# Enabling Packet Spraying over Commodity RNICs with In-Network Support

**Xiangzhou Liu**, Wenxue Li, Kai Chen

iSingLab
HKUST

# ECMP's Dilemma in AI Training Workloads

RDMA is essential for scale-out networks to meet AI training's high throughput demands.

ECMP, the de-facto standard for RDMA load balancing, is poorly suited for the unique traffic patterns of AI training workloads.

**AI workload traffic pattern**

**ECMP's Dilemma**

**Low Entropy Pattern:**
Small number of bursty elephant flows

High hash collision rate

**Coflow Pattern:**
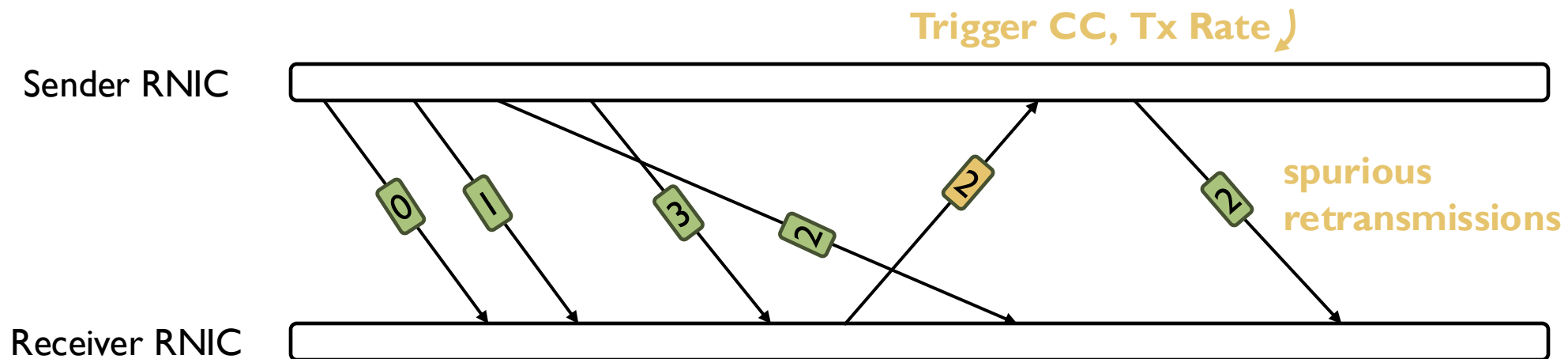One collective communication consists of multiple flow

Stragglers bottleneck the collective

**Packet spray is a promising solution to address the limitations of ECMP**

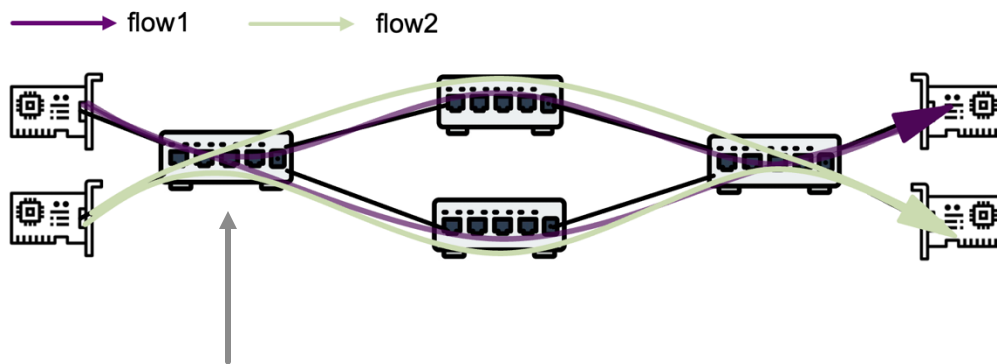# Packet Spray is Incompatible with Commodity RNIC Reliability Mechanisms

- Packet Spray inevitably results in **out-of-order (OOO) packet arrival**

- CX6, CX7, BF3 support OOO packet reception, but their RNIC-SR **treats OOO as packet loss signal and blindly generate NACKs,** causing:

  - **Unnecessary slow start**

  - **Retransmissions disrupt the RNIC TX data path**

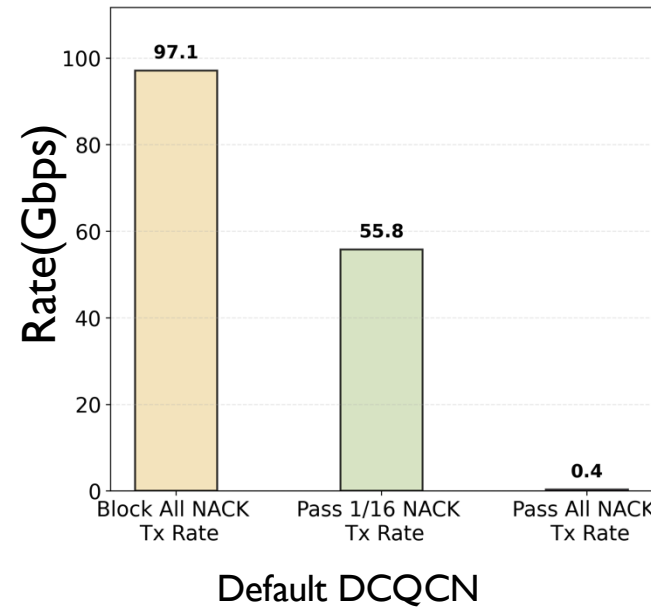  - **Bandwidth waste due to spurious retransmissions**

# Packet Spray is Incompatible with Commodity RNIC Reliability Mechanisms
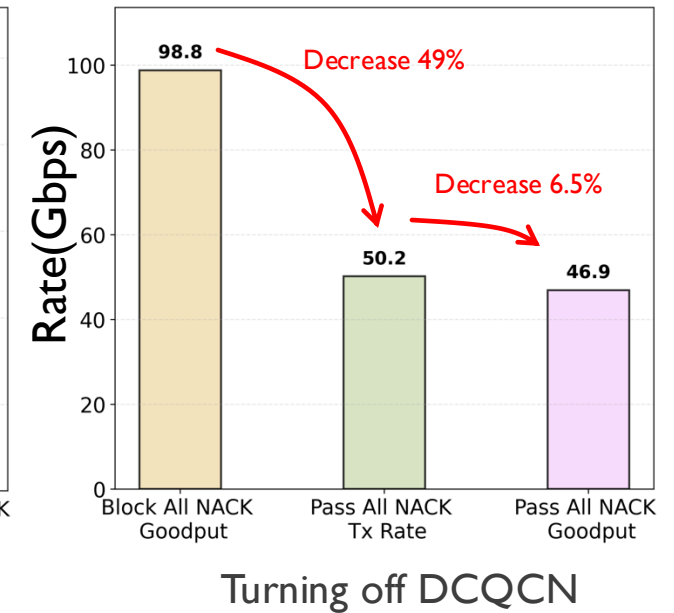
- **Settings:**
  - Two CX7 NIC pairs connected via dual paths
  - Link line rate = 100Gbps
  - Switch enable random packet spray



