

Multi-Path Transport for RDMA in Datacenters

Yuanwei Lu^{1,2}, Guo Chen³, Bojie Li^{1,2}, Kun Tan⁴, Yongqiang Xiong², Peng Cheng², Jiansong Zhang², Enhong Chen¹, Thomas Moscibroda⁵

¹USTC, ²Microsoft Research, ³Hunan University,

⁴Huawei Technologies, ⁵Microsoft Azure

Remote Direct Memory Access

Kernel bypass

- Low and stable base latency: $<10\mu\text{s}$



Transport hardware offloading

- High throughput with low CPU overhead

Remote Direct Memory Access

RDMA (RoCEv2) is deployed at scale in Microsoft Datacenter

- e.g. Half of the traffic in a online service datacenter is RDMA [1]

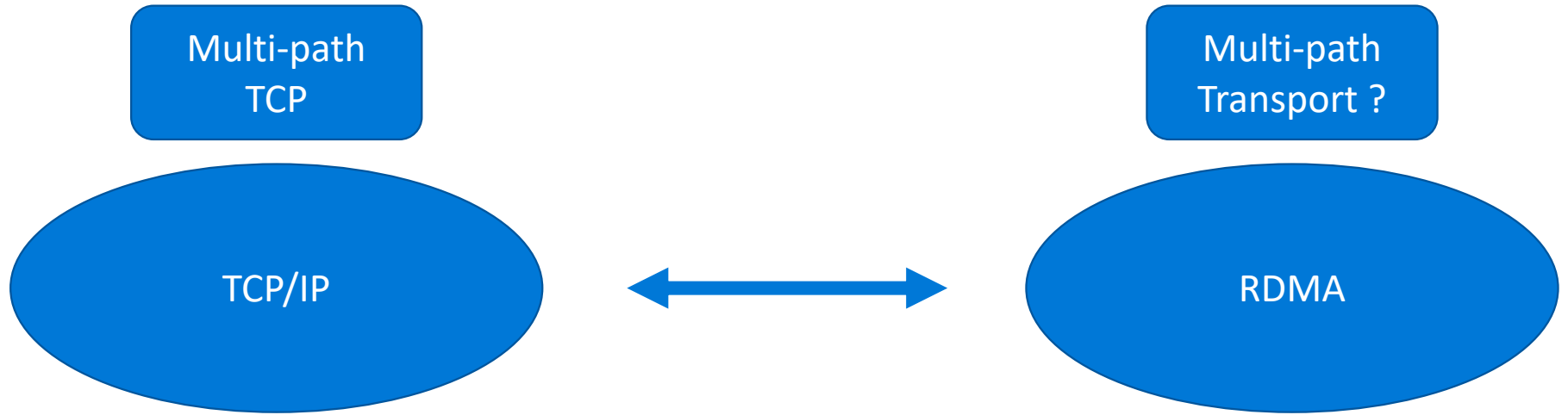
Datacenter topologies: FatTree, BCube, VL2

- Lots of parallel paths between any pair of nodes

Current RDMA is single-path with ECMP

- Hot spot caused high latency and low throughput
- Not robust to failure

[1] Guo, Chuanxiong, et al. "RDMA over Commodity Ethernet at Scale." *Proceedings of the 2016 conference on ACM SIGCOMM 2016 Conference*. ACM, 2016.



Goal: design a high performance multi-path transport for RDMA. (MP-RDMA)

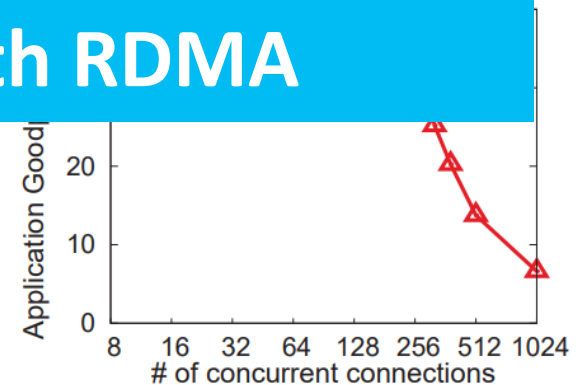
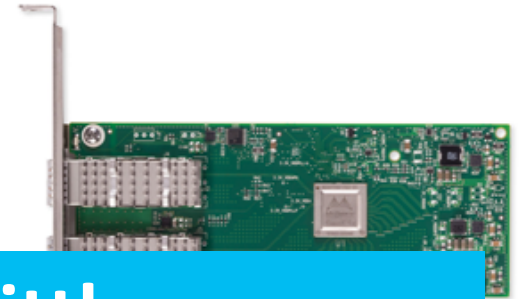
Uniqueness of RDMA Environment

RDMA must be implemented in hardware

NIC on-chip memory is scarce

- **MP-RDMA design needs to add as little hardware memory footprint as possible compared with existing single-path RDMA**
- Large and unstable latency

Frequent cache miss results in poor performance^[1]



[1] Kalia, Anuj, Michael Kaminsky, and David G. Andersen. "Design Guidelines for High Performance RDMA Systems." *2016 USENIX Annual Technical Conference (USENIX ATC 16)*. 2016.

Challenges

#1: How to achieve congestion-aware traffic split without per-path states at NIC?

#2: How to track out-of-order packets with small on-chip memory footprint?

#3: How to provide message ordering without hurting performance?

Challenges

#1: How to achieve congestion-aware traffic split without per-path states at NIC?

#2: How to track out-of-order packets with small on-chip memory footprint?

#3: How to provide message ordering without hurting performance?

Challenge #1

Congestion-aware load distribution

- Tracking the transmission and congestion states on every paths

**Challenge #1:
How to achieve congestion-aware traffic split
without per-path states at NIC?**

- RoCEv2: 248B per connection

Multi-Path ACK Clocking

A transport without per-path states

- Insight: window-based protocol with ACK clocking

Algorithm:

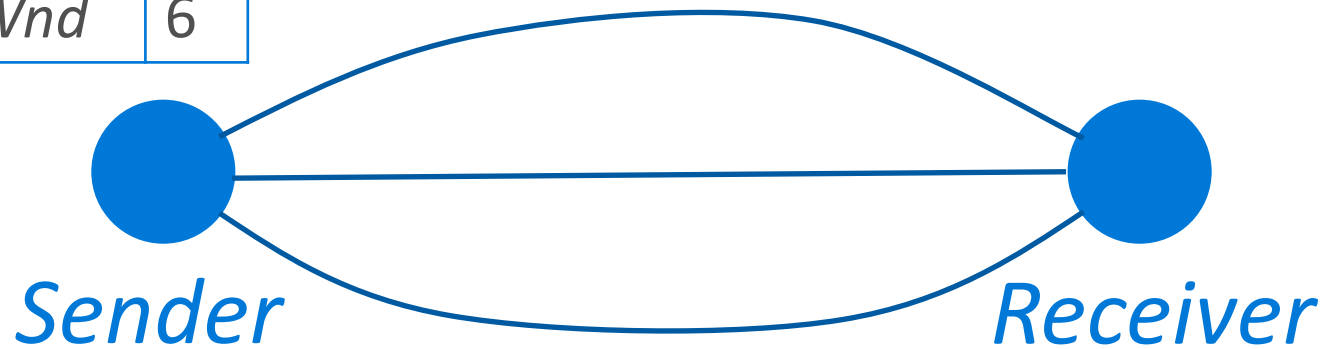
- At beginning, spread N packets randomly to N paths
- One *cwnd* for all paths, for each received ACK:

$$cwnd \leftarrow \begin{cases} cwnd + \frac{1}{cwnd}; & \text{if } ECN = 0 \\ cwnd - \frac{1}{2}; & \text{if } ECN = 1 \end{cases}$$

- Clock one packet to the path ACK returns when *cwnd* allows

An Example

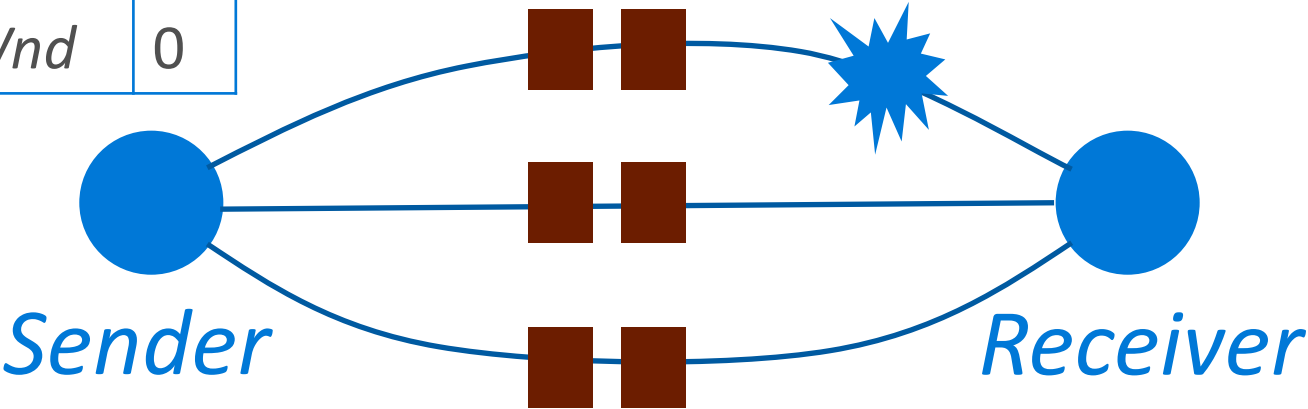
<i>cWnd</i>	6
<i>Inflight</i>	0
<i>aWnd</i>	6



