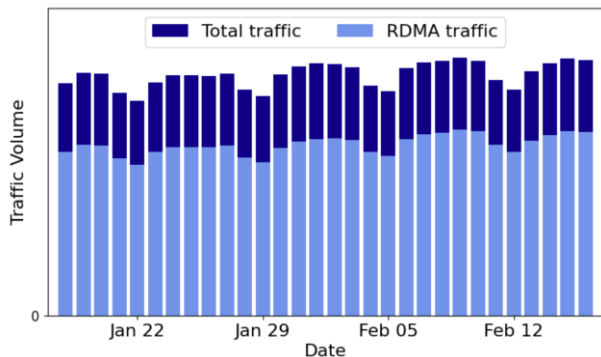# Revisiting Congestion Control for Lossless Etherent

Yiran Zhang, Qingkai Meng, Chaolei Hu, Fengyuan Ren

# Lossless Ethernet





Wide adoption of RDMA

*[1] Wei Bai, Shanim Sainul Abdeen, Ankit Agrawal et al. Empowering Azure Storage with RDMA. NSDI 2023*

**RDMA + PFC: lossless Ethernet**
No packets dropping
Full potential of RDMA

**But PFC comes with side effects!**
Head of line (HoL) blocking,
Deadlock, etc

**Congestion control is a key enabler for lossless Ethernet at scale**

# Congestion Control in Lossless Ethernet

|  | **Congestion Detection** | **Rate Control** |
|---|---|---|
| **DCQCN**[SIGCOMM'15]<br>**TIMELY**[SIGCOMM'15] | Following lossy networks (ECN, RTT) | Traditional heuristic rules |
| **HPCC** [SIGCOMM'19] | Advanced telemetry technique | Larger overhead |

# TCD: Ternary Signal for Lossless Networks



*Yiran Zhang, Yifan Liu, Qingkai Meng, Fengyuan Ren.*
*Congestion Detection in Lossless Networks.*
*SIGCOMM'21*

✓ **Congested flows (CE)**
✓ **Victim/Undetermined flows (UE)**
✓ **Uncongested flows (NO)**

# Congestion Control in Lossless Ethernet

| | **Congestion Detection** | **Rate Control** |
|---|---|---|
| **DCQCN**[SIGCOMM'15]<br>**TIMELY**[SIGCOMM'15] | Following lossy networks (ECN, RTT) | Traditional heuristic rules |
| **HPCC** [SIGCOMM'19] | Advanced telemetry technique | Larger overhead |
| **TCD** [SIGCOMM'21] | Precise ternary signal for lossless Ethernet | Traditional heuristic rules |

# Congestion Control in Lossless Ethernet

**DCQCN**[SIGCOMM'15]

**TIMELY**[SIGCOMM'15]

**HPCC** [SIGCOMM'19]

**TCD** [SIGCOMM'21]

Congestion Detection    Rate Control

Precise ternary signal
for lossless Ethernet

Traditional
heuristic rules

**Sub-optimal performance**

☹ **Suffering HoL blocking**
☹ **Hard to achieve low latency and high throughput, etc**

# A Desirable Congestion Control for Lossless Ethernet

**Congestion Detection**  **+**  **Rate Control**  **=**  **Congestion Control**

**Tailored for lossless Ethernet**  **Tailored for lossless Ethernet**  **High-performance for lossless Ethernet**

- ✓ Fast convergence to alleviate HoL blocking, deadlock etc.
- ✓ Low latency
- ✓ High throughput

*Can we rethink congestion control for lossless Ethernet by* <span style="color:red">*taking full advantage of its intrinsic properties?*</span>

# Revisiting the Impact of PFC

Host

Host

*Injected :
15 packets*

*Pipe capacity :
5 packets*

# Revisiting the Impact of PFC



Host    Switch    PAUSE    Switch    PAUSE    Switch    Host

*5 packets flying, 10 packets queueing*

*Injected : 15 packets*

*Finally ejected: 15 packets*

**An end-to-end lossless network path**

**Packet Conservation Property**
- Number of injected packets = number of ejected (acked) packets
- All injected packets are either flying or queueing