→ flow1  → flow2

Configure ToR to drop a specified proportion of NACK packets



Default DCQCN

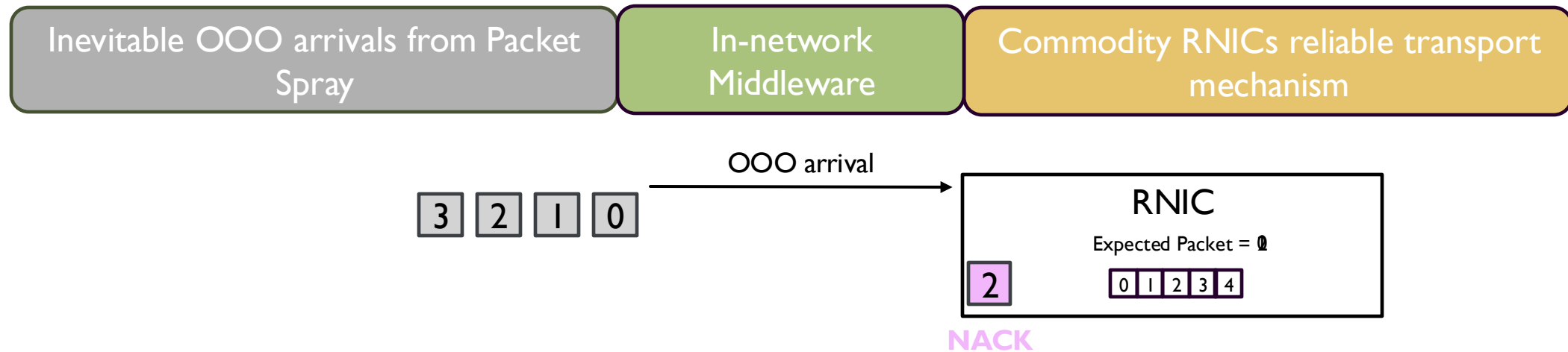**Unnecessary slow start**



Turning off DCQCN

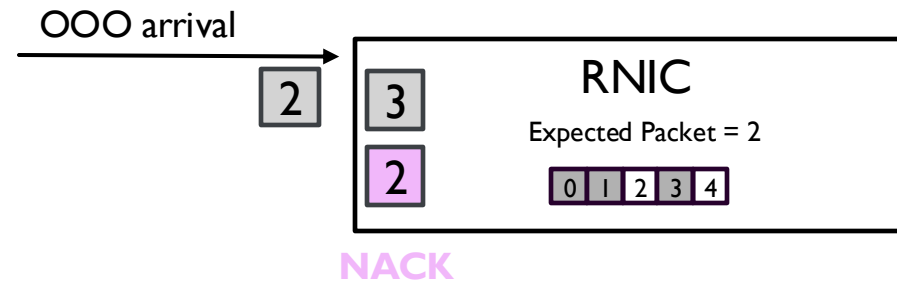Decrease 49%

Decrease 6.5%

**Retransmissions disrupt the RNIC TX data path**

**Bandwidth waste due to spurious retransmissions**

# Themis: Middleware at ToR switch for NACK validation & blocking

| Inevitable OOO arrivals from Packet Spray | In-network Middleware | Commodity RNICs reliable transport mechanism |
|---|---|---|

OOO arrival

3  2  1  0

**RNIC**

Expected Packet = 0

0  1  2  3  4

2

NACK

Themis operates as middleware on **off-the-shelf programmable ToR switches** to identify and block unnecessary NACKs, reconciling the gap between Packet Spray and RNIC-SR

# Themis: Key Method



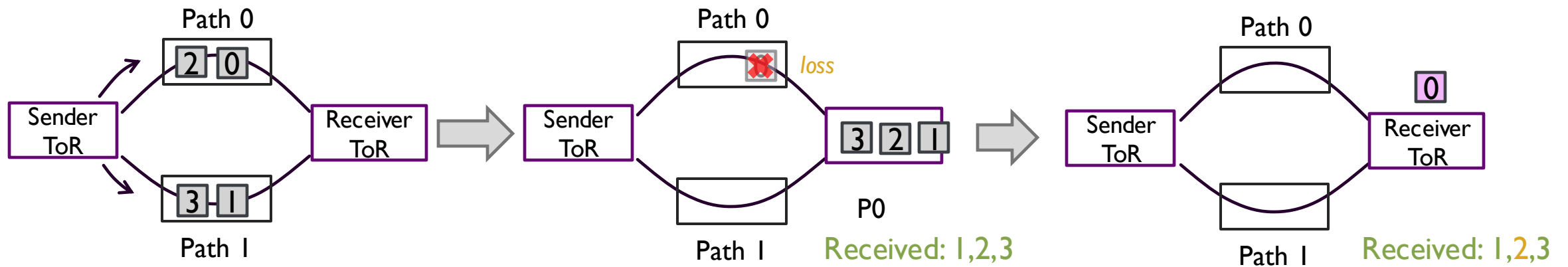Themis exploits that out-of-order arrival on the same path definitively indicates packet loss:

If the OOO packet traverses the same path as the unreceived expected packet, the corresponding NACK is valid

The receiver **only knows its sequence number (PSN)**

The receiver can **only use PSN to infer the expected packet's path**

**Themis use PSN as path selection entropy,** Enabling receiver identify expected packet's path before receiving it

# Themis Design Overview



- The Sender-ToR applies **PSN-based Packet Spray**

- The receiver-side ToR maintains PSN **records** for packets transmitted to RNIC