- Data Packet
- ▲ ACK Packet

An Example

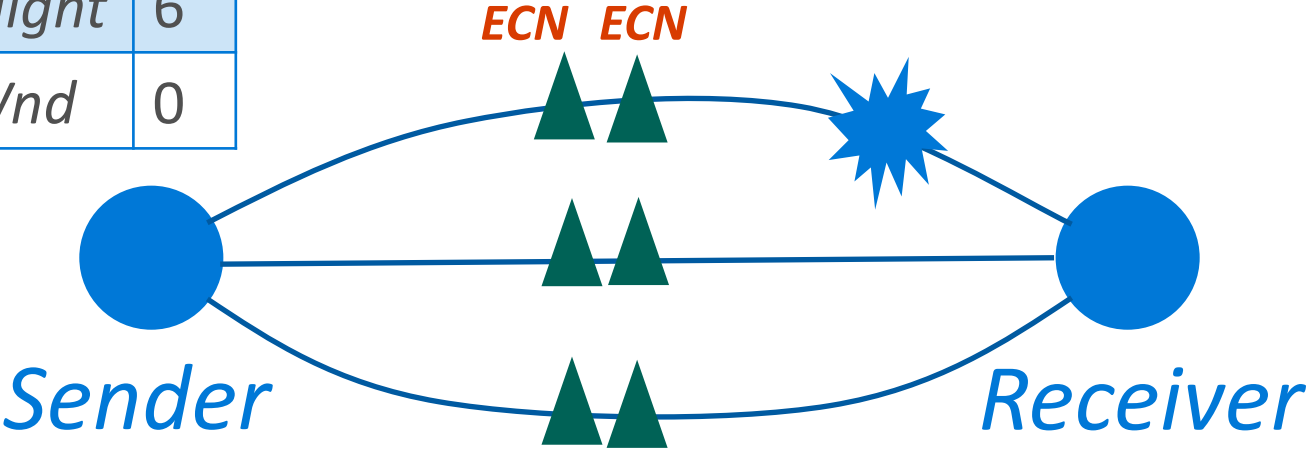
<i>cWnd</i>	6
<i>Inflight</i>	6
<i>aWnd</i>	0



- Data Packet
- ▲ ACK Packet

An Example

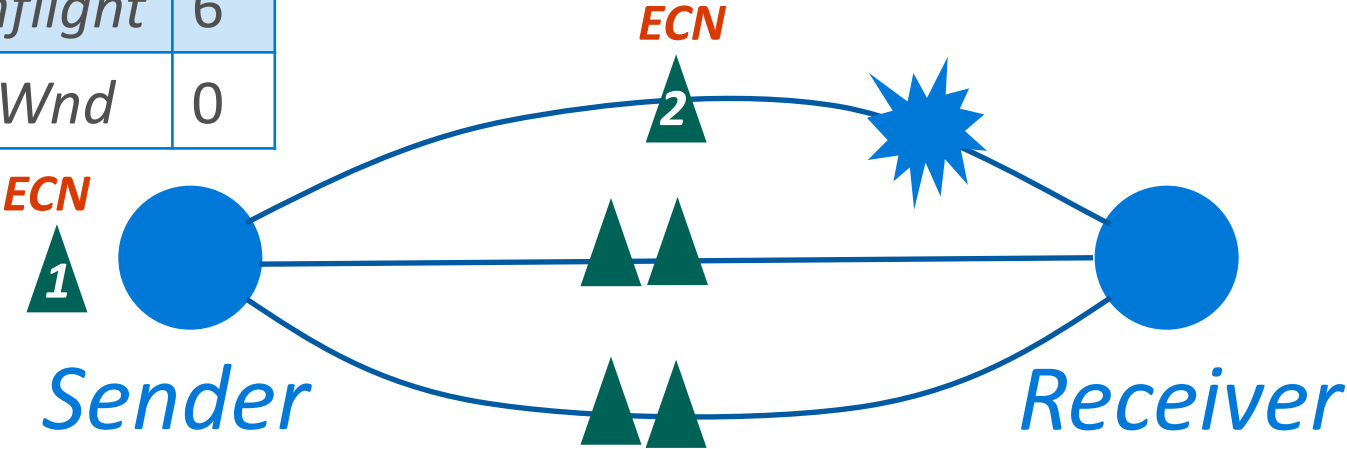
<i>cWnd</i>	6
<i>Inflight</i>	6
<i>aWnd</i>	0



- Data Packet
- ▲ ACK Packet

An Example

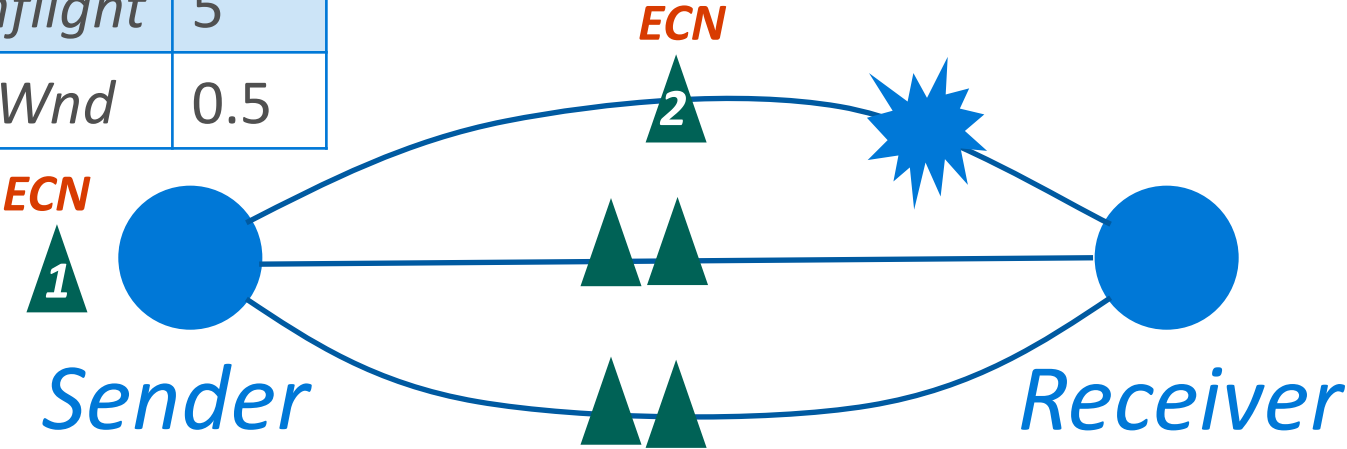
<i>cWnd</i>	6
<i>Inflight</i>	6
<i>aWnd</i>	0



- Data Packet
- ▲ ACK Packet

An Example

<i>cWnd</i>	5.5
<i>Inflight</i>	5
<i>aWnd</i>	0.5



- Data Packet
- ▲ ACK Packet

An Example

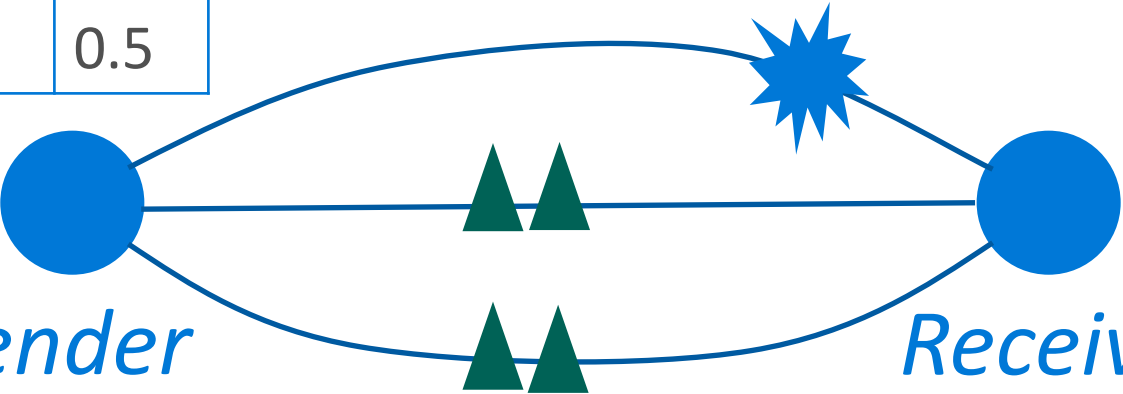
<i>cWnd</i>	5.5
<i>Inflight</i>	5
<i>aWnd</i>	0.5

