

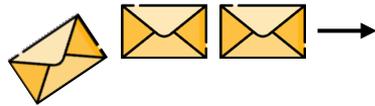
Masking Corruption Packet Losses in Datacenter Networks with Link-local Retransmission

Raj Joshi, Cha Hwan Song, Xin Zhe Khooi, Nishant Budhdev, Ayush Mishra,
Mun Choon Chan, Ben Leong



NUS
National University
of Singapore

Packet Loss in Datacenter Networks



Packet Loss



**Degrades Application
Performance**



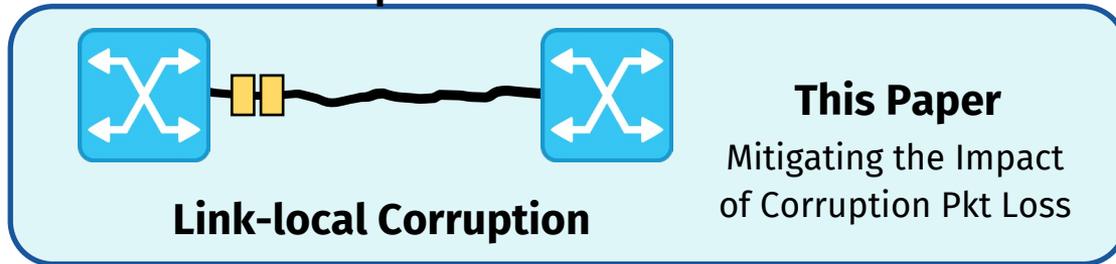
**Impact on
Revenue**



Congestion



Significant Prior Work



Little Attention

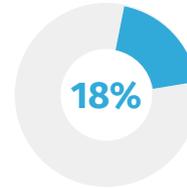
Why do we care about corruption packet loss?

Corruption Packet Loss – Significant



Comparable to Congestion Loss

Large-scale study by Microsoft
(350K links, 15 datacenters)
[Zhuo et al. SIGCOMM'17]



Packet drops affected customers

due to corruption – study by
Alibaba Cloud
[Zhou et al. SIGCOMM'20]



Increase in tail FCTs



Drop in throughput

How can we fix corruption packet loss?



Physical Repair

Several hours to days

Until then



Mitigate/Mask

Effects of corruption pkt loss

Prior Work – Mitigating Corruption Losses

Current state-of-the-art

Disable corruption

[e.g. NetPilot, etc.]

(* more related work)



Wireless Networks

Not revisited in the context of datacenter networks

Classical loss-recovery strategy



Challenges unique to DC context

Our Work

Link-local Retransmission

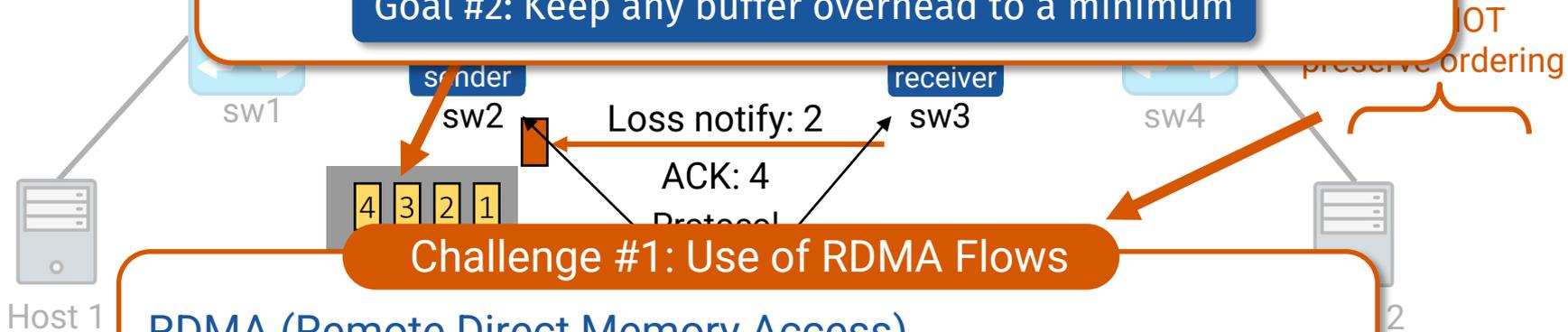
[LinkGuardian]