- The receiver-ToR determines whether a NACK is valid based on **records** and **PSN-based Packet Spray policy**

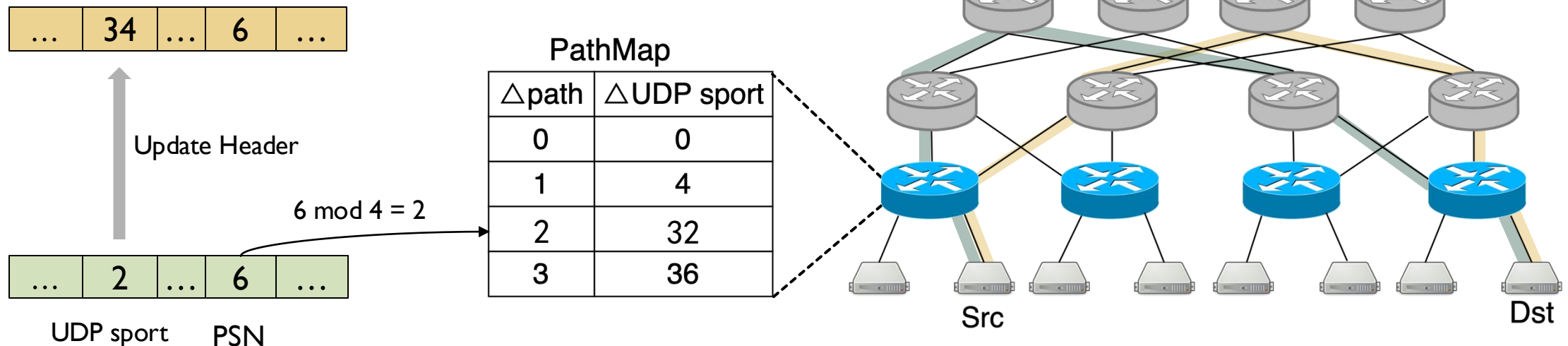  *2 is received and traverse the same path as 0*

  ➡ ☐0☐ *is lost* ➡ ☐0☐ *is valid*

# PSN-based Packet Spray

$$Path_i = (PSN_i \bmod N + Path_{ECMP}) \bmod N$$

Offset

$N$ is number of path from source to destination

$Packet_i$ traverse the same path as $Packet_j$ ⬌ $\boxed{PSN_i \bmod N == PSN_j \bmod N}$

This can be achieved through the **Relative Path Control**[1,2].

$Path_{ECMP}$     $Path_{PSN=6}$



PathMap

| △path | △UDP sport |
|-------|------------|
| 0 | 0 |
| 1 | 4 |
| 2 | 32 |
| 3 | 36 |

... 34 ... 6 ...

Update Header

6 mod 4 = 2

... 2 ... 6 ...

UDP sport    PSN

path 0    path 1    path 2    path 3

Src       Dst

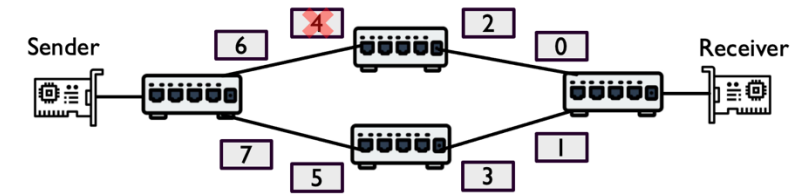[1] Hashing Linearity Enables Relative Path Control in Data Center. ATC, 2021.
[2] Unlocking ECMP Programmability for Precise Traffic Control. NSDI, 2025.

# Maintaining PSN Records to Filter Invalid NACKs

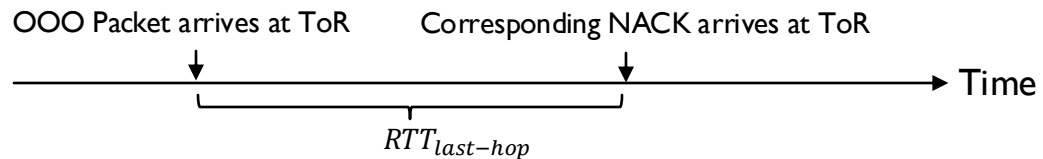NACKs generated by commodity RNICs only contain receiver's expected PSN