# Packet Conservation Empowers ACK-Driven

Simulation:
$H_1$-$H_N$: concurrent burst

F0: victim flow
F1: congested flow



**Queue size**

**Total queue size**

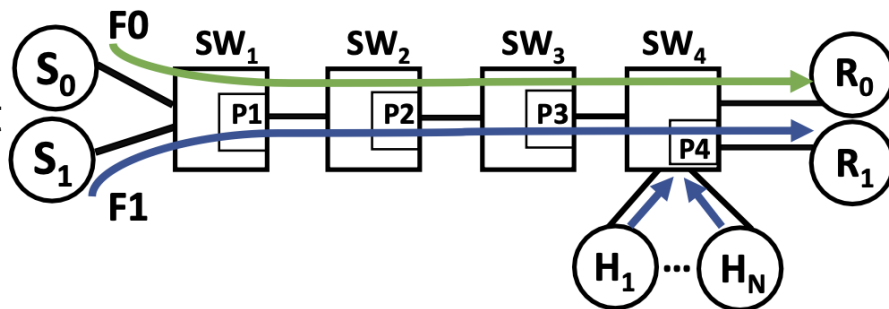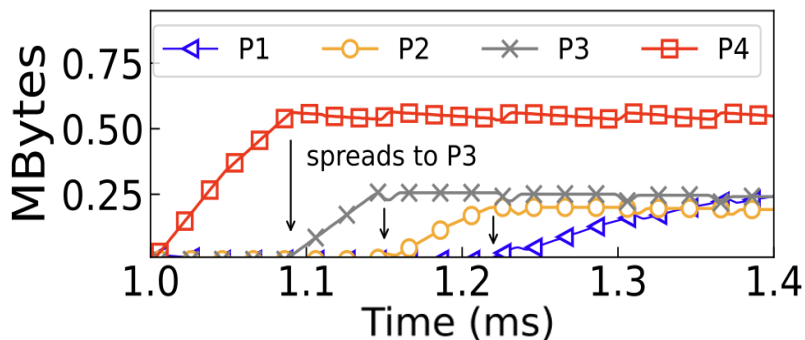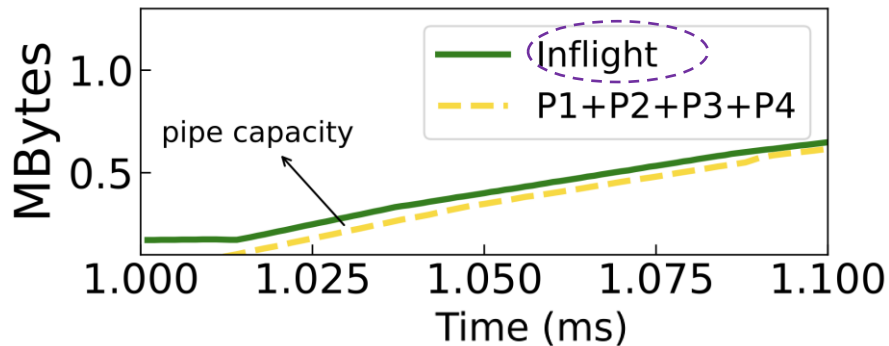**Excessive packets queueing in switches = inflight packets - network pipe capacity**

# Packet Conservation Empowers ACK-Driven

Simulation：
$H_1$-$H_N$: concurrent burst

F0: victim flow
F1: congested flow

**ACK arriving pattern**

**ACK arriving rate**

**The ACK arrival rate can imply the available bandwidth**

# Handling HoL Blocking

**Queue occupancy at a HoL blocked port.**
F0: victim flow. F1: congested flow

**Stopping only congested flow F1**
enough can eliminate HoL blocking



**Stopping congested flows sufficiently long can empty buffers as soon as possible**

# Handling HoL Blocking

**If HoL blocking is more severe:**
F0: victim flow. F1: congested flow

**Only stopping congested flow F1**
can **not** eliminate HoL blocking:



**Whether should throttling victim flows depends on the extent of congestion**

# Summary of Principles

✓ **The ACK-driven paradigm should be renewed**:
   Inferring the proper <span style="color:red">throttled rate</span> and the <span style="color:red">precise number of excessive packets</span> for congested flows