Basic Link-local Retransmission

Challenge #2: Shallow Buffered Switches

Switches in DC networks very constrained on the packet buffer

Goal #2: Keep any buffer overhead to a minimum



Challenge #1: Use of RDMA Flows

RDMA (Remote Direct Memory Access)

- Pkt reordering → degrades FCT performance of RDMA flows

Goal #1: Preserve packet ordering (at high link speeds)



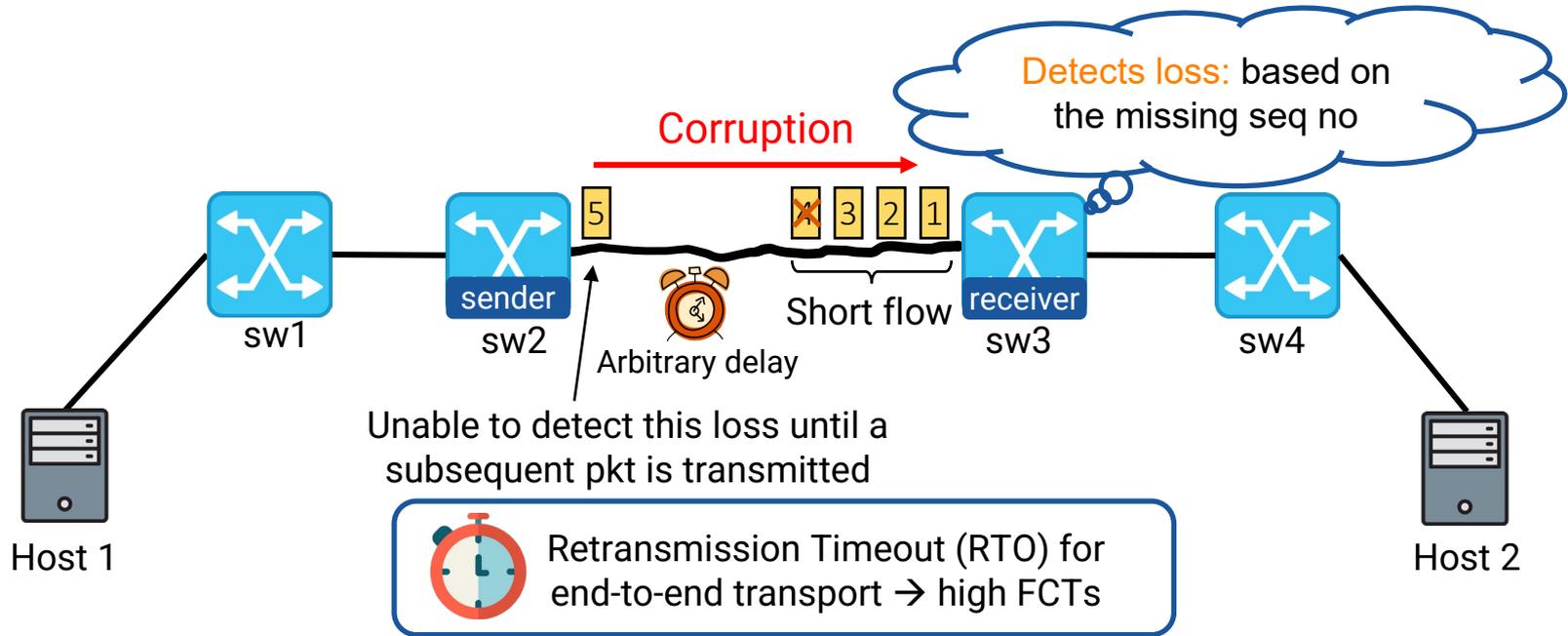
Works



next of
datacenter networks!

