# Reshaping Transport and Traffic Engineering in Reconfigurable Data Center Networks

**Shawn Shuoshuo Chen**

Workshop on Reconfigurable Networks
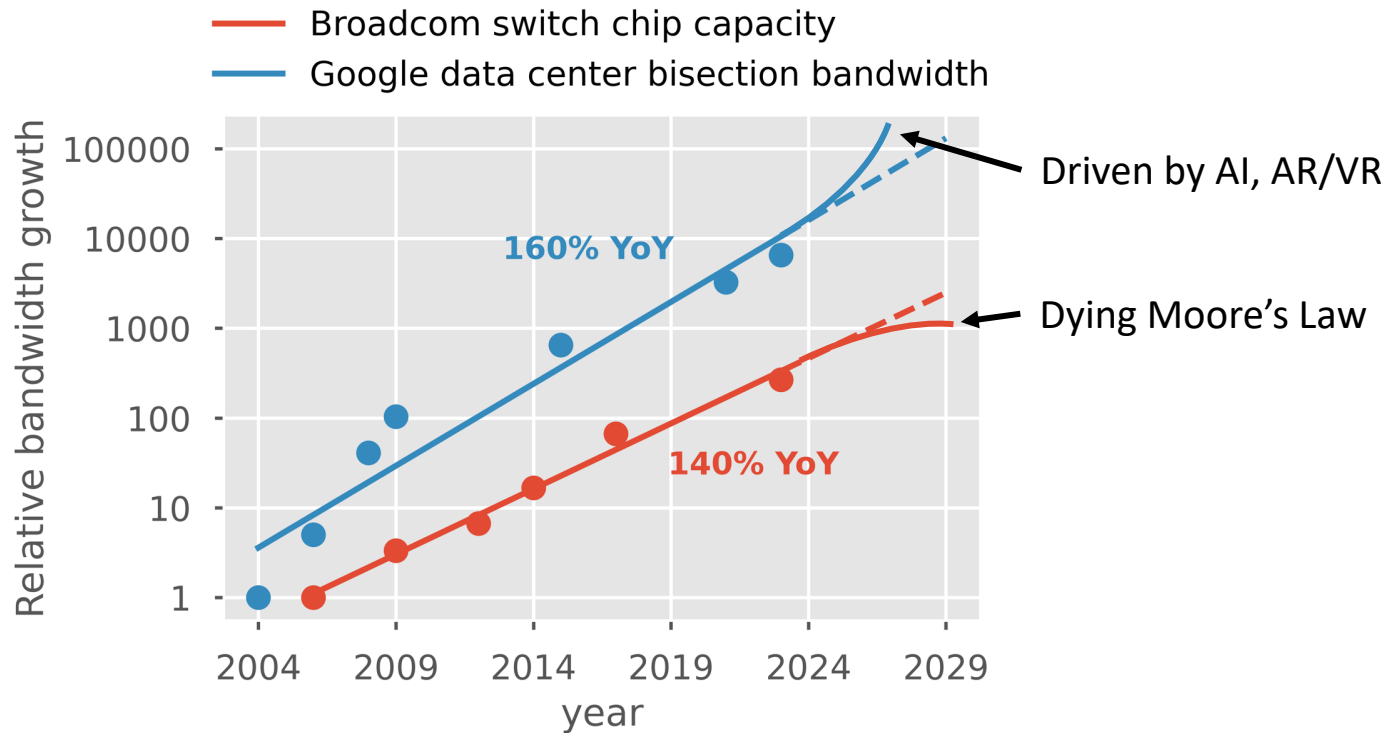June 2, 2025

