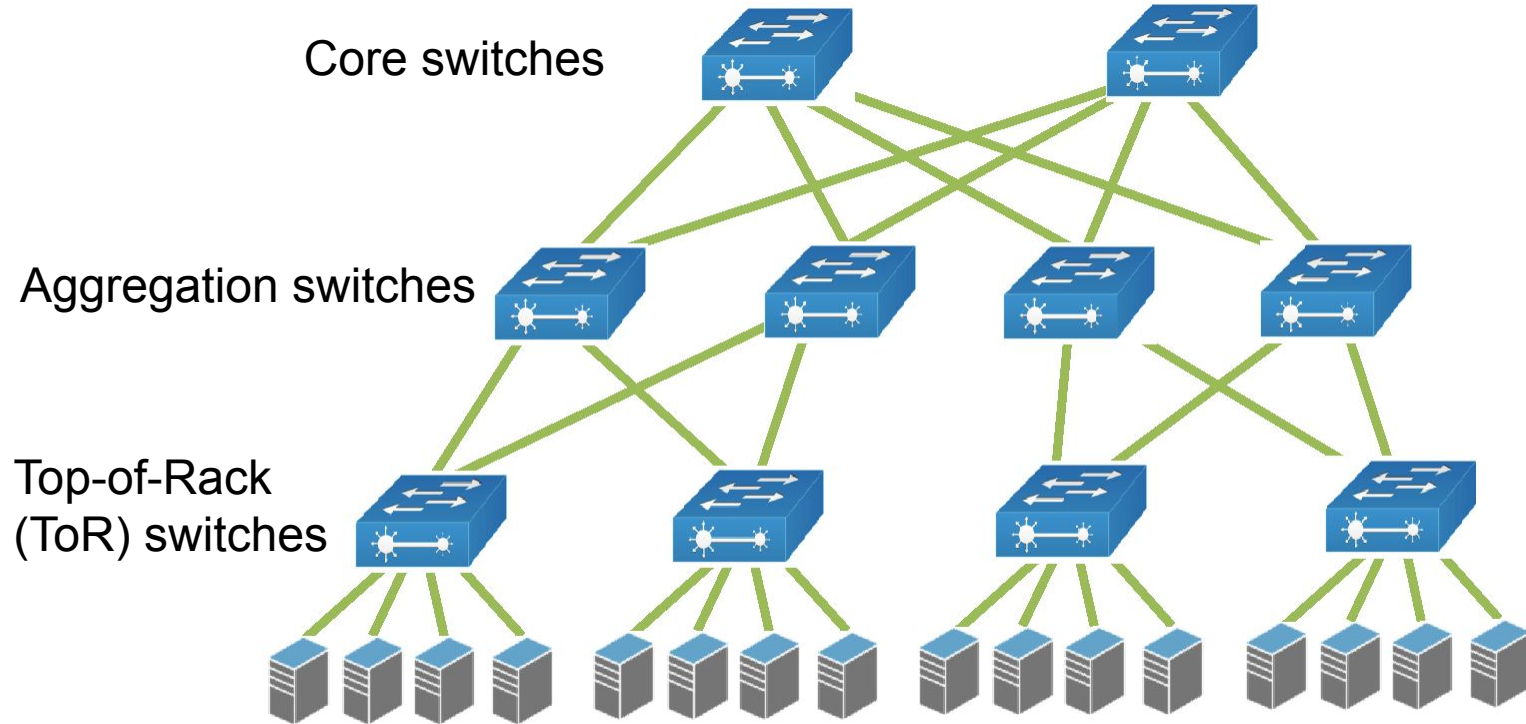# Hop-On Hop-Off Routing

## A Fast Tour across the Optical Data Center Network for Latency-Sensitive Flows
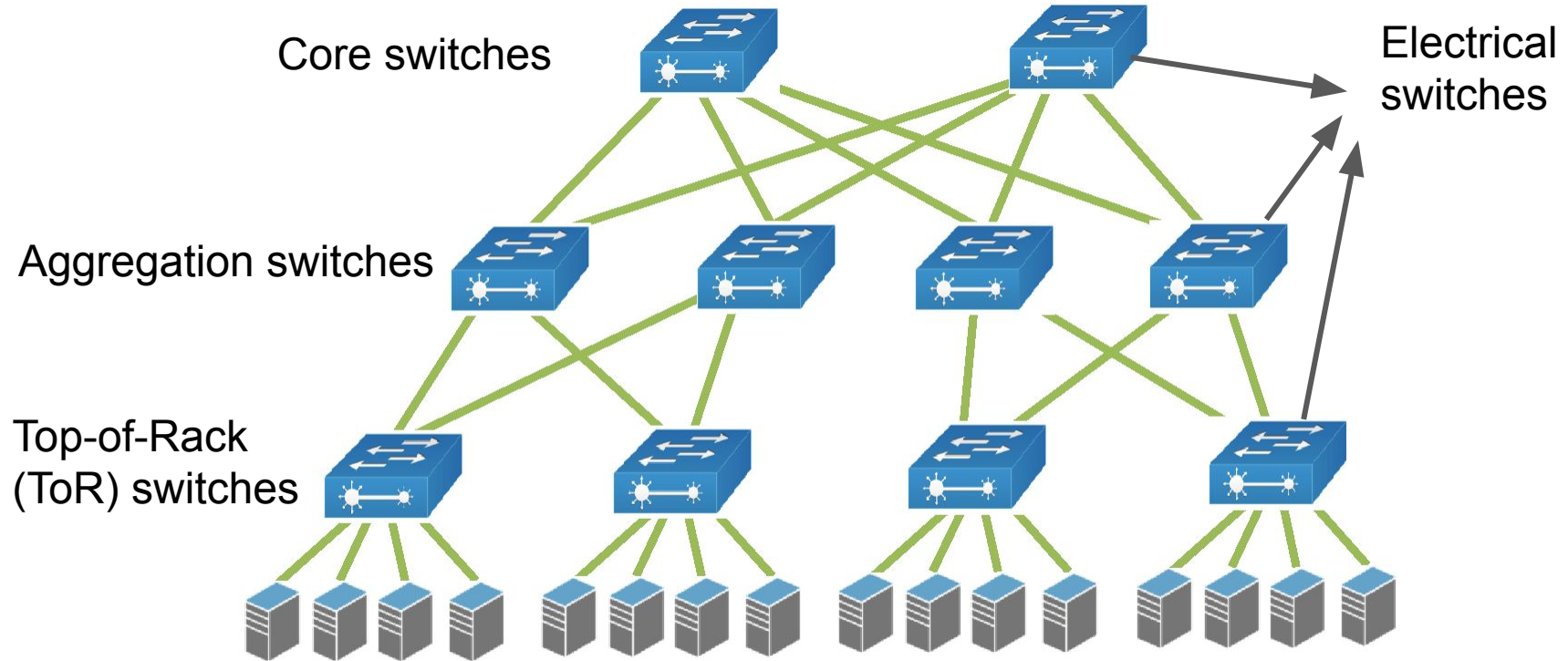
Jialong Li[1], Yiming Lei[1], Federico de Marchi[2]