ECN
2

Sender

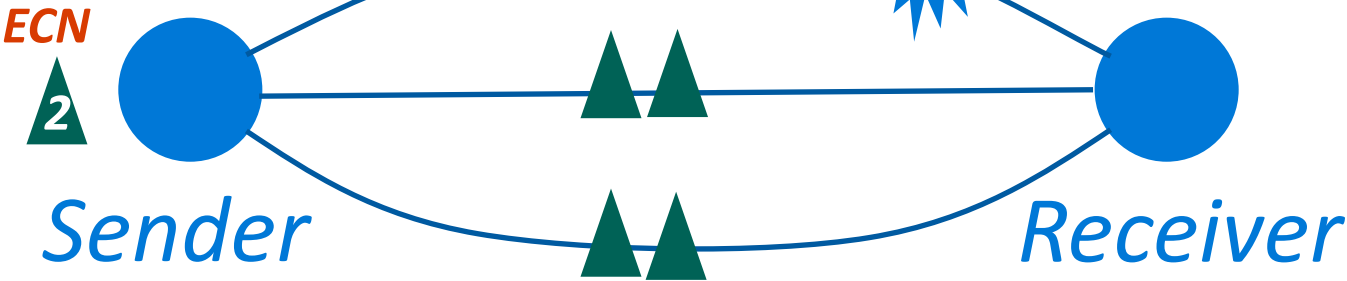
Receiver

- Data Packet
- ▲ ACK Packet



An Example

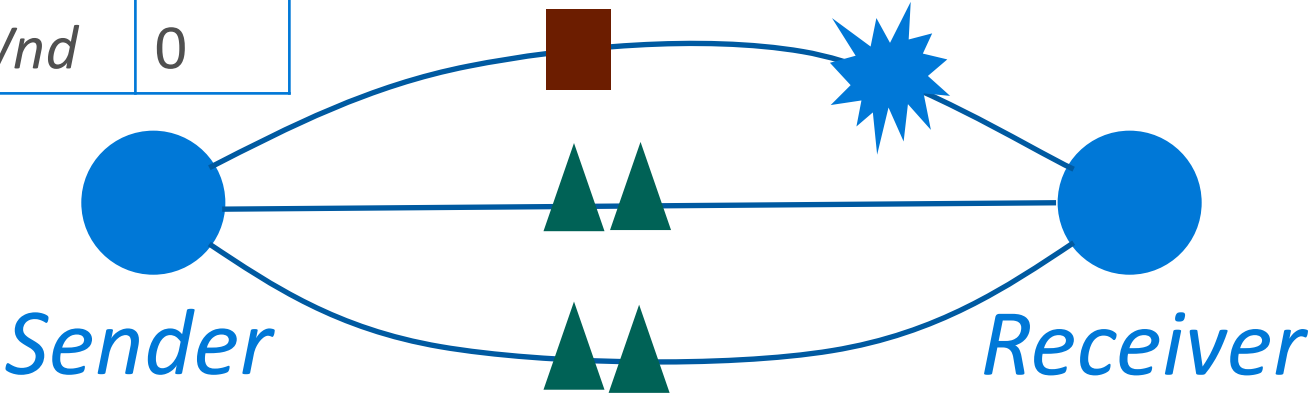
<i>cWnd</i>	5
<i>Inflight</i>	4
<i>aWnd</i>	1



- Data Packet
- ▲ ACK Packet

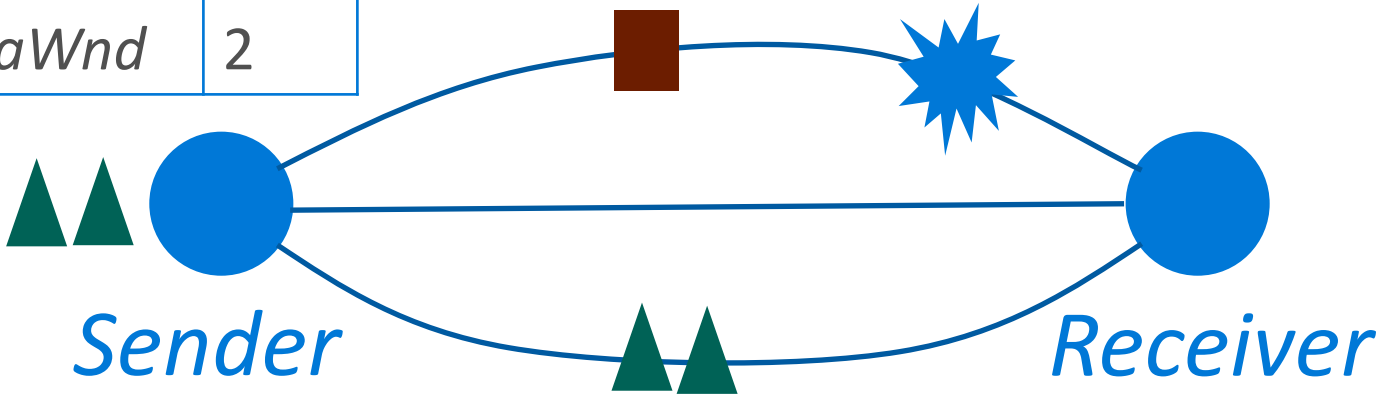
An Example

<i>cWnd</i>	5
<i>Inflight</i>	5
<i>aWnd</i>	0



An Example

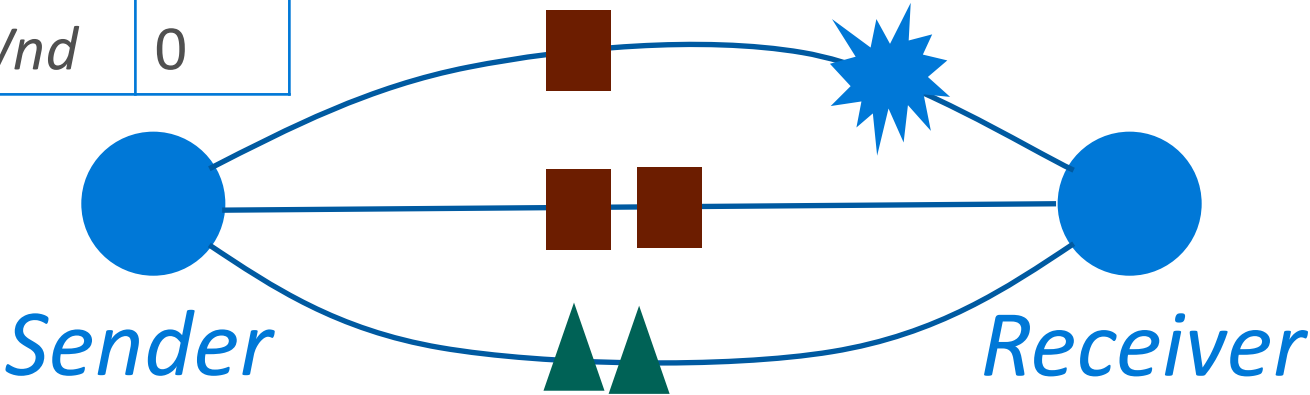
<i>cWnd</i>	5
<i>Inflight</i>	3
<i>aWnd</i>	2



