

Do Switches Still Need to Deliver Packets in Sequence?

Ufuk Usubütün, Fraida Fund, Shivendra S. Panwar



NYU

TANDON SCHOOL
OF ENGINEERING



NYU

WIRELESS



Conventional wisdom:

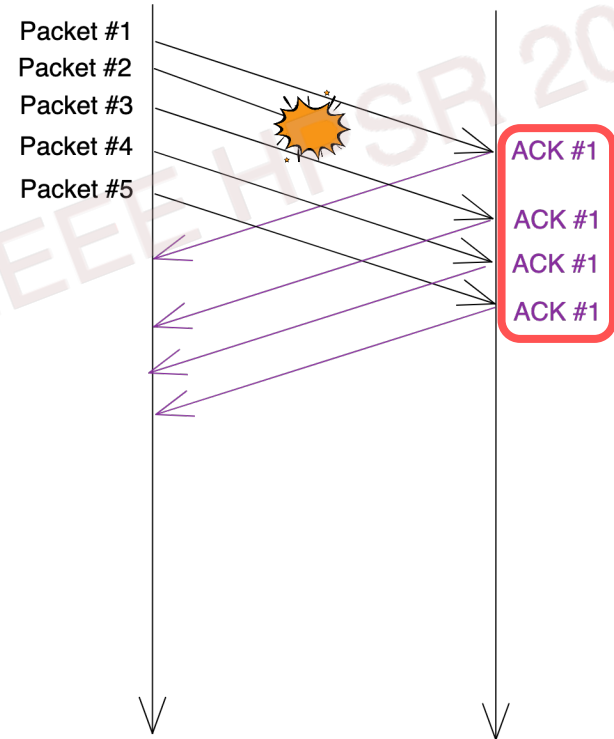
Packet reordering should be avoided
wherever possible

Ufuk Usubutun et. al. - IEEE HPSR 2023

How did Classical TCP Detect Packet Loss?

Triple Duplicate ACK Rule

- Counts repeated ACKs



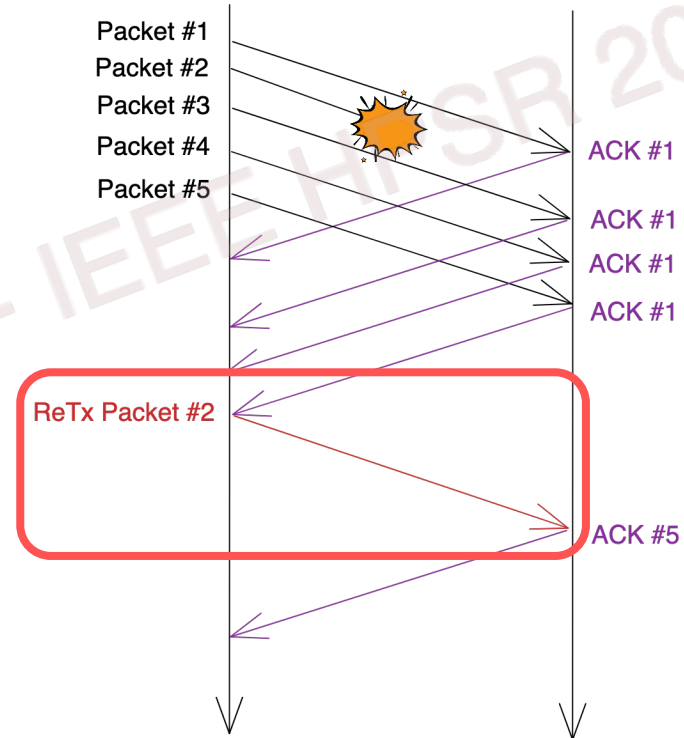
How did Classical TCP Detect Packet Loss?