Challenge #3: Flows in datacenter are very short!

Short flows + corruption → higher probability tail (last) packet is lost



Goal #3: Detect and recover tail packet losses (as quickly as possible)

LinkGuardian

LinkGuardian – 3 Key Techniques



Preserve Pkt Ordering

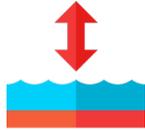


Ordering Buffer

At receiver switch



This talk!



Shallow Buffered Switches



Fast & Efficient ACK Mechanism

At receiver switch



Tail Pkt Loss



Self-Replenishing Queue of Dummy Pkts

At sender switch

Basic Link-Local Retransmission Scheme

Handle Tail Packet Loss

Sender Switch

Corruption →

Receiver Switch

Decreasing Strict Priority



ACK: 3

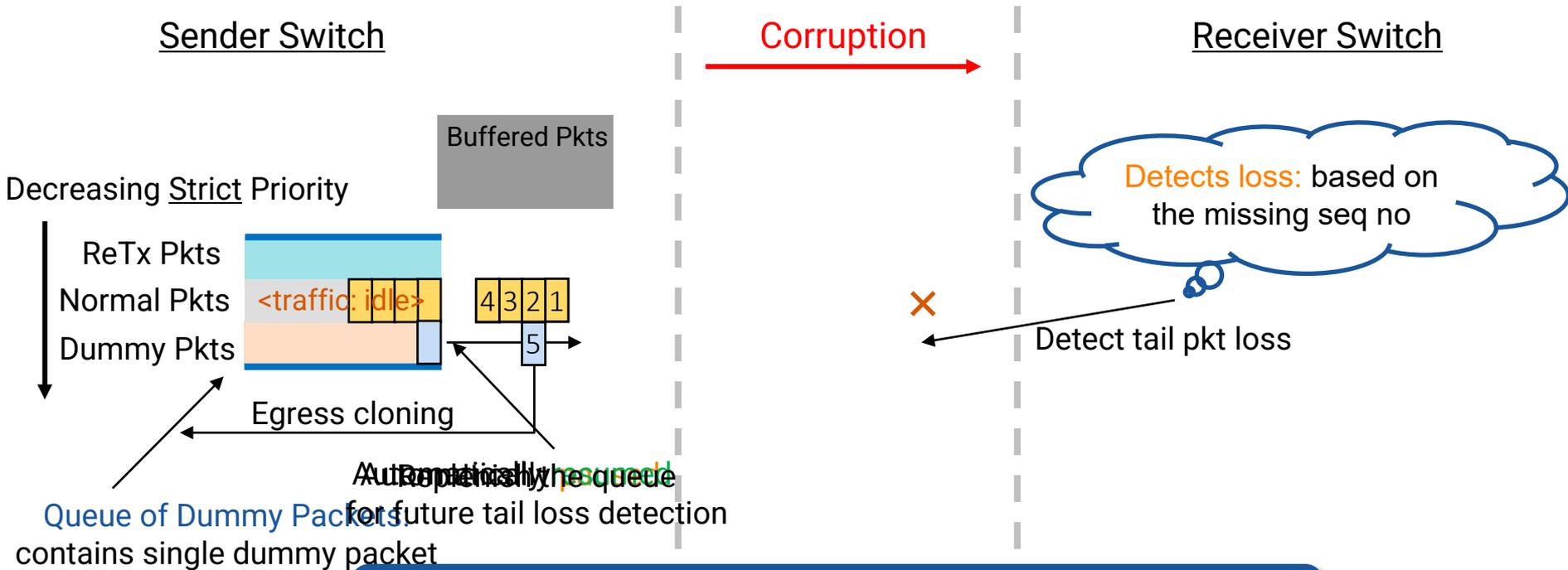
Standard: ReTx after a timeout
LinkGuardian: Timeout-less approach!
Risk spurious retransmissions

Key Idea



Send a "dummy" link-local packet when traffic is idle → help receiver switch detect the tail pkt loss

Handle Tail Packet Loss



LinkGuardian Implementation and Evaluation



Intel Tofino programmable switch (~1800 lines of P4)



Source code available!