**Carnegie Mellon University**

# The **scaling crisis** of data center networks

# The **scaling crisis** of data center networks

Broadcom switch chip capacity
Google data center bisection bandwidth

Switch
Host



Relative bandwidth growth

100000
10000
1000
100
10
1

2004   2009   2014   2019   2024   2029
year

Gap covered
switches & h

Not sustainable
- High power consumption
- Expensive

# What is optical circuit switch (OCS)?



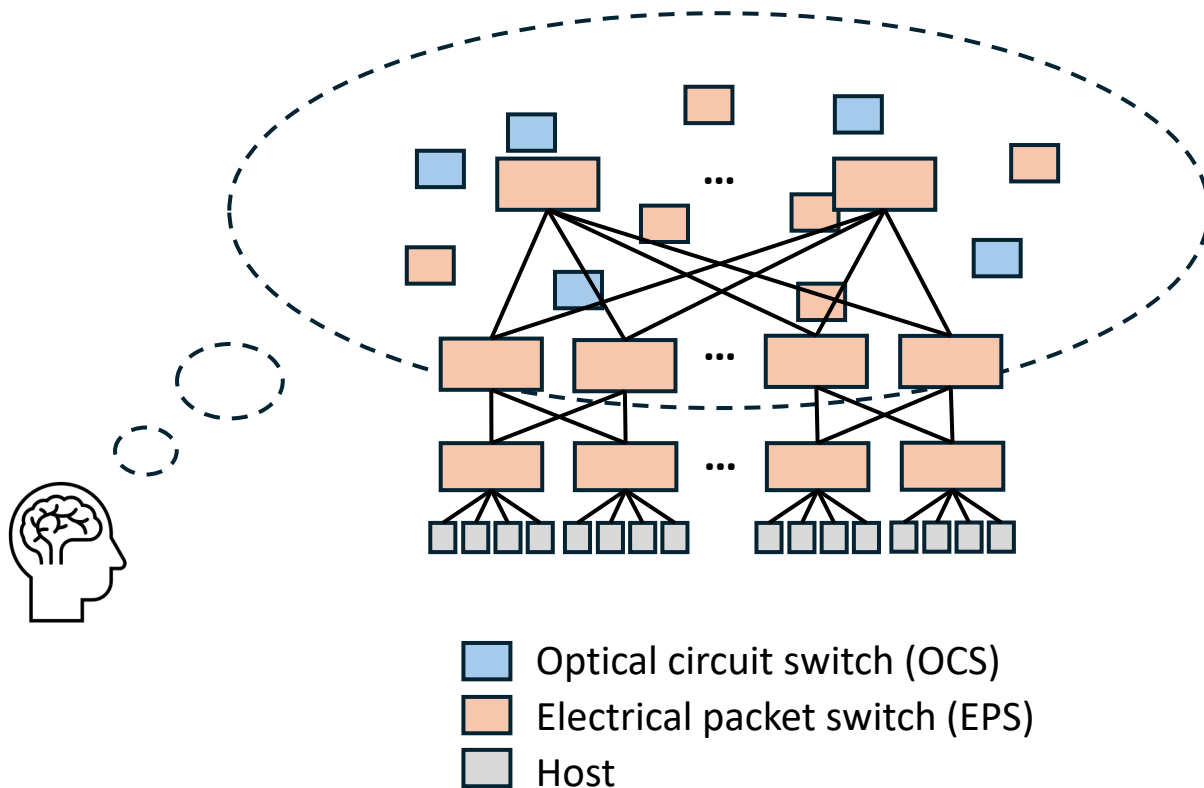| Optical circuit switch |
| --- |
| 1:1 in/out mapping |
| Down during **reconfiguration** (μs-ms) |
| Much higher b/w, lower latency |
| Data rate agnostic |

| Electrical packet switch |
| --- |
| Packet-level multiplexing |
| No down period |
| Lower b/w, higher latency |
| Fixed rate per generation |

# Reconfigurable Data Center Networks (RDCNs)



Optical circuit switch (OCS)
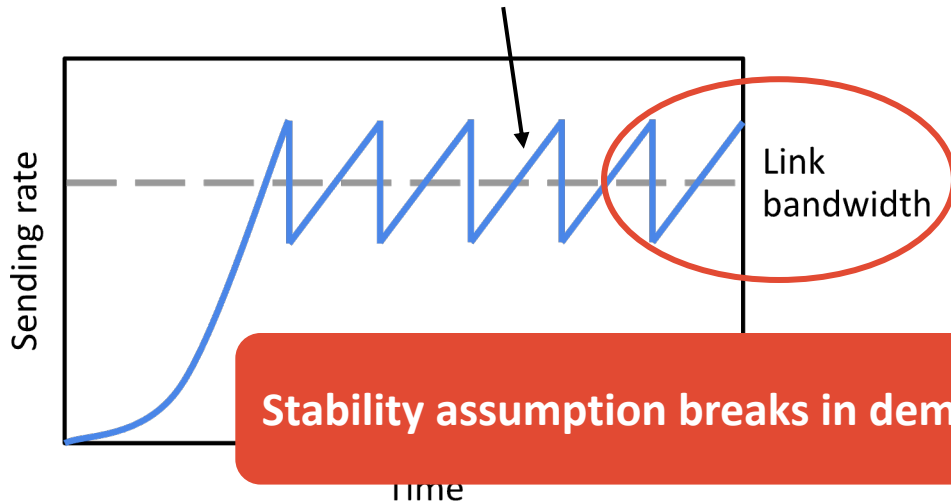Electrical packet switch (EPS)
Host

Today's talk

- **Transport**
  - Time-division TCP for demand-oblivious RDCNs *[SIGCOMM'22]*

- **Traffic engineering**
  - Precise traffic engineering for demand-aware RDCNs *[NSDI'24]*

# Existing TCP's assumption: stable network path

TCP's goal: match sending rate to available bandwidth.

