# **LoWAR**: Enhancing RDMA over Lossy WANs with Transparent Error Correction

Tianyu Zuo（左天宇）, Tao Sun, Shuyong Zhu*,
Wenxiao Li, Lu Lu, Zongpeng Du, and Yujun Zhang*

# Background: RDMA

**Remote Direct Memory Access (RDMA)**

- The de-facto standard for datacenter network (DCN)

- Leverage kernel bypass and transport offloading

- Achieve high throughput, low latency and CPU overhead

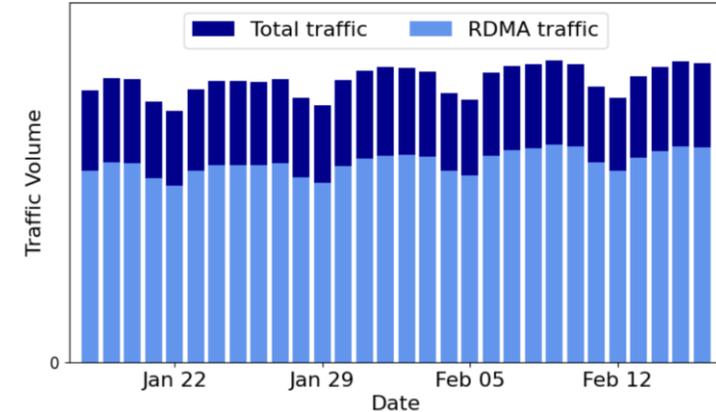**Can RDMA extend to WAN? Yes!**

- Increased transmission demand and capacity over WAN

- Unified interface, low CPU overhead, high throughput

- Recent studies on inter-region/cross-DC RDMA [1-4]

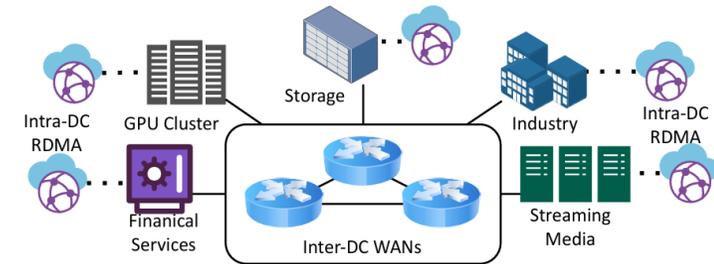[1] **NSDI`23**, *Empowering Azure Storage with RDMA.*
[2] **TPDS`23**, *Swing: Providing Long-Range Lossless RDMA via PFC-Relay.*
[3] **INFOCOM`24**, *BiCC: Bilateral Congestion Control in Cross-datacenter RDMA Networks.*
[4] **IWQoS`24**, *LSCC: Link-Segmented Congestion Control for RDMA in Cross-DC Networks.*

Wide adoption of RDMA in DCN [1]



Emerging scenarios for RDMA over WAN

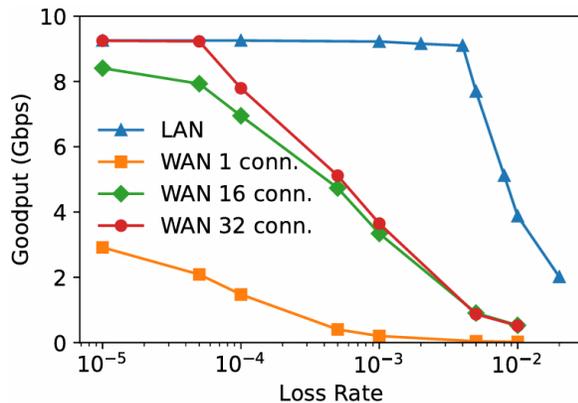RDMA's success in DCN has the potential to extend to WAN!