Hardware Testbed

- Intel Tofino switches, Intel Xeon servers w/ 100 Gbps NICs

LinkGuardian Variants

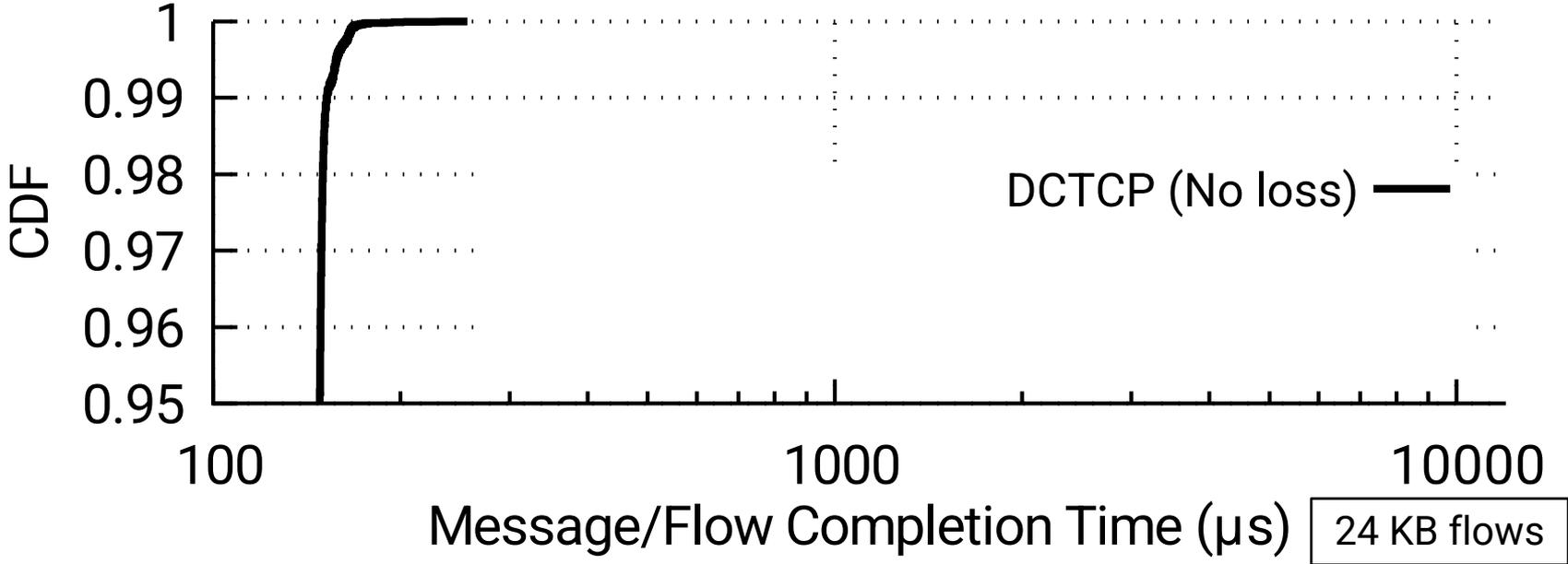
- LinkGuardian (default): preserves packet ordering
- LinkGuardian**NB** (non-blocking): “lite” version → everything, except preserving ordering