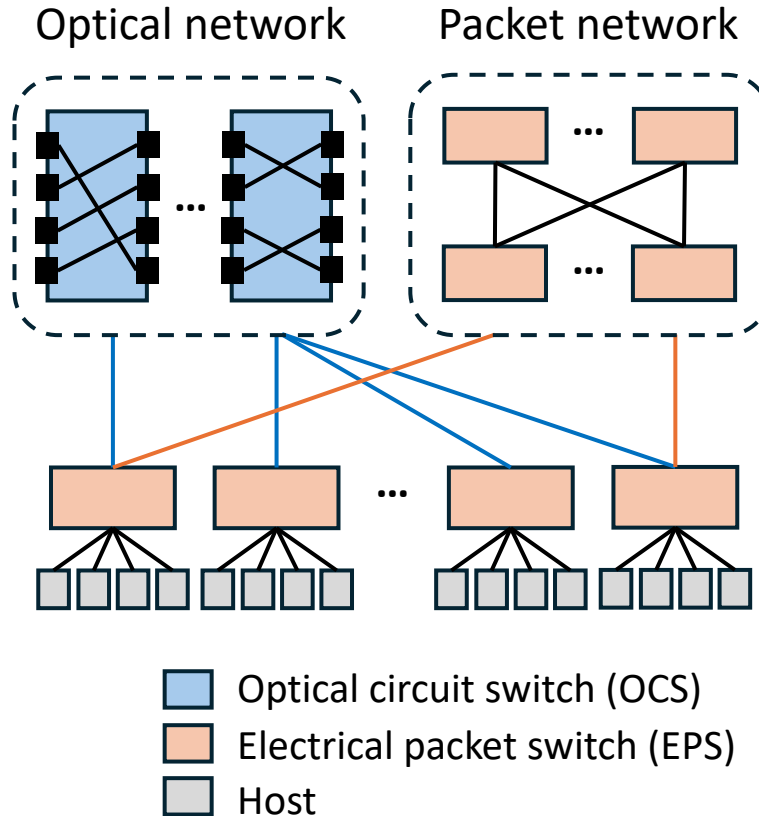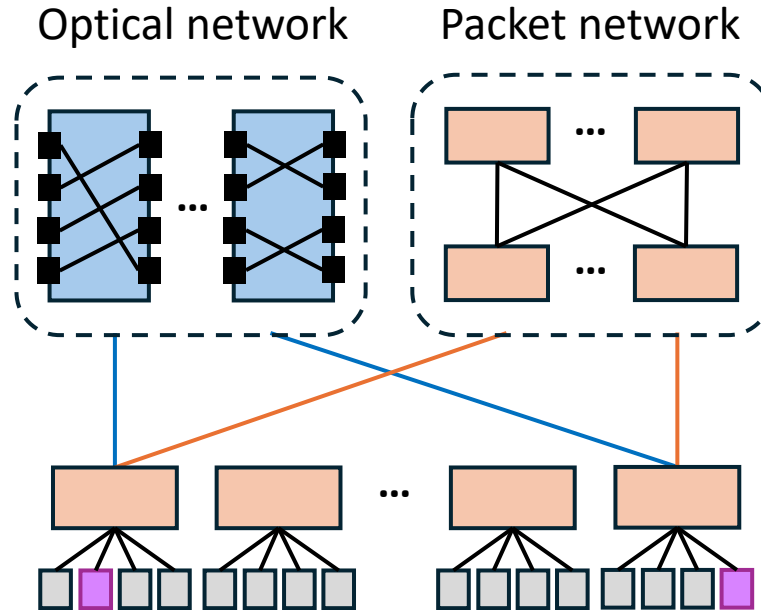Sending rate

Time

Link bandwidth

TCP's mechanism
- Probe & converge
  - On round-trip time (RTT) scale
- Model path characteristics
  - *cwnd*: sending rate