Raj Joshi[3], Balakrishnan Chandrasekaran[4], Yiting Xia[1]

1. Max Planck Institute for Informatics   2. Saarland University
3. National University of Singapore   4. Vrije Universiteit Amsterdam
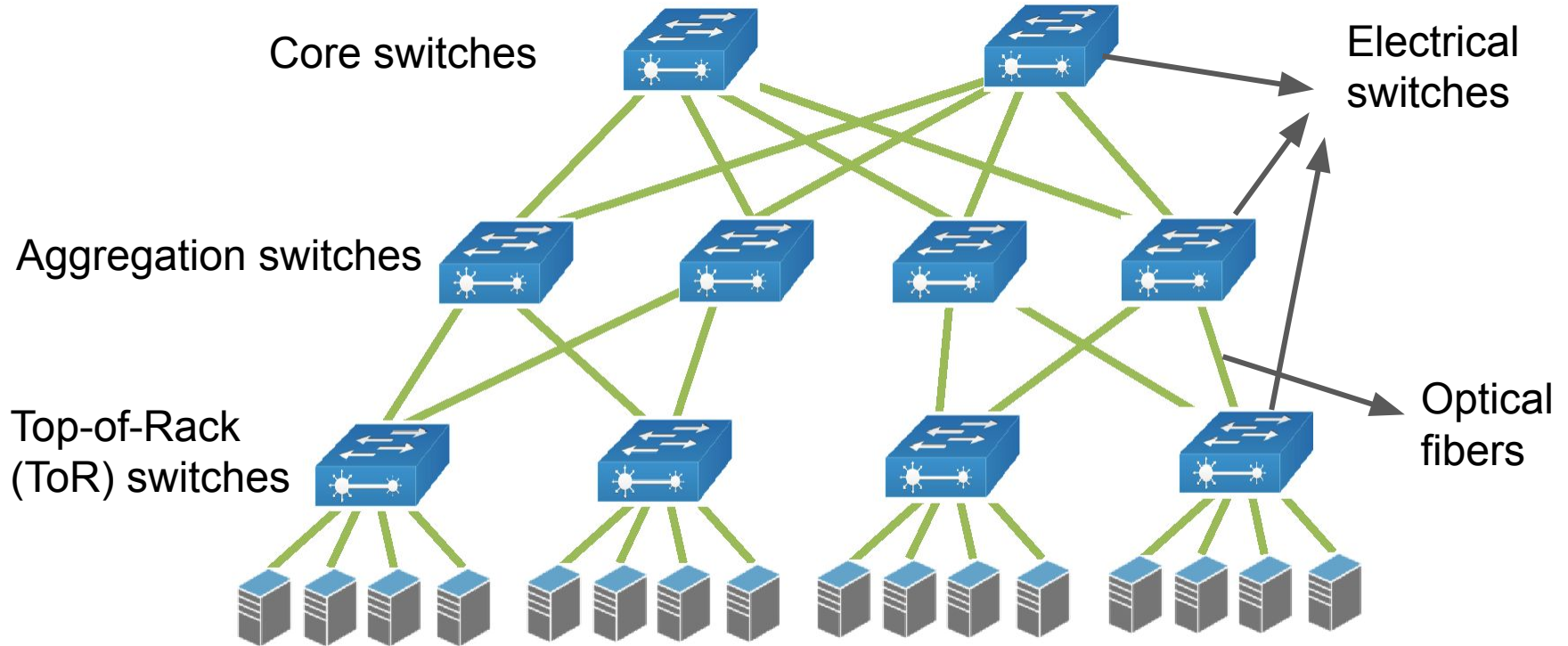
# Today's Data Center Networks