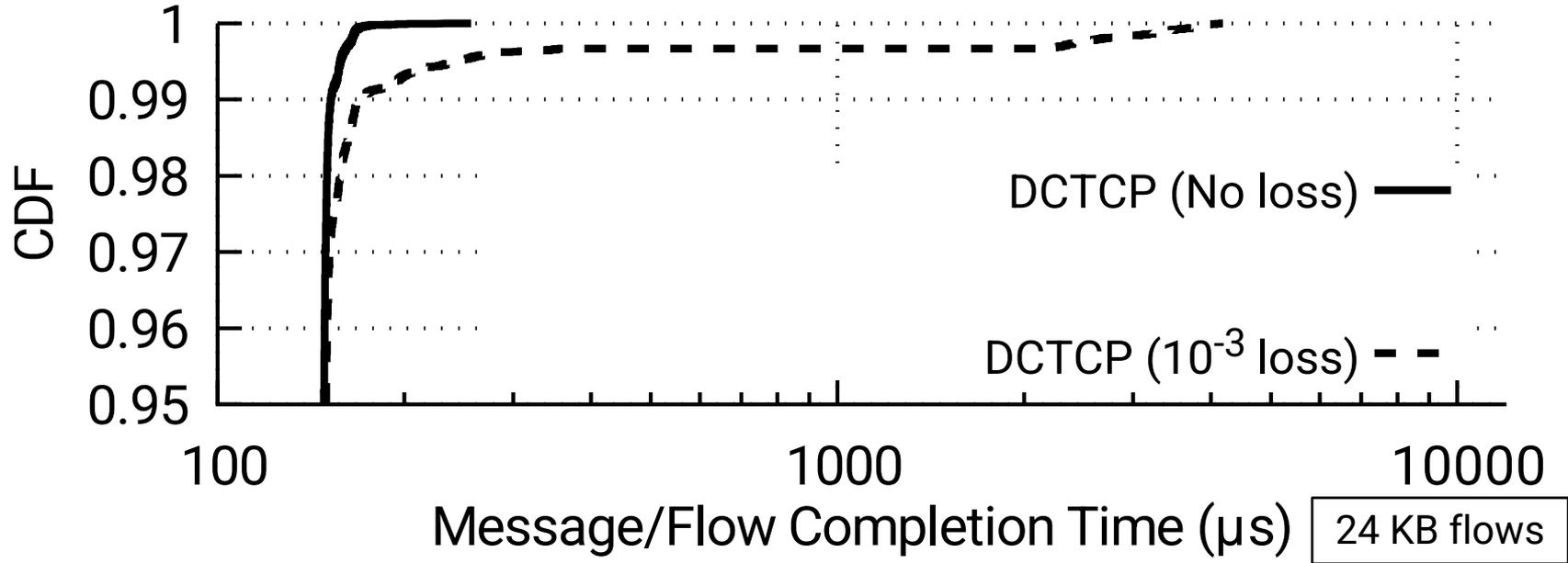
Variable Optical Attenuator

Corruption Loss Rate: 10^{-3}
(more loss rates in the paper)

