# INTRODUCING OPTICAL SWITCHING INTO DATACENTER NETWORKS

George Porter (on behalf of many co-authors!)
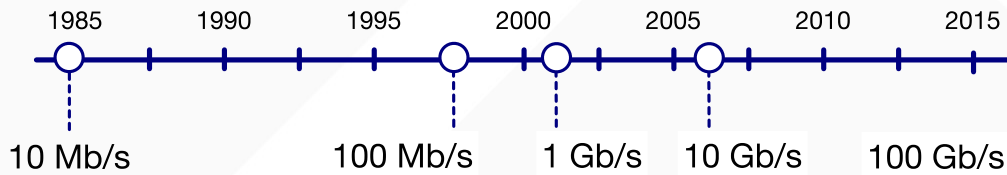
UC San Diego

University of Washington
Feb 8, 2021

# GROWTH OF "HYPERSCALE" DATACENTERS

Network problem:
connecting >100,000 servers

| 1985 | 1990 | 1995 | 2000 | 2005 | 2010 | 2015 |

10 Mb/s     100 Mb/s   1 Gb/s   10 Gb/s   100 Gb/s

Can't buy
sufficiently fast

Petabit line card
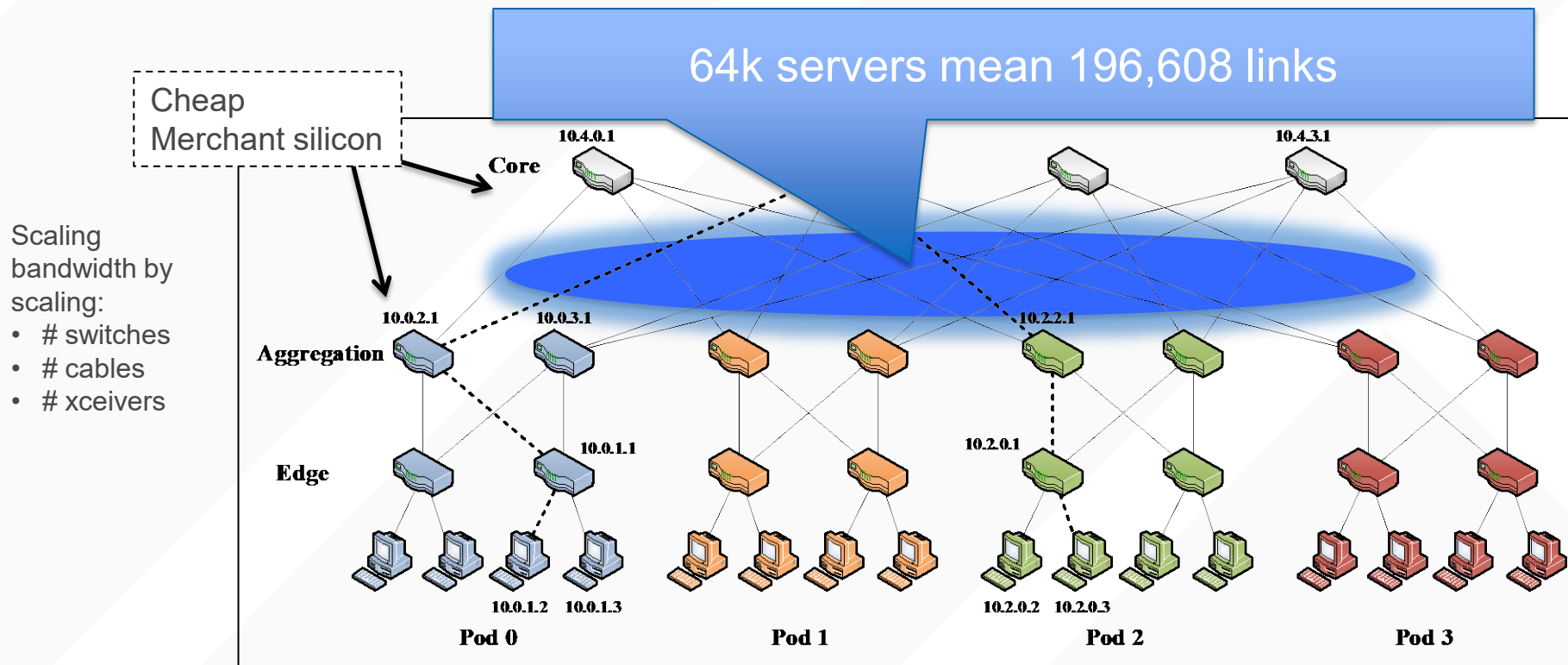Search

Your search - **Petabit line card** - did not match any documents.

Did you mean: *Gigabit* line card
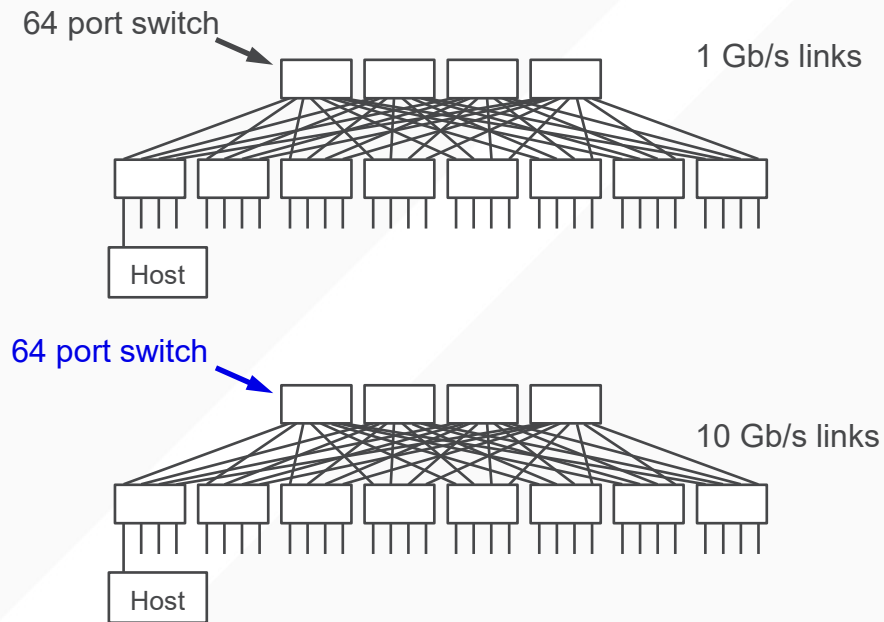
100,000 x 10 Gb/s = 1 Pb/s

# 2009: RISE OF "SCALE OUT" NETWORKS
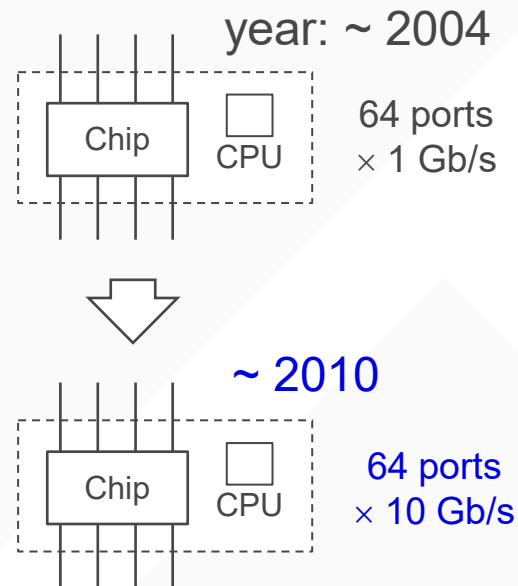
64k servers mean 196,608 links
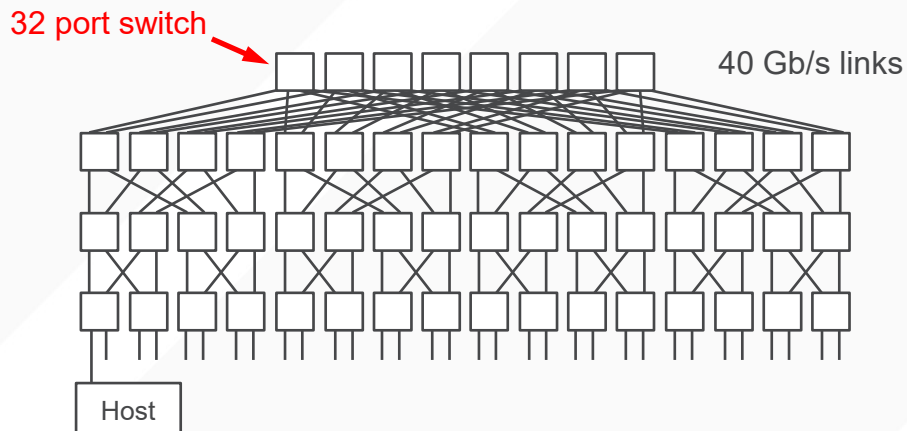
Cheap
Merchant silicon

Scaling bandwidth by scaling:
- # switches
- # cables
- # xceivers

Core

10.4.0.1          10.4.3.1

Aggregation

10.0.2.1     10.0.3.1          10.2.2.1

Edge

10.0.1.1

10.2.0.1

10.0.1.2   10.0.1.3          10.2.0.2   10.2.0.3

Pod 0          Pod 1          Pod 2          Pod 3

*M. Al-Fares, A. Loukissas, and A. Vahdat, "A Scalable, Commodity Data Center Network Architecture," SIGCOMM 2008.*

# SCALING "TRADITIONAL" FATTREES IS BECOMING EXPENSIVE

64 port switch

1 Gb/s links

Host

64 port switch

10 Gb/s links

Host

Merchant Silicon

year: ~ 2004

Chip    CPU

64 ports
× 1 Gb/s

~ 2010

Chip    CPU

64 ports
× 10 Gb/s
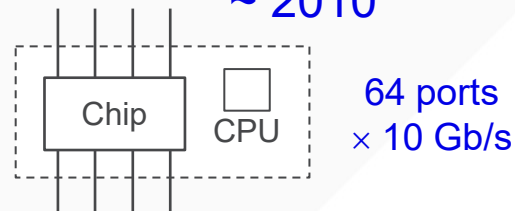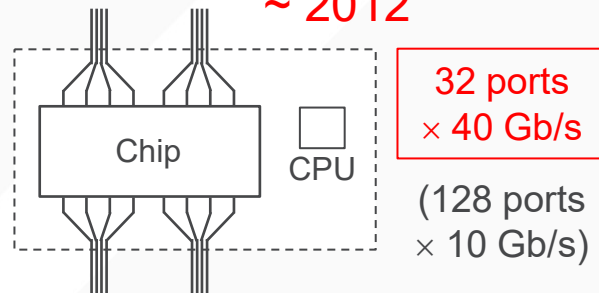
# SCALING "TRADITIONAL" FATTREES IS BECOMING EXPENSIVE
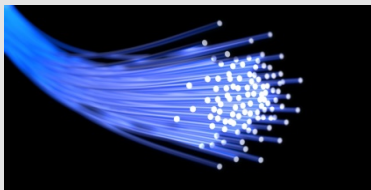
# PROLIFERATION OF FAT TREE LAYERS A GROWING PROBLEM