✓ **Handling HoL blocking needs individual rules**:
   <span style="color:red">Stopping congested flows sufficiently long</span> is the foremost means to suppressing HoL blocking
   <span style="color:red">Victim flows</span> should <span style="color:red">adapt to the severity</span> of congestion

# ACK-Driven Congestion Control (**ACC**)

**Ternary signal provided by TCD**

Receiving ACK

CE → **Congested**

ACK sequence

ACK arrival rate

→ Source Halt → *Rate decrease* → Sending

NO → **Uncongested** → *Rate increase* → Sending

UE → **Undetermined** → *Rate keep or decrease* → Sending

ACC  State Machine

# Halting and Throttling Congested Flows

- *First* halting to wait for the excessive packets to drain out
- *Then* matching the rate to the pipe capacity
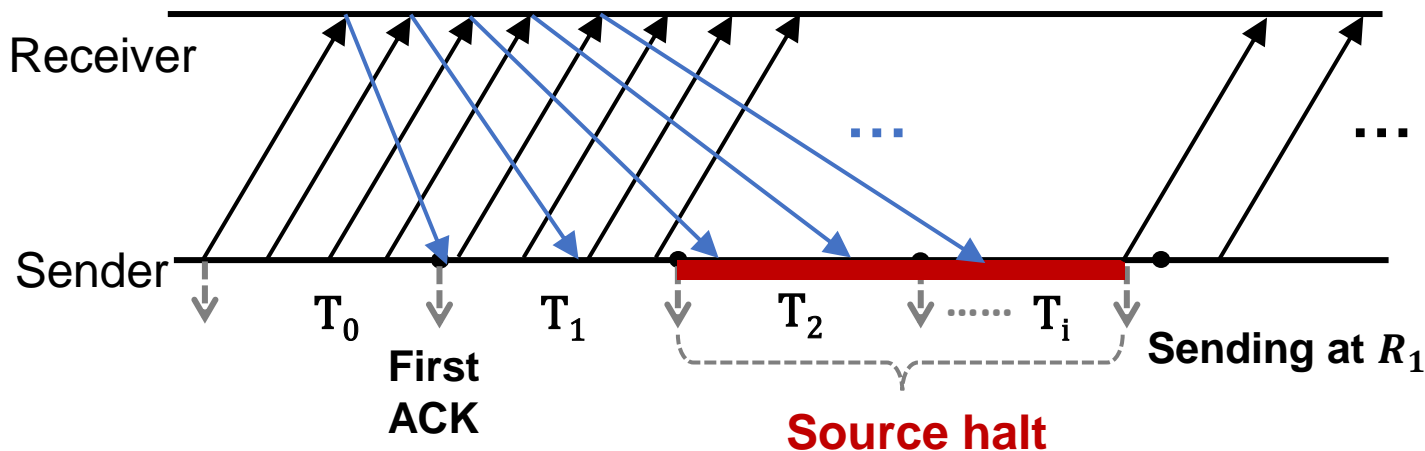
Receiver

