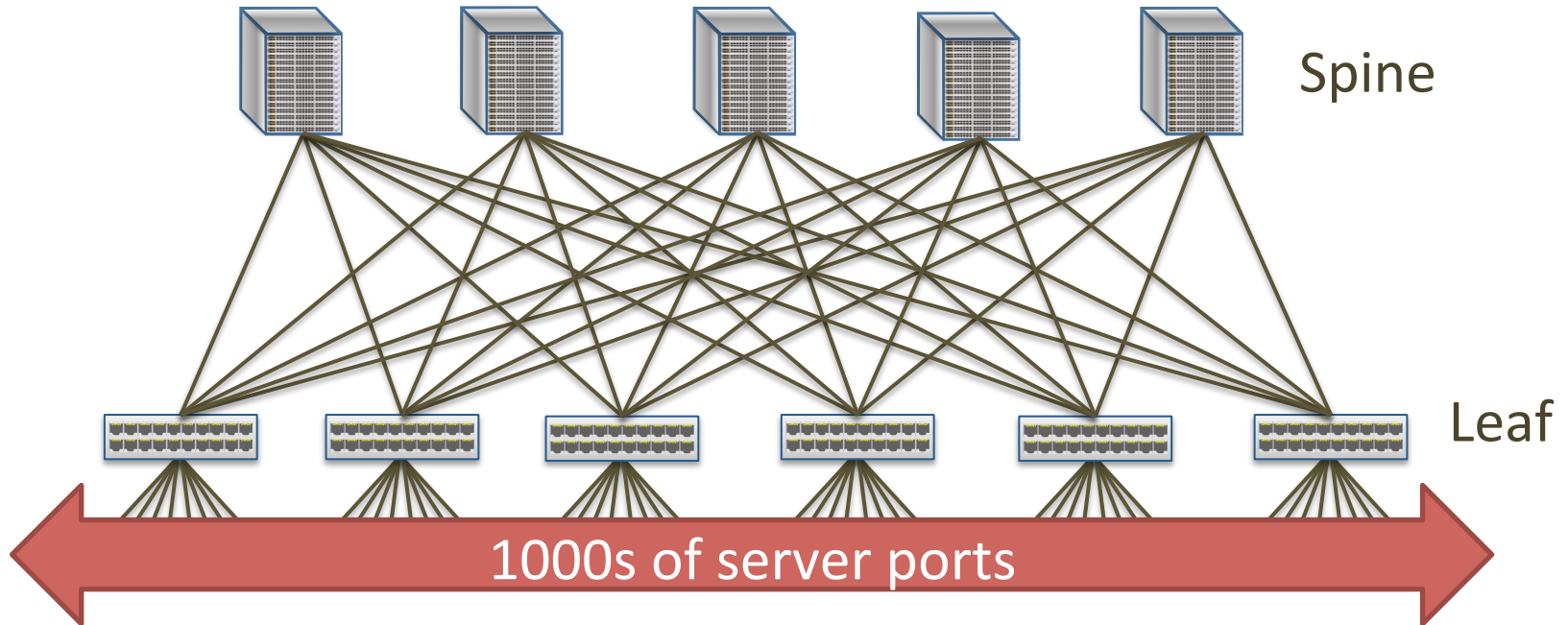
6.888
Lecture 9:
Wireless/Optical Datacenters

Mohammad Alizadeh and Dinesh Bharadia

✧ Many thanks to George Porter (UCSD) and Vyas Sekar (Berkeley)

Spring 2016

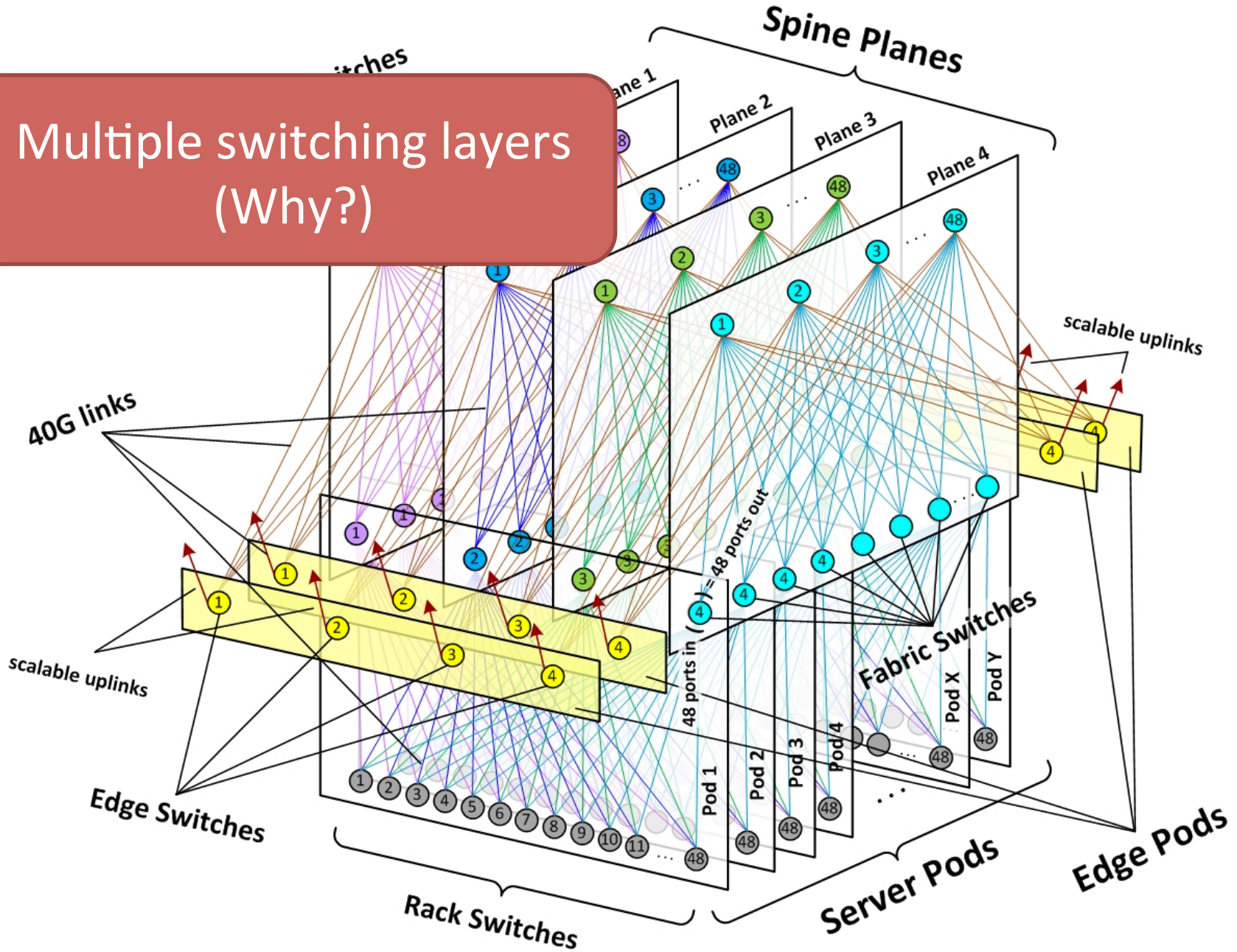
Datacenter Fabrics



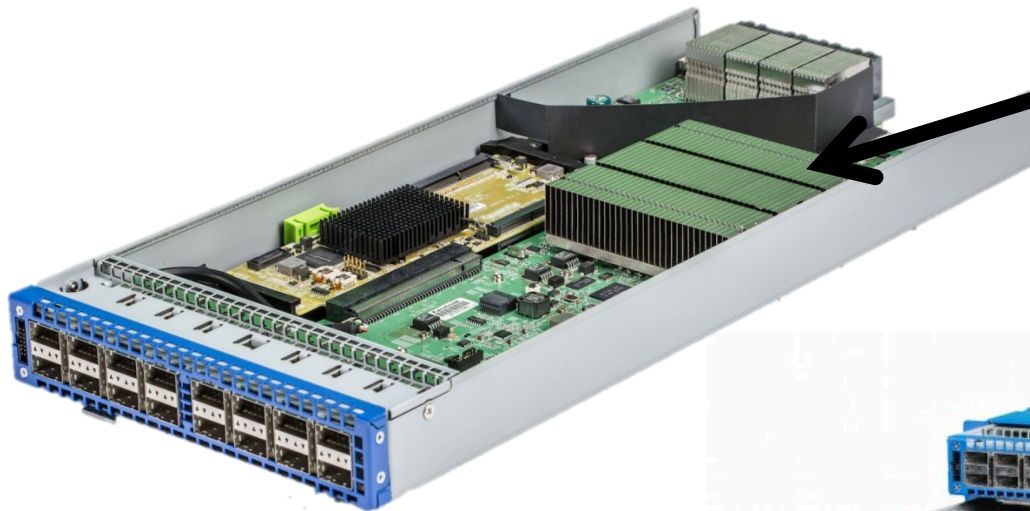
Scale out designs (VL2, Fat-tree)

- Little to no oversubscription
- Cost, power, complexity

Multiple switching layers (Why?)