Core switches

Aggregation switches

Top-of-Rack (ToR) switches

# Today's Data Center Networks



Core switches

Aggregation switches

Top-of-Rack (ToR) switches

Electrical switches

# Today's Data Center Networks



Core switches

Aggregation switches

Top-of-Rack (ToR) switches

Electrical switches

Optical fibers

2

# From Electrical Switching to Optical Switching



Electrical Switching

Packet buffer, switching…

# From Electrical Switching to Optical Switching



Electrical Switching

Packet buffer, switching…

Optical switch

Optical circuit

ToR A   ToR B   ToR C   ToR D

Optical Switching

# Why Optical Switching

- No queuing delay along the circuit

# Why Optical Switching

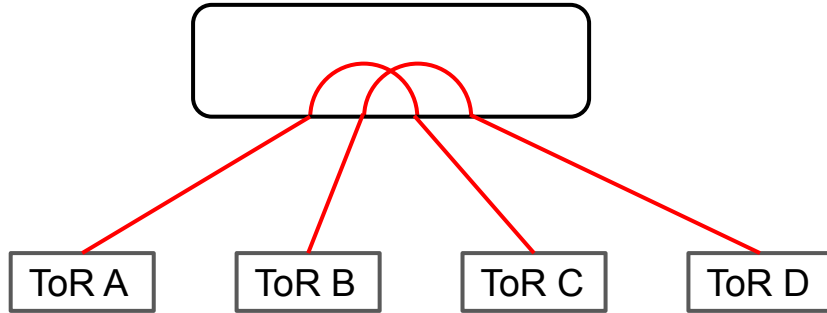- No queuing delay along the circuit

- Less power consumption

# Why Optical Switching

- No queuing delay along the circuit

- Less power consumption

- Independent to link bandwidth

# Optical Switching



Time slot 0 (time 0 - T)
A ⇔ C
B ⇔ D

# Optical Switching



Time slot 0 (time 0 - T)
  A ⇔ C
  B ⇔ D

Time slot 1 (time T - 2T)
  A ⇔ B
  C ⇔ D

# Optical Switching



Time slot 0 (time 0 - T)
    A ⇔ C
    B ⇔ D

Time slot 1 (time T - 2T)
    A ⇔ B
    C ⇔ D

Time slot: When and how long the circuit exists

# Optical Switching



- A packet arrives at time 0

- Src = A, dst = B

- Waiting at A for time T

# Optical Switching



- A packet arrives at time 0

- Src = A, dst = B

- Waiting at A for time T

Circuit waiting delay: delay at src ToR before required circuit comes

# Latency-Sensitive Flows

Latency-Sensitive Flows: Flow Completion Time (FCT) is critical

# Latency-Sensitive Flows

Latency-Sensitive Flows: Flow Completion Time (FCT) is critical

FCT: tens of *ns* to hundreds of *µs*[1]

Circuit waiting delay: tens of *µs* to *ms*

[1] Expanding across time to deliver bandwidth efficiency and low latency, NSDI' 20

# Latency-Sensitive Flows

Latency-Sensitive Flows: Flow Completion Time (FCT) is critical

FCT: tens of *ns* to hundreds of *µs*[1]

Circuit waiting delay: tens of *µs* to *ms*

Reduce circuit waiting delay: Using multi-hop paths

[1] Expanding across time to deliver bandwidth efficiency and low latency, NSDI' 20

# Multi-Hop Path

src ToR

dst ToR

t = *i-th* time slot

t = 5

Direct path

# Multi-Hop Path



ToR A

Multi-hop path

t = 3

t = 4

src ToR

dst ToR

t = *i-th* time slot

Direct path

t = 5

# Multi-Hop Path



ToR A

Multi-hop path

t = 3

t = 4

src ToR

dst ToR

Multi-Hop path = cascade of several circuits

t = *i-th* time slot

Direct path

t = 5

# State-of-the-art: Opera (NSDI' 2020)

Opera
- A multi-hop path is available at any moment between any ToR pairs
- This multi-hop path is non-stop path

# State-of-the-art: Opera (NSDI' 2020)

Opera

- A multi-hop path is available at any moment between any ToR pairs
- This multi-hop path is non-stop path
- Non-stop path ≠ fastest path

# Hop-On Hop-Off (HOHO) Routing

- Search for the fastest path, instead of non-stop paths

- Packets could wait at ToRs, so they can "hop on" a circuit, "hop off", and "hop on" another circuit

# HOHO Routing: Intuition

- The latency of a path is only determined by the time slot of the last-intermediate ToR and dst ToR

# HOHO Routing: Intuition

- The latency of a path is only determined by the time slot of the last-intermediate ToR and dst ToR