- Data Packet
- ▲ ACK Packet

An Example

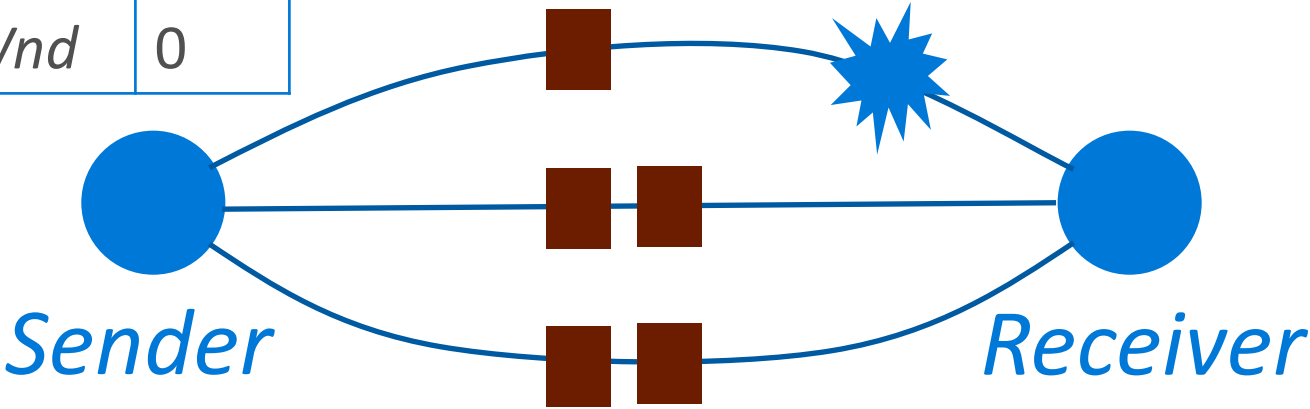
<i>cWnd</i>	5
<i>Inflight</i>	5
<i>aWnd</i>	0



- Data Packet
- ▲ ACK Packet

An Example

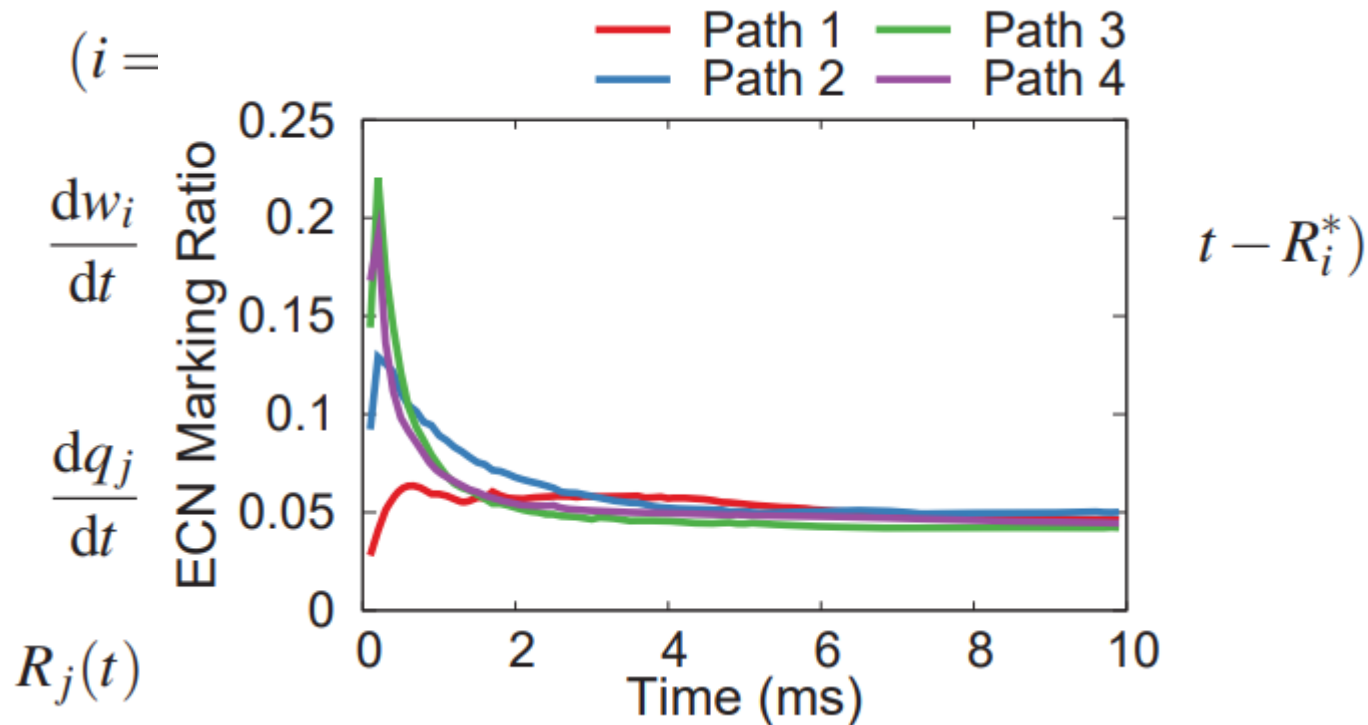
<i>cWnd</i>	5
<i>Inflight</i>	5
<i>aWnd</i>	0



- Data Packet
- ▲ ACK Packet

Fluid Model

N flows, M_v virtual paths, M_p physical paths



MP-RDMA can balance the ECN marking ratio among all sending paths

Challenges

#1: How to achieve congestion-aware traffic split without per-path states at NIC?

#2: How to track out-of-order packets with small on-chip memory footprint?

#3: How to provide message ordering without hurting performance?

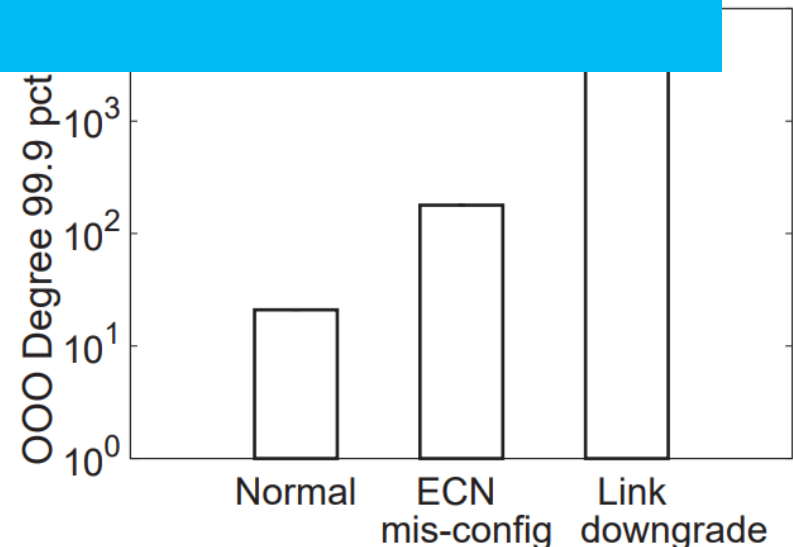
Challenge #2

Multi-path leads to packets out-of-order

- Need meta data to track out-of-order

Challenge #2: How to track out-of-order packets with small on-chip memory footprint?