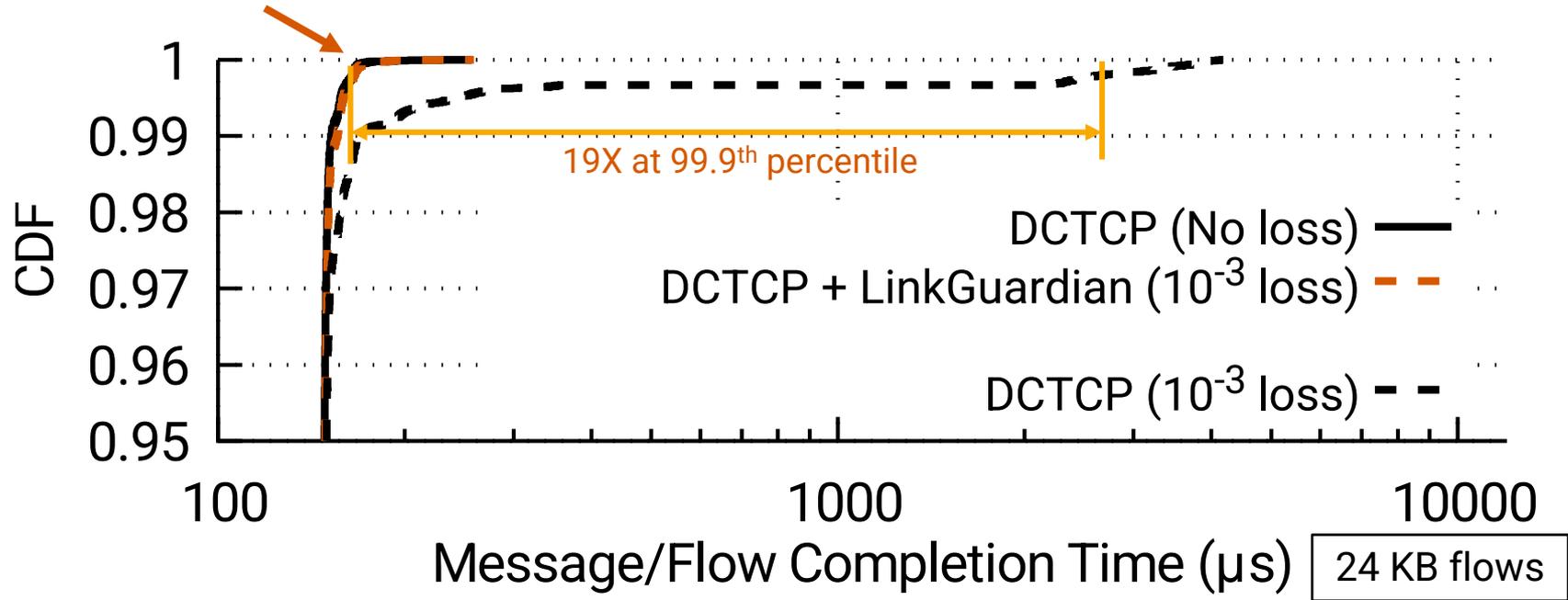
LinkGuardian improves tail FCTs for TCP (DCTCP)



LinkGuardian improves tail FCTs for TCP (DCTCP)