Path 1: src => A => dst

Path 2: src => B => C => dst

t = *i-th* time slot
a bus = circuit

# HOHO Routing: Intuition

- The latency of a path is only determined by the time slot of the last-intermediate ToR and dst ToR
- Backtracking: Search from dst ToR to src ToR



Path 1: src => A => dst

Path 2: src => B => C => dst

t = *i-th* time slot
a bus = circuit

# HOHO Routing

| Time Slot | Circuits |
|:---:|:---|
| 1 | S – H |
| 2 | A – F |
| 3 | H – E, F – A, S – B |
| 4 | S – G, A – D |
| 5 | B – D, G – B |
| 6 | S – D |
| 7 | S – A |
| 8 | C – D |

Fixed, periodic optical schedule

- Time slot index indicating when two ToRs have a circuit

Input

- A fixed, periodic optical schedule
- Src and dst ToR
- Packet arrival time slot
- Maximum hops (optional)

# HOHO Routing

| Time Slot | Circuits |
|:---:|:---|
| 1 | S – H |
| 2 | A – F |
| 3 | H – E, F – A, S – B |
| 4 | S – G, A – D |
| 5 | B – D, G – B |
| 6 | S – D |
| 7 | S – A |
| 8 | C – D |

Fixed, periodic optical schedule

- Time slot index indicating when two ToRs have a circuit

Input
- A fixed, periodic optical schedule
- Src and dst ToR
- Packet arrival time slot
- Maximum hops (optional)

Output
- A fastest path between src and dst ToR

Rerun per ToR pair, per packet arrival time slot

# HOHO Routing

Send a packet from S to D in the minimal time within max hops, max = 3 here

Tree level 0



| Time Slot | Circuits |
|-----------|----------|
| 1 | S – H |
| 2 | A – F |
| 3 | H – E, F – A, S – B |
| 4 | S – G, A – D |
| 5 | B – D, G – B |
| 6 | S – D |
| 7 | S – A |
| 8 | C – D |

# HOHO Routing

Send a packet from S to D in the minimal time within max hops, max = 3 here



Tree level 0

Tree level 1

t = 6

| Time Slot | Circuits |
|-----------|----------|
| 1 | S – H |
| 2 | A – F |
| 3 | H – E, F – A, S – B |
| 4 | S – G, A – D |
| 5 | B – D, G – B |
| 6 | S – D |
| 7 | S – A |
| 8 | C – D |

# HOHO Routing

Send a packet from S to D in the minimal time within max hops, max = 3 here

Tree level 0

Tree level 1

t = 8

t = 6

Later than S->D

| Time Slot | Circuits |
|-----------|----------|
| 1 | S – H |
| 2 | A – F |
| 3 | H – E, F – A, S – B |
| 4 | S – G, A – D |
| 5 | B – D, G – B |
| 6 | S – D |
| 7 | S – A |
| 8 | C – D |

# HOHO Routing

Send a packet from S to D in the minimal time within max hops, max = 3 here

Tree level 0

Tree level 1

t = 6

| Time Slot | Circuits |
|-----------|----------|
| 1 | S – H |
| 2 | A – F |
| 3 | H – E, F – A, S – B |
| 4 | S – G, A – D |
| 5 | B – D, G – B |
| 6 | S – D |
| 7 | S – A |
| 8 | C – D |

# HOHO Routing

Send a packet from S to D in the minimal time within max hops, max = 3 here

Tree level 0

Tree level 1



| Time Slot | Circuits |
|-----------|----------|
| 1 | S – H |
| 2 | A – F |
| 3 | H – E, F – A, S – B |
| 4 | S – G, A – D |
| 5 | B – D, G – B |
| 6 | S – D |
| 7 | S – A |
| 8 | C – D |

# HOHO Routing

Send a packet from S to D in the minimal time within max hops, max = 3 here

| Time Slot | Circuits |
|-----------|----------|
| 1 | S – H |
| 2 | A – F |
| 3 | H – E, F – A, S – B |
| 4 | S – G, A – D |
| 5 | B – D, G – B |
| 6 | S – D |
| 7 | S – A |
| 8 | C – D |

Tree level 0

Tree level 1

t = 4
t = 5
t = 6

A    B    S

t = 7

Tree level 2

S

Miss t = 4 slot
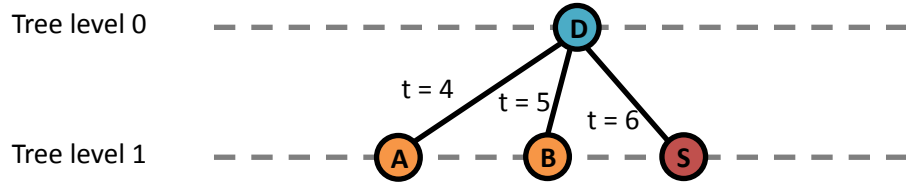
# HOHO Routing

Send a packet from S to D in the minimal time within max hops, max = 3 here

| Time Slot | Circuits |
|-----------|----------|
| 1 | S – H |
| 2 | A – F |
| 3 | H – E, F – A, S – B |
| 4 | S – G, A – D |
| 5 | B – D, G – B |
| 6 | S – D |
| 7 | S – A |
| 8 | C – D |