# Motivation:

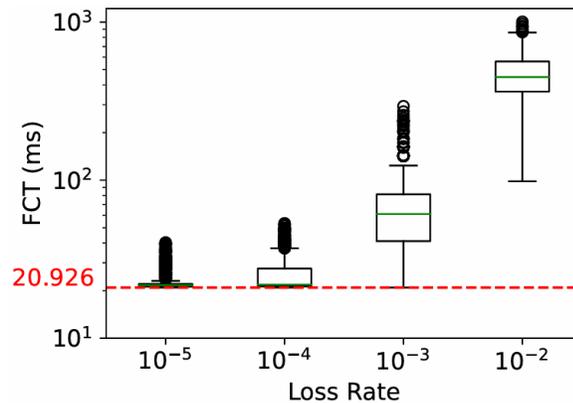## WAN and DCN has diverse characteristics:

|  | Latency | Reliability | Non-congestion Loss Rate |
|------|---------|-------------|--------------------------|
| DCN | ~10us | Lossless (PFC) | $< 4 \times 10^{-9}$ [5] |
| WAN | 1ms ~ 100ms | Lossy | $10^{-5} \sim 10^{-3}$ |

[5] Torsten Hoefler et al, *Datacenter Ethernet and RDMA: Issues at Hyperscale.*

## RDMA's performance drops in lossy WAN:



(a) Average Goodput      (b) FCT

- **Experiments Setting:** Mellanox ConnectX-5; BW = 10Gbps, RTT = 40ms, Loss Rate = $[10^{-5}, 10^{-2}]$; 1 MB flows for flow completion time (FCT)

## Observation:

- Goodput and FCT degradation
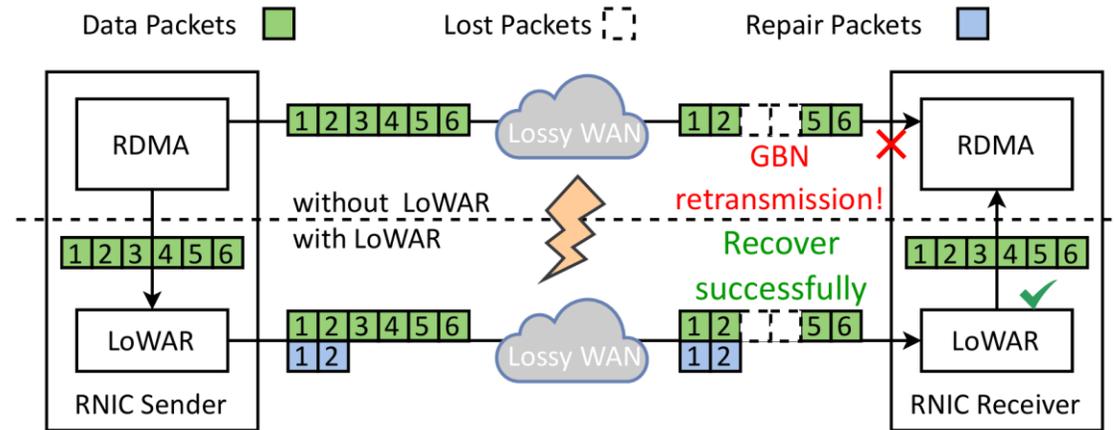- *e.g.*: 0.001% loss, goodput drops from 9.21 Gbps to 2.92Gbps

## Why:

- Go-Back-N retransmission
- High latency & loss rate

❌ Great bandwidth is wasted when GBN is frequently triggered.

> GBN limits RDMA's performance in lossy WAN!

# Methodology

Incorporate forward error correction (FEC) to protect RDMA messages from packet loss, thus minimizing the inefficiency of GBN.
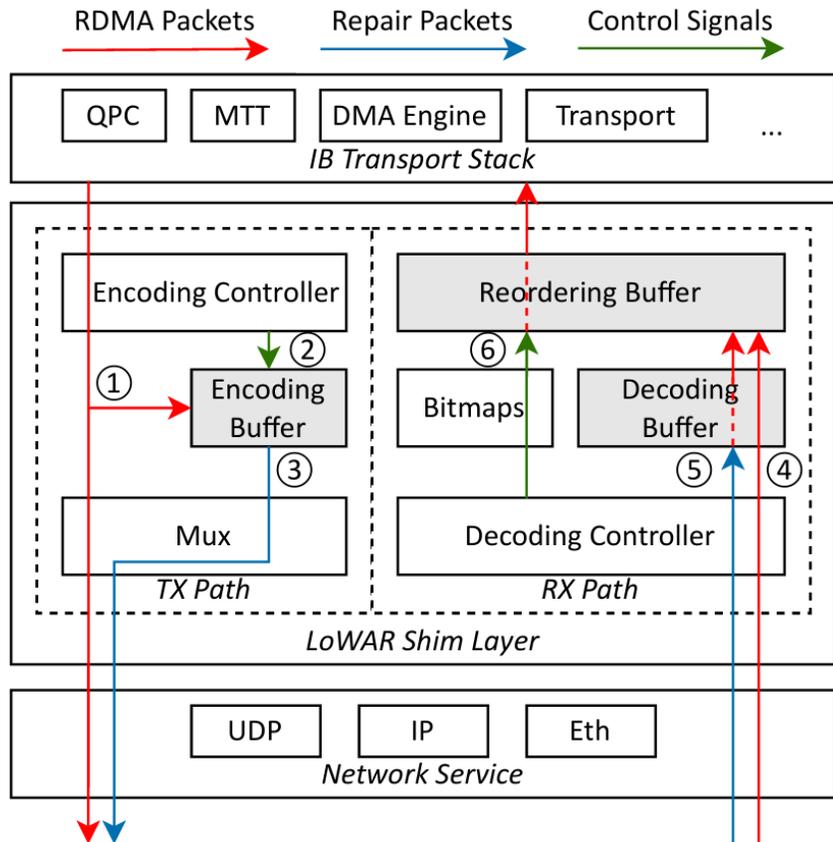