## Optical links

1,000 + Gb/s – 1,000 + meters

10 Gb/s – 2,000 meters

Single mode fiber

SFP+ transceiver

## Electrical links

1 Gb/s – 100m

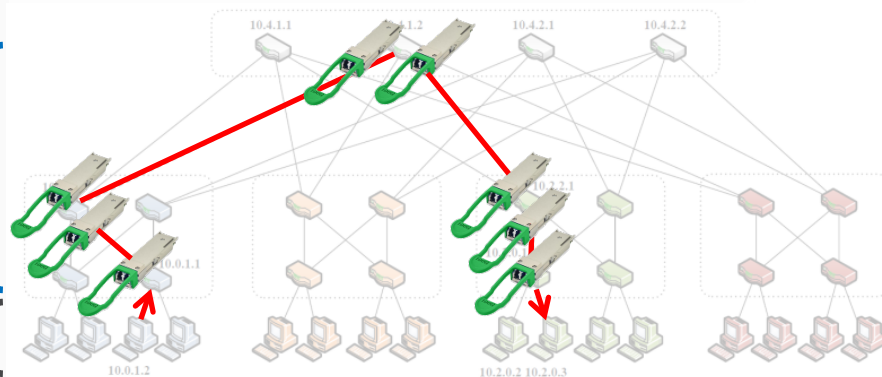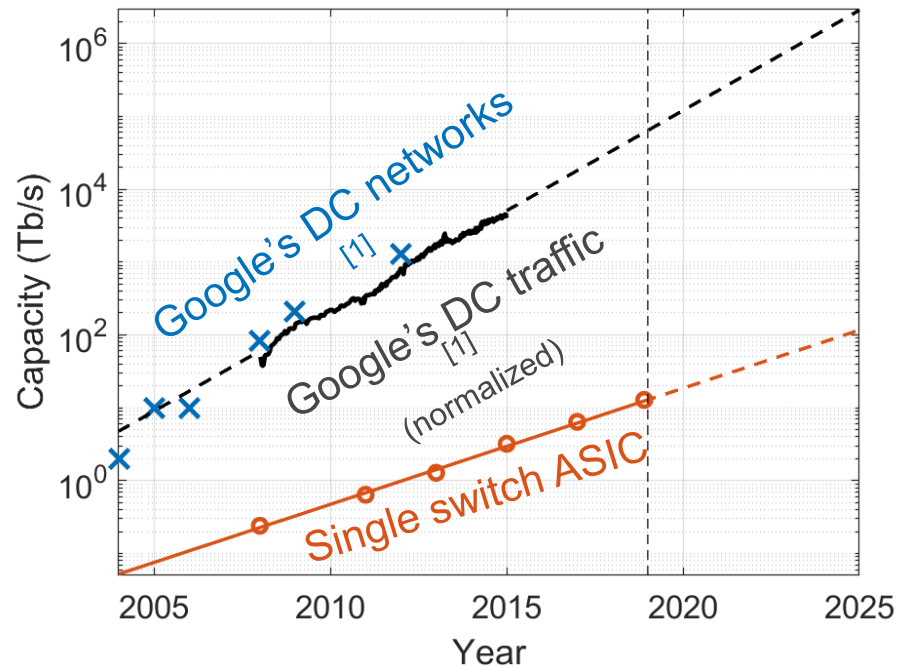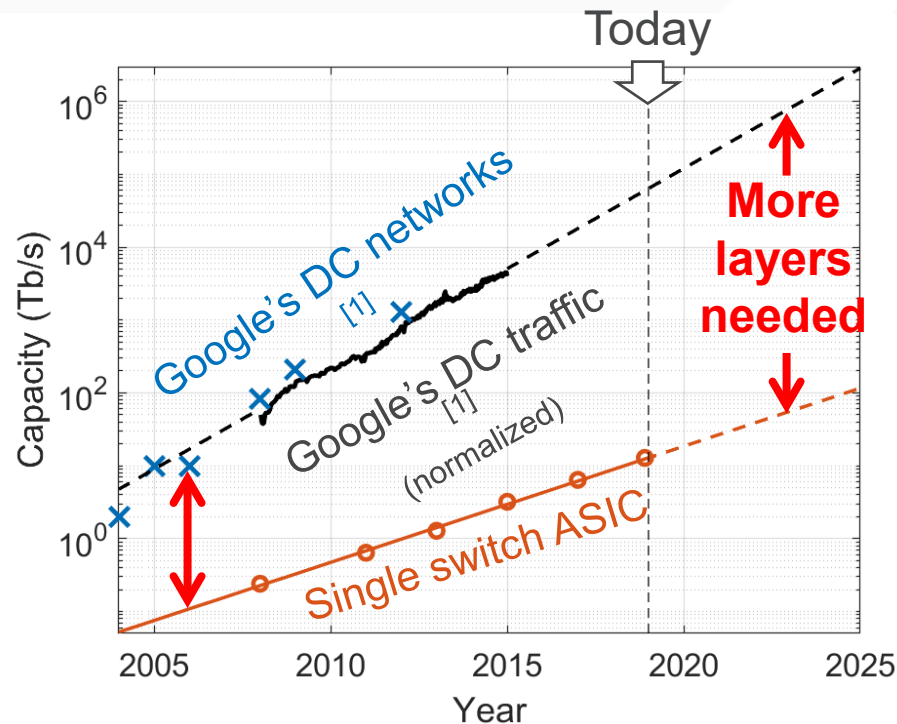10 Gb/s – 10 meters

CAT 5

10G DAC

## Datacenter Network

For every device attached to the network, there are multiple transceivers in the network

100k nodes: O(100kW) and O($$$)
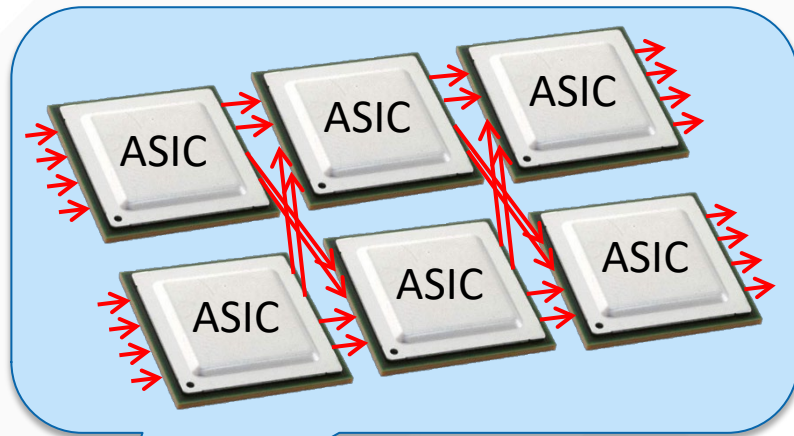
# SCALING LIMITATIONS OF CMOS-BASED PACKET SWITCH CHIPS

- Increasing difficulty getting data in/out of the chip

- Divergence between *link* rate and *channel* rate

  - E.g. 100G vs 4x25G

- More fabric layers = higher cost & power



"Hiding" layers

| | 0.64 TB/s | 5.12 TB/s | 12.8 TB/s |
|---|---|---|---|
| Max. chip radix = | **64** x 10G | **128** x 10G | **128** x 25G |
| Used chip radix = | **16** x 40G | **32** x 40G | **32** x 100G |

# MOVE TO "CHASSIS" BASED FAT TREES (FACEBOOK, GOOGLE)

"Traditional" packet switch

Multistage chassis switch

Fully-provisioned network — 8,192 end hosts

| Architecture | # Tiers | # Hops | # Transceivers | # Switch chips | # Switch boxes | # Fibers |
|---|---|---|---|---|---|---|
| Traditional | 3 | 5 | 49 k | 1,280 | 1,280 | 25 k |
| Multistage Chassis | 2 | 9 | 33 k | 2,304 | 192 | 16 k |

Improvement:                             1.5× (cost)                        6.7× (cost)   1.6× (cost)

Penalty:            1.8× (latency)                    1.8× (power)