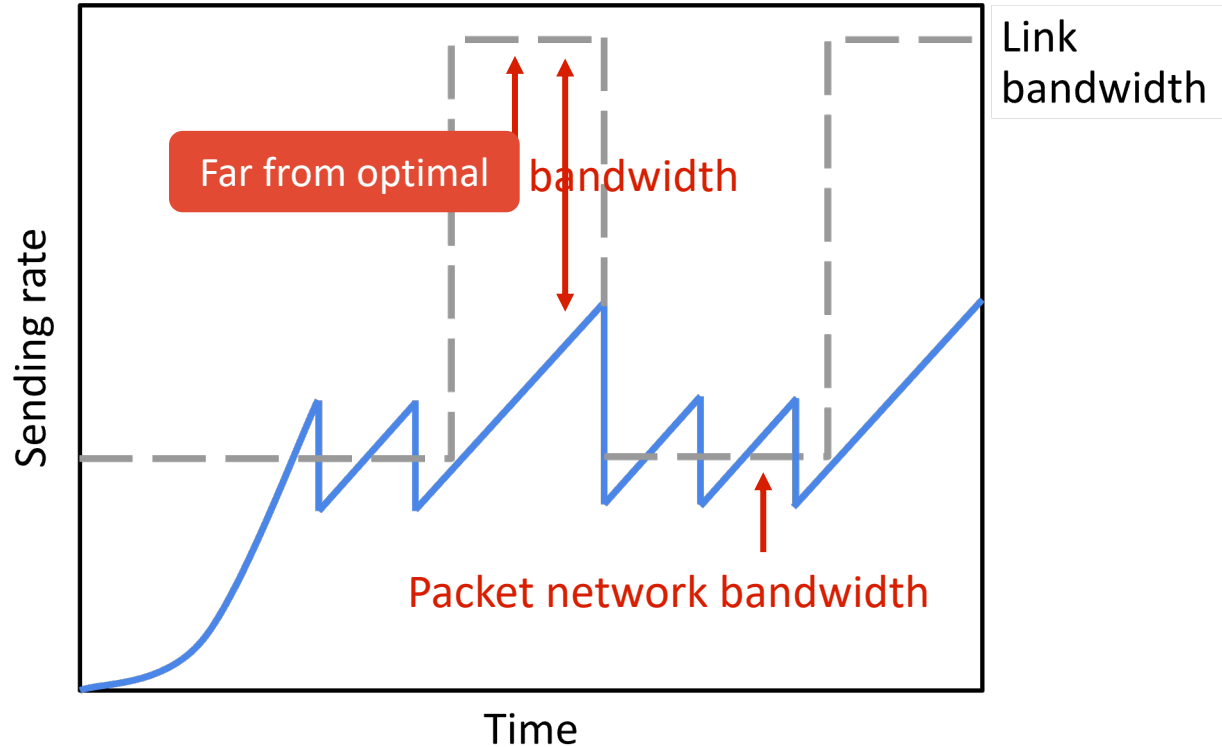**Stability assumption breaks in demand-oblivious RDCN.**
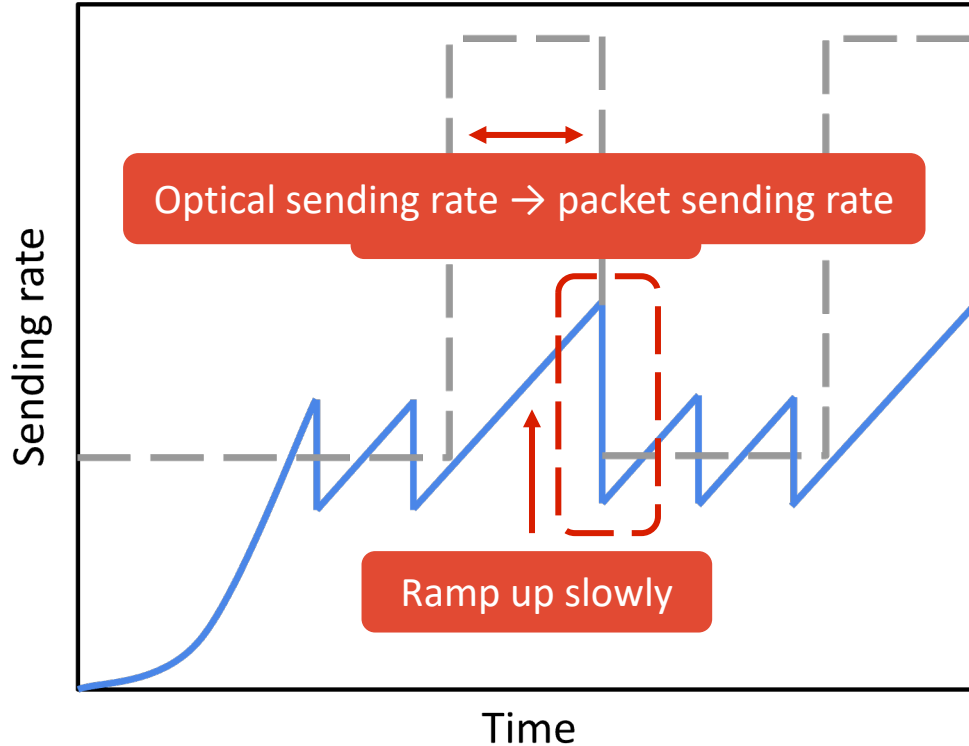
# Demand-oblivious RDCN