Workflow of end-to-end packet-level FEC

**By integrating packet-level FEC into RDMA, it can:**

- Enables independent loss recovery with less bandwidth consumption than retransmission.
- Deploys cost-effectively without modification on any WAN infrastructure.
- Leverages the hardware offload capability of SmartNIC for FEC calculations.

# Lossy Wide Area RDMA (LoWAR)

A high-goodput, high-reliability RDMA solution optimized for lossy WAN



LoWAR Architecture

**Hardware-offloaded FEC shim layer:**

- **TX** : RDMA transport layer packets ➡ Repair packets

- **RX** : Broken message ✚ Repair packets ➡ Full message

**Design Requirements:**

① **Compatible** with RDMA's message-based pattern

② **Real-time** packet loss recovery

③ **Minimal** storage and computing latency

④ **Transparent** to applications, upper layer RDMA transport stacks, and network infrastructures
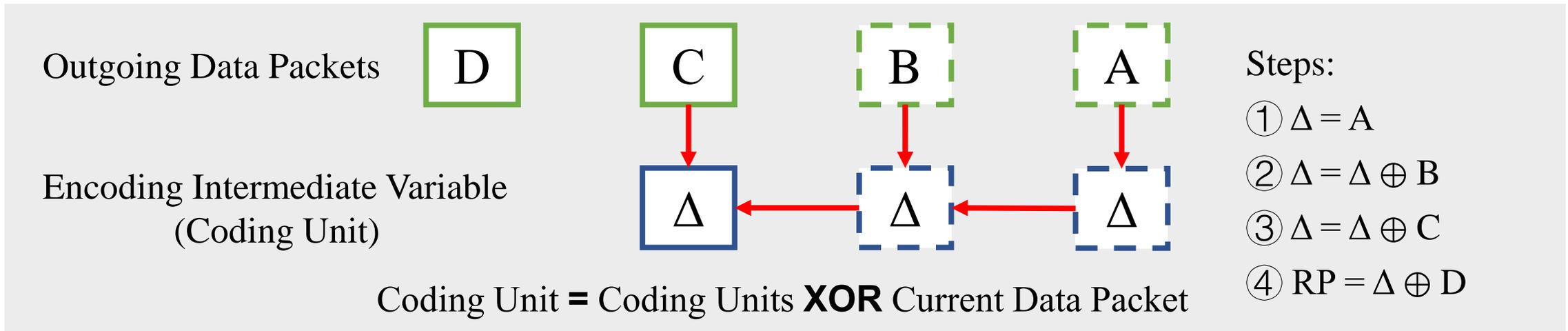
☑ Fully hardware-offloaded shim layer design ④

# Packet Alignment Coding Model

**Packet Alignment**: each new message re-initiate new coding process

**Coding block**: $r$ data packets to $c$ interleaved repair packets

☑ Compatible with RDMA's message-base communication. (①)

# Buffer-based Update-style Calculation

**XOR-based code:** calculation in steps. Coding finishes once all data sent/received

☑ Real-time repair packets decoding and lost packets recovery (②)

**Coding Units:** save the intermediate encoding/decoding results for subsequent unfinished calculation, instead of the whole coding block