Sender

···   ···

$T_0$

First ACK

$T_1$

$T_2$   ······ $T_i$
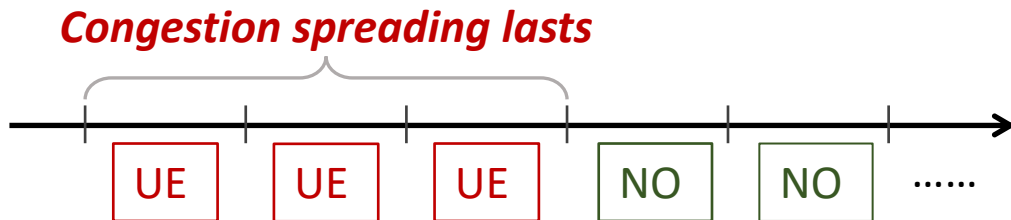
Sending at $R_1$

**Source halt**

$$t_{Halt} = (\Delta Sent_0 - \Delta Ack_1) / R_1$$

*Refer to paper for more details*
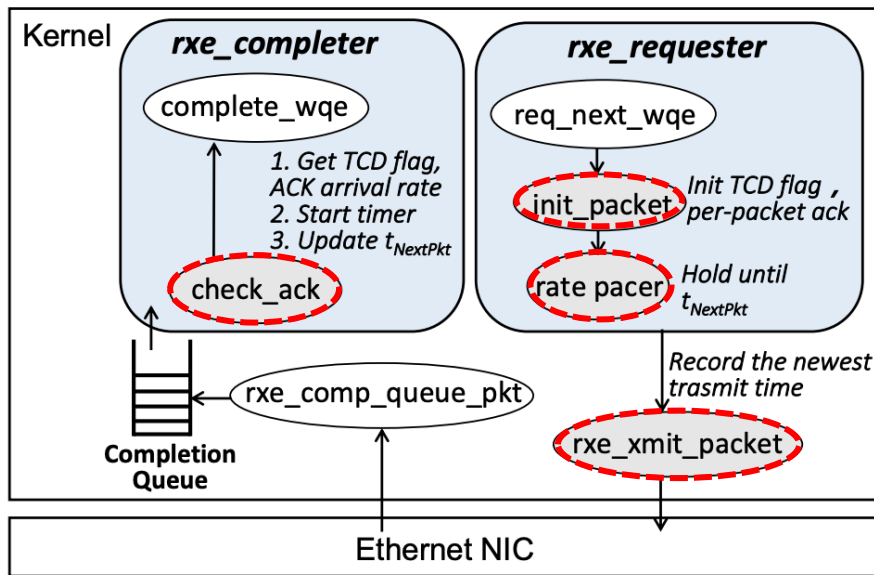
# Victim Flows & Uncongested Flows

- **Victim flows: adapt to the severity of congestion spreading**
  *Indicator: the number of consecutive periods with UE marks*
  *Exceeding a threshold $P_{thresh}$: decreasing the injection rate*

  **Congestion spreading lasts**
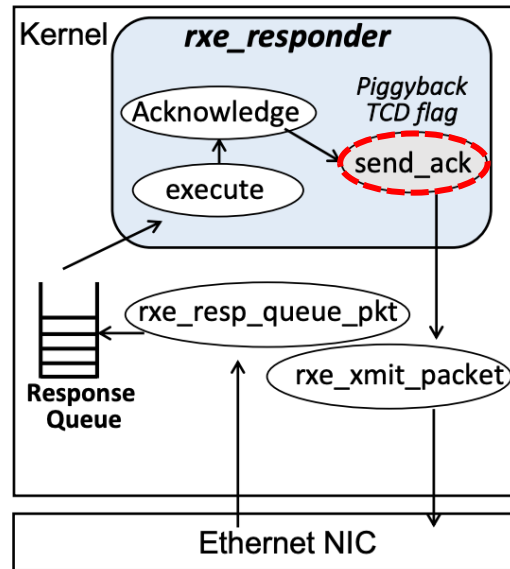
  | UE | UE | UE | NO | NO | ...... |

- **Uncongested flows:**
  *increasing gradually at first and then aggressively*

# Implementation

SoftRoCE:  software implementation of RDMA



ACC sender

ACC receiver

# Implementation

**Testbed**
- 242 and 3 lines of code added in SoftRoCE sender and receiver
- 119 lines of code added to SoftRoCE common library
- 5 hosts with Intel 82599ES 10GbE NIC + 1 Tofino switch