Building Block: Merchant Silicon Switching Chips



Switch ASIC

6 pack

Facebook Wedge

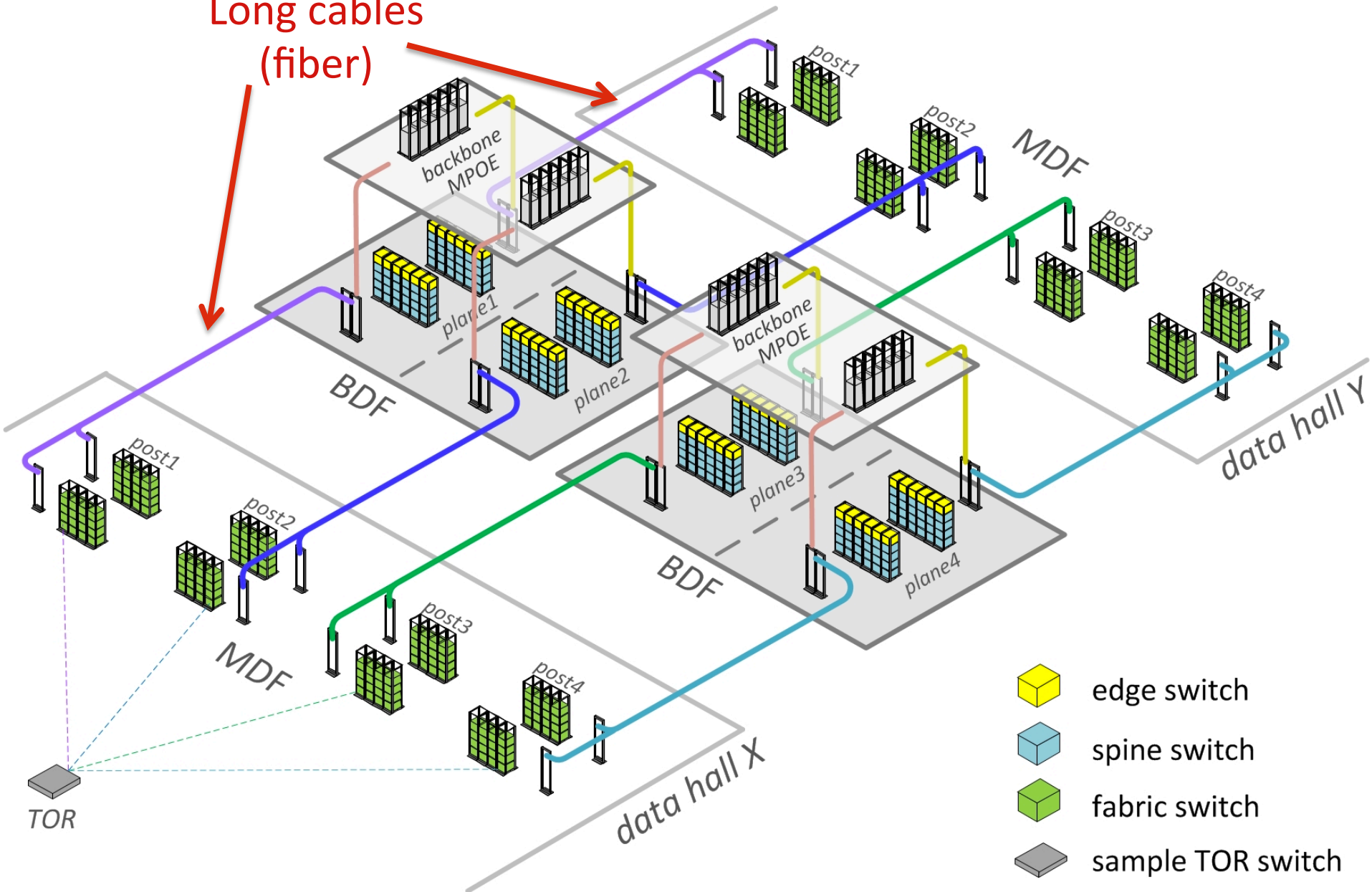
Limited radix: 16x40Gbps

High power: 17 W/port



✧ Image courtesy of Facebook

Long cables
(fiber)



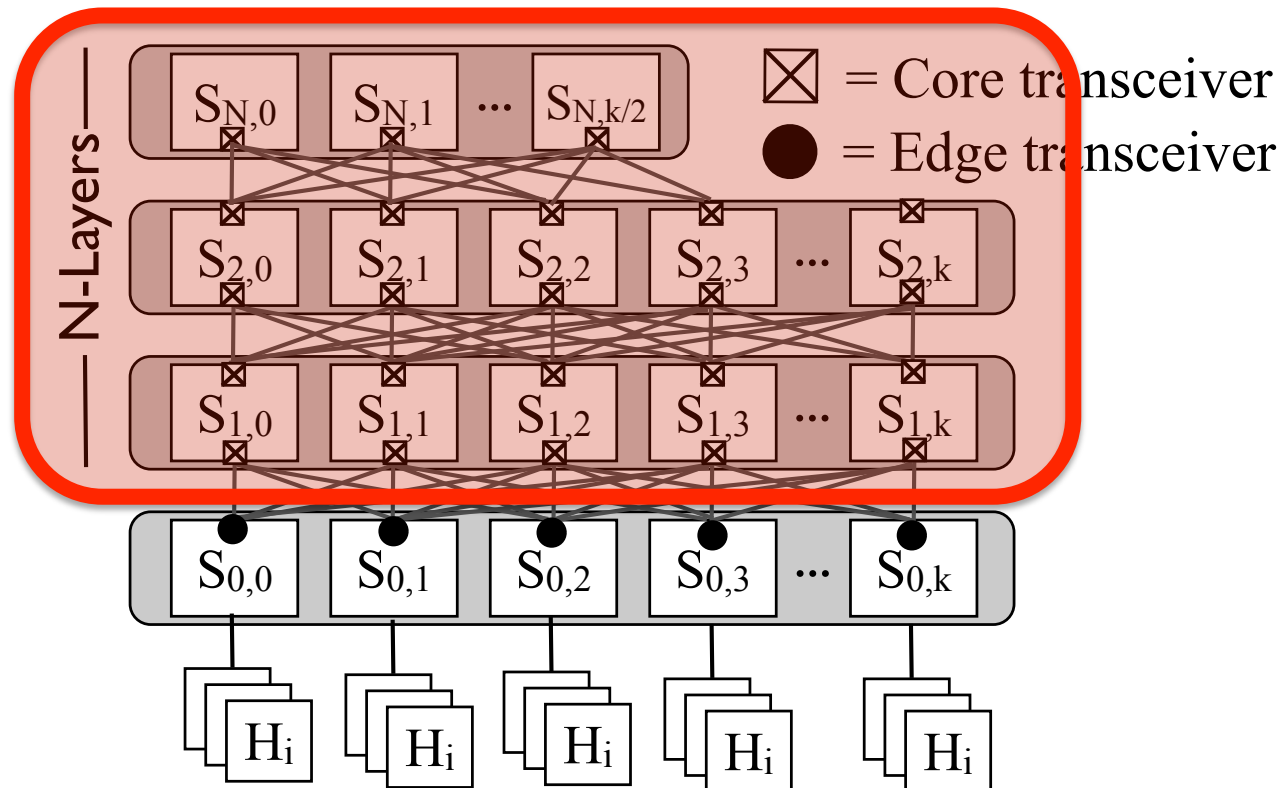
✧ <https://code.facebook.com/posts/360346274145943/introducing-data-center-fabric-the-next-generation-facebook-data-center-network/>