Host

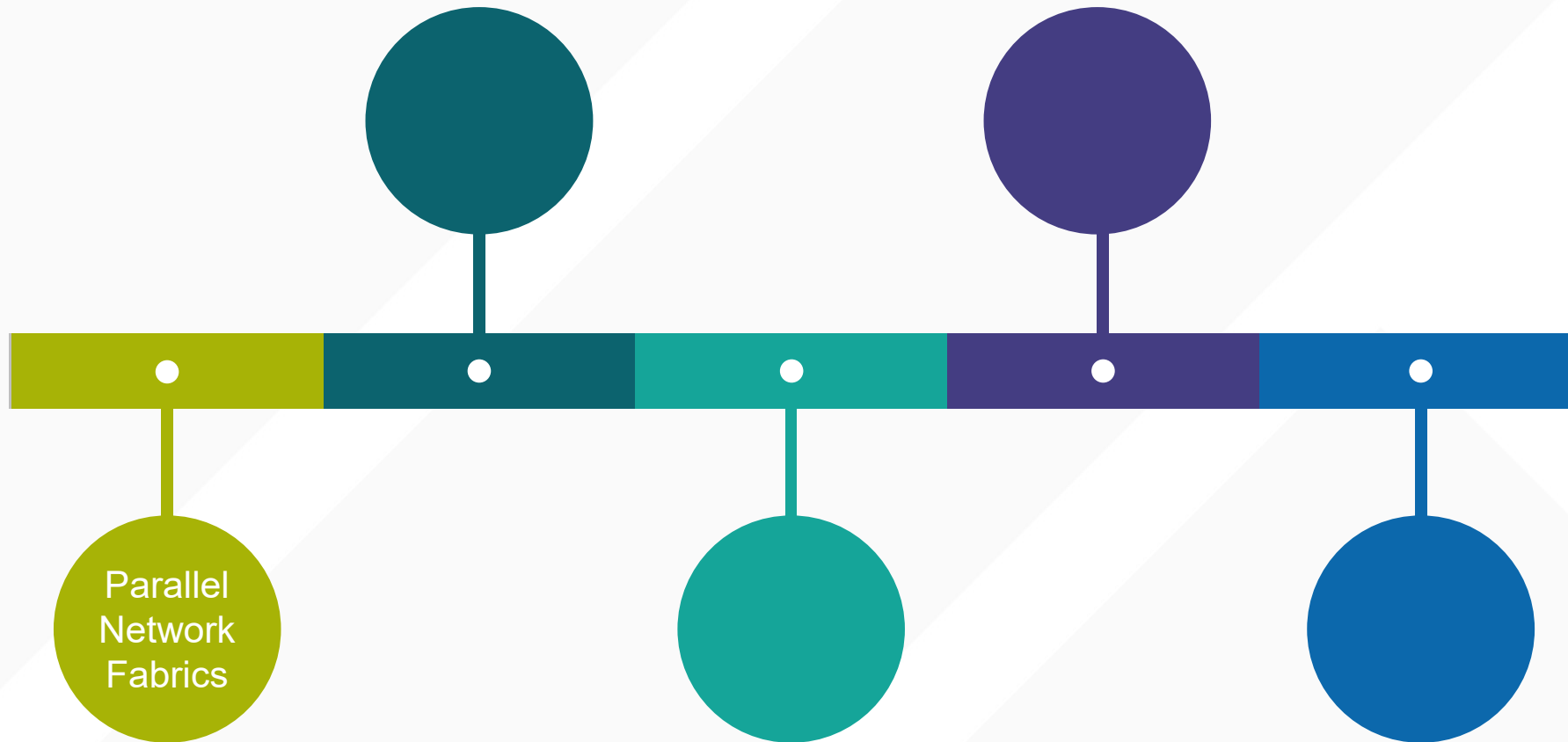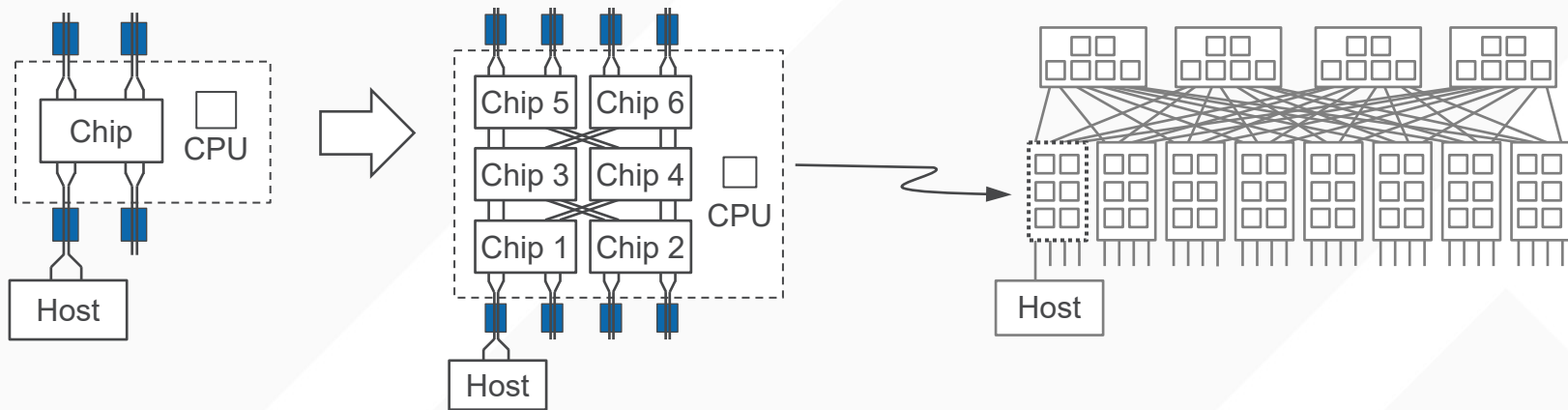Host

# TRENDS

- Conventional datacenter networks facing scaling limitations

  - Largely due to scaling limits of underlying packet switch chips


- Direction 1: Parallel network fabrics

  - Adopted thus far by Facebook and LinkedIn

  - Short-to-medium timeframe


- Direction 2: Replace packet switches with optical switches / circuit switches

  - Medium-to-longer timeframe

Conventional architectures:



Chassis architecture still scaling the network *up*… just hiding the tiers in switch chassis.
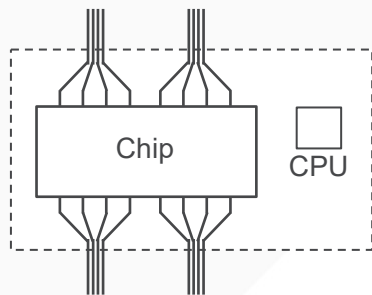
Alternative: Scale out via separate physical data planes

- Benefits: Reduced cost, power, and latency
- Tradeoff: Give up a single "fast" network abstraction
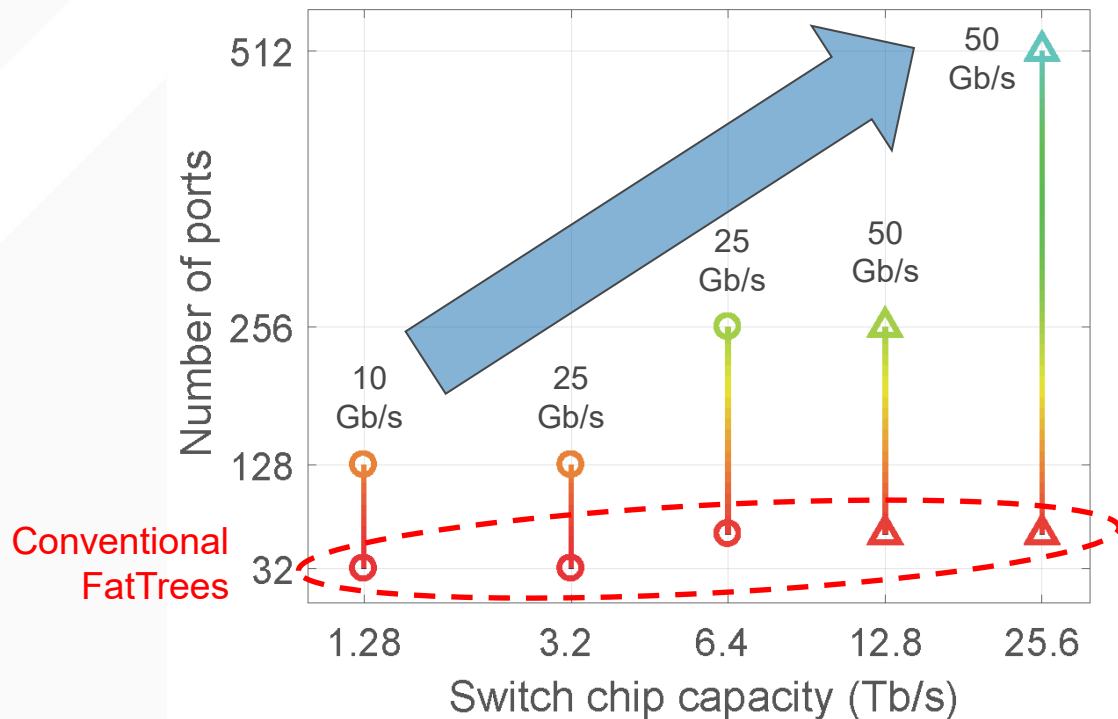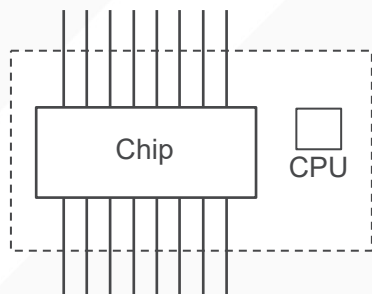
# UNDERLYING SWITCH RADIX IS INCREASING

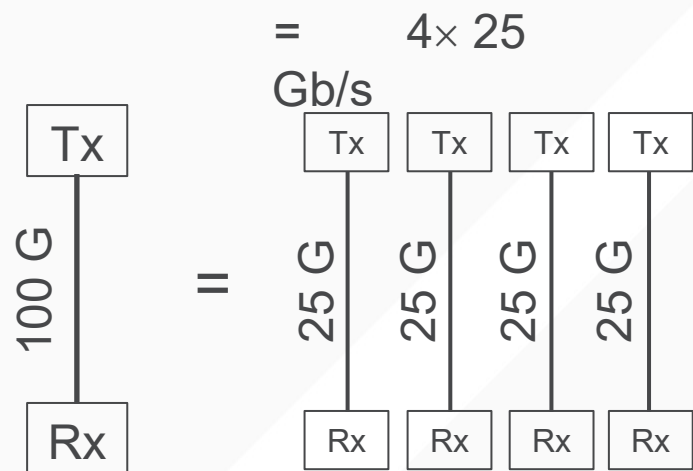Ex. Broadcom's Tomahawk switch:
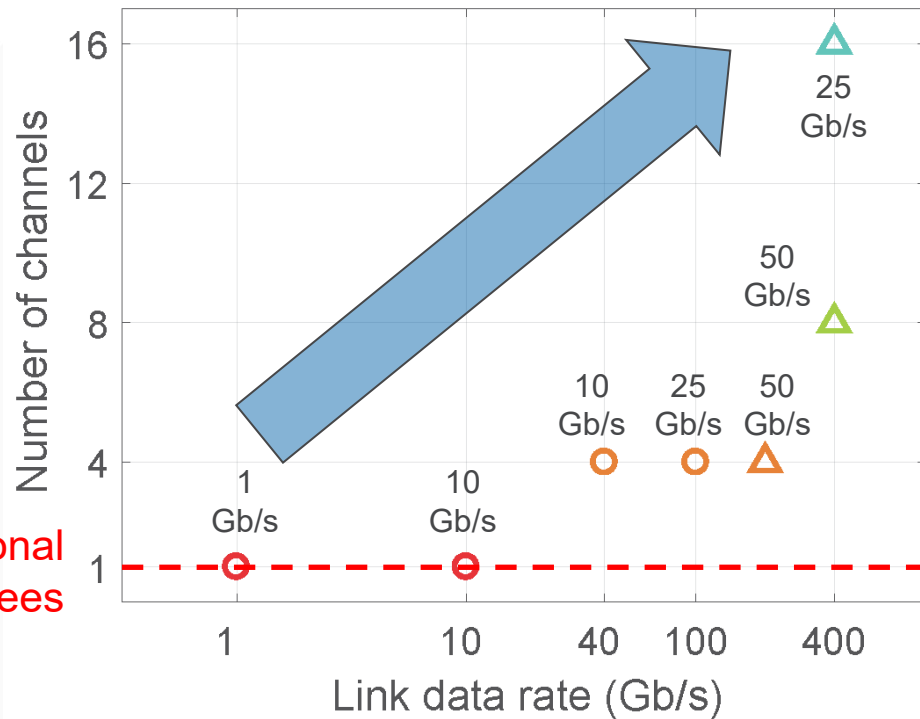
32 ports @ 100 Gb/s