**Simulator**
- Customized NS3 packet simulator
- Fat-tree network with 320 servers in 20 racks
- 100Gbps/400Gbps

# Evaluation Summary
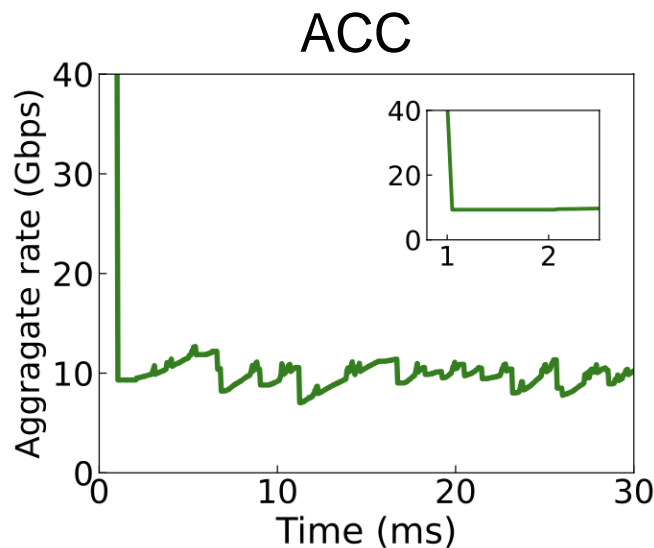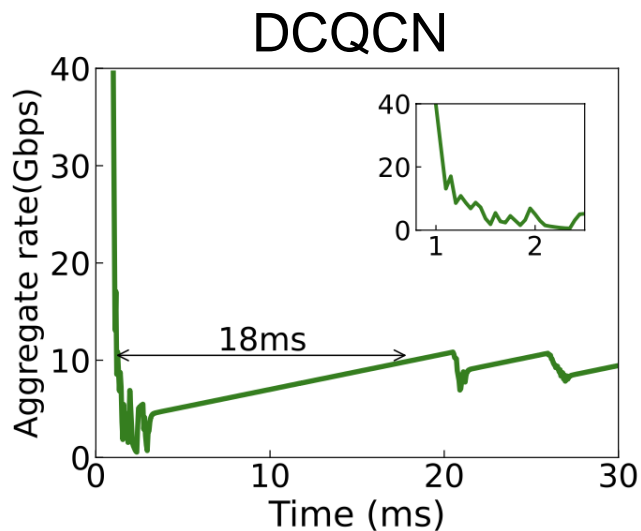
- **Basic Properties**
  - ✓ **25X faster** Convergence
  - ✓ **Full** link utilization
  - ✓ **1.35X faster** emptying the queues and **suppressing HoL blocking**
  - ✓ Effectively **prevent deadlocks**
  - ✓ Good fairness
  - ✓ Proper parameters (UE periods threshold for victim flows, etc)
- **FCT Performance**
  - ✓ **1.3~3.3X** and **1.4~11.5X** better FCT (avg and P99) of small flows
  - ✓ Not sacrificing throughput of large flows
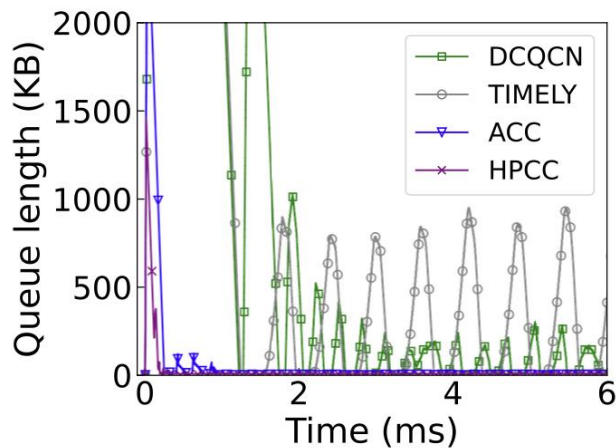  - ✓ Source halt greatly benefits low latency and reduces PFC PAUSEs