Scale-out packet-switch fabrics

Large number of switches, fibers, **optical transceivers**

Power hungry

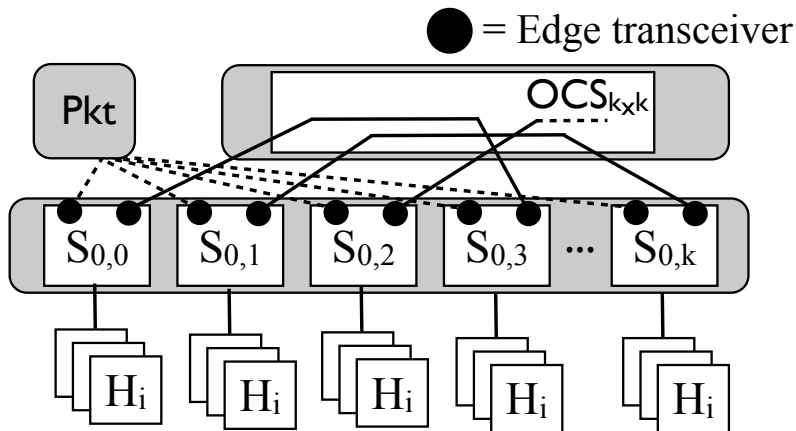
Hard to expand



Beyond Packet-Switched DC Fabrics

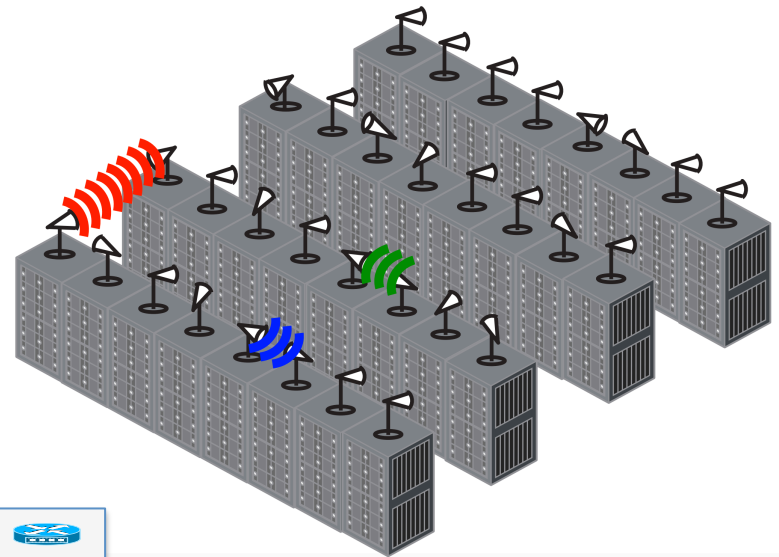
Optical circuit switching

[Helios, cThrough, Mordio, Reactor, ...]



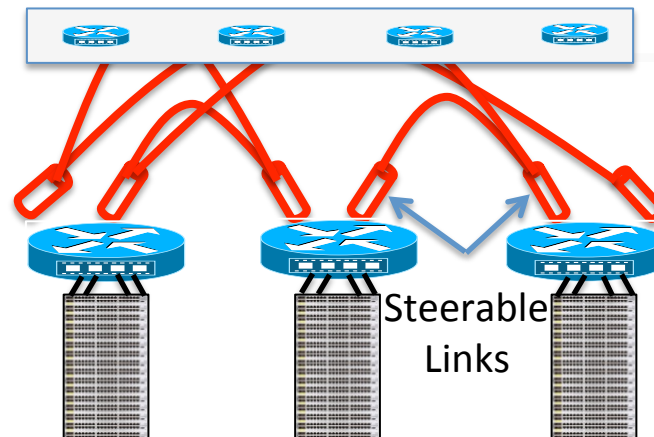
60 GHz RF

[Flyways, MirrorMirror]



Free-space Optics

[FireFly]

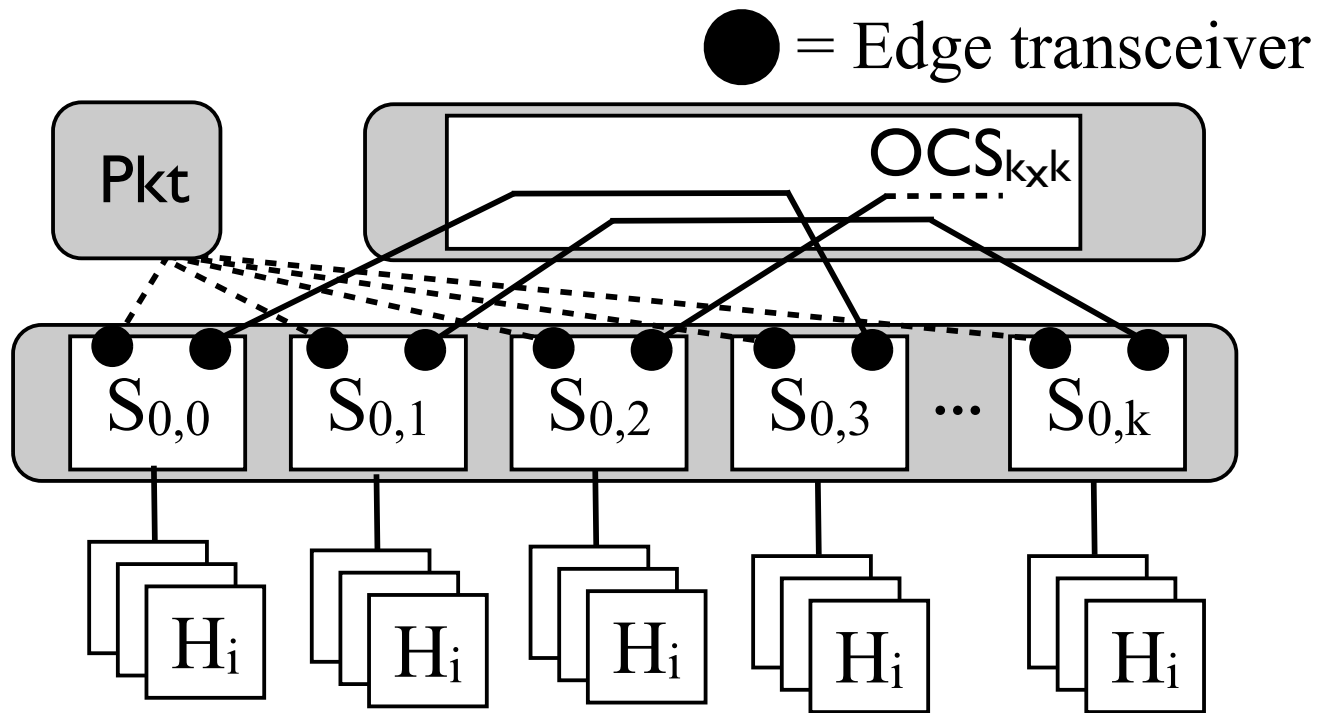


✧ Fig. from presentation by Xia Zhou

Integrating Microsecond Circuit Switching into the Data Center

✧ Slides based on presentation by George Porter (UCSD)

Key idea: *Hybrid Circuit/Packet Networks*



Why build hybrid switch?

Circuit vs. Packet Switching

Observation: Correlated traffic → Circuits

Electrical Packet



\$500/port

10 Gb/s fixed rate

12 W/port

Transceivers (OEO)

Buffering

Per-packet switching

In-band control

Optical Circuit



\$500/port

Rate free (10/40/100/400/+)

240 mW/port

No transceivers

No buffering

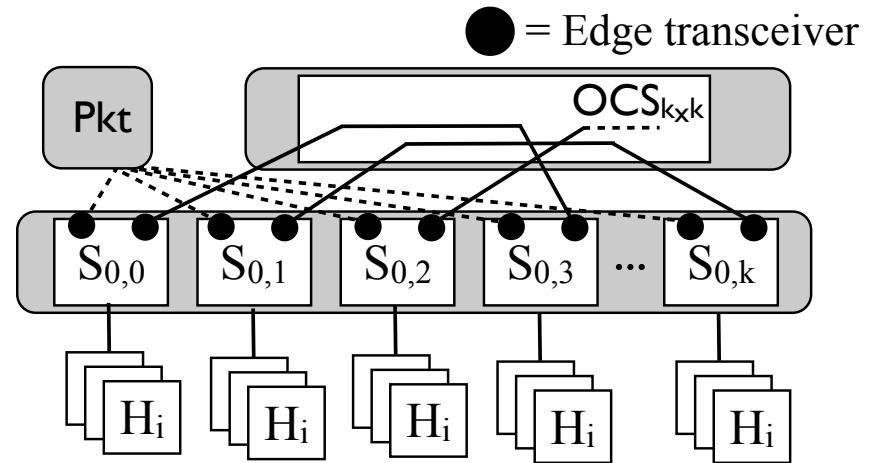
Duty cycle overhead

Out-of-band control

Disadvantages of Circuits

Despite advantages, circuits present different service model:

- Point-to-point connectivity
- Must wait for circuit to be assigned
- Circuit “down” while being reconfigured



affects throughput, latency

*affects network duty cycle;
overall efficiency*

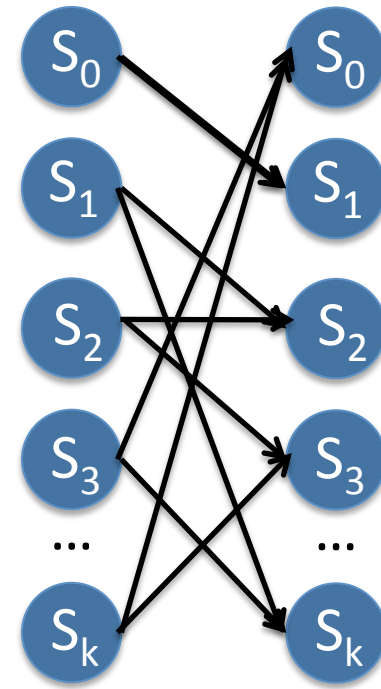
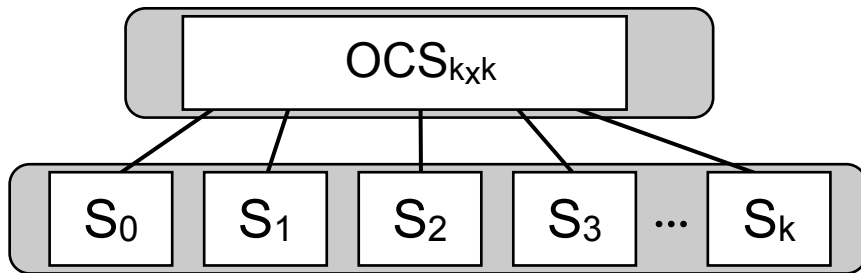
Stability Increases with Aggregation

Inter-Data Center
Inter-Pod
Inter-Rack
Inter-Server
Inter-Process
Inter-Thread

Where is the Sweet Spot?