☑ Minimal storage and computing latency (③)



Outgoing Data Packets

Encoding Intermediate Variable (Coding Unit)

Coding Unit **=** Coding Units **XOR** Current Data Packet

Steps:
① $\Delta = A$
② $\Delta = \Delta \oplus B$
③ $\Delta = \Delta \oplus C$
④ $RP = \Delta \oplus D$

# Repair Header (RH) Extension

**Negotiation Channel:**

- Synchronize parameters of FEC (redundancy rate, interleaving depth, packet length…)

**Repair Packets:**

- Repair Header + XOR Payload

**Negotiation Packets:**

- Repair Header Only

- Receiver to proactively change parameters

☑ Control path transparency to application and CPU(④)

Negotiation Packet format:

| Eth Header | IP Header | UDP Header | **Repair Header** | ICRC | FCS |
|---|---|---|---|---|---|

Repair Packet format:

| Eth Header | IP Header | UDP Header | **Repair Header** | XOR Payload | ICRC | FCS |
|---|---|---|---|---|---|---|

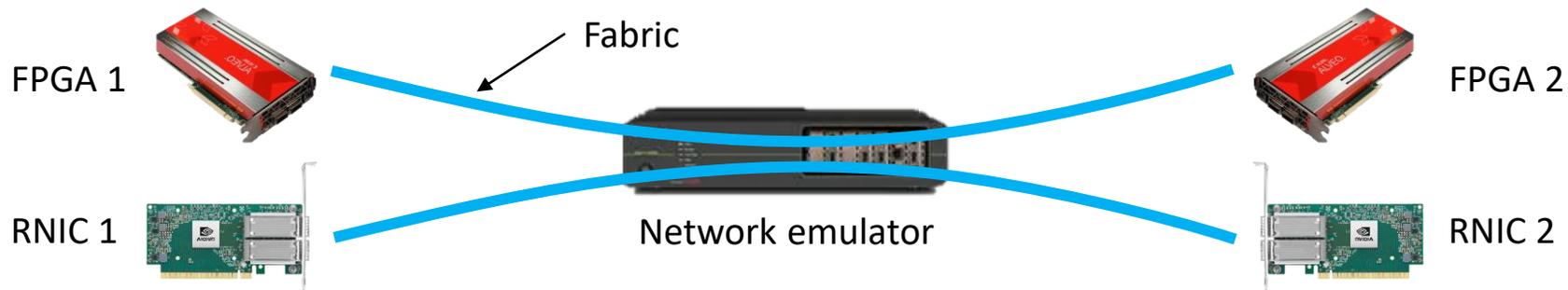| Bytes/bits | 31-24 | 23-16 | 15-8 | 7-0 |
|---|---|---|---|---|
| 0-3 | 0x1F | Destination QPN | | |
| 4-7 | Type | Type-specified Fields | | |
| … | | … | | |