*OR*

128 ports @ 25 Gb/s

# PARALLEL NETWORKS IN INDUSTRY: FACEBOOK



➤ from 4 x 128p multi-chip 400G fabric switches

FSW1    FSW2    FSW3    FSW4

**48 FSW ASICs + Control Planes per Pod**

**4 x 400G = 1.6T**
uplink per rack

➤ to 16 x 128p **single-chip 100G** fabric switches

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16

**16 FSW ASICs + Control Planes per Pod**

Sample Server Pod

**16 x 100G = 1.6T**
uplink per rack

*https://engineering.fb.com/data-center-engineering/f16-minipack/*

## Simpler and Flatter

→ Over **3X less** switch ASICs and control planes in fabric

→ **2.25X less** tiers of chips in the topology

→ Up to **2X less** host-to-host network hops intra-fabric

→ Up to **3X less** host-to-host network hops intra-region

https://engineering.fb.com/data-center-engineering/f16-minipack/

RESEARCH TIMELINE: DIRECTION 2: OPTICAL NETWORKS

Removing Transceivers

Parallel Network Fabrics

# MOTIVATION FOR OPTICAL NETWORKING

The faster the data rate of a cable, the shorter it has to be

**OPTICAL SWITCHES**

Output 1
Output 2
Input 1

Glass Fiber
Bundle

Lenses

Fixed
Mirror

te Mirror

Mirrors on Motors

✓ **No transceivers needed**
✓ **Supports unlimited bandwidth**
- **Different service model**
- **Not a drop-in replacement**
- **Reconfiguration delay δ**

# REMOVING TRANSCEIVERS

Performance

Practicality

Time

2009

2019

Millisecond switching
**Helios**, Sigcomm '10

Microsecond switching
**Mordia**, Sigcomm '13

Hybrid network & scheduling
**ReacToR**, NSDI '14
**Solstice**, CoNEXT '15

Simple control
**RotorNet**, Sigcomm '17

Low-latency for RotorNet
**Opera**, *NSDI'20*

# 2009 – USING 3D MEMS TO REMOVE TRANSCEIVERS



- Model: 15ms switch time

- Observed: 12.1ms switch time

- Technology: telecom-grade 3D-MEMs

- Scalability: 100s of ports

- Target: Inter-"pod"

*N. Farrington et al., SIGCOMM 2010*

# BOTTLENECKS IN NON-SWITCH COMPONENTS

- Telecom not designed for rapid reconfiguration

- Many non-switch bottlenecks in optical components

2-second link flap prevention



Packet-switch baseline

Observed circuit-switched bw

Transceiver Electronic Dispersion Compensation

# CONTROL PLANE 100X SLOWER THAN SWITCH TIME



Source Pods  Destination Pods

Hedera demand estimator
+
Edmund's Algorithm

Demand Matrix i

Circuit Switch i

1. Collect counters from packet switches

2. Estimate "true" demand

3. Calculate max-weighted matching

4. Reconfigure packet and optical switches

- One cycle ≈ one second

- Circuits try to "match" current network conditions

- Stateless in between assignments

# APPLICABILITY LIMITED BY SLOW SWITCH TIME & CONTROL PLANE

- Model: 15ms switch time

- Reality: 1000ms control plane

- To "capture" more of the traffic in optics, need a **faster switch** and **faster control plane**



*"Hardware Requirements for Optical Circuit Switched Data Center Networks", Farrington et al., OFC 2011*

Applicability of circuit switching determined by switch time δ

# REMOVING TRANSCEIVERS

Performance (vertical axis)

Practicality / Time (axes)

Hybrid network & scheduling
**ReacToR**, NSDI '14
**Solstice**, CoNEXT '15

Microsecond switching
**Mordia**, Sigcomm '13

Millisecond switching
**Helios**, Sigcomm '10

Simple control
**RotorNet**, Sigcomm '17

Low-latency for RotorNet
**Opera**, *NSDI'20*

2009

2019

# USING 2D MEMS TO "CHASE MICE"



- Needed a faster switch

- 2D MEMS very fast…

  - 2 µs switch time + ringing

  - Approx 11.5 µs total

- …but not scalable (~24 ports)

  - Lots of ports → slow

  - Few ports → fast

# 2011 - MORDIA – A 2D-MEMS 24-PORT MICROSECOND SWITCH



*Porter et al., "Integrating Microsecond Circuit Switching into the Data Center", Sigcomm'13.*

- Microsecond switching prevents scheduling with "fresh" data

  - Collecting demand a bottleneck!

- Insight: amortize series of switch configurations across a single demand estimate:

$$TM' = \sum_{i}^{N} t_i P_i$$

- Embodied by *Solstice* and *Eclipse* algorithms

- Result: "Chasing" demand

  - Reactive and responsive

Step 1. Gather traffic matrix TM     Step 2. Scale TM into TM′

$$\text{TM} \quad \begin{bmatrix} \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{bmatrix} \quad \longrightarrow \quad \text{TM}' \quad \begin{bmatrix} \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{bmatrix}$$

Step 3. Decompose TM′ into schedule

$$P_1 \quad P_2 \quad P_N$$

$$t_1 \begin{bmatrix} \vdots & \vdots & \vdots \end{bmatrix} + t_2 \begin{bmatrix} \vdots & \vdots & \vdots \end{bmatrix} + \cdots + t_N \begin{bmatrix} \vdots & \vdots & \vdots \end{bmatrix}$$

Step 4. Execute schedule in hardware

$$t_1 \qquad t_2 \qquad t_N$$

# NON-CROSSBAR NETWORKS



Performance

Practicality

Time

2009

2019

Hybrid network & scheduling
**ReacToR**, NSDI '14
**Solstice**, CoNEXT '15

Microsecond switching
**Mordia**, Sigcomm '13

Millisecond switching
**Helios**, Sigcomm '10

Simple control
**RotorNet**, Sigcomm '17

Low-latency for RotorNet
**Opera**, *NSDI'20*

# Toward 100+ Petabit/second datacenters

**Challenge:** deliver (very) low-cost bandwidth at scale



**Co-design:**

Protocol

Topology

Hardware

- **New protocols**
  Load balancing, congestion control, …

- **New topologies**
  Jellyfish, Longhop, Slimfly, …

- **New hardware**
  Optical circuit switching, RF/optical wireless, …

- ~~**Same switching model**~~

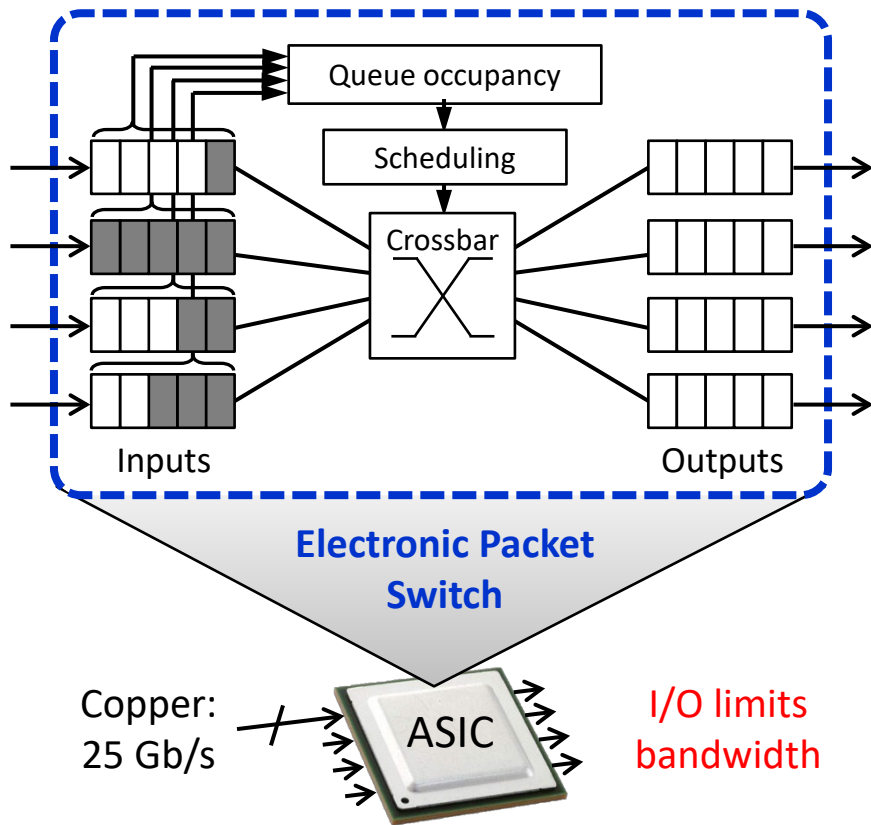  **New "Rotor" switching model**

**RotorNet** → **"**Future-proof" bandwidth (2× today)  +  simple control  +  …

# Optical switching – benefits & barriers



Data plane doesn't scale to entire datacenter!

**Electronic Packet Switch**

Inputs — Queue occupancy — Scheduling — Crossbar — Outputs

Copper: 25 Gb/s — ASIC — I/O limits bandwidth

Sending racks/hosts — Queue occupancy — Scheduling — Crossbar — Receiving racks/hosts

**Optical Circuit Switch**

Fiber: > 1 Tb/s — Cheap, future-proof bandwidth

Queue occupancy

Scheduling

Crossbar

# Rotor switching model simplifies control

**Crossbar model:**

**Real-time schedule**

Matching 1

*N* input ports

*N* output ports

Queue occupancy

Scheduling

Crossbar

→ **No (central) control**

**Rotor switch model:**

*N* − 1 matchings, **Fixed schedule**

$1 \rightarrow 2$   $1 \rightarrow 3$   $1 \rightarrow 4$

*N* input ports

*N* output ports

Rotor switch

→ **Bounded reduction in throughput**

*N* input
ports

*N* output
ports

# Rotor switches have a simpler implementation

Optical Crossbar:

**Optical Rotor switch:**



*N* input ports

*N* output ports

**Hard-wired matchings**

Mirror

Mirror

- Cost and complexity scale with:

Ports

Matchings (<< Ports)

**Ex.  2,048 ports:**     4,096 mirrors
2,048 directions

2 mirrors
16 directions

# RotorNet architecture overview

- **Forwarding?**

- **Topology?**

✔ - Optical Rotor switch → **More scalable**

✔ - Rotor switching model → **Simpler control**

# 1-hop forwarding over Rotor switch

- ## Wait for direct path:



Node 1 → 2, 3, 4
Node 2 → 3, 4, 1
Node 3 → 4, 1, 2
Node 4 → 1, 2, 3

Matching cycle 1        Matching cycle 2

Time

**Uniform traffic → 100% throughput**

- ## But datacenter traffic can be sparse …

41

# 1-hop forwarding & sparse traffic = low throughput

- Wait for direct path:



Matching cycle 1    Matching cycle 2
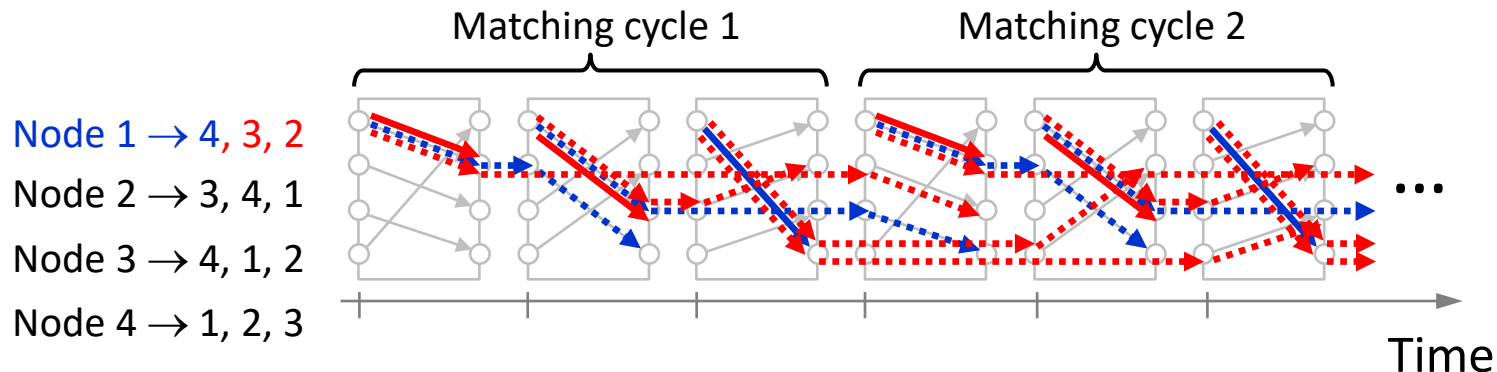
Node 1 → 4

Time

**Problem: single flow → 33% throughput**

- Hint at improvement: network is underutilized

# 2-hop forwarding better for sparse traffic

- Not new: Valiant ('82) & Chang et al. ('02)



**Throughput:** **Single flow** **33% (1-hop)** → **100% (2-hop)**
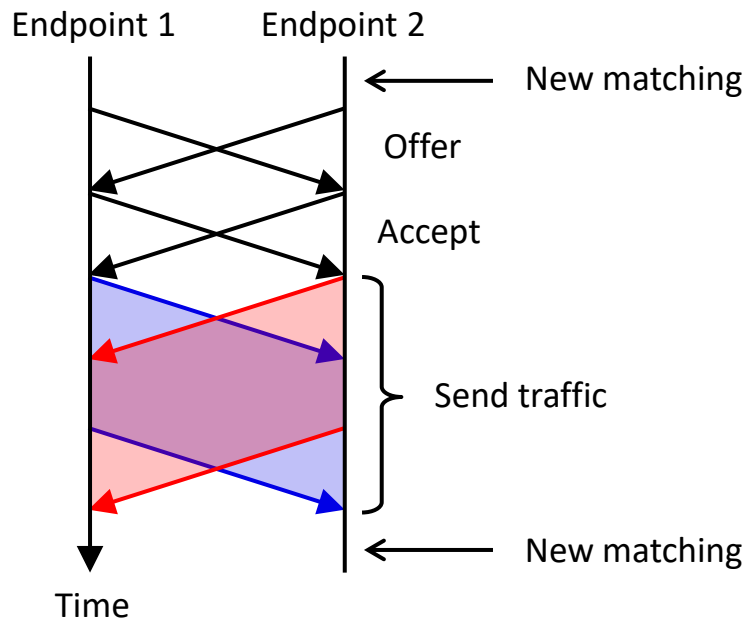
**Uniform traffic** **100% (1-hop)** → **50% (2-hop)**

- Optimization: can we adapt between **1-hop** and **2-hop** forwarding?

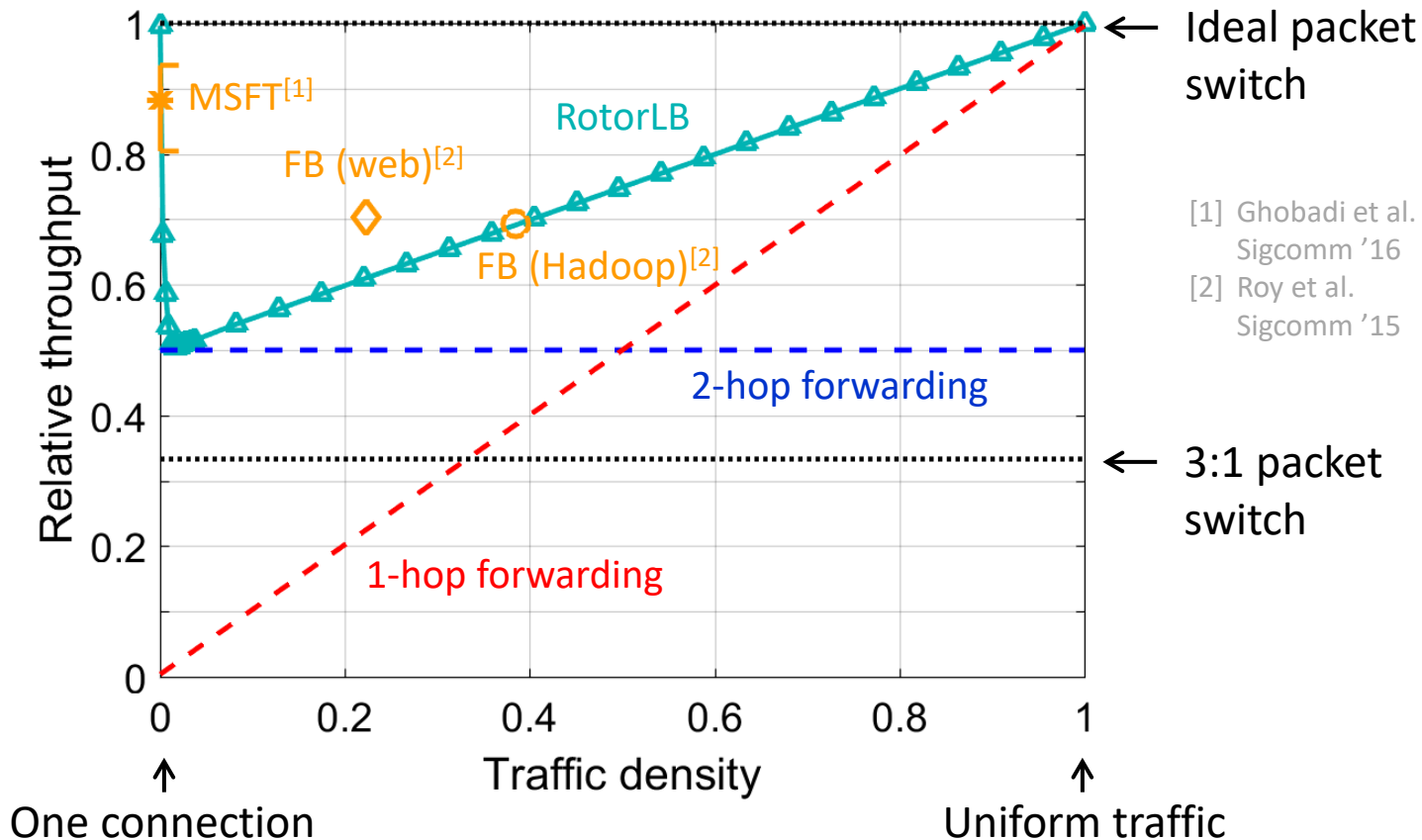# RotorLB: adapting between 1 & 2-hop forwarding

RotorLB (Load Balancing) overview:

- Default to 1-hop forwarding

- Send traffic over 2 hops only when there is extra capacity
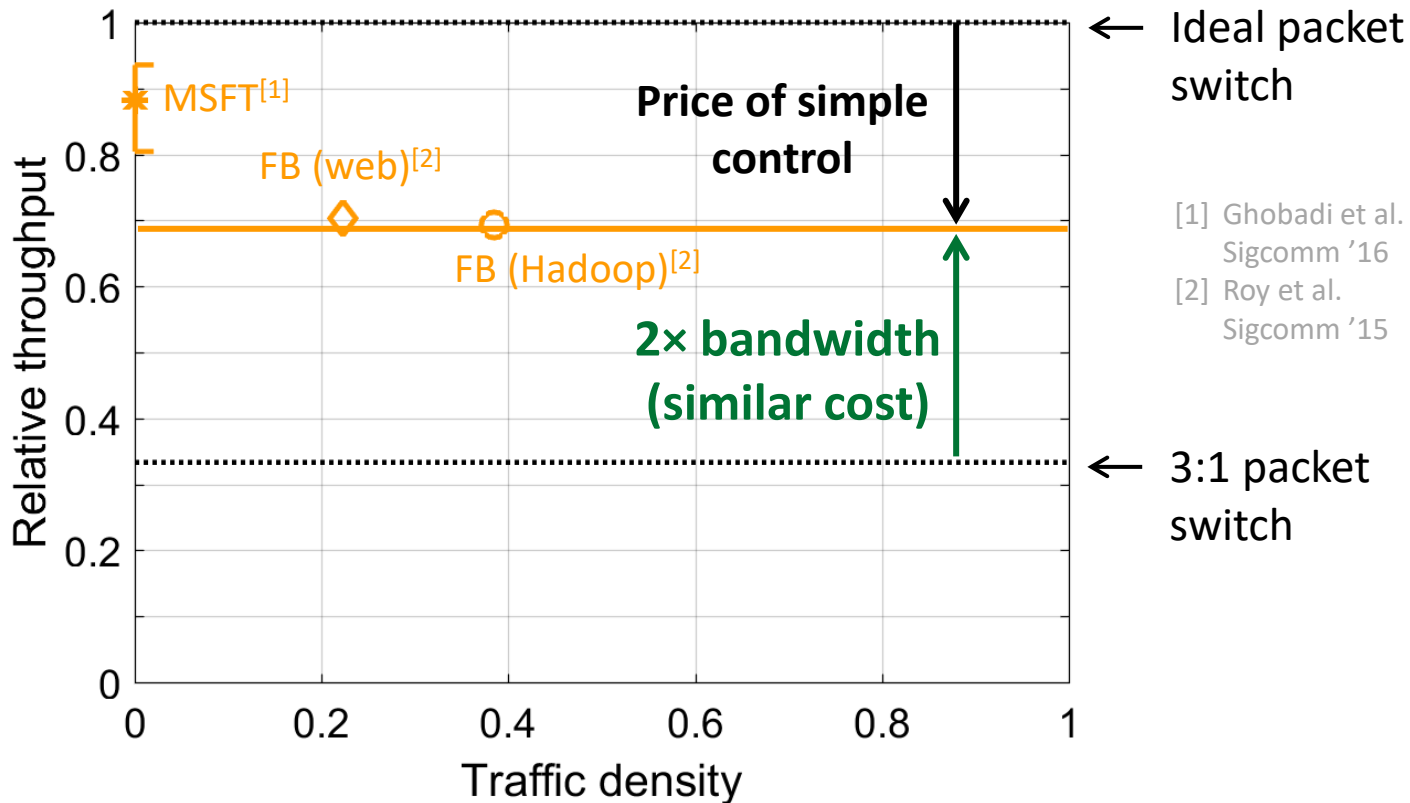
- Discover capacity using in-band pairwise protocol:



→ **RotorLB is fully distributed**

# Throughput of forwarding approaches (256 ports)

- Ideal packet switch
- MSFT[1]
- FB (web)[2]
- RotorLB
- FB (Hadoop)[2]
- 2-hop forwarding
- 3:1 packet switch
- 1-hop forwarding