- 1. Enough Stability**
- 2. Enough Traffic**

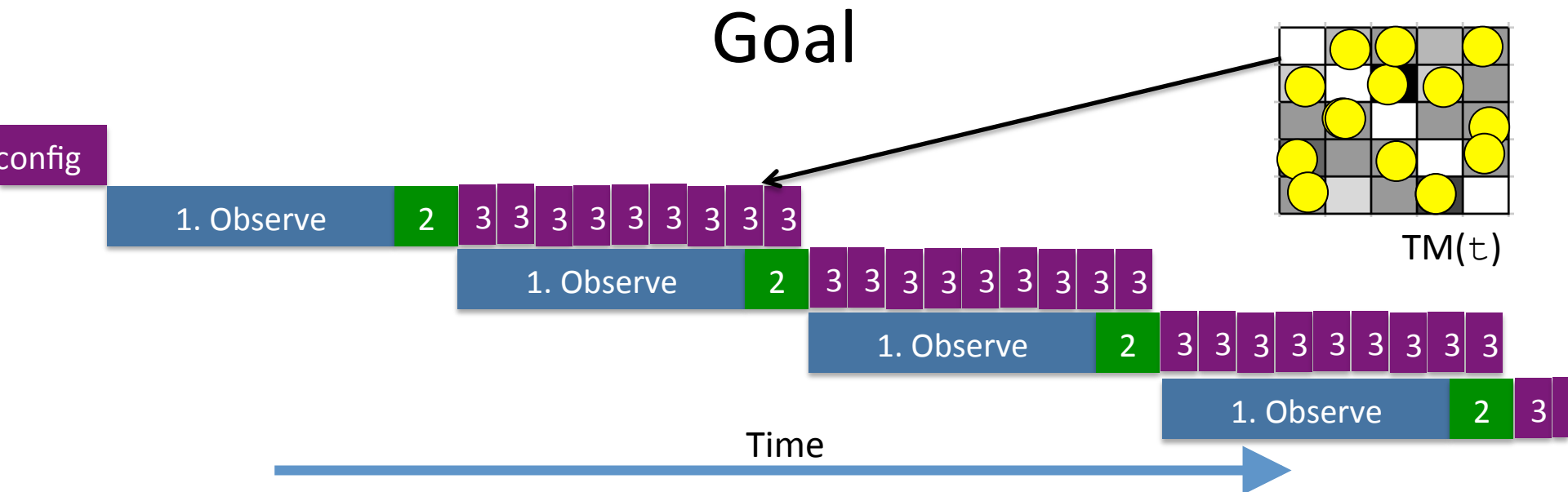
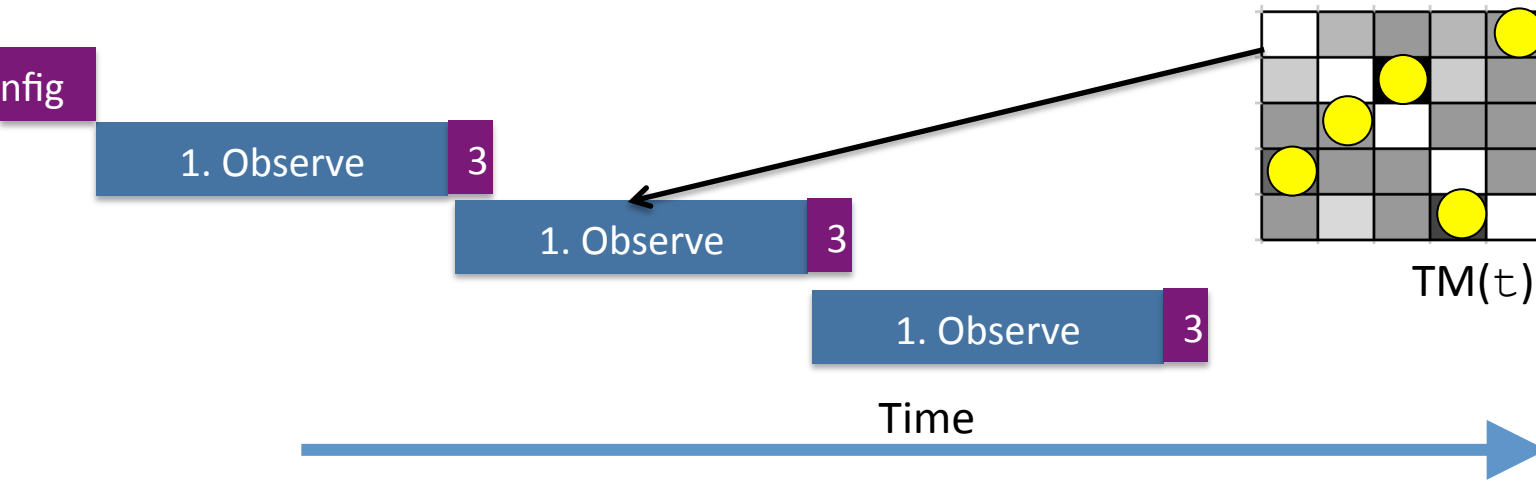
Mordia OCS model



Bi-partite graph

- Directly connects inputs to outputs
- Reconfiguration time: 10us
 - “Night” time (T_n): no traffic during reconfiguration
 - “Day” time (T_d): circuits/mapping established
- Duty cycle: $T_d / (T_d + T_n)$

Limitations of Hotspot Scheduling



Traffic Matrix Scheduling

Step 1. Gather traffic matrix TM

Step 2. Scale TM into TM'

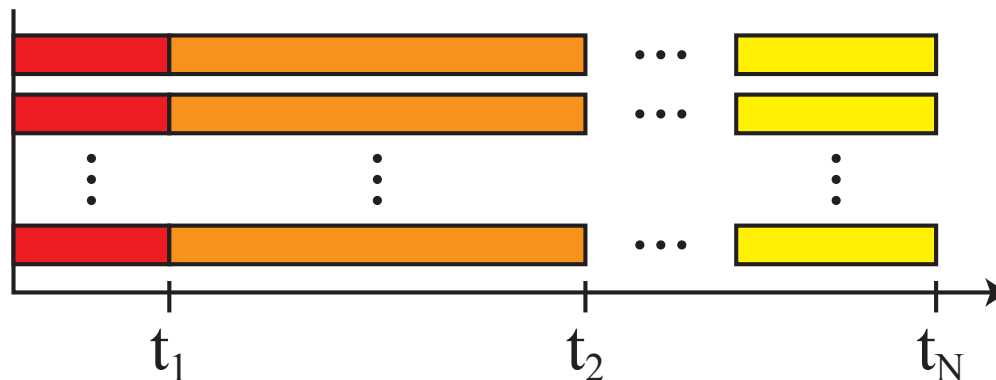


Step 3. Decompose TM' into schedule

Birkhoff von-Neumann
Decomposition

$$t_1 \begin{matrix} P_1 \\ \left[\begin{array}{ccc} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{array} \right] \end{matrix} + t_2 \begin{matrix} P_2 \\ \left[\begin{array}{ccc} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{array} \right] \end{matrix} + \dots + t_N \begin{matrix} P_N \\ \left[\begin{array}{ccc} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{array} \right] \end{matrix}$$

Step 4. Execute schedule in hardware



BvN Decomposition

$\exists (\alpha_1, P_1), (\alpha_2, P_2), \dots, (\alpha_{k'}, P_{k'})$ s.t.

$$T = \alpha_1 P_1 + \alpha_2 P_2 + \dots + \alpha_{k'} P_{k'}$$

T has to be
doubly-stochastic

k' could be large
($\Omega(n^2)$ in worst case)