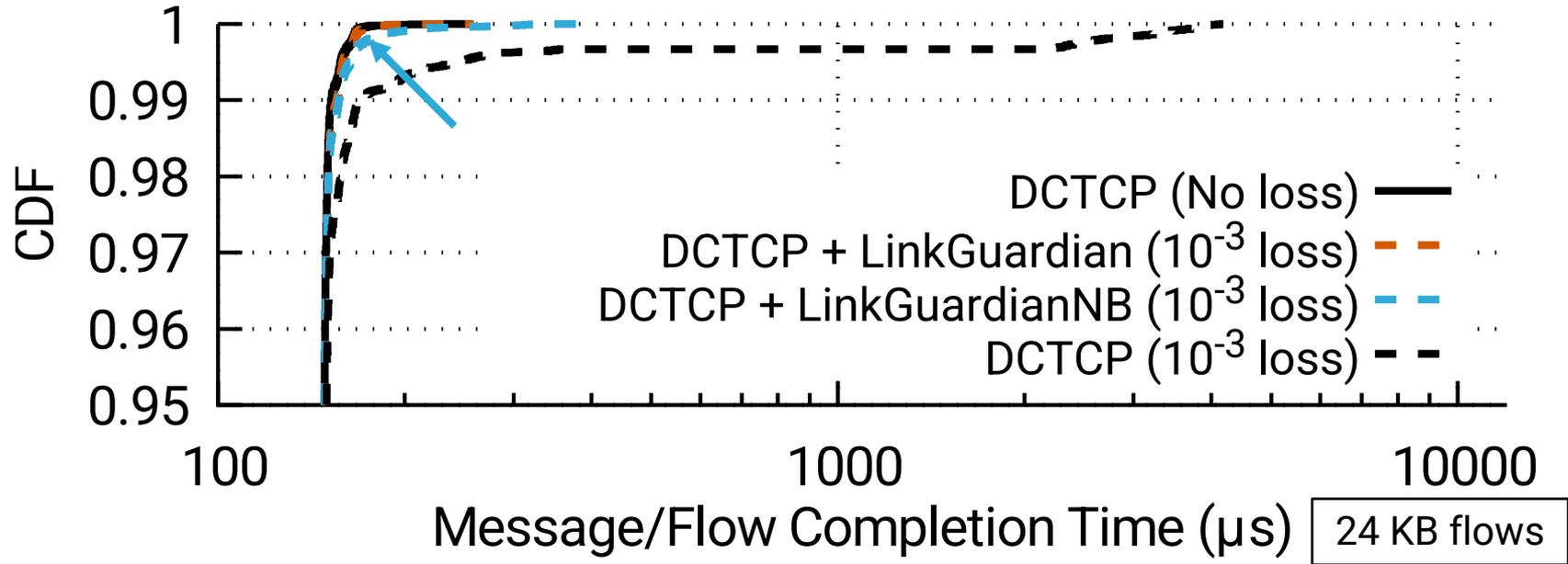


LinkGuardian improves tail FCTs for TCP (DCTCP)



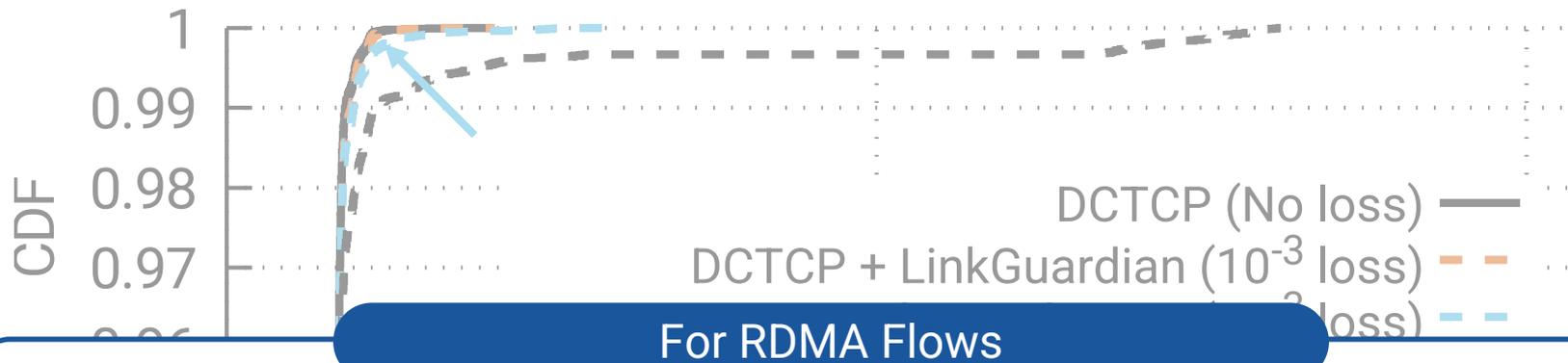
LinkGuardian: improves 99.9th percentile FCT by 19 times. Matches no loss case!

LinkGuardian improves tail FCTs for TCP (DCTCP)



LinkGuardianNB: nearly the same improvement, even without maintaining pkt order!
TCP flows are short: reduction in cwnd (due to out-of-order) has minimal impact!
(more detailed analysis in the paper)

LinkGuardian improves tail FCTs for TCP (DCTCP)