Triple Duplicate ACK Rule

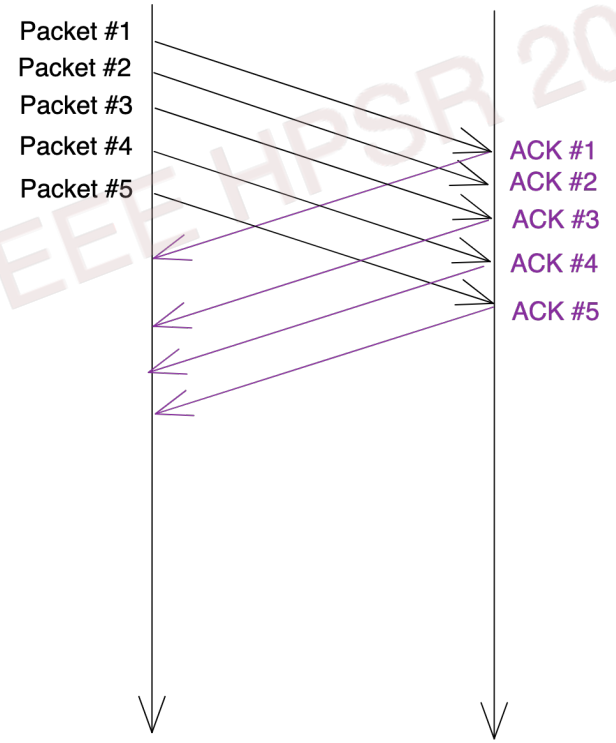
- Counts repeated ACKs

Loss detected:

- Retransmits lost packets
- Reduces *cwnd*



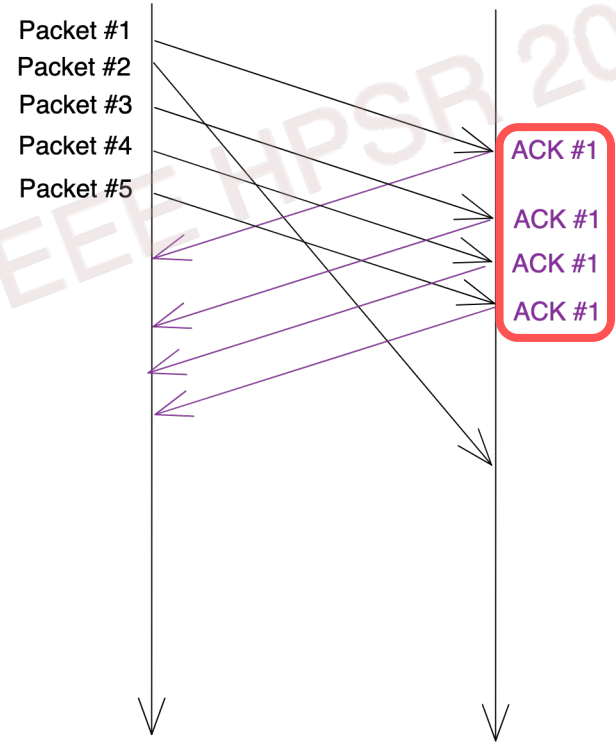
Assumption: Packets will arrive in sequence!



Ufuk Usubutun et. al. - IEEE HPSR 2023

Assumption: Packets will arrive in sequence!

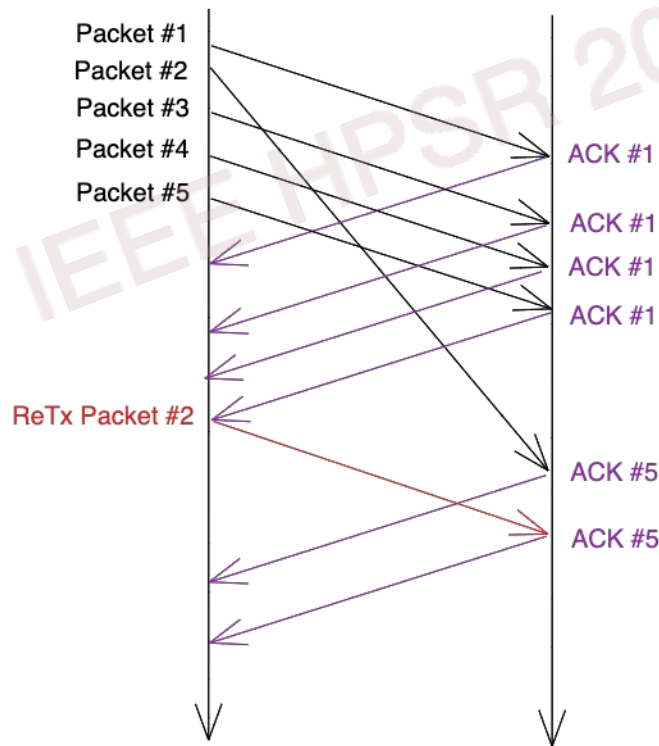
- What happens if they don't?



Assumption: Packets will arrive in sequence!

- What happens if they don't?