# Implementation

## LoWAR Prototype

- Xilinx Alveo U200 FPGA

- Based on open-source RoCE FPGA implementation

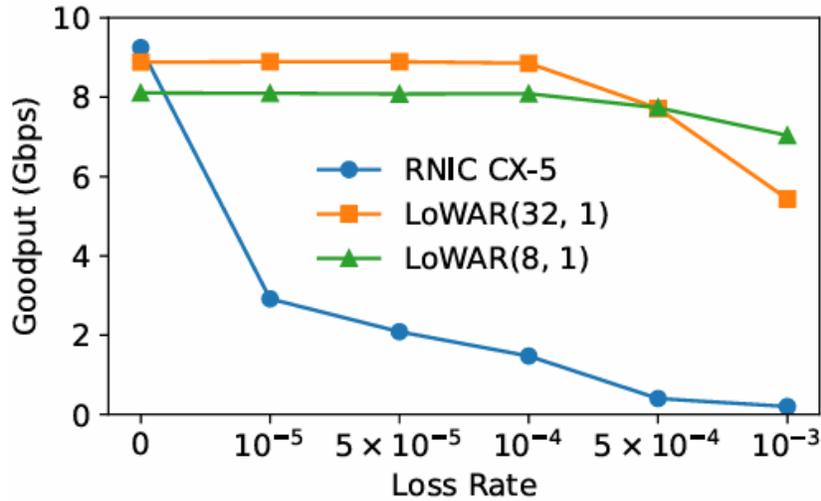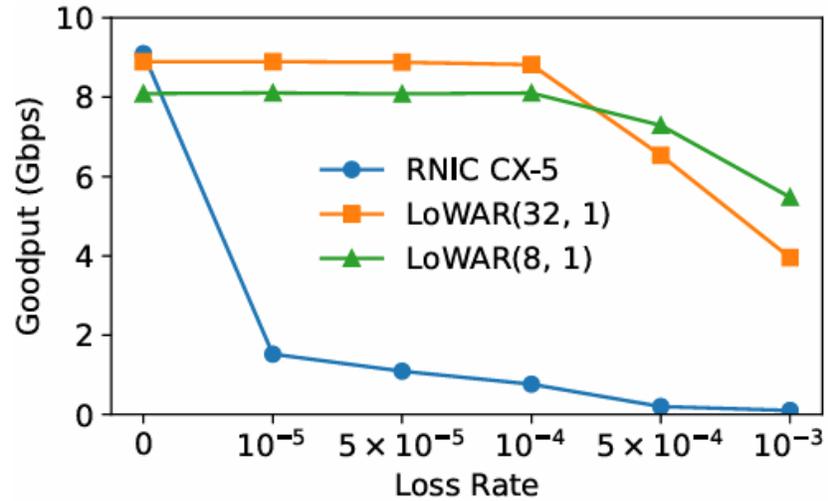- Place LoWAR between UDP and IB Transport



## End-to-End Testbed Settings

- Two hosts with Xilinx U200 FPGA and Mellanox Connect-X 5 RNIC

- Spirent SNE-X network emulator to  emulate bidirectional lossy wide-are links

- Applications based on Xilinx driver for LoWAR, and OFED *perftest* for CX-5
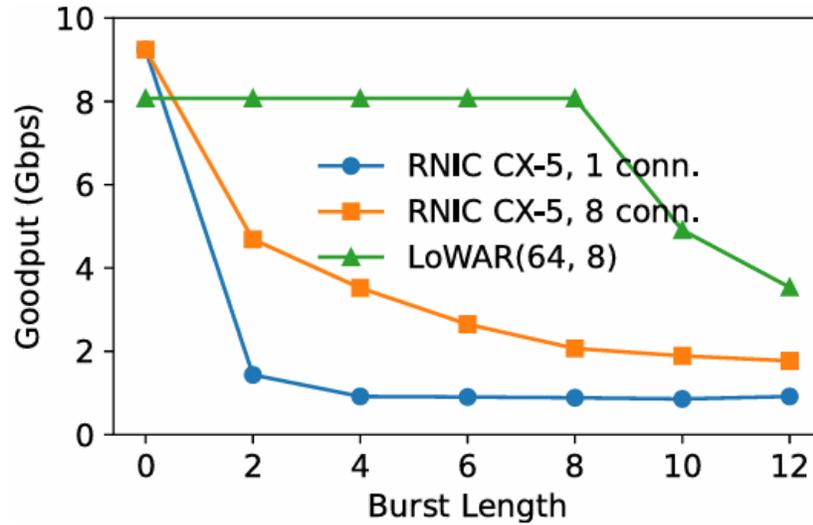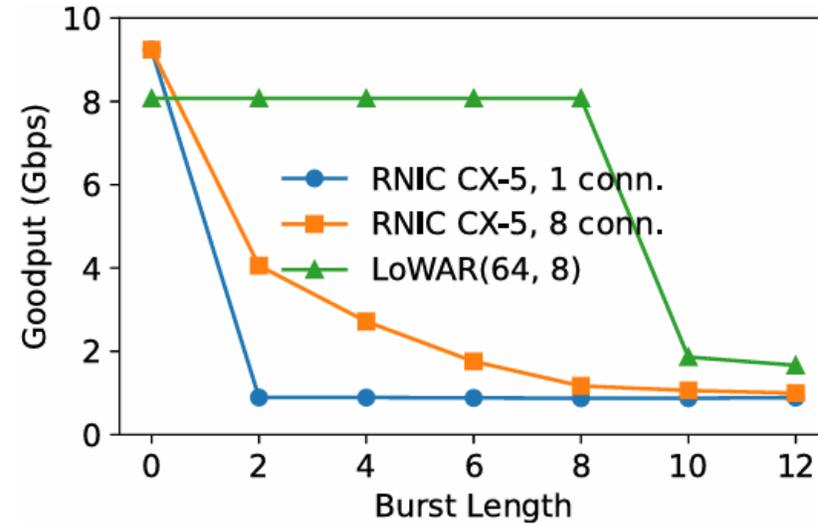
# Goodput Improvement