# Demand-oblivious RDCN

# TCP performs poorly under invalid assumption



Link bandwidth

Far from optimal bandwidth

Sending rate

Packet network bandwidth

Time

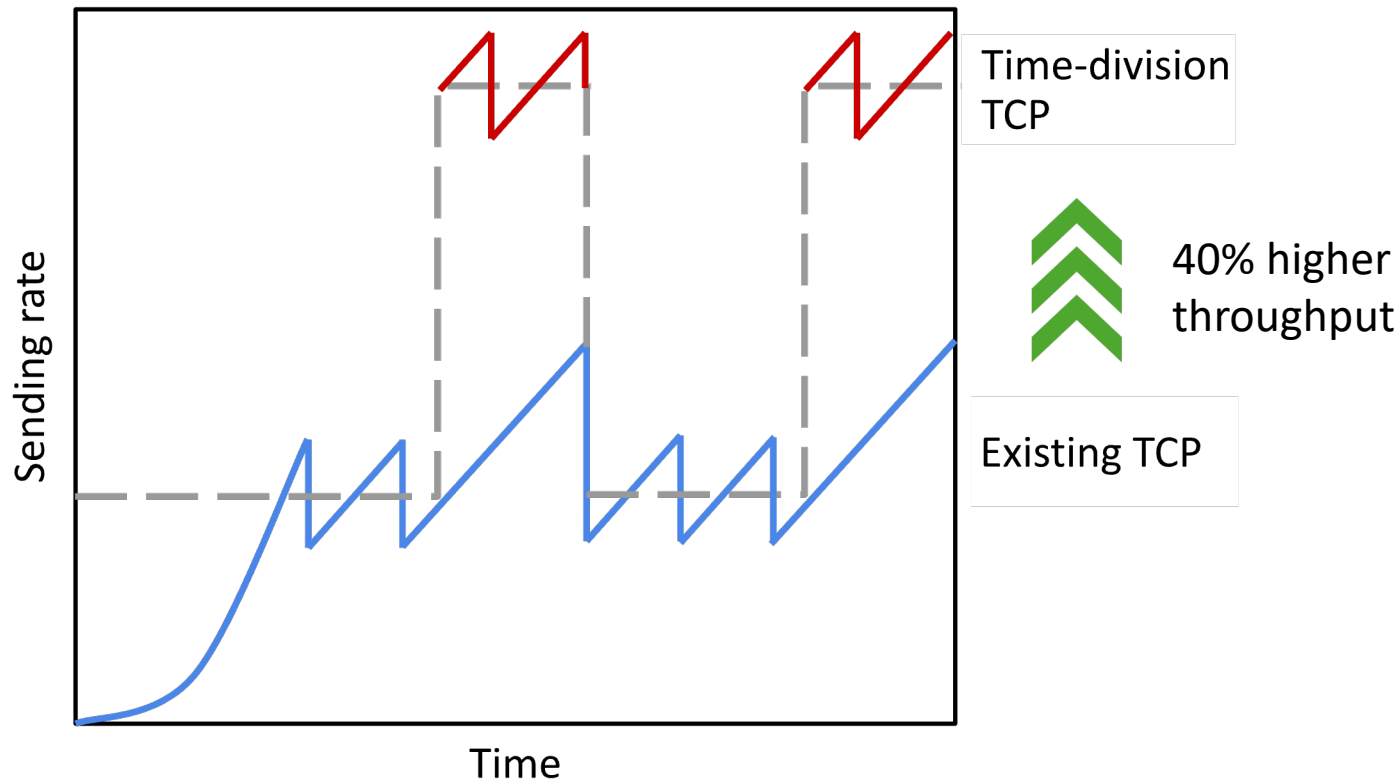# TCP performs poorly under invalid assumption



What happened?
1) Reactive probing
2) Insufficient time to converge
3) Overwritten states

# Our proposal: Time-division TCP

| | Existing TCP | Time-division TCP |
|---|---|---|
| Change discovery | Reactive<br>- in-band, probing | Proactive<br>- out-of-band, switch notification |
| Path modeling | One state:<br>- *cwnd, srtt* | 2 (N) states:<br>- *cwnd[], srtt[]* |