Tree level 0

Tree level 1

Tree level 2



t = 4

t = 5

t = 6

t = 3

t = 4

# HOHO Routing

Send a packet from S to D in the minimal time within max hops, max = 3 here

| Time Slot | Circuits |
|-----------|----------|
| 1 | S – H |
| 2 | A – F |
| 3 | H – E, F – A, S – B |
| 4 | S – G, A – D |
| 5 | B – D, G – B |
| 6 | S – D |
| 7 | S – A |
| 8 | C – D |



Tree level 0

Tree level 1

Tree level 2

Tree level 3

t = 4
t = 5
t = 6
t = 3
t = 4
t = 2

Repetitive A

# HOHO Routing

Send a packet from S to D in the minimal time within max hops, max = 3 here

| Time Slot | Circuits |
|-----------|----------|
| 1 | S – H |
| 2 | A – F |
| 3 | H – E, F – A, S – B |
| 4 | S – G, A – D |
| 5 | B – D, G – B |
| 6 | S – D |
| 7 | S – A |
| 8 | C – D |



Tree level 0 — D

Tree level 1 — A, B, S (t = 4, t = 5, t = 6)

Tree level 2 — E (t = 4)

Tree level 3

# HOHO Routing

Send a packet from S to D in the minimal time within max hops, max = 3 here



Tree level 0

Tree level 1

Tree level 2

Tree level 3

| Time Slot | Circuits |
|-----------|----------|
| 1 | S – H |
| 2 | A – F |
| 3 | H – E, F – A, S – B |
| 4 | S – G, A – D |
| 5 | B – D, G – B |
| 6 | S – D |
| 7 | S – A |
| 8 | C – D |

# HOHO Routing

Send a packet from S to D in the minimal time within max hops, max = 3 here

| Time Slot | Circuits |
|-----------|----------|
| 1 | S – H |
| 2 | A – F |
| 3 | H – E, F – A, S – B |
| 4 | S – G, A – D |
| 5 | B – D, G – B |
| 6 | S – D |
| 7 | S – A |
| 8 | C – D |



Tree level 0

t = 4   t = 5   t = 6

Tree level 1   A   B   S

t = 4

Tree level 2   E

t = 3

Tree level 3   H

Reach max hop count   t = 1

Tree level 4   S

# HOHO Routing

Send a packet from S to D in the minimal time within max hops, max = 3 here

| Time Slot | Circuits |
|-----------|----------|
| 1 | S – H |
| 2 | A – F |
| 3 | H – E, F – A, S – B |
| 4 | S – G, A – D |
| 5 | B – D, G – B |
| 6 | S – D |
| 7 | S – A |
| 8 | C – D |

Tree level 0

Tree level 1

Tree level 2

Tree level 3



t = 5

t = 6

# HOHO Routing

Send a packet from S to D in the minimal time within max hops, max = 3 here

| Time Slot | Circuits |
|-----------|----------|
| 1 | S – H |
| 2 | A – F |
| 3 | H – E, F – A, S – B |
| 4 | S – G, A – D |
| 5 | B – D, G – B |
| 6 | S – D |
| 7 | S – A |
| 8 | C – D |



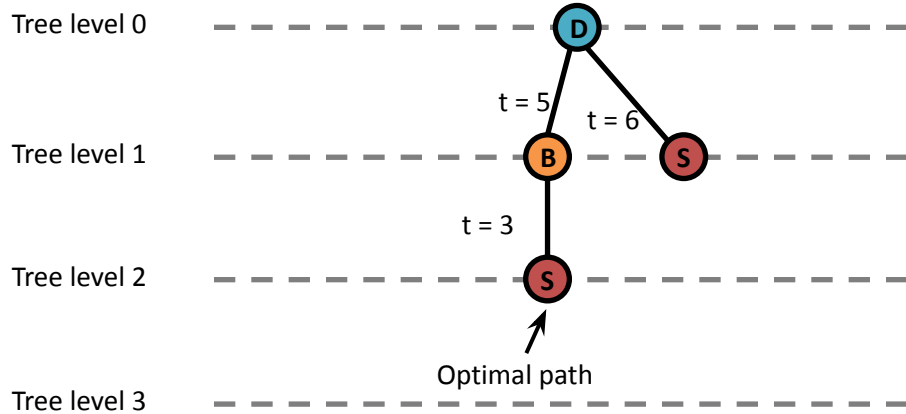Tree level 0

Tree level 1

t = 5

t = 6

Tree level 2

t = 3

Optimal path

Tree level 3

# HOHO Routing

Send a packet from S to D in the minimal time within max hops, max = 3 here

| Time Slot | Circuits |
|-----------|----------|
| 1 | S – H |
| 2 | A – F |
| 3 | H – E, F – A, S – B |
| 4 | S – G, A – D |
| 5 | B – D, G – B |
| 6 | S – D |
| 7 | S – A |
| 8 | C – D |