- e.g. Link downgrade: 1.2KB bitmap
- e.g. PFC pause: infinite

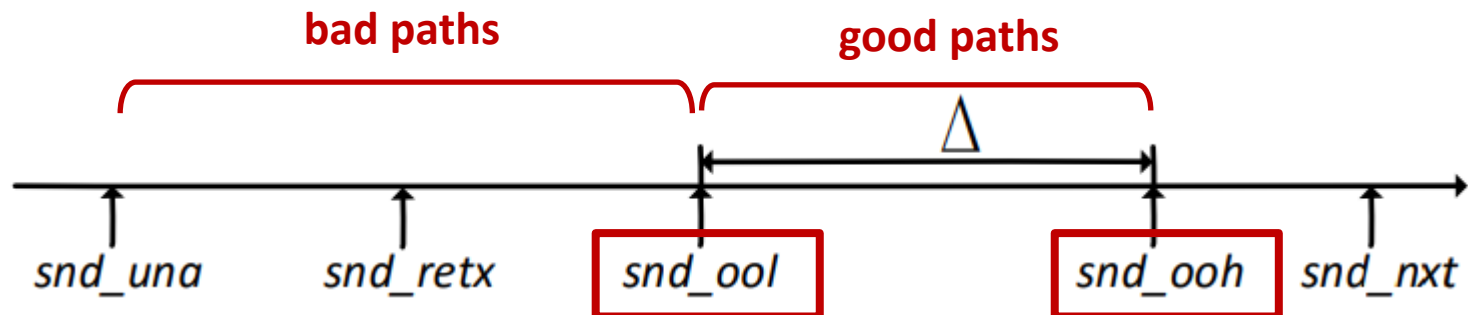


Out-of-order Aware Path Selection

MP-RDMA chooses paths with similar delay

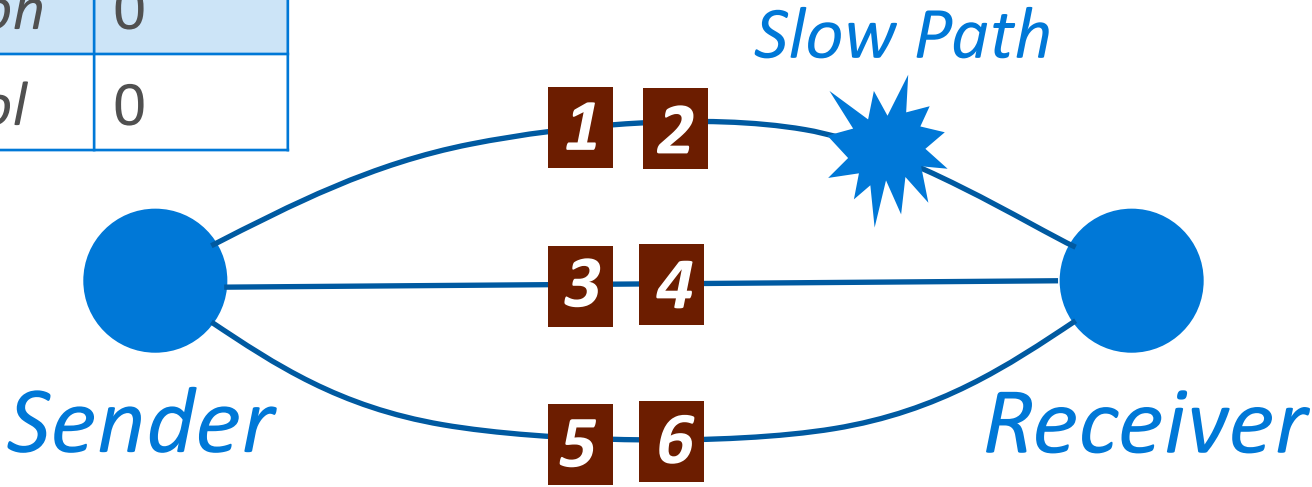
Sender maintains two pointers:

- *snd_ooh*: the highest PSN successfully received at receiver (fast paths)
- *snd_ooh*: the smallest ACK PSN that can clock out packet
- $snd_ool = snd_ooh - \Delta$