✧ Suppose: T is a scaled doubly-stochastic matrix

Scheduling

circuit switch configuration: bipartite graph matching

Traffic Matrix: T

$$\begin{bmatrix} \cdot & 1 & \cdot & 4 & \cdot \\ \cdot & \cdot & 1 & \cdot & 4 \\ \cdot & 4 & \cdot & 1 & \cdot \\ 4 & \cdot & \cdot & \cdot & 1 \\ 1 & \cdot & 4 & \cdot & \cdot \end{bmatrix}$$

n = 5 nodes



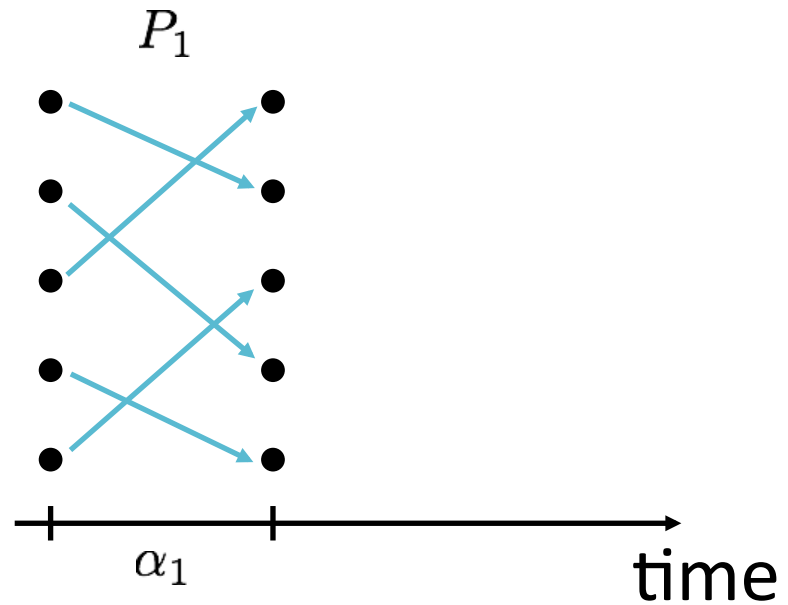
Scheduling

configuration of circuit switch modeled as bipartite graph matching

Traffic Matrix: T

$$\begin{bmatrix} \cdot & 1 & \cdot & 4 & \cdot \\ \cdot & \cdot & 1 & \cdot & 4 \\ \cdot & 4 & \cdot & 1 & \cdot \\ 4 & \cdot & \cdot & \cdot & 1 \\ 1 & \cdot & 4 & \cdot & \cdot \end{bmatrix}$$

n = 5 nodes



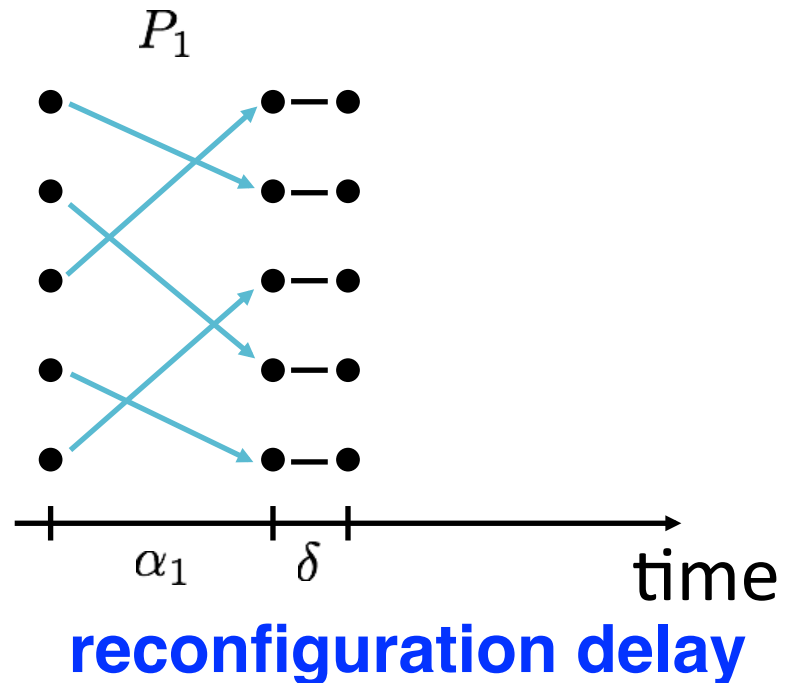
Scheduling

configuration of circuit switch modeled as bipartite graph matching

Traffic Matrix: T

$$\begin{bmatrix} \cdot & 1 & \cdot & 0 & \cdot \\ \cdot & \cdot & 1 & \cdot & 0 \\ \cdot & 0 & \cdot & 1 & \cdot \\ 0 & \cdot & \cdot & \cdot & 1 \\ 1 & \cdot & 0 & \cdot & \cdot \end{bmatrix}$$

n = 5 nodes



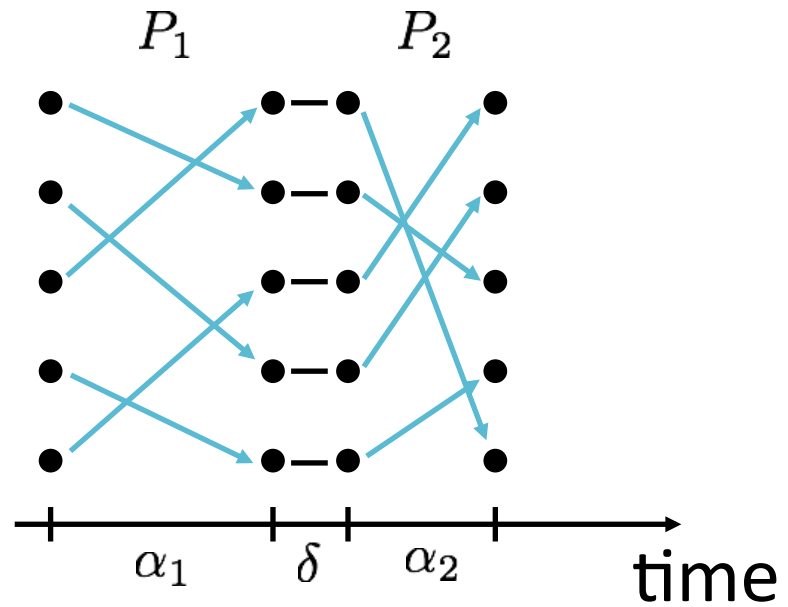
Scheduling

configuration of circuit switch modeled as bipartite graph matching

Traffic Matrix: T

$$\begin{bmatrix} \cdot & \mathbf{1} & \cdot & 0 & \cdot \\ \cdot & \cdot & \mathbf{1} & \cdot & 0 \\ \cdot & 0 & \cdot & \mathbf{1} & \cdot \\ 0 & \cdot & \cdot & \cdot & \mathbf{1} \\ \mathbf{1} & \cdot & 0 & \cdot & \cdot \end{bmatrix}$$

n = 5 nodes



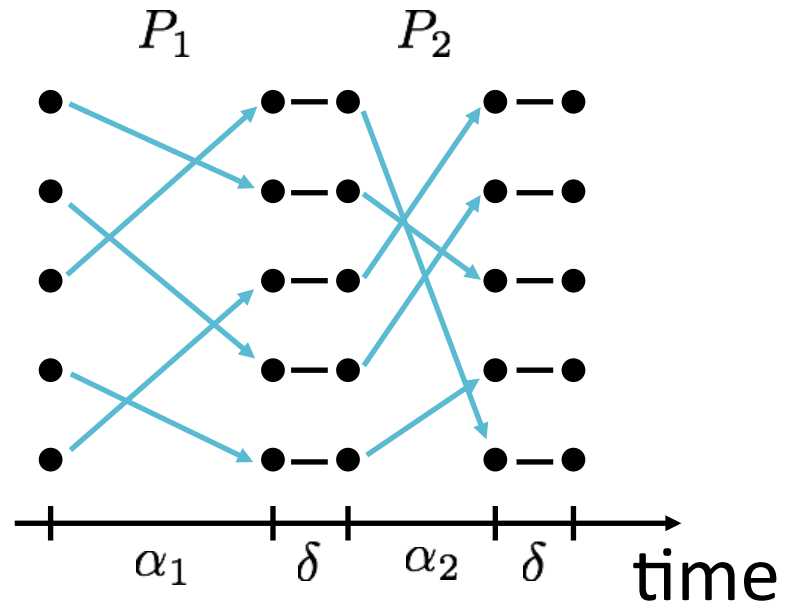
Scheduling

configuration of circuit switch modeled as bipartite graph matching

Traffic Matrix: T

$$\begin{bmatrix} \cdot & 0 & \cdot & 0 & \cdot \\ \cdot & \cdot & 0 & \cdot & 0 \\ \cdot & 0 & \cdot & 0 & \cdot \\ 0 & \cdot & \cdot & \cdot & 0 \\ 0 & \cdot & 0 & \cdot & \cdot \end{bmatrix}$$