# Fast Convergence

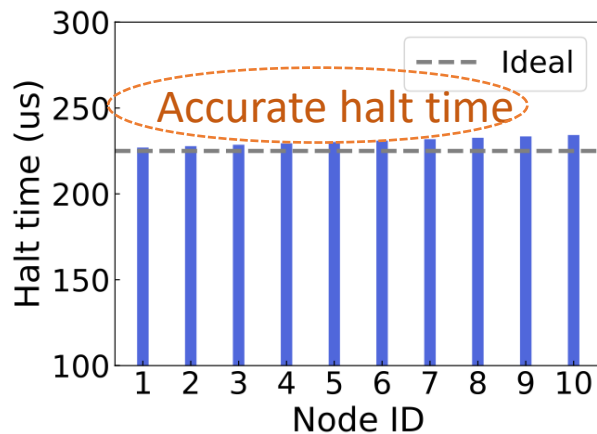## Testbed results

DCQCN

ACC



**25X** faster convergence than DCQCN
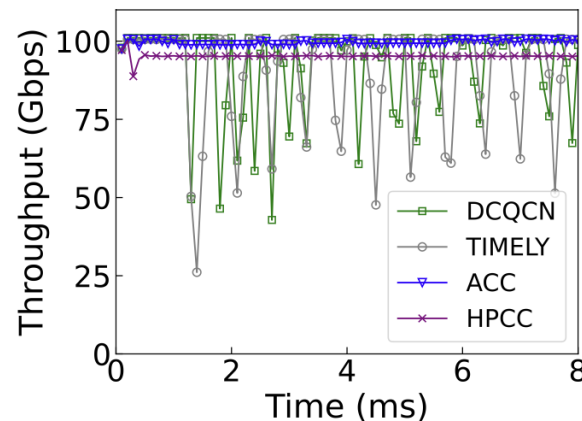
# High Link Utilization and Low Queues
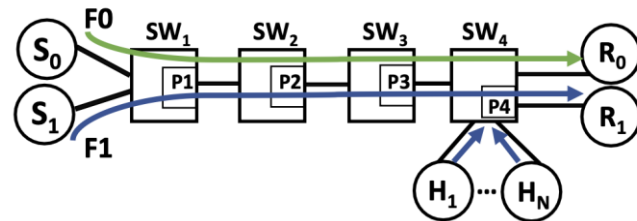
## Simulation results



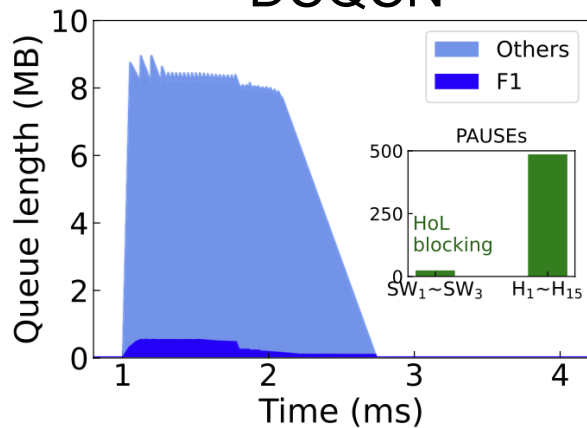(a) Queue length  (b) Halt time in ACC  (c) Bottleneck link utilization

ACC can quickly eliminate congestion and maintain
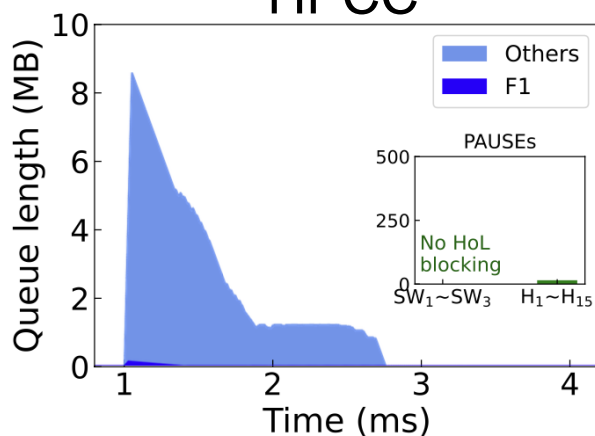near full link utilization