An Example

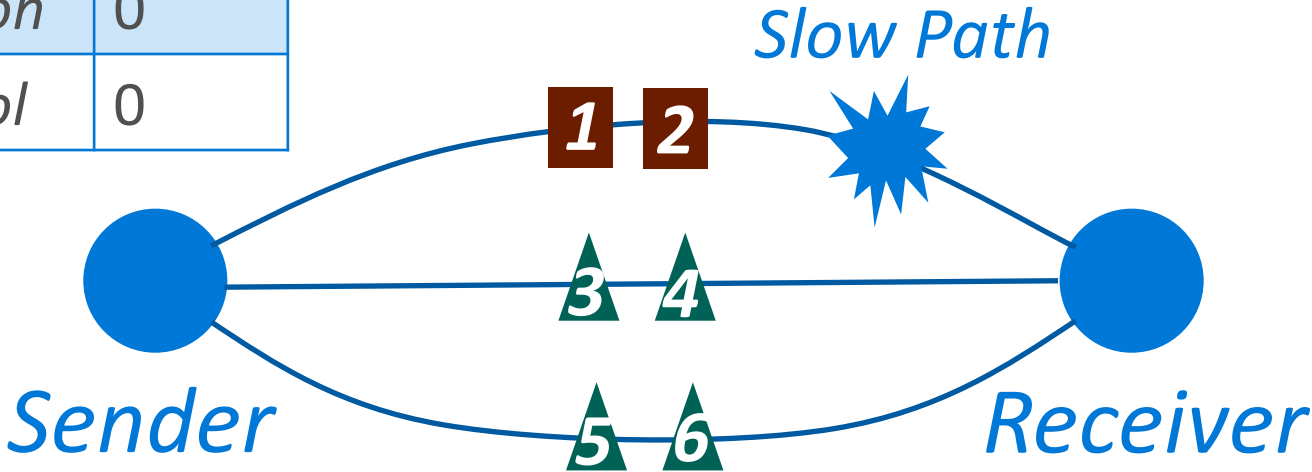
Δ	3
<i>snd_ooh</i>	0
<i>snd_ool</i>	0



- Data Packet
- ▲ ACK Packet

An Example

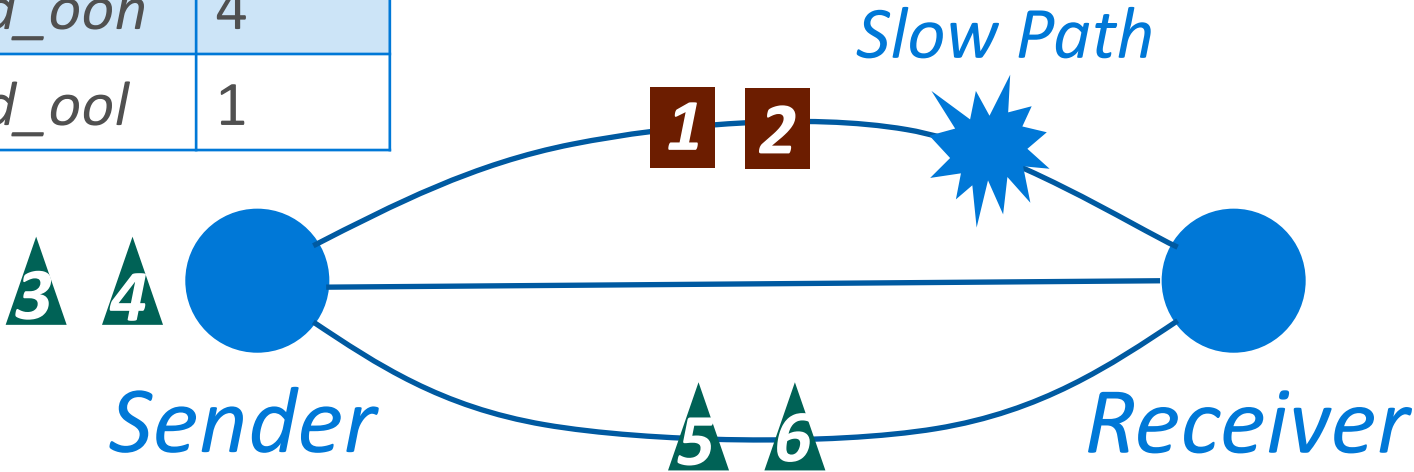
Δ	3
<i>snd_ooh</i>	0
<i>snd_ool</i>	0



- Data Packet
- ▲ ACK Packet

An Example

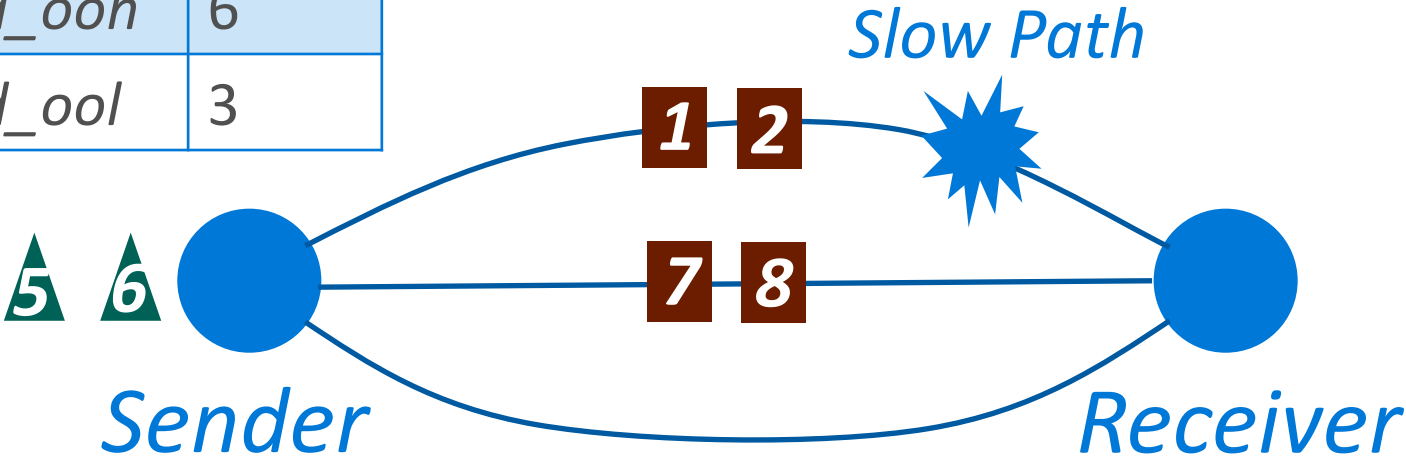
Δ	3
<i>snd_ooh</i>	4
<i>snd_ool</i>	1



- Data Packet
- ▲ ACK Packet

An Example

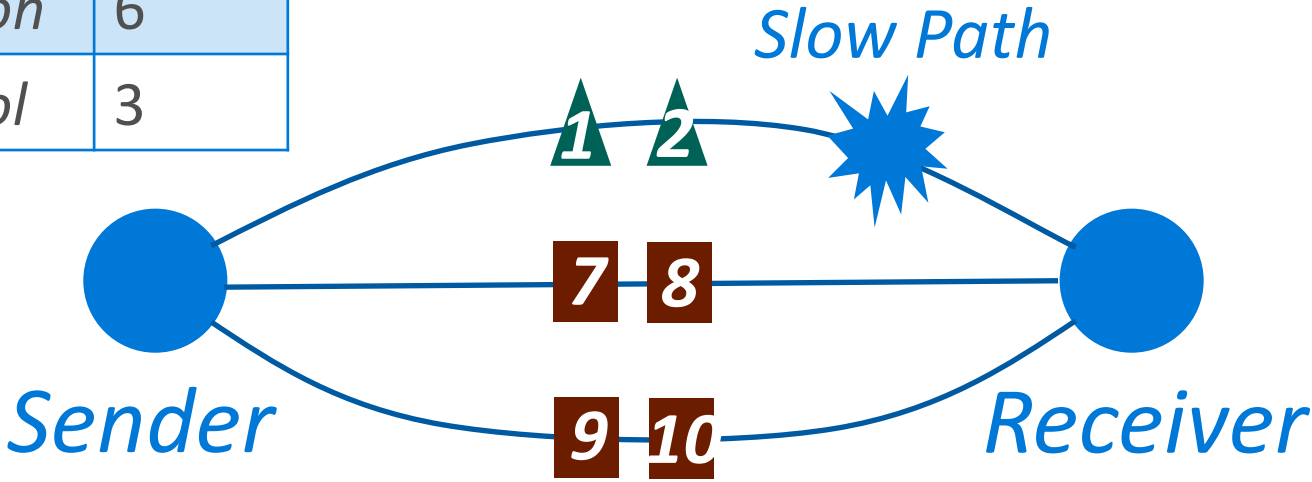
Δ	3
<i>snd_ooh</i>	6
<i>snd_ool</i>	3



- Data Packet
- ▲ ACK Packet

An Example

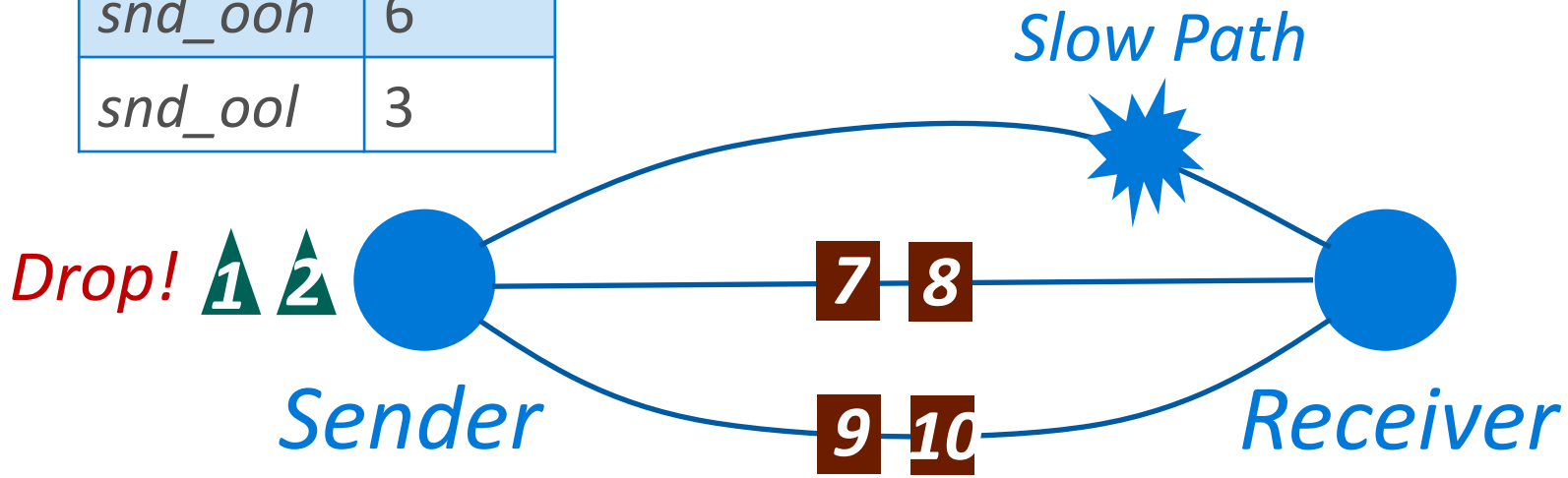
Δ	3
<i>snd_ooh</i>	6
<i>snd_ool</i>	3



- Data Packet
- ▲ ACK Packet

An Example

Δ	3
<i>snd_ooh</i>	6
<i>snd_ool</i>	3

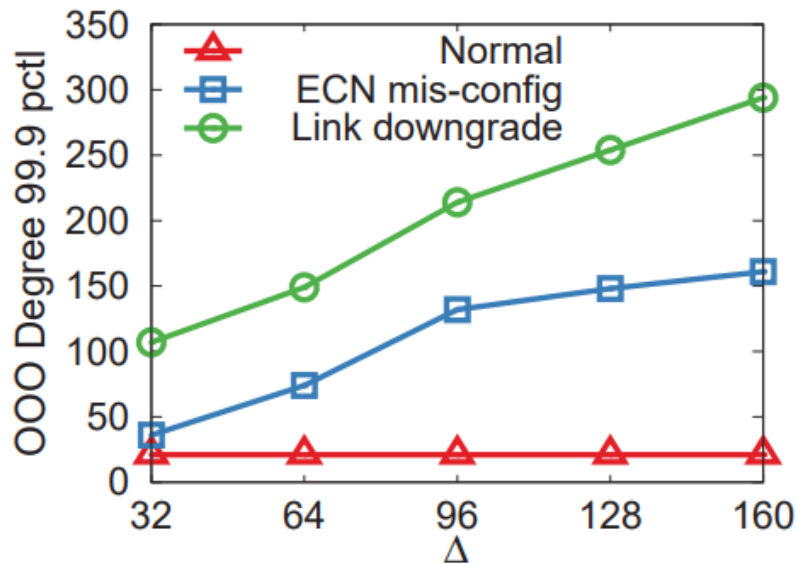


- Data Packet
- ▲ ACK Packet

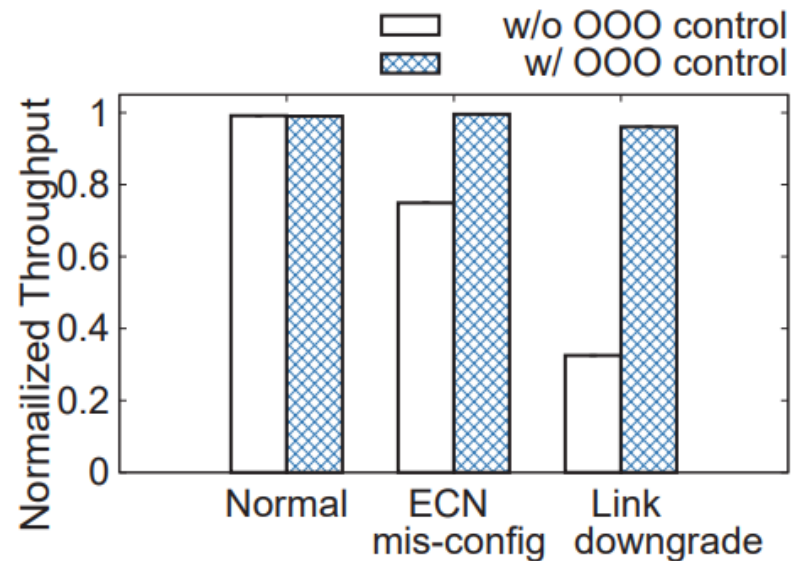
Experiment Results

Adjust Δ can control out-of-order degree

Small out-of-order improves throughput



(a) OOO control.



(b) Goodput.

MP-RDMA can track OOO with a tiny bitmap

Challenges

#1: How to achieve congestion-aware traffic split without per-path states at NIC?

#2: How to track out-of-order packets with small on-chip memory footprint?

#3: How to provide message ordering without hurting performance?