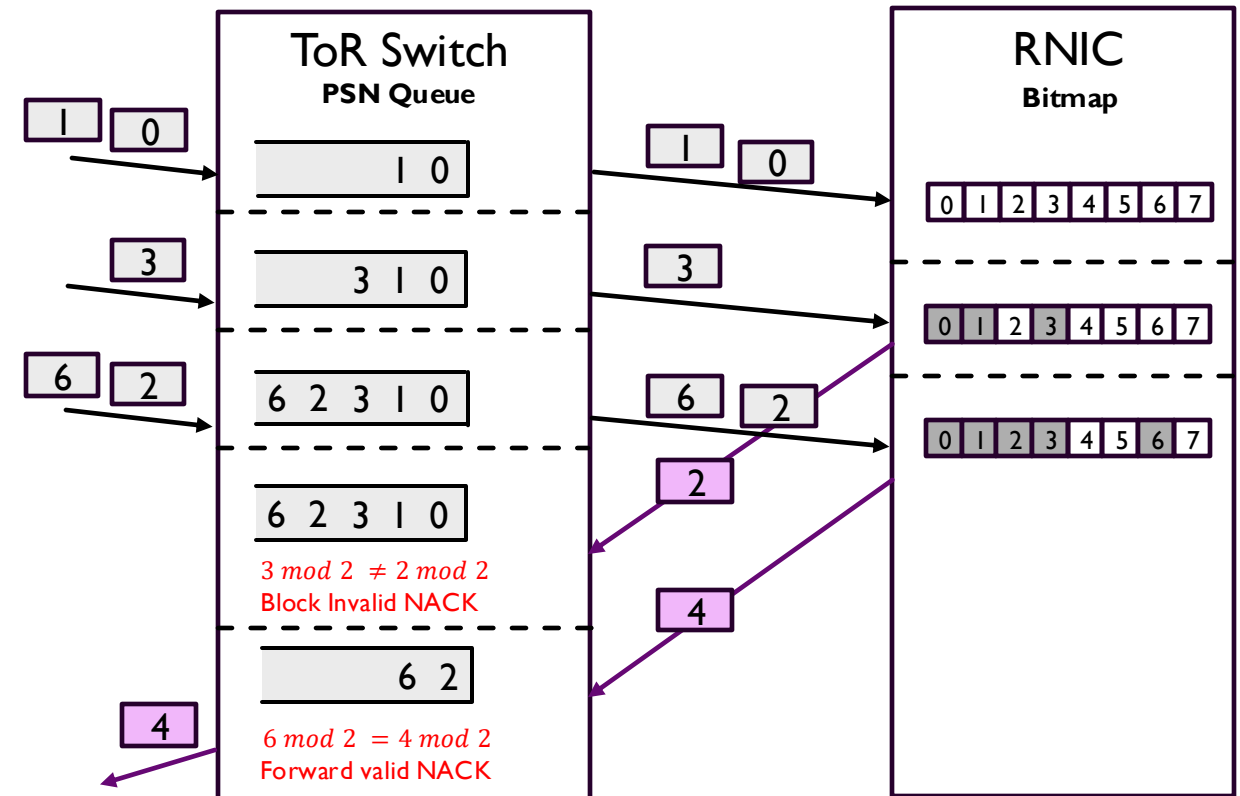How can the ToR switch identify which OOO packet triggered the NACK?