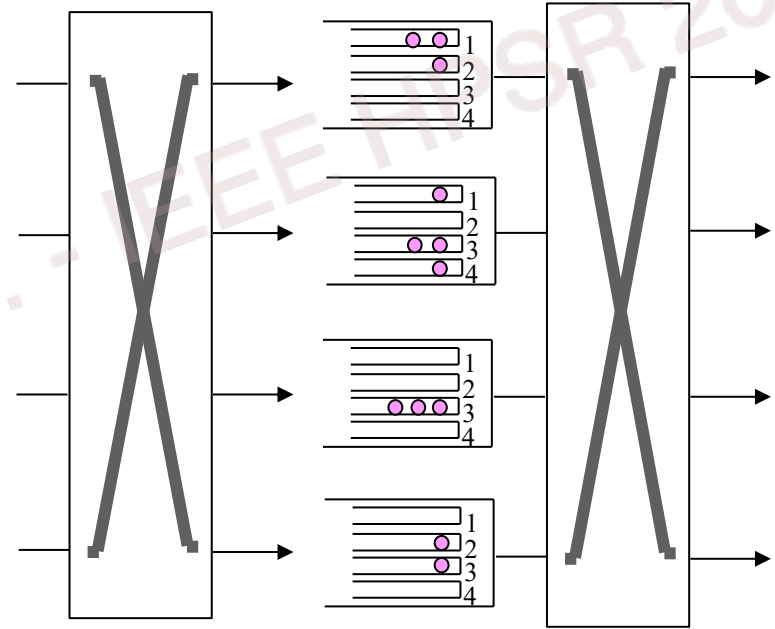
Reordering is misinterpreted as loss!



Switches today are expected to deliver packets in-sequence

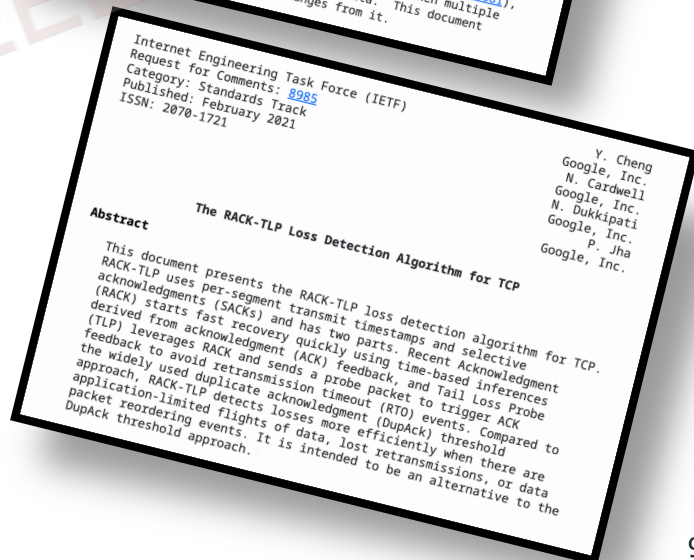
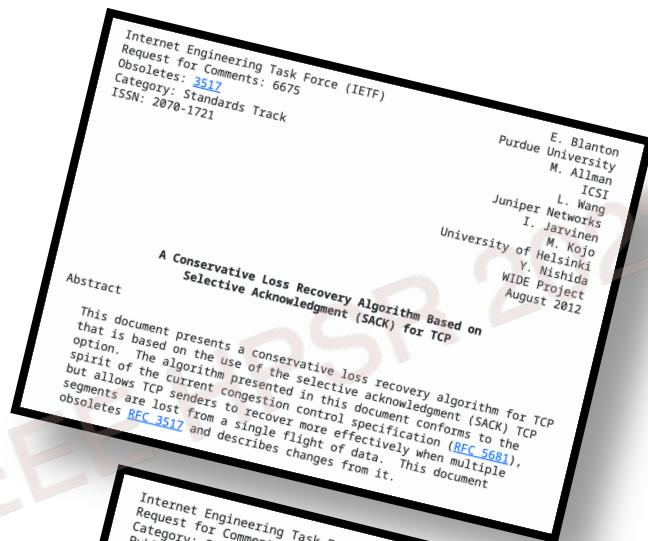
E.g.: Load-Balanced Birkhoff–von Neumann Switch:

Was not adopted since it caused packet reordering



In the past two decades: Two major changes

- Advanced loss detection algorithms for TCP widely deployed

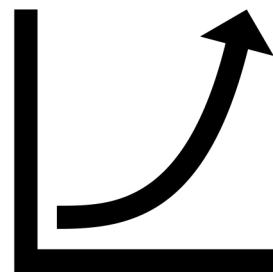


RFC 8985, "The RACK-TLP Loss Detection Algorithm for TCP", Y. Cheng et. al.