# Time-division TCP outperforms existing TCP



Sending rate

Time

Time-division TCP
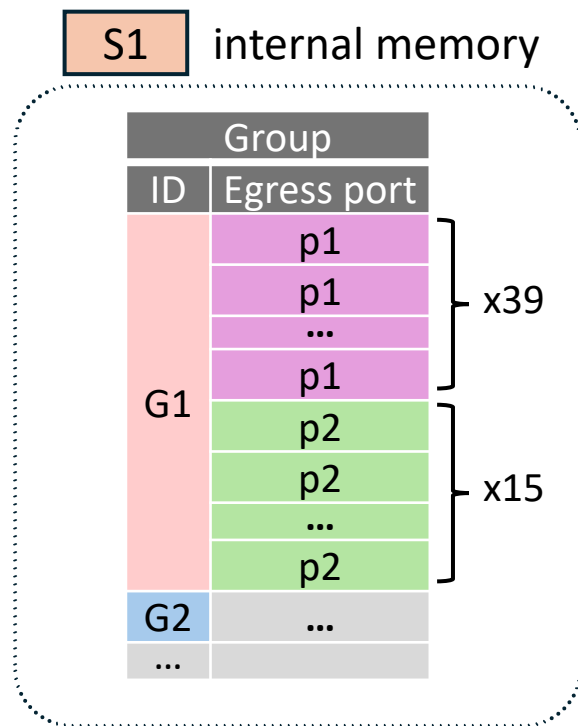
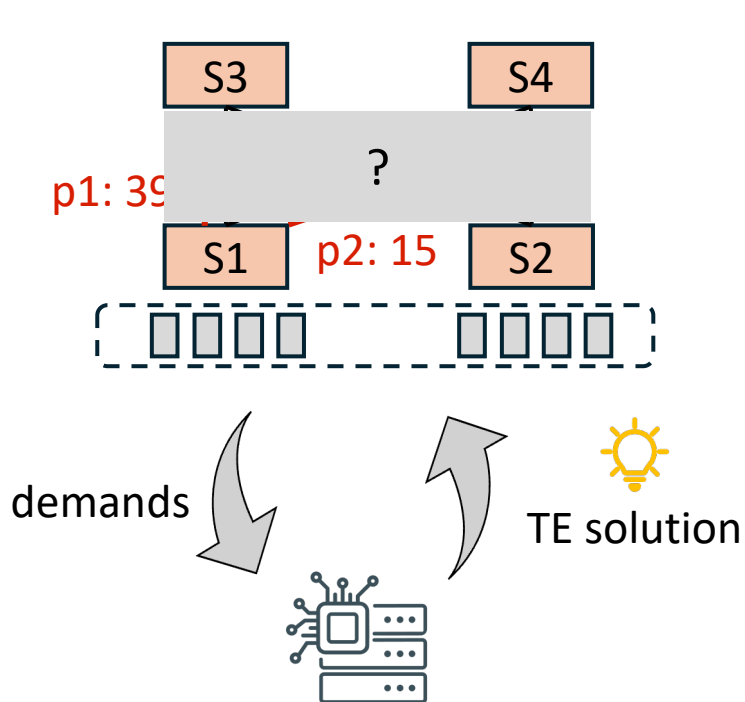40% higher throughput

Existing TCP

Today's talk

- **Transport**
  - Time-division TCP for demand-oblivious RDCNs *[SIGCOMM'22]*

- **Traffic engineering**
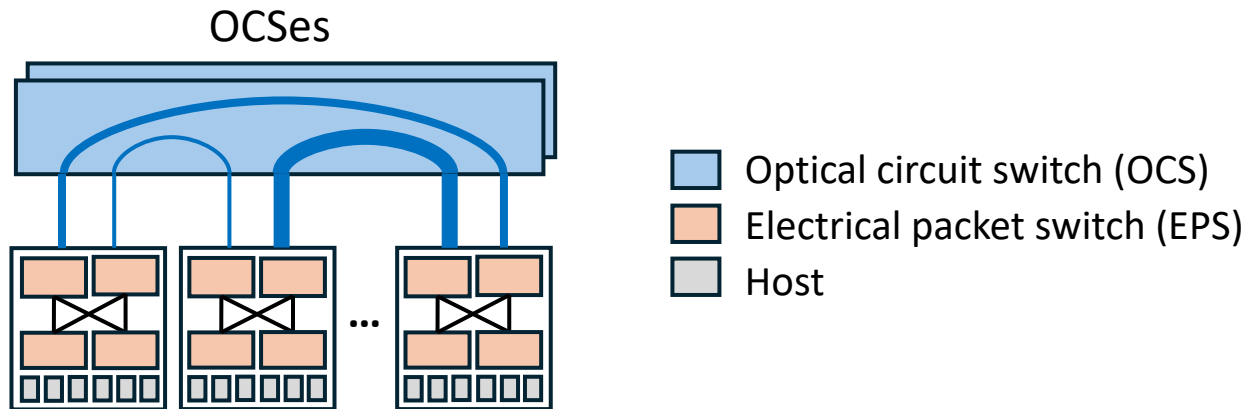  - Precise traffic engineering for demand-aware RDCNs *[NSDI'24]*