(a) RTT = 40ms, 1 connection    (b) RTT = 80ms, 1 connection

- Increase goodput by 2.05x to 5.01x with 40ms RTT, and 11.55x to 19.07x with 80ms
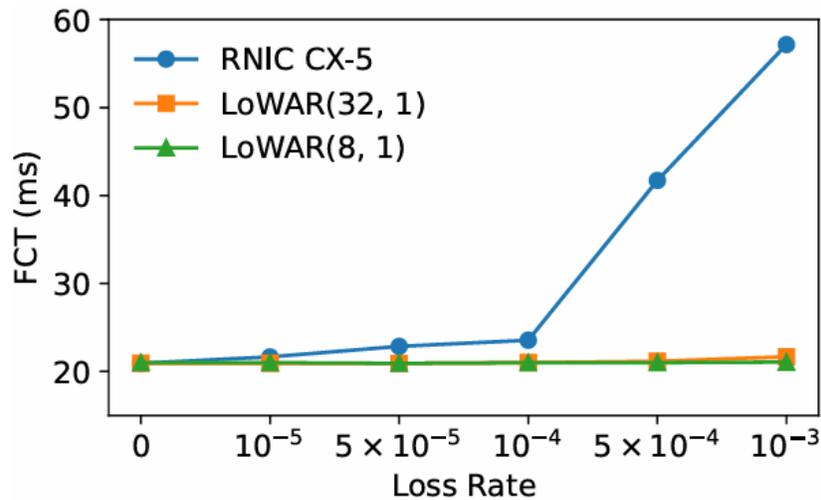- A low redundancy rate (32 data -> 1 repair) suffices in most scenarios

# Burst Tolerance



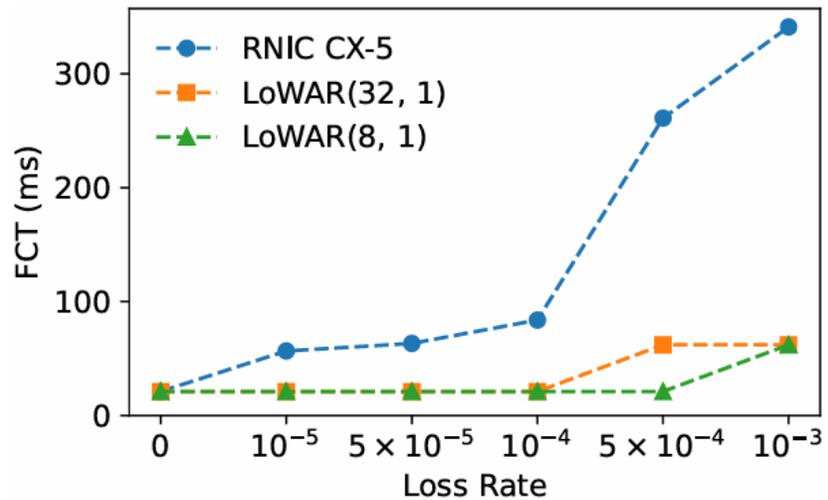(a) $p_{high} = 0.3$  (b) $p_{high} = 0.7$