Tree level 0

Tree level 1

Tree level 2

Tree level 3

t = 5

t = 6

t = 3
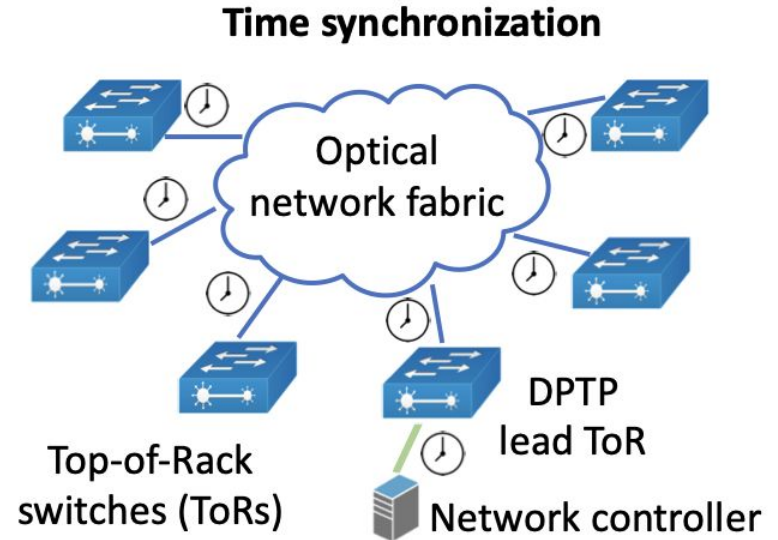
t = 5

t = 4

Optimal path

Feasible path

# HOHO Routing: Features

- Backtracking

    - Search from dst to src ToR

- Optimal

    - Generate the fastest and shortest path

- Offline

    - Calculate offline and build routing table

- Decouple algorithm and run-time system designs

    - Assume no queuing delay at offline calculation and consider this at run-time system

# System Design: Time Synchronization

- Every ToR needs to know when to send packet
- Leverage existing protocols
  - DPTP[1], nanosecond-level synchronization precision



**Time synchronization**

Optical network fabric

Top-of-Rack switches (ToRs)

DPTP lead ToR

Network controller

[1] Precise time-synchronization in the data-plane using programmable switching ASICs, SOSR' 19
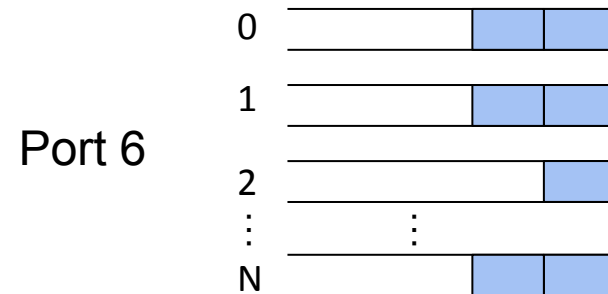
# System Design: Routing Lookup

- HOHO generates a routing table

  - Current time slot: when the packet arrives

  - Dst ToR: final destination ToR

  - Egress port

  - Send time slot: when to send out the packet

| Current Time Slot | Dst ToR | Egress Port | Send Time Slot |
|---|---|---|---|
| 0 | 1 | 6 | 0 |
| 0 | 2 | 23 | 4 |
| 1 | 1 | 6 | 2 |
| ... | | | |
| 7 | 0 | 12 | 10 |
| ... | | | |

# System Design: Routing Lookup

| Current Time Slot | Dst ToR | Egress Port | Send Time Slot |
|---|---|---|---|
| 0 | 1 | 6 | 0 |
| 0 | 2 | 23 | 4 |
| 1 | 1 | 6 | 2 |
| ... | | | |
| 7 | 0 | 12 | 10 |
| ... | | | |

- Look up current time slot and dst

  ToR to get egress port and send time

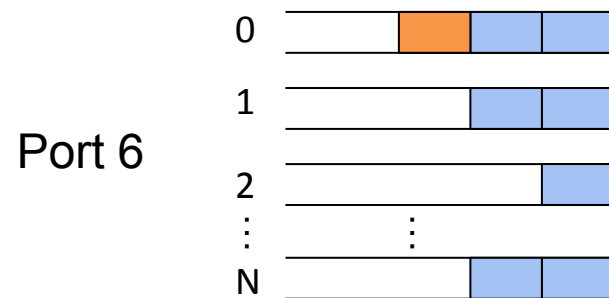  slot

  - Packet with current time slot 0

    and dst ToR 1

Port 6

0

1

2
⋮
N

14

# System Design: Routing Lookup

| Current Time Slot | Dst ToR | Egress Port | Send Time Slot |
|---|---|---|---|
| 0 | 1 | 6 | 0 |
| 0 | 2 | 23 | 4 |
| 1 | 1 | 6 | 2 |
| ... | | | |
| 7 | 0 | 12 | 10 |
| ... | | | |

- Look up current time slot and dst ToR to get egress port and send time slot

  - Packet with current time slot 0 and dst ToR 1

  - Buffer at queue 0

Port 6

# System Design: Routing Lookup

| Current Time Slot | Dst ToR | Egress Port | Send Time Slot |
|---|---|---|---|
| 0 | 1 | 6 | 0 |
| 0 | 2 | 23 | 4 |
| 1 | 1 | 6 | 2 |
| ... | | | |
| 7 | 0 | 12 | 10 |
| ... | | | |

- If missing the planned send time slot due to the queue occupancy, look up from next current time slot

Queue 0 is full

Port 6

0

1

2

...