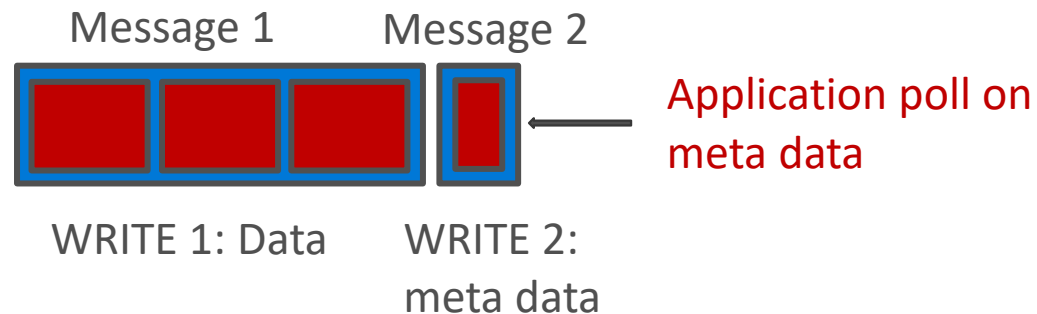
Challenge #3

Out-of-order data is directly placed to application buffer

- Problem: application may assume in-order messages delivery

Challenge #3:

How to provide message ordering without hurting the application performance?



Synchronise Mechanism

MP-RDMA provides a “*Synchronise*” interface

- A synchronise message is guaranteed to arrive after all previous messages

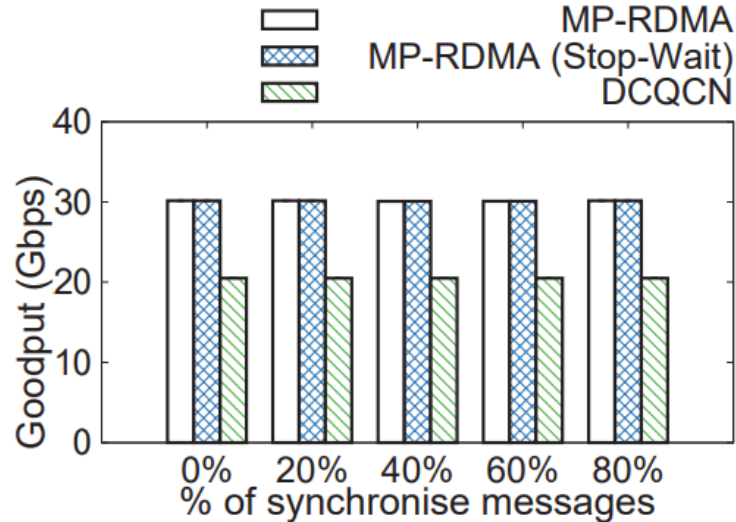
Optimistically transmit a *synchronise* message

- A synchronise message just wait Δt before previous messages
- $\Delta t = \alpha * \Delta / (cwnd / RTT)$, large enough to avoid out-of-order

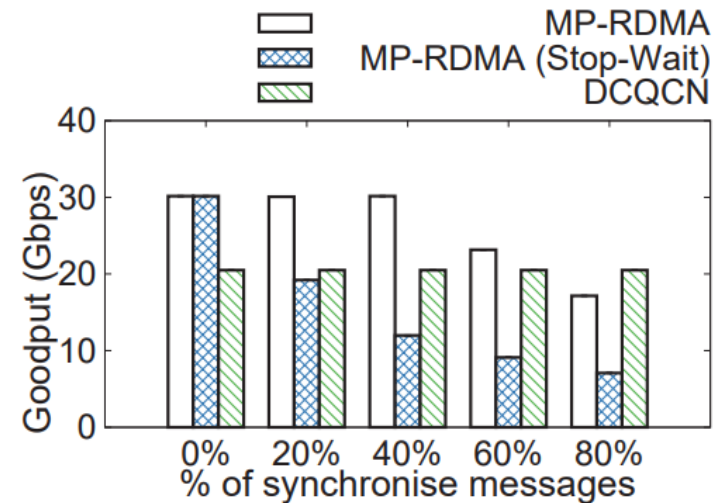
Experiment Results

No performance sacrifice for large messages

Little downgrade for small messages



(a) 512KB message.



(b) 32KB message.

Synchronise mechanism provides message ordering without sacrificing throughput

Implementation

ClickNP_[1] FPGA platform

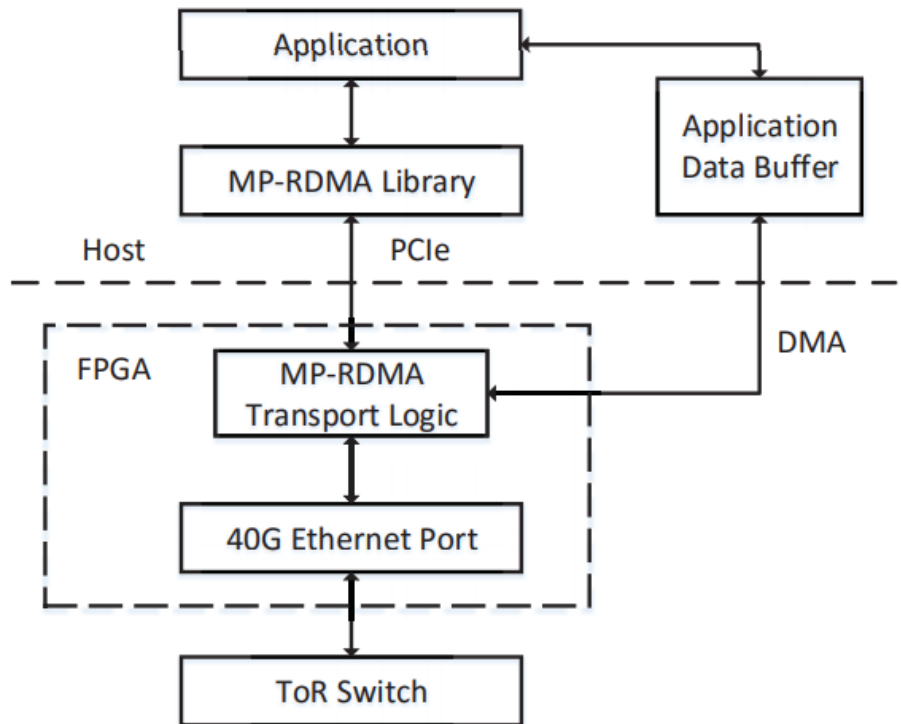
- Altera Stratix V D5 FPGA board with two 40Gb Ethernet ports
- 14 elements, ~2K lines of OpenCL code

Send logic:

- Dependency: 9 cycles
- Clock Frequency: ~206MHz

Receive logic:

- Dependency: 3 cycles
- Clock Frequency: ~196MHz



Performance

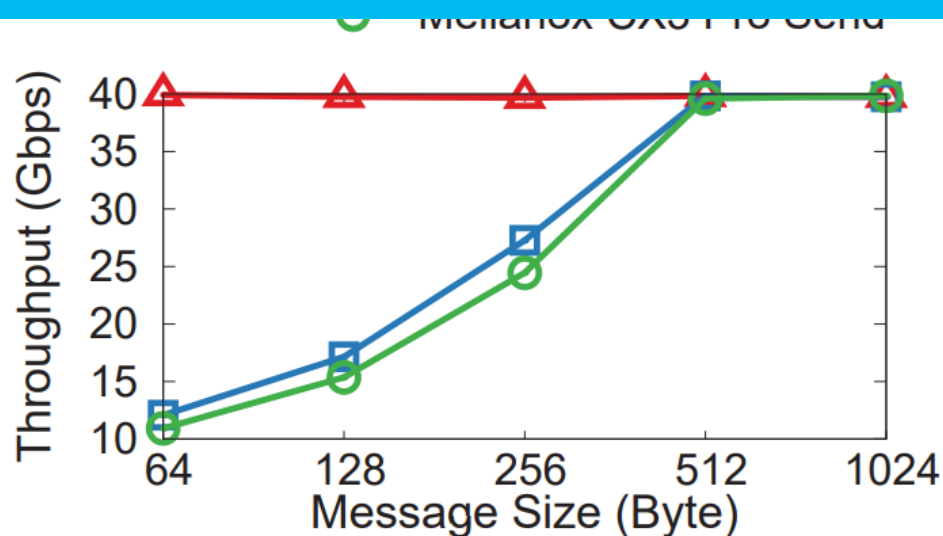
Send and receive at line rates

Latency:

- Send: 0.54 μ s

- **Conclusion:**

- **MP-RDMA can be implemented in hardware with high performance**

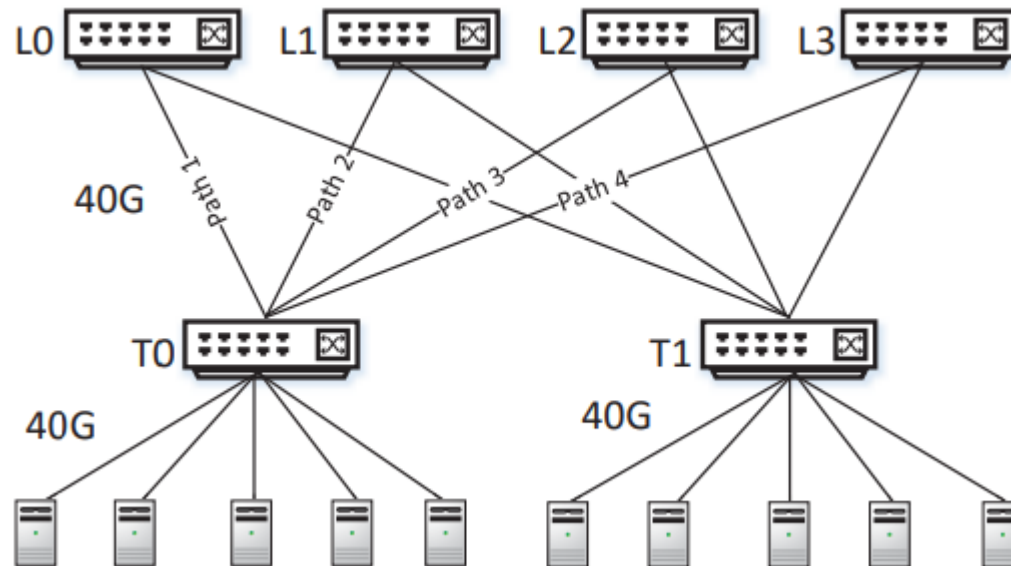


Testbed

10 FPGA boards

Leaf-Spine topology: 4 Paths

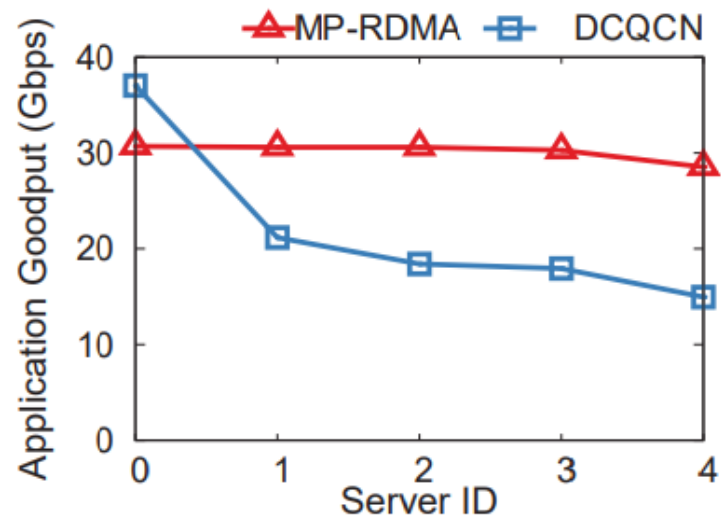
Base RTT: 16 μ s



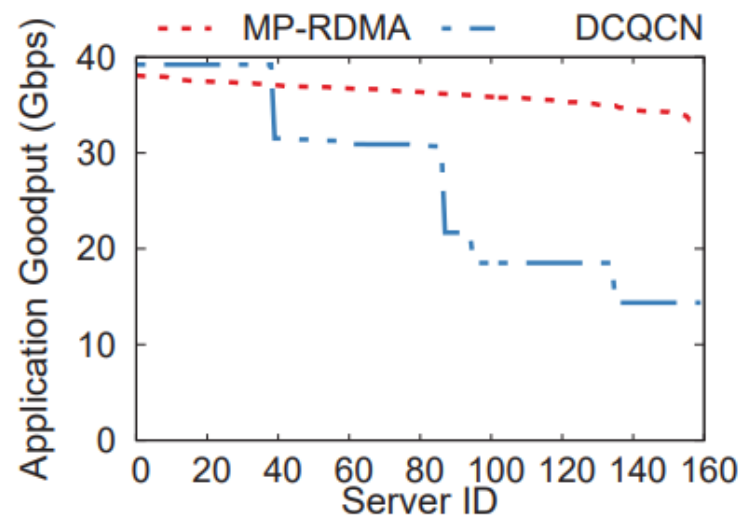
Key Results

Throughput

- 47.05% higher application throughput v.s. DCQCN



(a) Small-scale testbed.

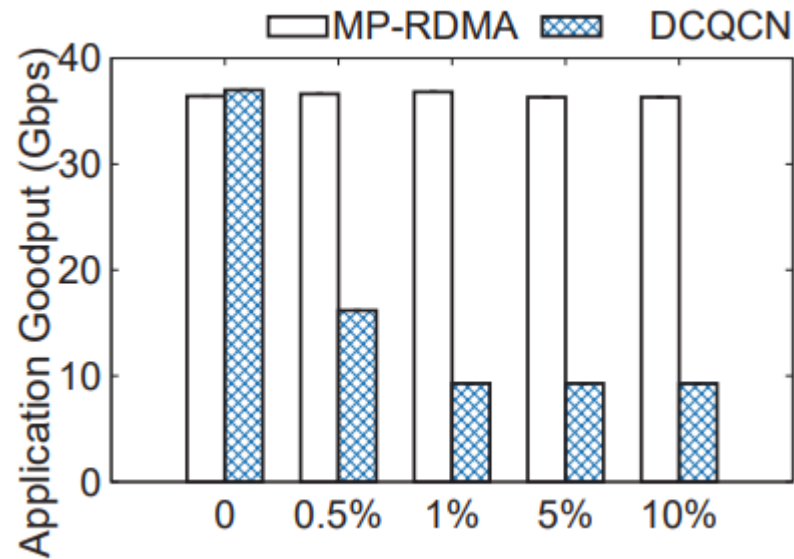


(b) Large-scale simulation.

Key Results

Robustness

- Under loss: > 3x throughput improvement



Summary

First effort to design a multi-path transport for RDMA environment

Minimize on-chip memory footprint

- No per-path states
- Out-of-order aware path selection
- Synchronise mechanism to ensure message ordering

Improve throughput, reduce latency and is robust to failure

Thank you !

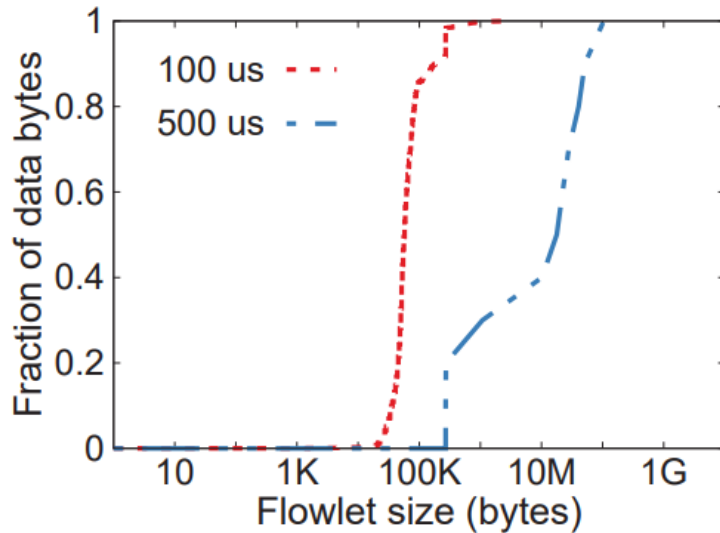
Q& A

Backup slides

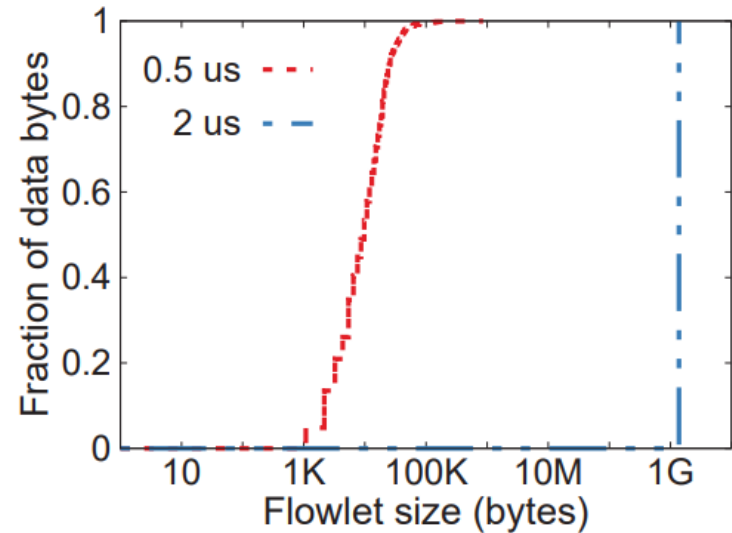
Flowlet in RDMA

TCP: inactive gap $100\mu\text{s}$, flowlet size $\sim 60\text{KB}$

RDMA: inactive gap $2\mu\text{s}$, flowlet size $\sim 2\text{GB}$



(a) TCP.



(b) RDMA.