n = 5 nodes



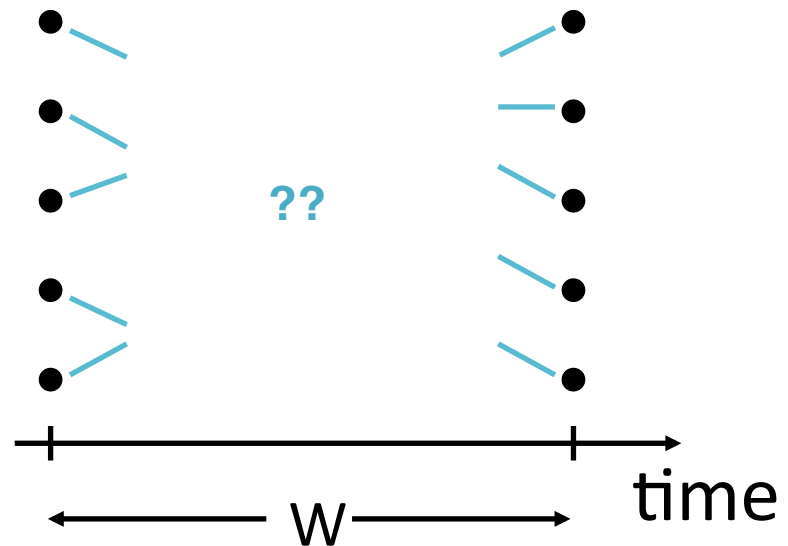
Scheduling

maximize throughput in time-window W

Traffic Matrix: T

$$\begin{bmatrix} \cdot & 1 & \cdot & 4 & \cdot \\ \cdot & \cdot & 1 & \cdot & 4 \\ \cdot & 4 & \cdot & 1 & \cdot \\ 4 & \cdot & \cdot & \cdot & 1 \\ 1 & \cdot & 4 & \cdot & \cdot \end{bmatrix}$$

$n = 5$ nodes



Problem Statement

$$\text{maximize } \left\| \min \left(\sum_{i=1}^k \alpha_i P_i, T \right) \right\|_1$$

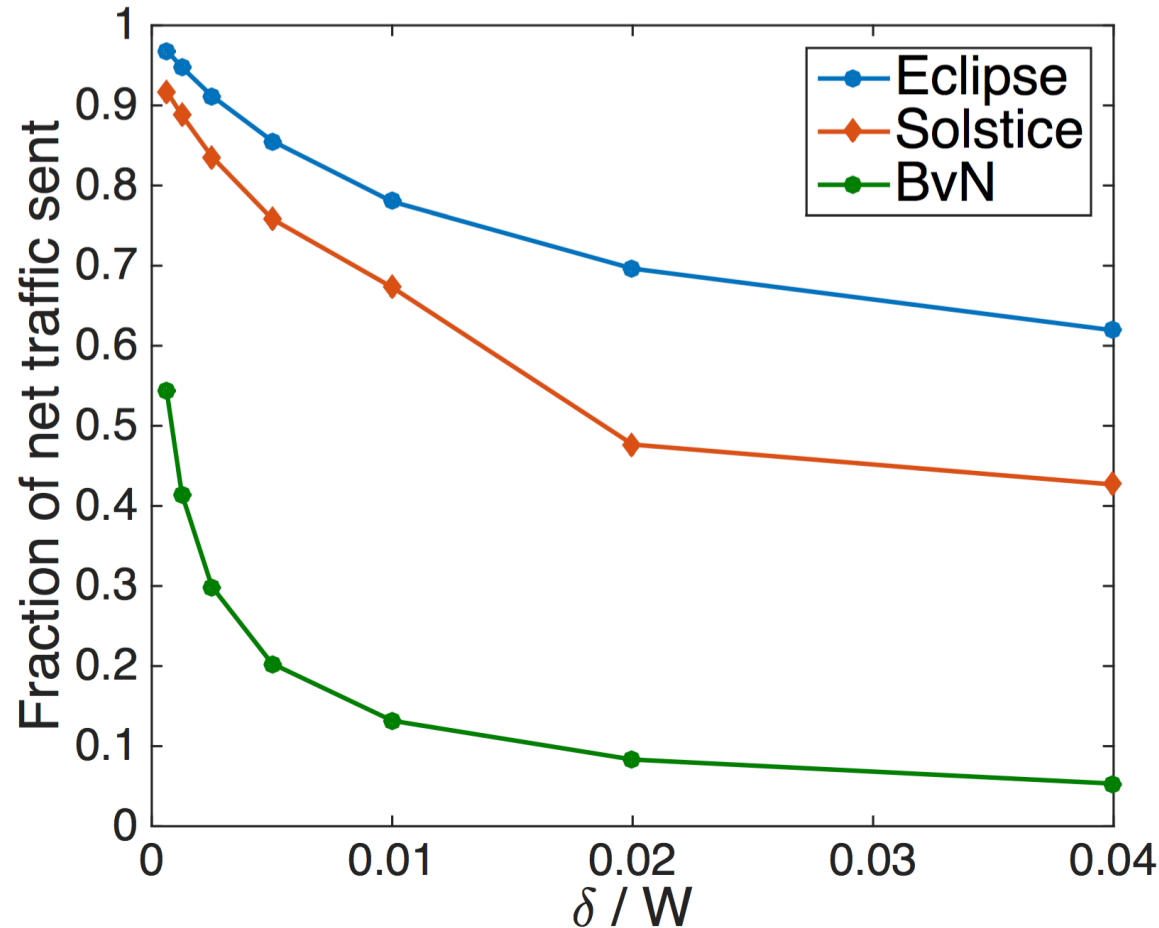
$$\text{s.t. } \alpha_1 + \alpha_2 + \dots + \alpha_k + k\delta \leq W$$

number of matchings $\longrightarrow k \in \mathbb{N}$

permutation matrices $\longrightarrow P_1, \dots, P_k \in \mathcal{P}$

duration $\longrightarrow \alpha_1, \dots, \alpha_k \geq 0$

Eclipse: Greedy Algorithm (with provable guarantees)



✧ Venkatakrisnan et al., “Costly Circuits, Submodular Schedules, Hybrid Switch Scheduling for Data Centers”, To appear in SIGMETRICS 2016.

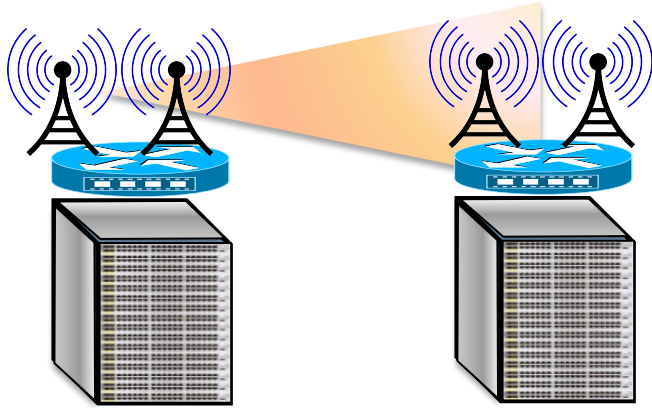
Discussion

Firefly

✧ Slides based on presentation by Vyas Sekar (CMU)

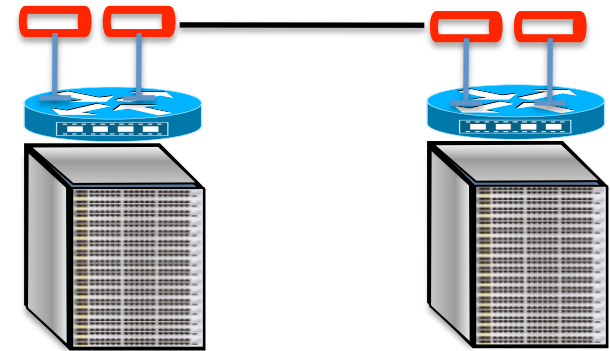
Why FSO instead of RF?