N

14

# System Design: Routing Lookup

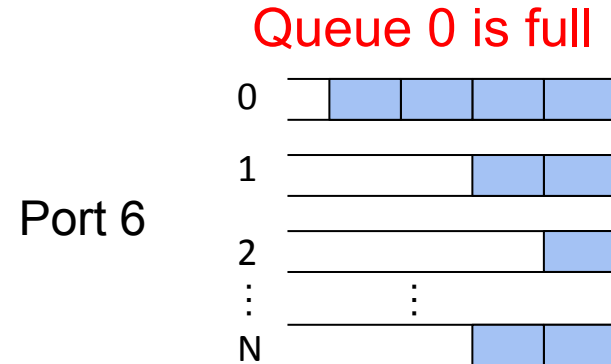| Current Time Slot | Dst ToR | Egress Port | Send Time Slot |
|:---:|:---:|:---:|:---:|
| 0 | 1 | 6 | 0 |
| 0 | 2 | 23 | 4 |
| 1 | 1 | 6 | 2 |
| ... | | | |
| 7 | 0 | 12 | 10 |
| ... | | | |

- If missing the planned send time slot due to the queue occupancy, look up from next current time slot

  - Move to current time slot 1

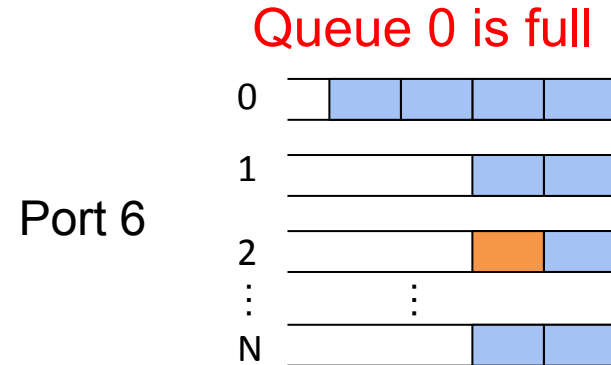Queue 0 is full

Port 6

# System Design: Routing Lookup

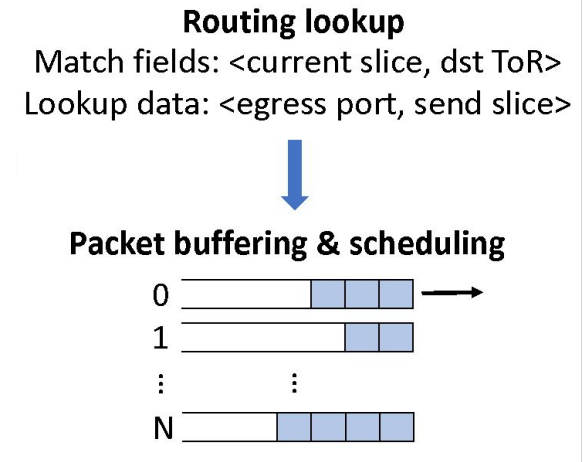| Current Time Slot | Dst ToR | Egress Port | Send Time Slot |
|:---:|:---:|:---:|:---:|
| 0 | 1 | 6 | 0 |
| 0 | 2 | 23 | 4 |
| 1 | 1 | 6 | 2 |
| ... | | | |
| 7 | 0 | 12 | 10 |
| ... | | | |

- If missing the planned send time slot due to the queue occupancy, look up from next current time slot

  - Move to current time slot 1
  - Buffer at queue 2

Queue 0 is full

Port 6

0

1

2
:
N

# System Design: Buffering

- Packet is buffered at the queue before its sending time slot
- Could be realized by queue pause



**Routing lookup**
Match fields: <current slice, dst ToR>
Lookup data: <egress port, send slice>

**Packet buffering & scheduling**

# Simulation

## Setup

- Reused the setup in Opera paper[1]

    - Topology: 108 ToRs and 648

      servers, each ToR with six 10G

      downlinks to servers and six 10G

      uplinks to optical fabric

    - Workload: 1%~40% data-mining
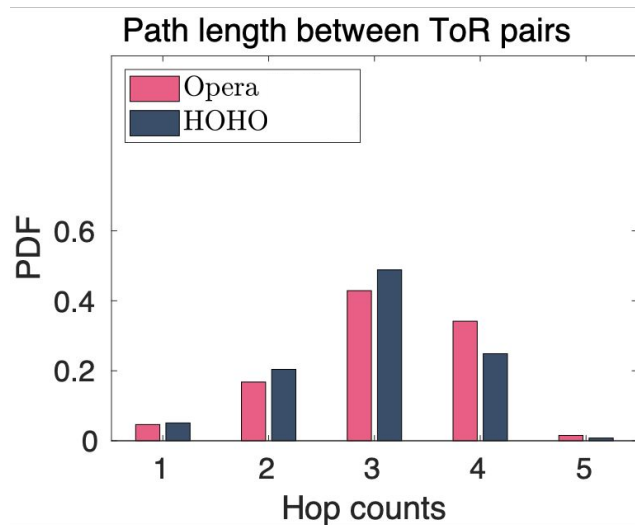
      traffic from Microsoft