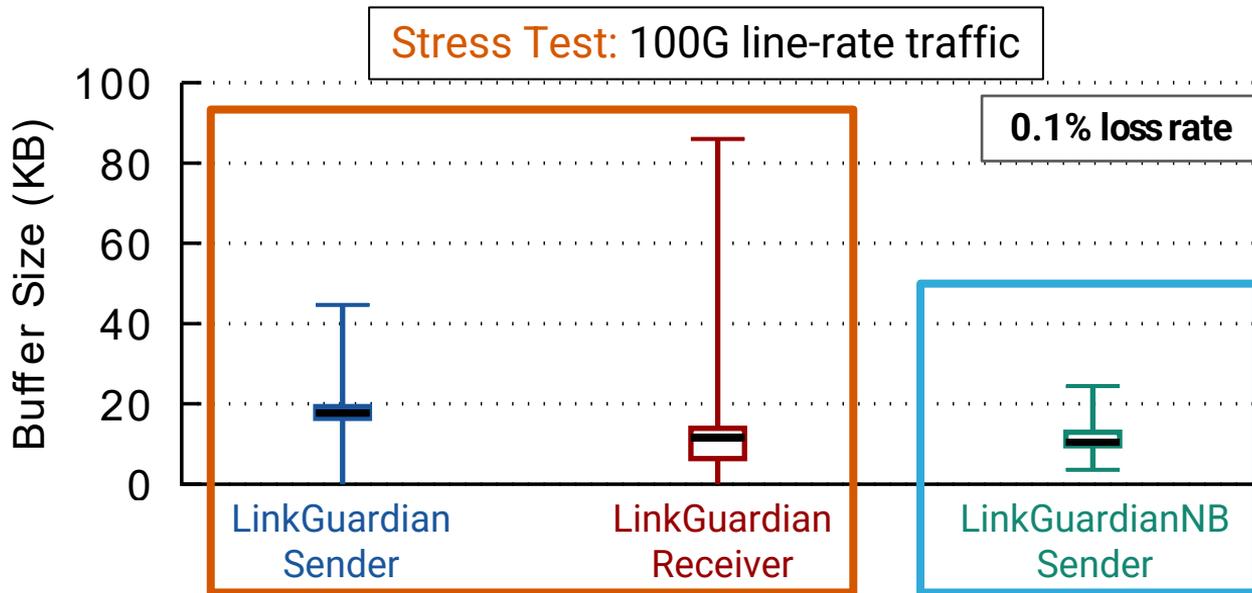


LinkGuardian: improves 99.9th percentile FCT by **~39 times**.
(detailed RDMA results in the paper)

Message/Flow Completion Time (μs) [24 KB flows]

LinkGuardianNB: nearly the same improvement, even without maintaining pkt order!
TCP flows are short: reduction in cwnd (due to out-of-order) has minimal impact!
(more detailed analysis in the paper)

LinkGuardian's Packet Buffer Overhead



LinkGuardian: < 90KB for both sender and receiver switches

LinkGuardianNB: < 25KB for sender switch. No ordering buffer at receiver!

Put this result in context: Switches have 10's of MBs of packet buffer

Summary

LinkGuardian → mitigate corruption pkt loss



Selective In-network Recovery (sub-RTT)

pkt caching in dataplane



Handles Tail Pkt Loss

Without employing any timeouts

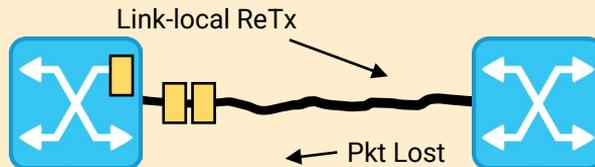


Preserves Pkt Ordering

For reordering sensitive RDMA flows

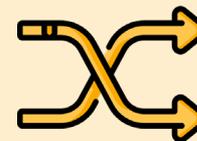


Key Takeaways



Link-local ReTx

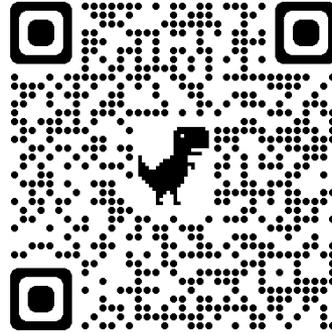
Both *practical* and *effective* in today's DC networks



Out-of-Order ReTx

Sufficient for short TCP flows

Thank you!



<https://github.com/NUS-SNL/linkguardian>

Credits: This presentation template was adapted from [Slidesgo](#), including icons by [Flaticon](#).

Flaticon icons by: Becris, Vectors Market, Flat Icons, Maxim Basinski Premium, Darius Dan, noomtah, iconixar, Kiranshastry, Bombasticon Studio, fjstudio, Vectorslab, Chattapat, Freepik, Andy Horvath, itim2101, Pavel Kozlov, Pixel perfect, Cuputo, Icon home, Fajri rama, Good Ware, Vitaly Gorbachev, Eucalyp, pongsakornRed, and Octopocto.