RF (e.g. 60GHz)



Wide beam →
Faster steering of beams
High interference
Limited active links
Limited Throughput

FSO (Free Space optical)



Narrow beam →
Slow steering of beams
Zero interference
No limit on active links
High Throughput

Today's FSO



Cost: \$15K per FSO

Size: 3 ft³

Power: 30w

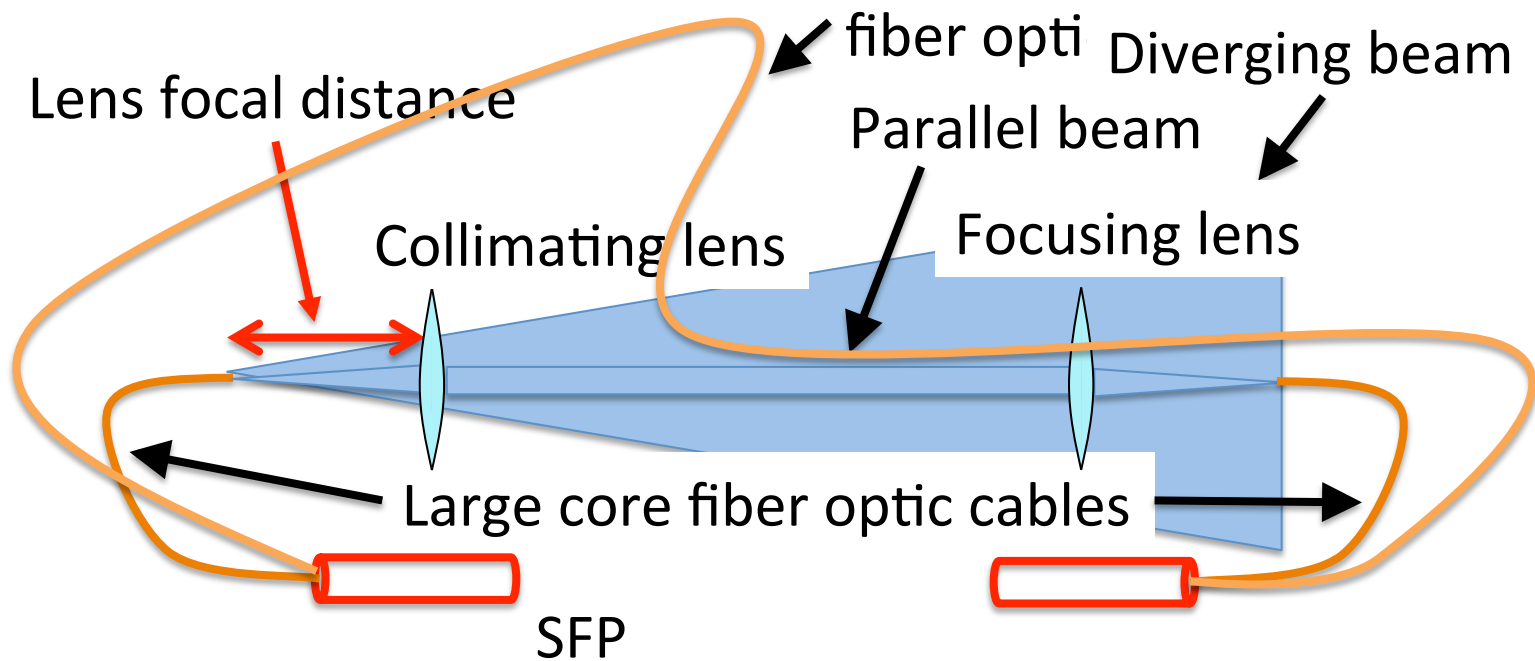
Non steerable

- Current: bulky, power-hungry, and expensive
- Required: small, low power and low expense

Why Size, Cost, Power Can be Reduced?

- Traditional use : outdoor, long haul
 - High power
 - Weatherproof
- Data centers: indoor, short haul
- Feasible roadmap via commodity fiber optics
 - E.g. Small form transceivers (Optical SFP)

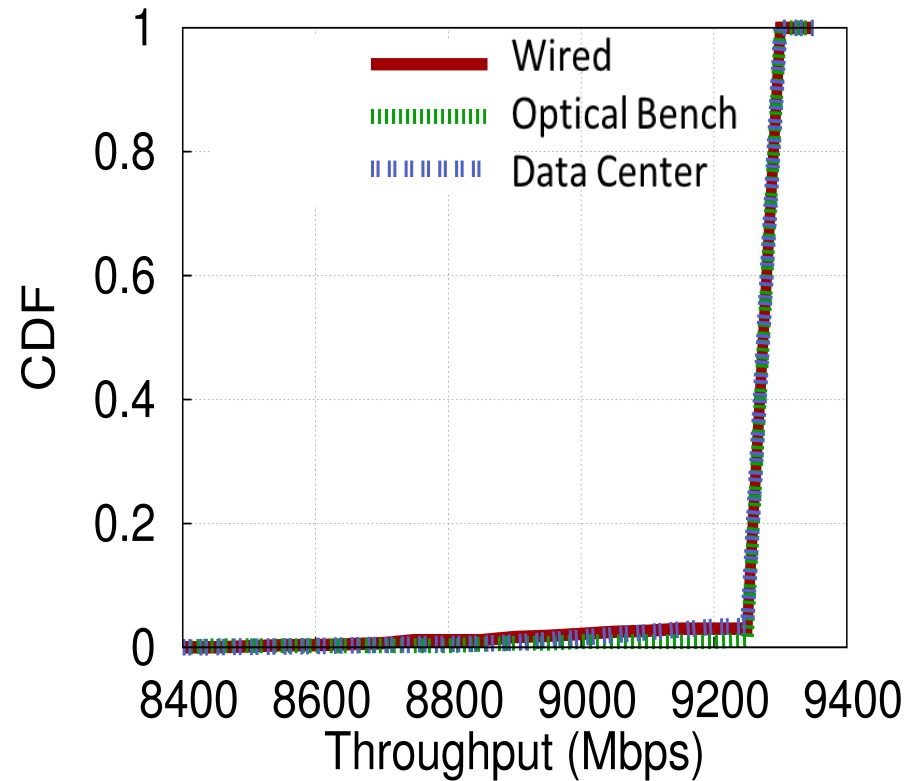
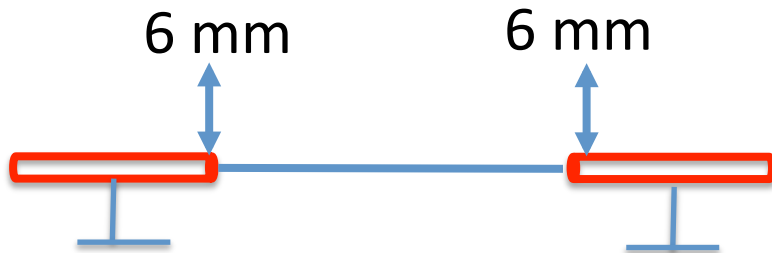
FSO Design Overview