Relative throughput

Traffic density

One connection

Uniform traffic

[1] Ghobadi et al. Sigcomm '16
[2] Roy et al. Sigcomm '15

# Throughput of forwarding approaches (256 ports)

- Ideal packet switch
- **Price of simple control**
- [1] Ghobadi et al. Sigcomm '16
- [2] Roy et al. Sigcomm '15
- **2× bandwidth (similar cost)**
- 3:1 packet switch

MSFT[1]

FB (web)[2]

FB (Hadoop)[2]

Relative throughput

Traffic density

# RotorNet architecture overview

- RotorLB → **Distributed, bounded throughput**

- **Topology?**

- Optical Rotor switch → **More scalable**

- Rotor switching model → **Simpler control**

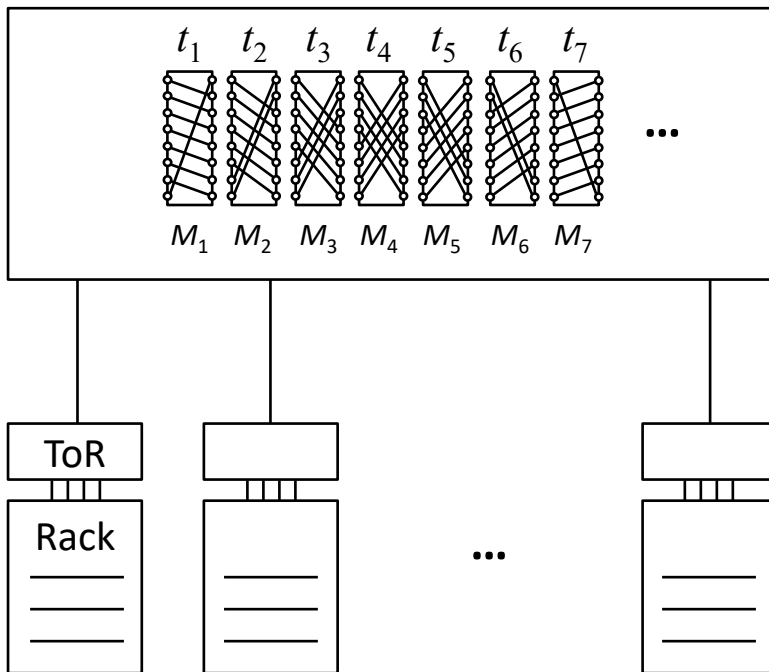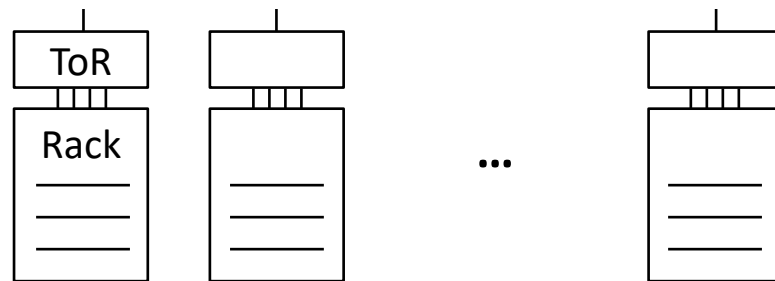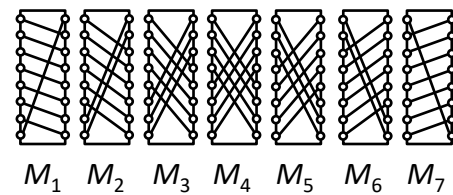# How should we build a network from Rotor switches?

At large scale:

- **High latency:**
  Sequentially step through many matchings

- **Fabrication challenge:**
  Monolithic Rotor switch with many matchings

- **Single point of failure**

Rotor switch

$M_1$  $M_2$  $M_3$  $M_4$  $M_5$  $M_6$  $M_7$

ToR

Rack

...

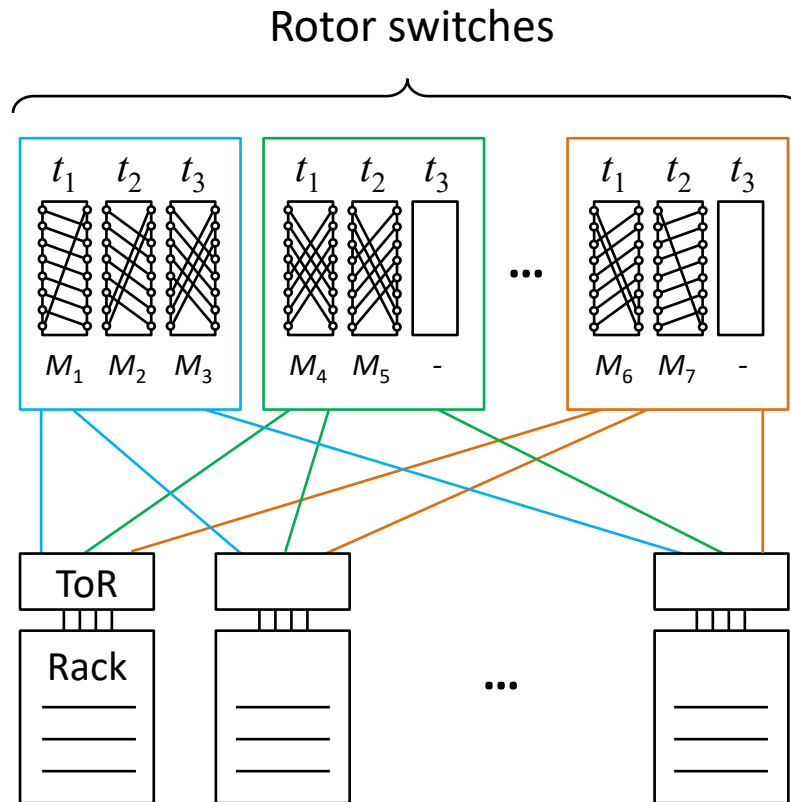# Distributing Rotor matchings = lower latency

**Fault tolerant**

**Reduced latency:**
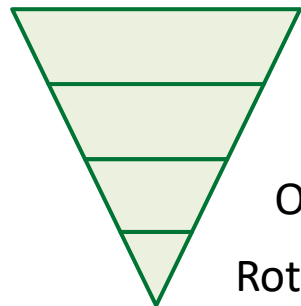
- Access matchings in parallel

**Simplifies Rotor switches:**

- Matchings << ports
- More scalable, less expensive



Rotor switches

# Rotor switching is feasible today

**Validated feasibility of entire architecture:**

(8 endpoints)

RotorLB
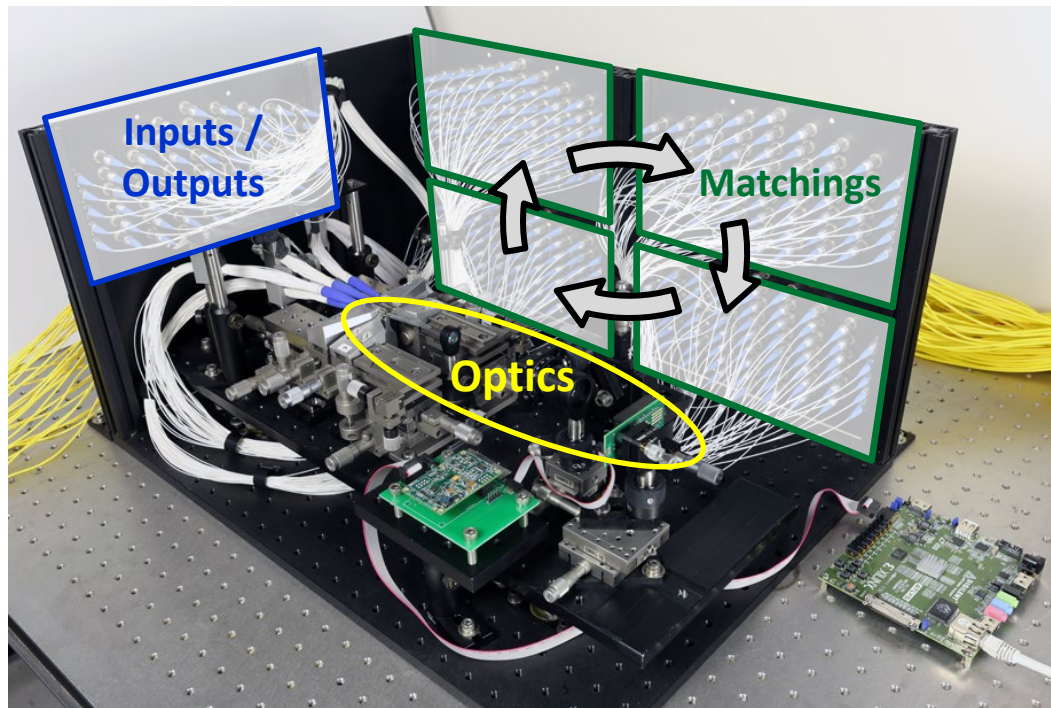
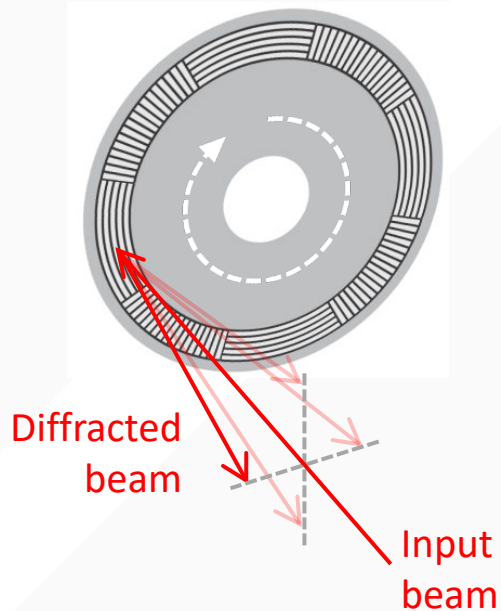RotorNet topology

Optical Rotor switch

Rotor switch model

100× faster switching than crossbar

Prototype Rotor switch



Inputs / Outputs

Matchings

Optics

"Pinwheel" sequential beam deflector



Diffracted beam

Input beam

=

High-speed spindle

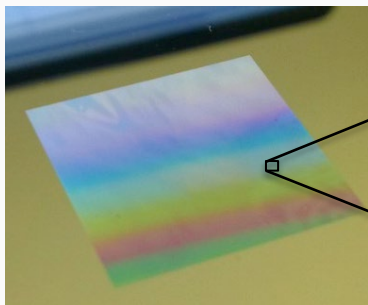(e.g. commercial 3.5" 7200 RPM drive)