[1] Expanding across time to deliver bandwidth efficiency and low latency, NSDI' 20

# Simulation

## Setup

- Reused the setup in Opera paper[1]

- Topology: 108 ToRs and 648 servers, each ToR with six 10G downlinks to servers and six 10G uplinks to optical fabric

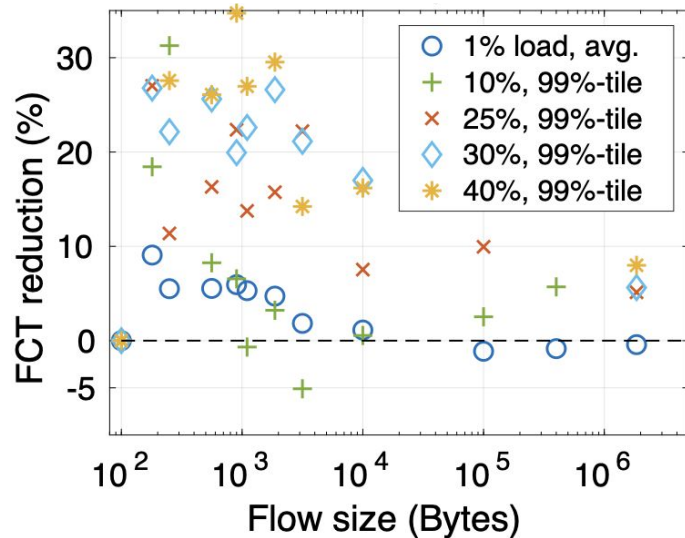- Workload: 1%~40% data-mining traffic from Microsoft

## Shorter paths



Path length between ToR pairs

- Avg. hops:  3.11 => 2.80
- >= 4 hops:  37% => 25%

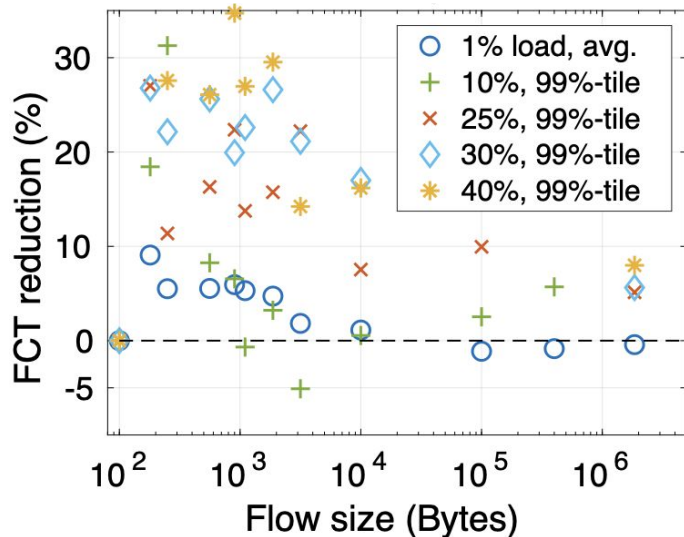[1] Expanding across time to deliver bandwidth efficiency and low latency, NSDI' 20

# Simulation

## Lower latency



- FCT reduction: up to <span style="color:red">35%</span>
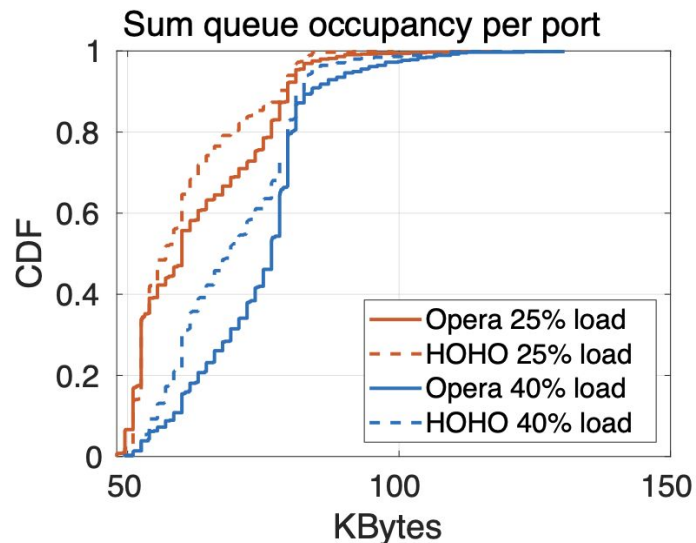
# Simulation

## Lower latency



- FCT reduction:  up to <span style="color:red">35%</span>

## Lower queue occupancy



- Queue occupancy reduction:  5% (25% load) and 10% (40% load)

# Summary

- Hop-On Hop-Off optical circuits

  - Allow packets to "wait" at ToRs

- HOHO routing algorithm

  - Works on any optical schedule

  - Optimal (fastest and shortest path)

- A system sketch

  - Time synchronization + routing lookup + buffer management

- Simulation

  - Shorter paths, lower FCT, lower queue occupancy

THANKS!