- Resilient to burst loss (with two-state Markov model)
- Low loss rate = 0. High loss rate = 0.3/ 0.7. Transition rate (low to high) = 0.01%
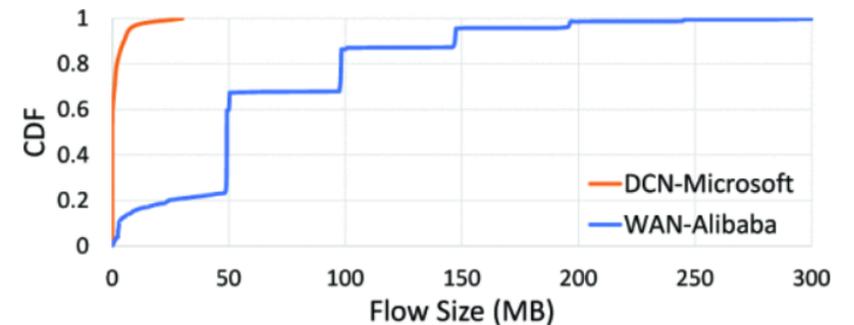
# Flow Completion Time
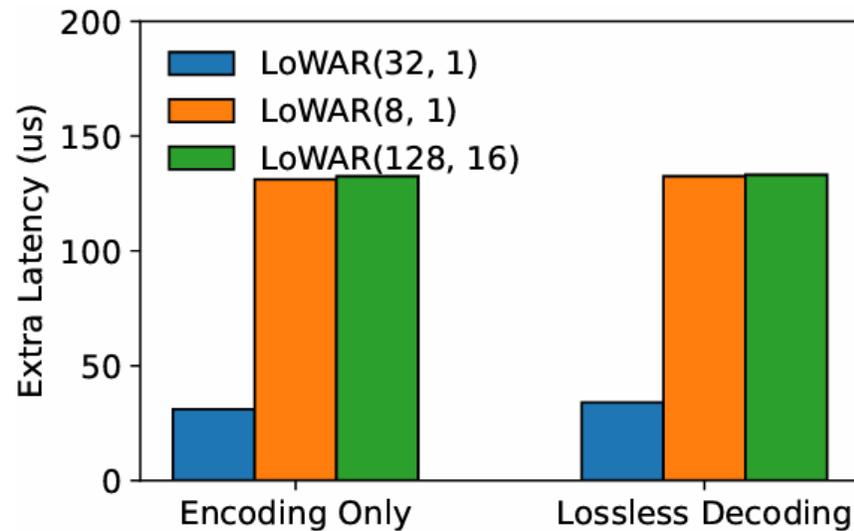


(a) Average, RTT = 40ms

(b) 99*th* tail, RTT = 40ms

- Use 1 MB message for FCT test.
  - Note: 1MB is actually a small flow for WAN[6]!

- 12% to 97% lower FCT with 0.01% to 0.05% loss rate
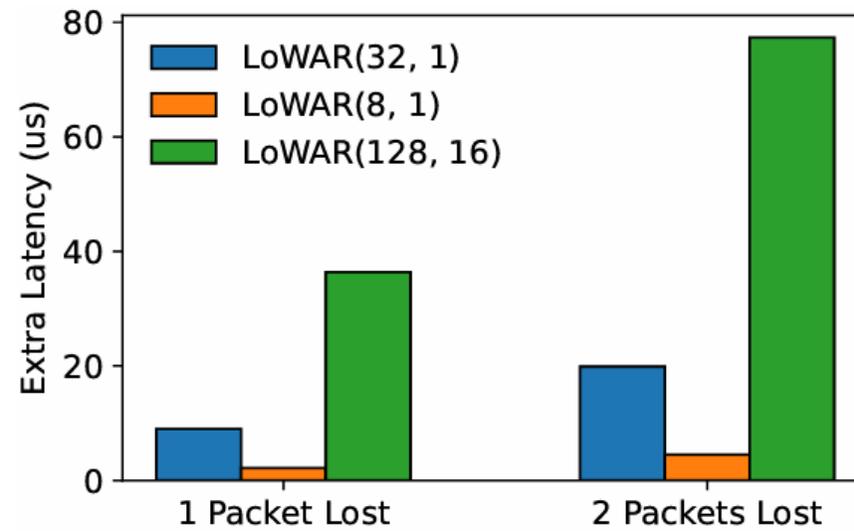
- Reduce long-tail FCT in most scenarios



[5] **ICNP`21**, G. Zeng, K. Chen, et al. FlashPass.

# Latency Overhead



(a) Lossless latency

(b) Recovering latency

- Encoding/Decoding: little latency when lossless
- Recovering: 10~100 us, $\propto$ coding block size and broken blocks number
- Latency comes from packet reordering and draining

# Conclusion

- **LoWAR provides high-goodput, high-reliability RDMA over lossy WANs by incorporating *end-to-end forward error correction***
  ✓ *Fully offloaded in RNIC as a shim layer, with minimal storage overhead and computational latency*
- LoWAR protects RDMA against packet loss and significantly mitigate the retransmission inefficiency.
- LoWAR works transparently and requires no modifications on both applications and networks.

We look forward to your questions and suggestions!