- large cores (> 125 microns) are more robust

FSO Link Performance

Effect of vibrations, etc.
6mm movement tolerance
Range up to 24m tested



FSO link is as robust as a wired link

Steerability

Shortcomings of current FSOs

✓ Cost

✓ Size

✓ Power

• Not Steerable



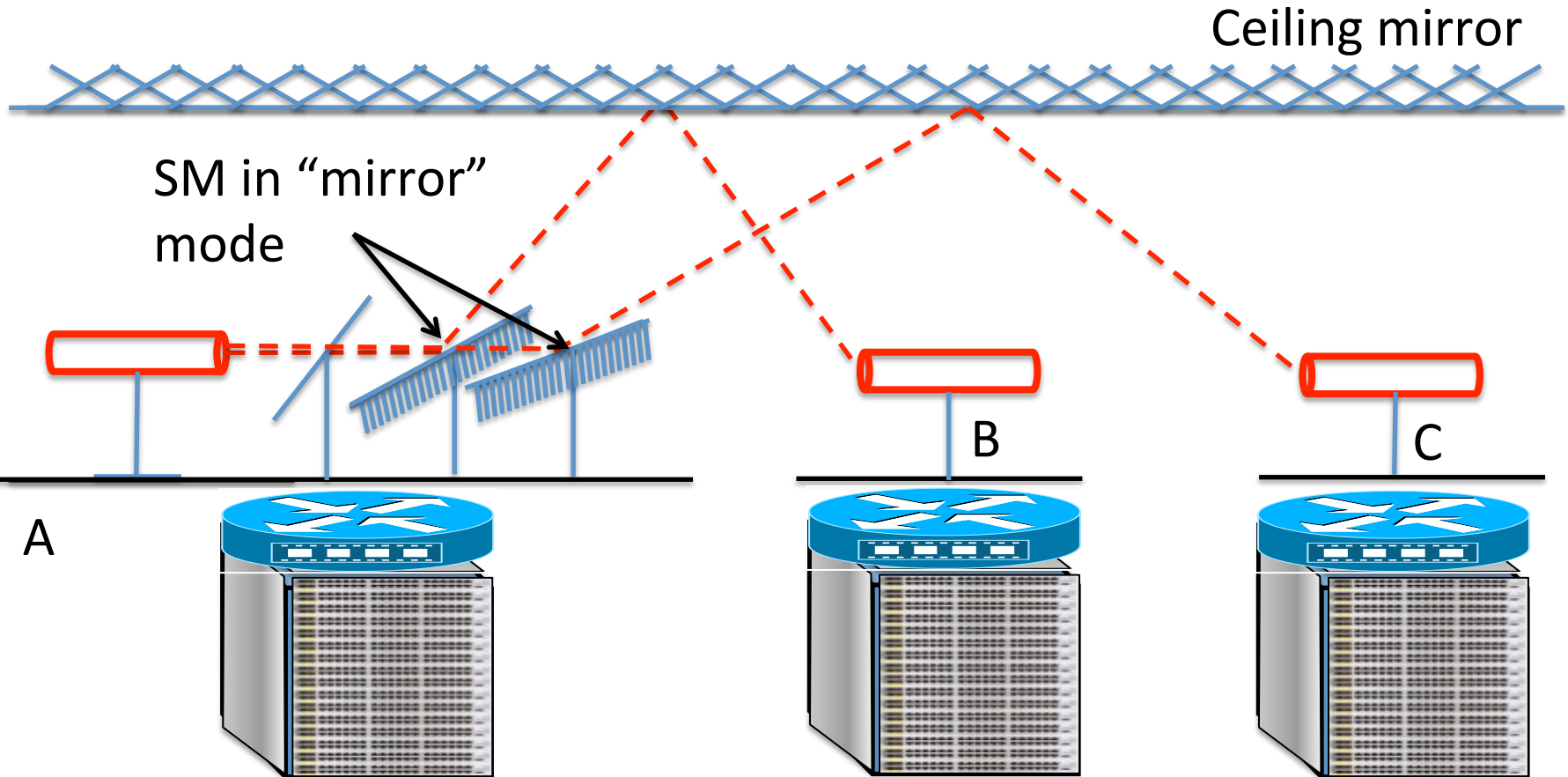
FSO design
using SFP

Via Switchable mirrors
or Galvo mirrors

Shortcomings of current FSOs

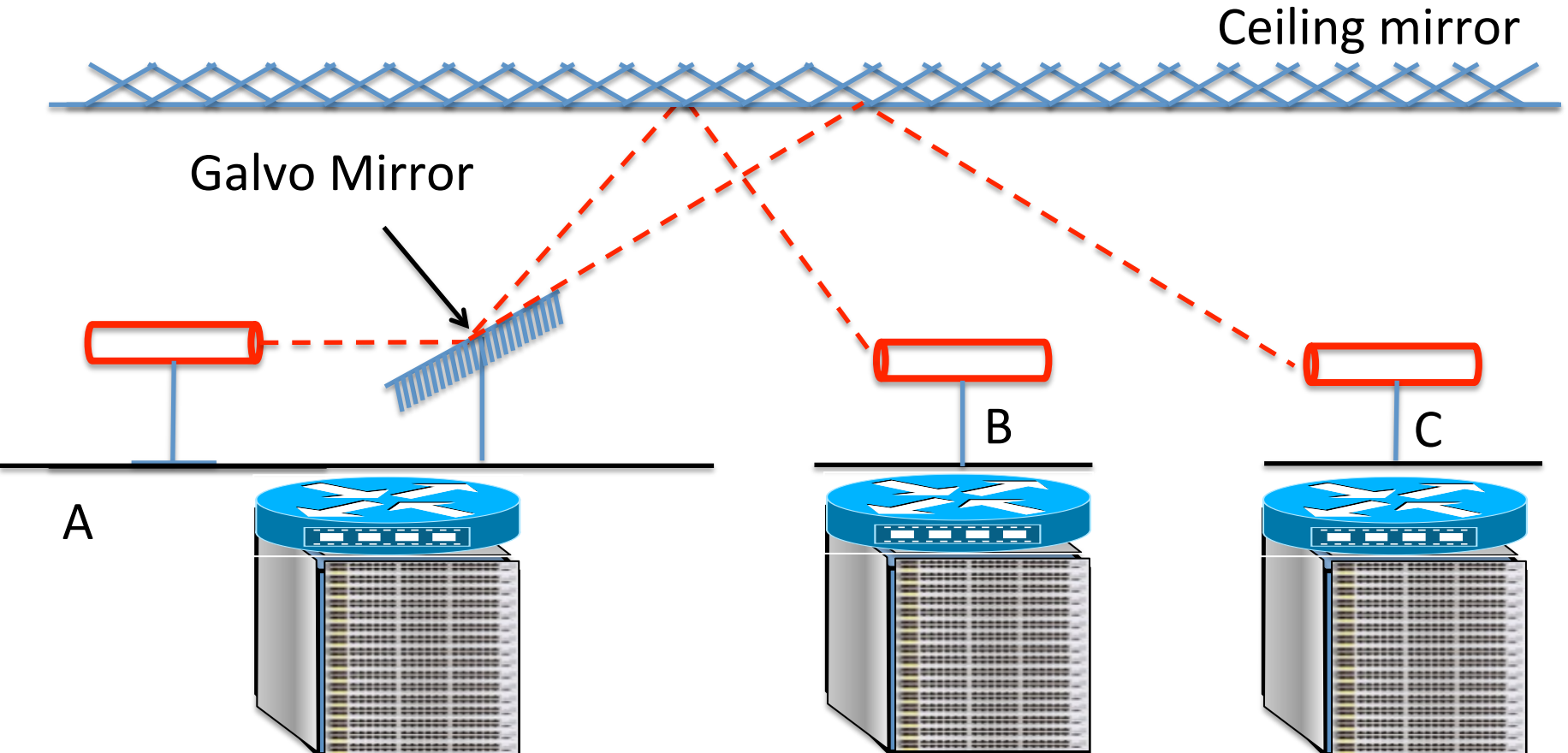
Steerability via Switchable Mirror

- Switchable Mirror: glass \leftrightarrow mirror
- Electronic control, low latency



Steerability via Galvo Mirror

- Galvo Mirror: small rotating mirror
- Very low latency



How to design FireFly network?

Goals: Robustness to current and future traffic

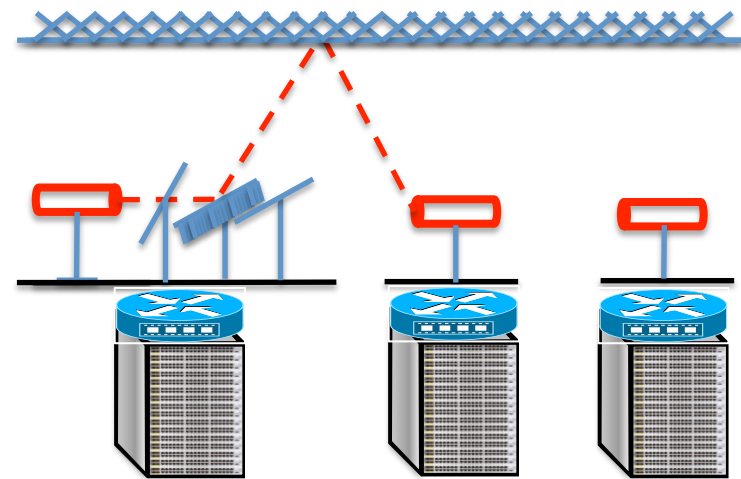
Budget & Physical Constraints

Design parameters

- Number of FSOs?
- Number of steering mirrors?
- Initial mirrors' configuration

Performance metric

- Dynamic bisection bandwidth



Discussion

Next Time: Rack-Scale Computing

R2C2: A Network Stack for Rack-scale Computers

Paolo Costa Hitesh Ballani Kaveh Razavi* Ian Kash
Microsoft Research

ABSTRACT

Rack-scale computers, comprising a large number of micro-servers connected by a direct-connect topology, are expected to replace servers as the building block in data centers. We focus on the problem of routing and congestion control across the rack's network, and find that high path diversity in rack topologies, in combination with workload diversity across it, means that traditional solutions are inadequate.

We introduce R2C2, a network stack for rack-scale computers that provides flexible and efficient routing and congestion control. R2C2 leverages the fact that the scale of rack topologies allows for low-overhead broadcasting to ensure that all nodes in the rack are aware of all network flows. We thus achieve rate-based congestion control without any probing; each node independently determines the sending rate for its flows while respecting the provider's allocation policies. For routing, nodes dynamically choose the routing protocol for each flow in order to maximize overall utility. Through a prototype deployed across a rack emulation platform and a packet-level simulator, we show that R2C2 achieves very low queuing and high throughput for diverse and bursty workloads, and that routing flexibility can provide significant throughput gains.

CCS Concepts

•Networks → Data center networks; Transport protocols; Cloud computing;

1. INTRODUCTION

While today's large-scale data centers such as those run by Amazon, Google, and Microsoft are built using commodity off-the-shelf servers, recently there has been an increasing trend towards server customization to reduce costs and improve performance [50, 54, 55, 58]. One such trend is the advent of "rack-scale computing". We use this term to refer to emerging architectures that propose servers or *rack-scale computers* comprising a large number of tightly integrated systems-on-chip, interconnected by a network fabric. This design enables thousands of cores per rack and provides high bandwidth for rack-scale applications. The consequent power, density and performance benefits means that racks are expected to replace individual servers as the basic building block of datacenters. Early examples of rack-scale computers include commercial (HP Moonshot [56], AMD SeaMicro [62], Boston Viridis [51], and Intel RSA [26, 59]) as well as research platforms [7, 9, 19, 34, 38].

A design choice that allows rack-scale computers to achieve high internal bandwidth and high density is to move away from a switched network fabric to a "distributed switch" architecture where each node functions as a small switch and forwards traffic from other nodes. This underlies many existing designs [19, 34, 38, 47, 51, 56, 59, 62], and results in a multi-hop direct-connect topology, with very high path diversity. This is a departure from today's data centers, which mostly use tree-like topologies. While direct-connect topologies have been used in high performance computing