# How does a traffic engineering (TE) system work?

# Demand-aware RDCN
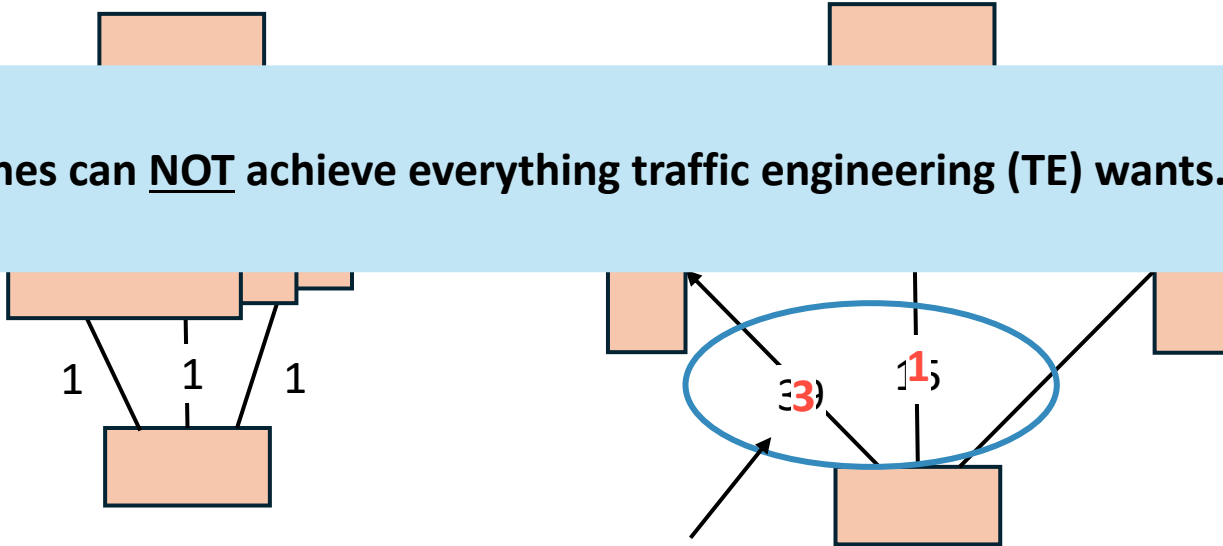
Demand-awareness introduces heterogeneity
- Skewed traffic distribution (weight ratios)
- Ratios have different actual impact on traffic

OCSes



Optical circuit switch (OCS)
Electrical packet switch (EPS)
Host

# Traffic engineering's assumption: omnipotent switches



Switches can **NOT** achieve everything traffic engineering (TE) wants.

Spine switches

1    1    1

3    1

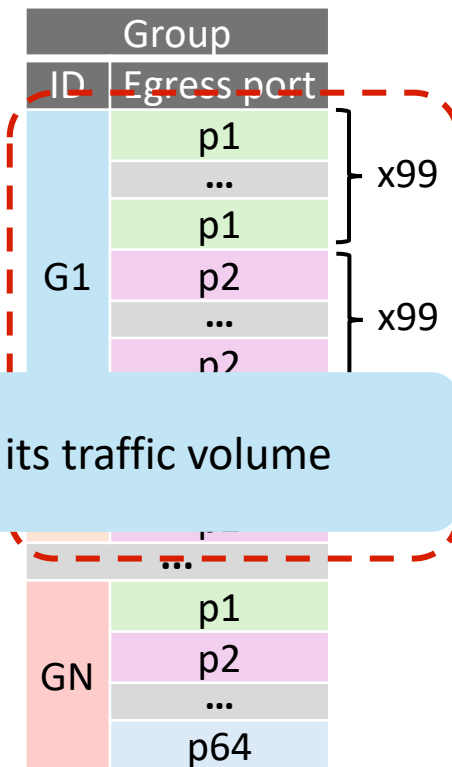Skewed ratio

# Heuristics to reduce group space usage

Insights

**Our heuristics**

| Group | |
|---|---|
| ID | Egress port |

Different groups contribute to the overall traffic imbalance differently

**Table Carving**

| G1 | p1 |
|---|---|
| | ... | } x99
| | p1 |
| | p2 |
| | ... | } x99
| | p2 |

**Table Carving**: allocate space to each group proportional to its traffic volume

...