**Tracking PSNs with a PSN Queue**

- **Receive Data Packet**: Enqueue arriving packets' PSN

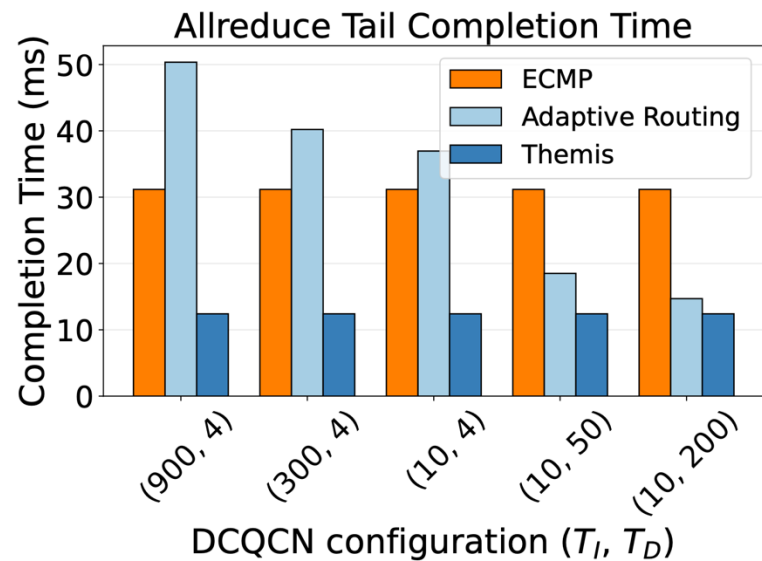- **Receive NACK**: Dequeues entries from the PSN queue until it finds the first PSN larger than the ePSN

PSN of the OOO packet that triggered the NACK

OOO Packet arrives at ToR    Corresponding NACK arrives at ToR

Time

$RTT_{last-hop}$

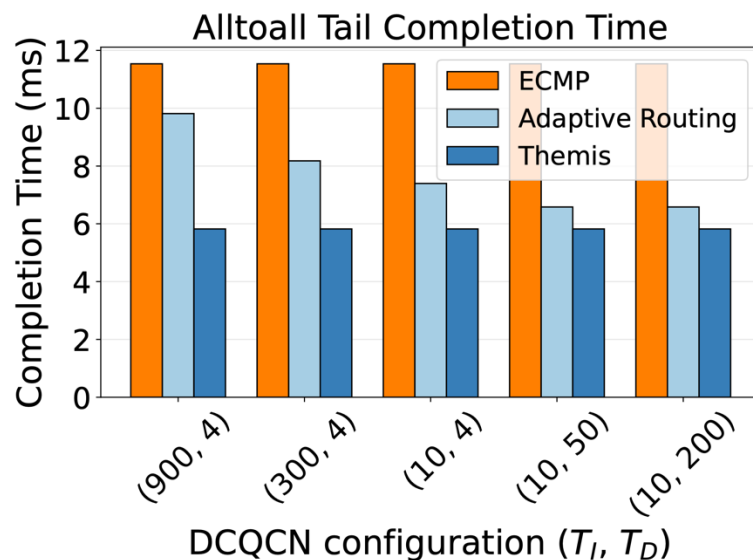$$PSN\ Queue\ Length \approx \left\lceil \frac{BW \times RTT_{last-hop}}{MTU} \right\rceil = \left\lceil \frac{400Gbps \times 3us}{4KB} \right\rceil = 38 \rightarrow 76B$$



ToR Switch
PSN Queue

1 0

3 1 0

6 2 3 1 0

6 2 3 1 0

3 mod 2 ≠ 2 mod 2
Block Invalid NACK

6 2

6 mod 2 = 4 mod 2
Forward valid NACK

RNIC
Bitmap

0 1 2 3 4 5 6 7

0 1 2 3 4 5 6 7

0 1 2 3 4 5 6 7

9

# Preliminary Simulation Results



Allreduce Tail Completion Time

15.6%~75.3% lower than AR



Alltoall Tail Completion Time

11.5%~40.7% lower that AR

**Settings:**
- 256 servers into 16 groups (16 servers each).
- Each group executes an AllReduce/AlltoAll operation, starting execution at the same time.
- DCQCN with different $(T_I, T_D)$ configuration
- $T_I$: rate increase interval(us)
- $T_D$: rate decrease interval(us)

Themis ensures **compatibility** with commodity RNICs and achieves **high-performance** packet-level load balancing by preventing unnecessary slow starts and spurious retransmissions.

# Summary

➢ Packet Spray inevitably results in out-of-order packet arrival, which is incompatible with commodity RNIC reliability mechanisms.

➢ We design Themis, a lightweight middleware deployed on programmble ToR switches that applies PSN-based packet spraying at the source ToR switch while identifying and blocking invalid NACKs at the destination ToR switch.

➢ By preventing spurious retransmissions and unnecessary slow starts, Themis enables effective packet spraying with commodity RNICs.

# Thank you!
contact email: xliugg@connect.ust.hk