+

Faceted disk

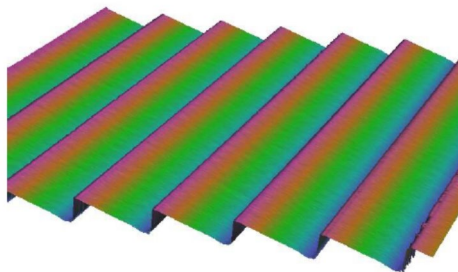(custom patterned with diffraction gratings)

# GRATING FABRICATION USING GREYSCALE LASER WRITING

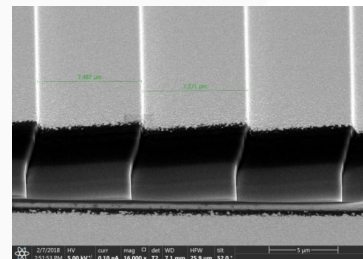Laser-written photoresist test grating
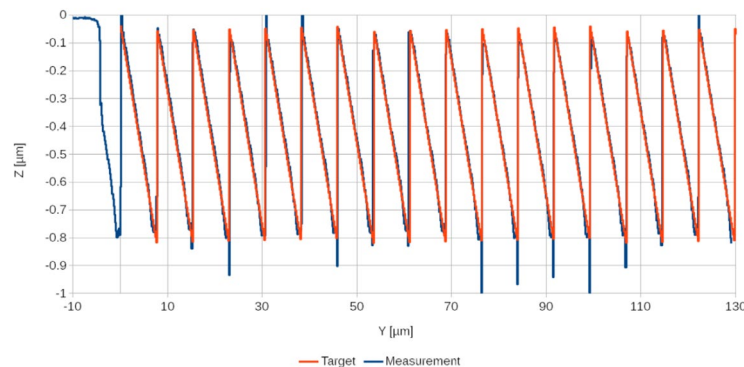(with gold coating)

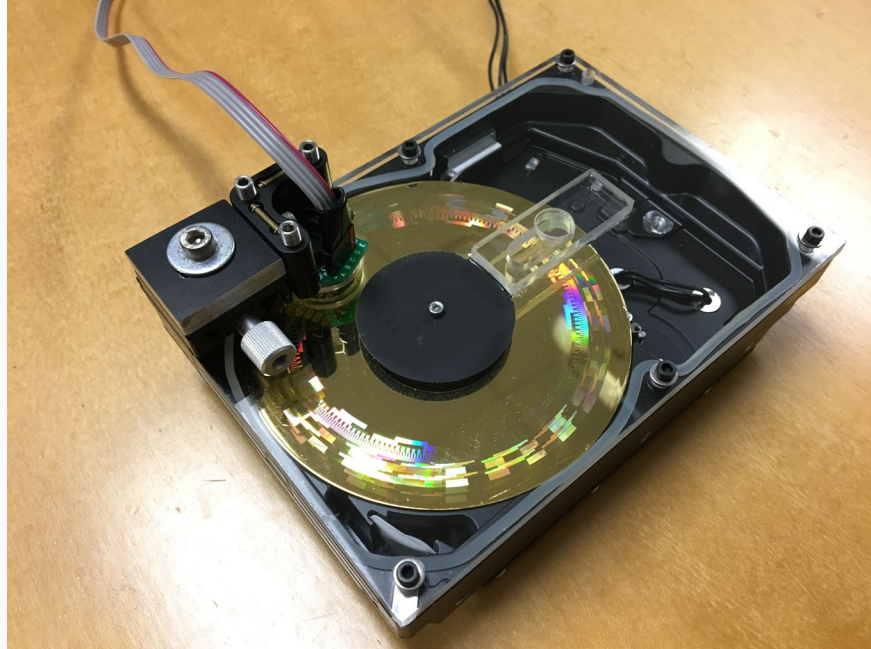

Surface profile of laser-written grating



SEM image
Pitch ≈ 6.67 μm, 150 lines / mm





Initial results indicated that laser writing can produce the features needed.

With encoder, encoder tracks, and clear cover

# ROTOR SWITCH PROTOTYPE

Optical layout:

Out 1

Input

Out 2

Grating pinwheel

95 mm

Laser-written grating pinwheel

Diffracted beam

Input beam

(WD) HGST Deskstar NAS drive

| Crosstalk: | < 30 dB |
| Operating spectrum: | > 120 nm |
| 2-pass insertion loss: | 5 – 8 dB* |

(*can be improved with better grating)

Image of fiber I/O array

Disk Sector 1

Disk Sector 2

Pinwheel rotation

Switching transitions

15 μs Rising

15 μs Falling

**15 μs reconfiguration @ 7200 RPM**
**(1,000 x faster than commercial MEMS OXC)**

# IMPROVED PERFORMANCE WITH NEW PROTOTYPE



**1st Prototype:** MEMS selector switch
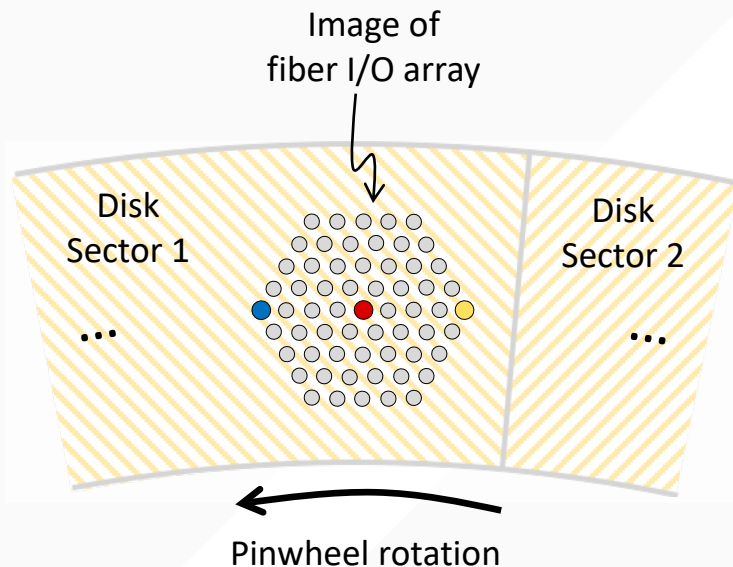- *Higher loss optics on enclosed ½" breadboard*
- *150 µs switching*
- *I/O to external connection patch panels*

**2nd Prototype:** "rotor" switch with pinwheel
- *Lower loss optics mounted on vibration-isolated rail*
- *15 µs switching (@ 7200 RPM)*
- *I/O with 4x internal connection patch panels*

Rotor Switch

9 Servers

Connections (back)

# RotorNet scales to 1,000s of racks

- Rotor switch design point:  2,048 ports,  1,000× faster switching than crossbar

  Details in:  W. Mellette et al., *Journal of Lightwave Technology* '16
                W. Mellette et al., *OFC* '16

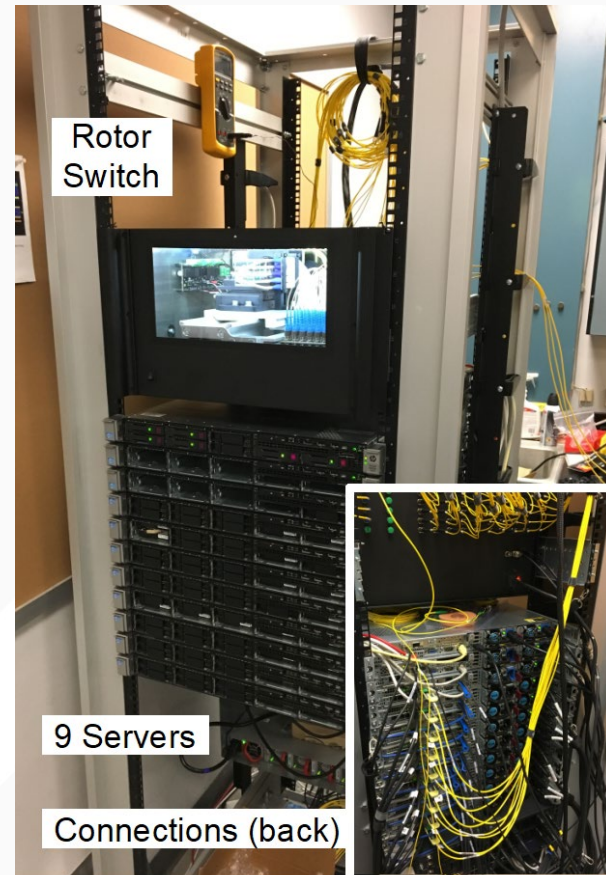- 2,048-rack data center:
  → **Latency (cycle time)**
     **= 3.2 ms**

- Faster than 10 ms crossbar <u>reconfiguration time</u>

- Hybrid network for low-latency applications

128 Rotor switches

Packet switches

ToR

Rack

# RotorNet component comparison

| Network | # Packet switches | # Transceivers | # Rotor switches |
|---------|-------------------|----------------|------------------|
| 3:1 Fat Tree | 2.6 k | 103 k | 0 |

RotorNet delivers:
- Today: Bandwidth 2× less expensive

- Future: Cost advantage grows with bandwidth

- **Benefits of optical switching without control complexity**

# REMOVING TRANSCEIVERS

✓ Similar hardware, cost, and power savings to an oversubscribed Fat Tree

✓ Improved throughput vs oversubscribed Fat Tree at low load

**"Bandwidth tax"** – Reduction in throughput at high traffic loads
– Proportional to average path length

# Bandwidth tax limits throughput in expander networks



Contention

Flow 1:  ~~100%~~ ~~70%~~ 50%     Flow 2:  ~~50%~~ ~~20%~~

Bandwidth tax = 2  →  Throughput = 50% at high load

→ *Is it possible to support high loads while reducing cost and power?*

# Reconfigurable networks enable higher throughput



Reconfigured *direct* links: **bandwidth tax = 1**

*Multi-hop* links
**bandwidth tax = 2**

Reconfigure links

Flow 1:  **100%**

Flow 2:  **100%**

✓ Reconfiguration permits high throughput at high load

*Added complexity: how do we decide which links to reconfigure and when?*

→ "RotorNet" (Sigcomm '17) – fixed schedule of direct circuits

*Today's circuit switching technologies reconfigure too slowly →high latency*

# Our contribution: we can have the best of static *and* reconfigurable

**Reconfigurable networks:** high latency high throughput

**Expander networks:** low latency low throughput

**Workload** = Short flows + Long flows

*Latency-bound* *Throughput-bound*

"Opera" – combining expanders and reconfiguration in a single, unified network

# Opera's design – part 1: providing low-bandwidth-tax connectivity



Full, direct inter-rack connectivity with *N* <u>matchings</u>:

Full, direct inter-rack connectivity with $N$ <u>matchings</u>:

$M_1$  $M_2$  $\cdots$  $M_N$

Time

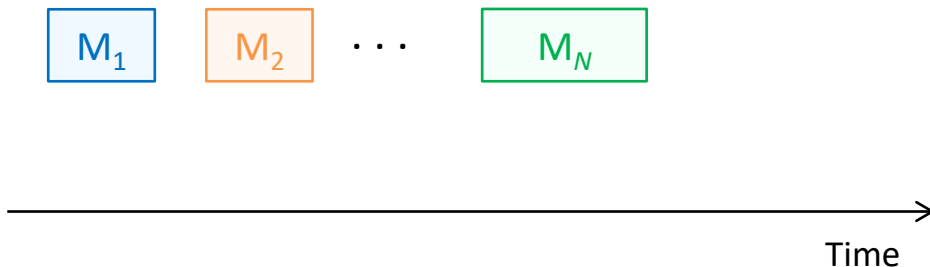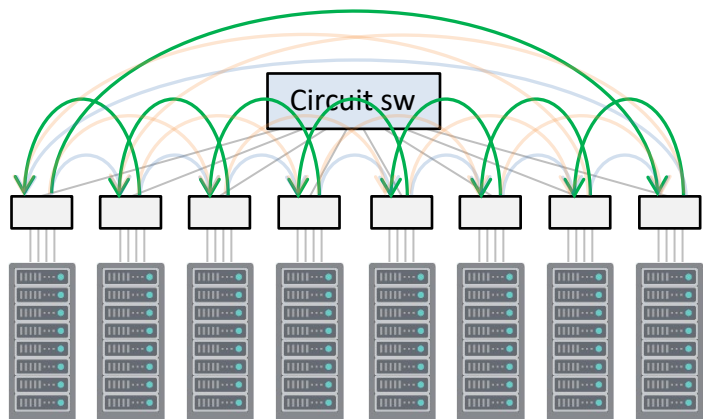# Opera's design – part 2: providing low-latency connectivity



Full, direct inter-rack connectivity with $N$ matchings:

- Short, latency-bound flows can be sent immediately over multi-hop paths  (high BW tax)
- Long, throughput-bound flows can wait for direct paths                   (low BW tax)

Key property: Opera only pays a bandwidth tax for short flows → **lower average tax**

# Choosing matchings

i. <u>Expansion</u>    Union of 3 or more <u>randomly-structured matchings</u> is an expander [1]

[1] N Alon, "Eigen values and expanders," Combinatorica, 6(2), 1986.

ii. <u>Direct connectivity between all racks over time</u>

Factor complete graph into *N* <u>randomly-structured</u> & <u>disjoint</u> matchings:



Complete graph  =  Factored complete graph  =  **M₁**  +  **M₂**  + ··· +  **Mₙ**

# Offsetting reconfigurations for continuous connectivity



Time to wait for direct path → cutoff between "short" & "long" flows

# Opera is well-suited for many published data center workloads

Published data center flow distributions:

## Quantifying the cutoff

For 10 Gb/s – 100 Gb/s links:

- Long flows ≥ 15-30 MB
  *can afford to wait for direct paths*
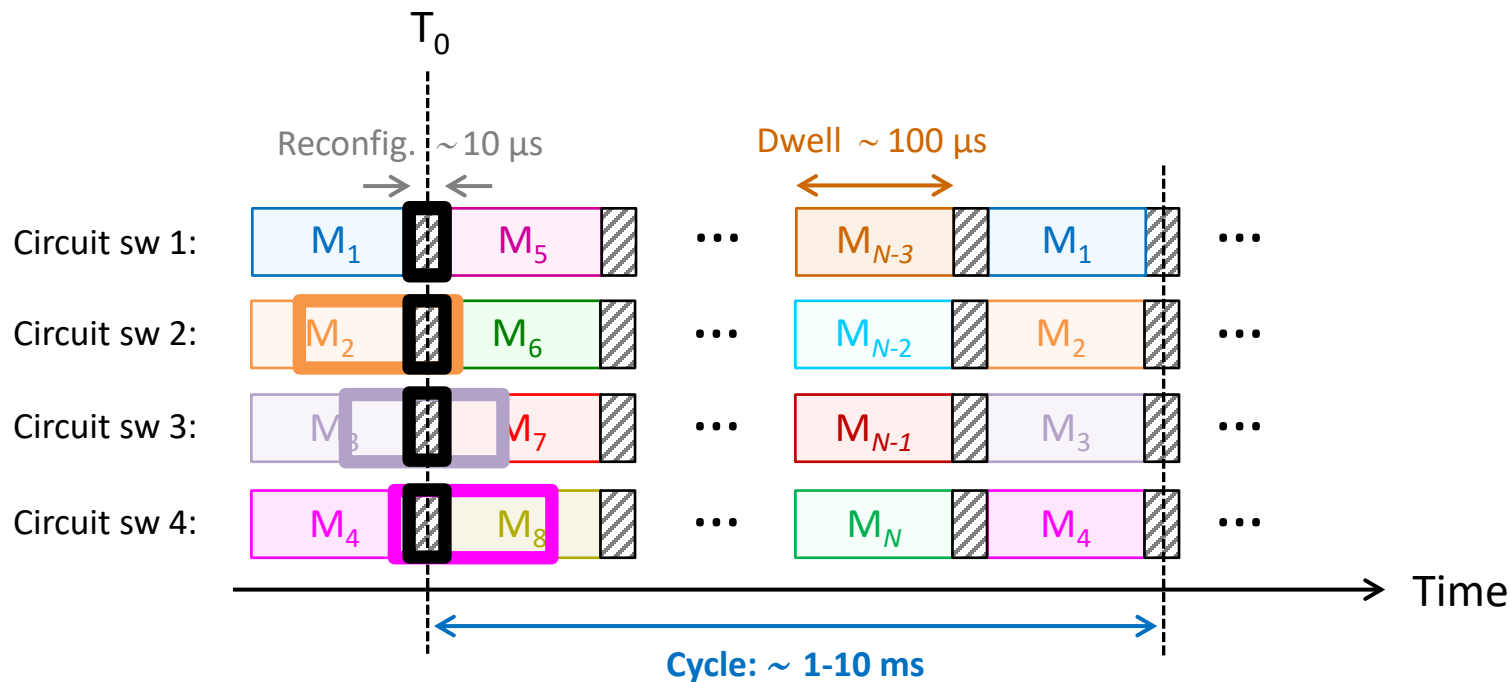
- Short flows < 15-30 MB
  *cannot wait for direct paths*



"Short" ← → "Long"

Microsoft websearch
Facebook Hadoop
Microsoft datamining

15-30 MB

CDF of flows
Flow size (bytes)

*Most **flows** must pay bandwidth tax*

✓ *Most **bytes** can avoid bandwidth tax*

> 90% of bytes

CDF of bytes
Flow size (bytes)

**Workload 1:** All-to-all shuffle
(favorable)



Opera
Expander graph
3:1 Fat Tree

→ 4x higher throughput & faster completion

**Workload 2:** Shuffle + MSFT websearch workload
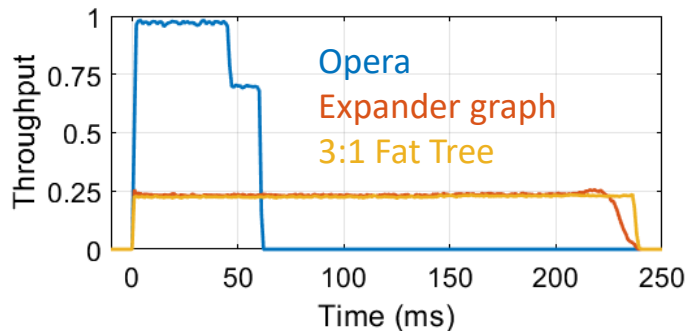(challenging)



Opera
Expander graph
3:1 Fat Tree

→ 2-4x higher throughput &
equivalent completion times for short flows

**Workload 3:** MSFT datamining (100 B – 1 GB flows)

→ 60% higher admissible load with equivalent FCTs

# Practical considerations

**Fault tolerance:**
- Full connectivity maintained with 4% of links, 7% of ToRs, or 40% of circuit switches failed
  (Better than oversubscribed Fat Tree, not as good as static expander)

- Failures detected and disseminated within $O$(10 ms)
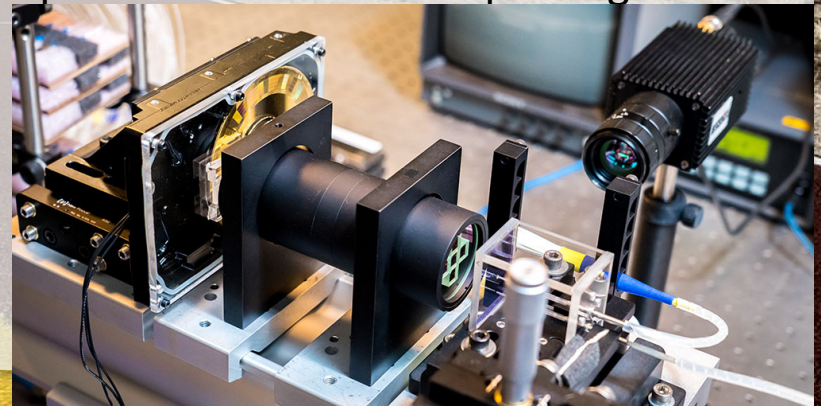
**Prototype implementation:**

- Time-synchronized routing implemented on programmable Barefoot Tofino switch with P4

- Opera scales to 1,000's of racks, 10,000's of servers with commodity switch table sizes

# CREDITS AND THANKS TO MY COLLEAGUES/STUDENTS

- Alex C. Snoeren
- Alex Forencich
- Amin Vahdat
- Andrew Grieco
- Anthony Lentine
- Arjun Roy
- Chang-Heng Wang
- Christopher DeRose
- Chunming Qiao
- Conglong Li
- David G. Andersen
- Douglas C. Trotter
- Feng Lu
- Geoffrey M. Voelker
- George C. Papen
- George Papen
- Glenn M. Schuster

- Guohui Wang
- Hamid Bazzaz
- He Liu
- Jordan Davis
- Joseph Ford
- Joshua Zhu
- Kai Chen
- Li Chen
- Malveeka Tewari
- Matthew K. Mukerjee
- Michael A. Kozuch
- Michael Kaminsky
- Minlan Yu
- Nathan Farrington
- Nicolas Feltman
- Pang-Chen Sun
- Rajdeep Das

- Richard Strong
- Rishi Kapoor
- Rob McGuinness
- Ryan Aguinaldo
- Shan Zhong
- Shaya Fainman
- Shayan Mookherjea
- Sivasankar Radhakrishnan
- Srinivasan Seshan
- Stefan Savage
- T. S. Eugene Ng
- Tajana Rosing
- Tara Javidi
- Vikram Subramanya
- William M. Mellette
- Yeshaiahu Fainman
- Yibo Guo

- Reaching the limits of CMOS-based packet switching
    - In terms of cost, power, performance…

- Direction 1: scale bandwidth by adding *parallel* dataplanes

- Direction 2: scale bandwidth by replacing packet switches with optical ones
    - Unique opportunity to incorporate novel optical devices such as spinning pinwheel/hard drive based switches

- Thank you for your time and attention!