# Suppressing HoL Blocking
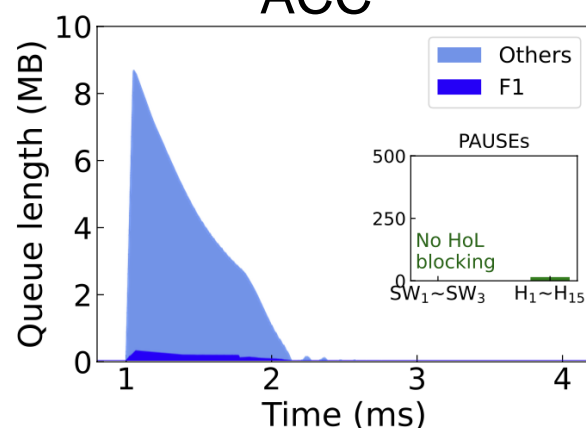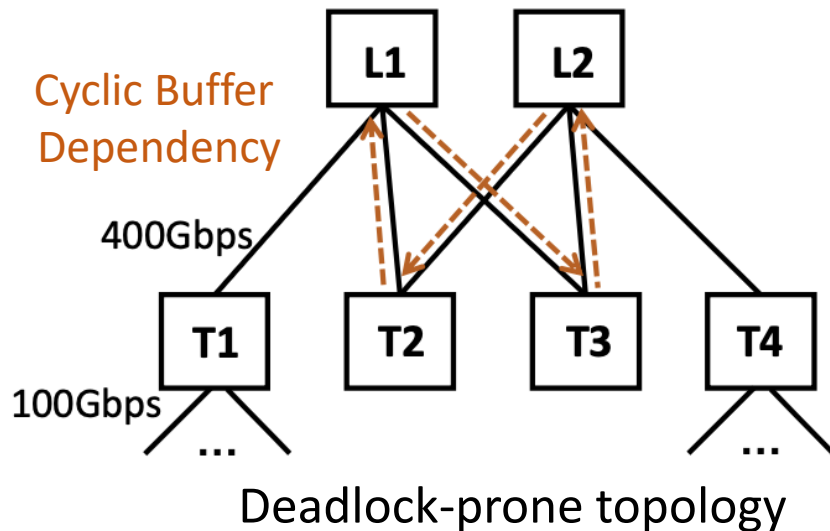


## Simulation results



**1.35X faster than HPCC**

ACC can effectively alleviate HoL blocking and congestion spreading under bursty traffic

# Resiliense to Deadlocks

## Simulation results



Cyclic Buffer Dependency

400Gbps

100Gbps

Deadlock-prone topology

Fraction of deadlock runs

| Scheme | Fraction |
|--------|----------|
| DCQCN  | 6%       |
| TIMELY | 74%      |
| HPCC   | 0%       |
| ACC    | 0%       |

**No deadlocks in 50 runs**

Fast convergence of CC does help prevent deadlocks

# FCT Performance Gains



Avg/P99 FCT(ms)

Size (Byte)

Web Search 80% load

**ACC vs. HPCC:** reduces avg FCT by **29%,** P99 FCT by **40%**

**ACC vs. DCQCN+TCD, TIMELY+TCD :** **3.9X** and **5.1X** better P99 FCT

# Conclusions

- **ACC pushes precise congestion control in lossless Ethernet by unlocking its intrinsic** *packet conservation property*
  - ✓ *Only utilizing ACKs to infer the throttled rate, excessive packets and conduct accurate source halting*
- ACC well alleviates thorny issues (HoL blocking , congestion spreading and deadlock) and achieves lower FCT
- ACC can inspire congestion control or traffic management in other lossless interconnects