| GN | p1 |
|---|---|
| | p2 |
| | ... |
| | p64 |

# Heuristics to reduce group space usage

Insights

**Our heuristics**

| Group | |
|---|---|
| ID | Egress port |
| G1 | p1 |
| | ... |
| | p1 |
| | p2 |
| | ... |
| | p2 |
| | p3 |
| G2 | p1 |
| | p2 |
| | ... |
| | p64 |

Different groups contribute to the overall traffic imbalance differently

**Table Carving**

Not all ports need to be preserved

**Group Pruning**

x99

x99

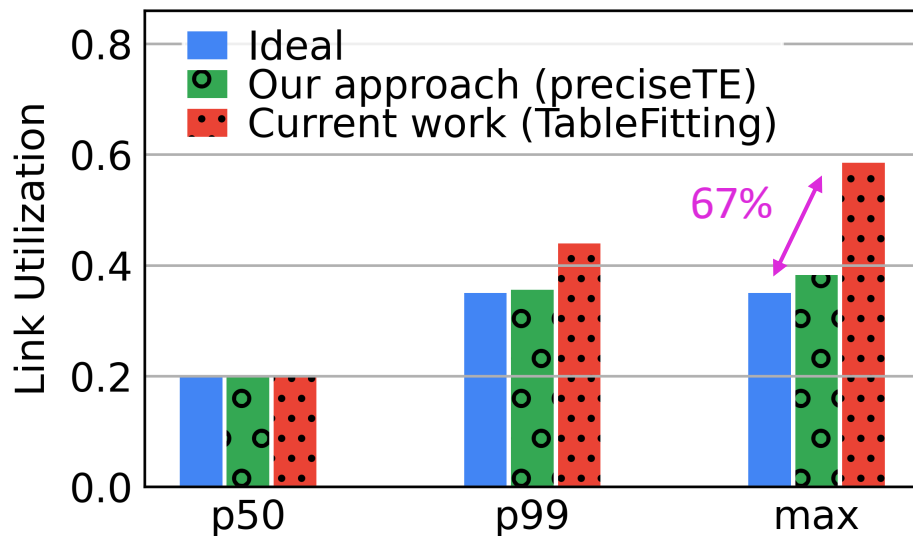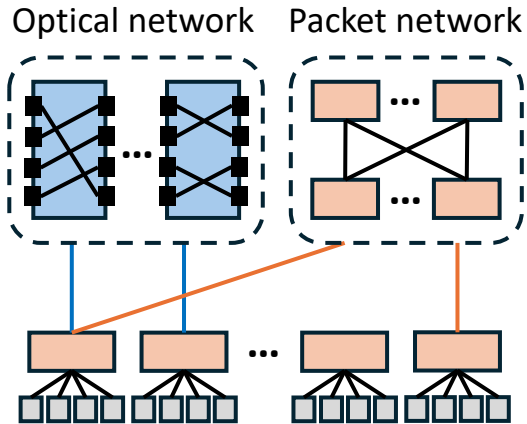**Group Pruning**: prune select ports from a group to enable size reduction

# Our approach is more precise than current work.

- preciseTE **7% error** vs. TableFitting **67% error**
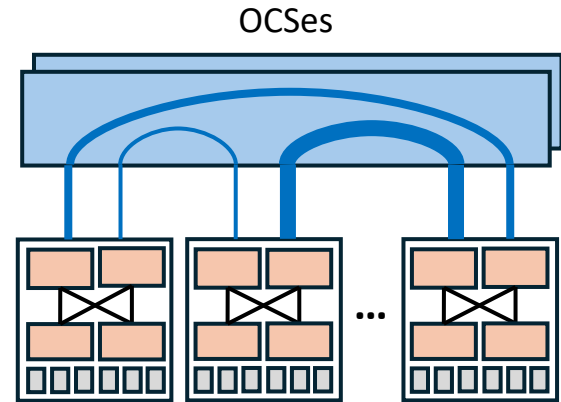- Being deployed at Google

# Summary



Demand-oblivious RDCN

Optical network    Packet network

Transport: coordination

Demand-aware RDCN

OCSes

TE: managing heterogeneity

# Future direction: All-optical RDCN

- Fast OCS (optical packet switching)
- Fully scheduled, source-routed network
- Scheduling challenge
  - "Incast" avoidance