RFC 6675, "A Conservative Loss Recovery Algorithm Based on Selective Acknowledgment (SACK) for TCP", E. Blanton et. al.

In the past two decades: Two major changes

- Advanced loss detection algorithms for TCP widely deployed
- Core network capacities grew from hundreds of Mbps to hundreds of Gbps



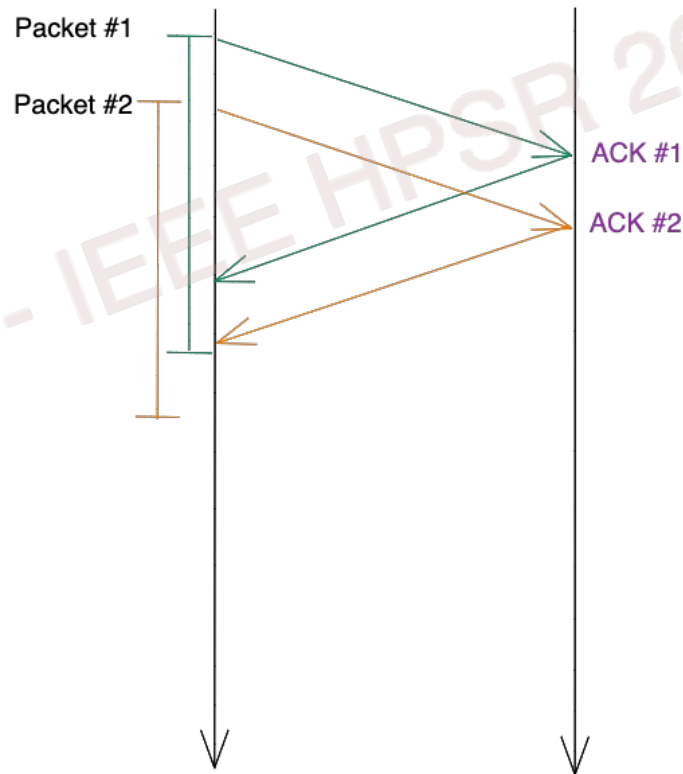
Advanced Loss Detection for TCP

Ufuk Usubutun et. al. IEEE HPSR 2023

Temporal Loss Detection is now the default way!

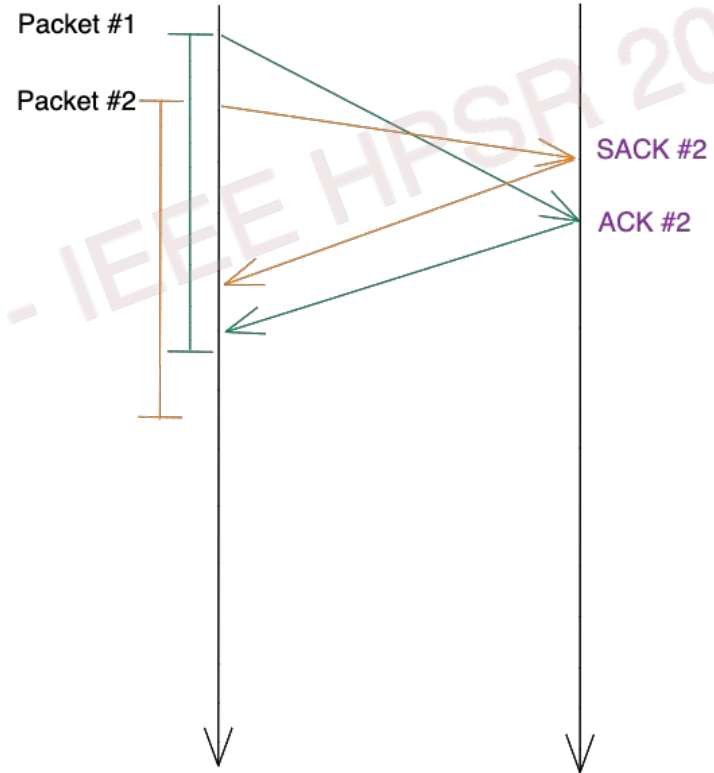
RFC 8985 - RACK:

- Replaces ACK counting
- Lateness triggers loss detection.



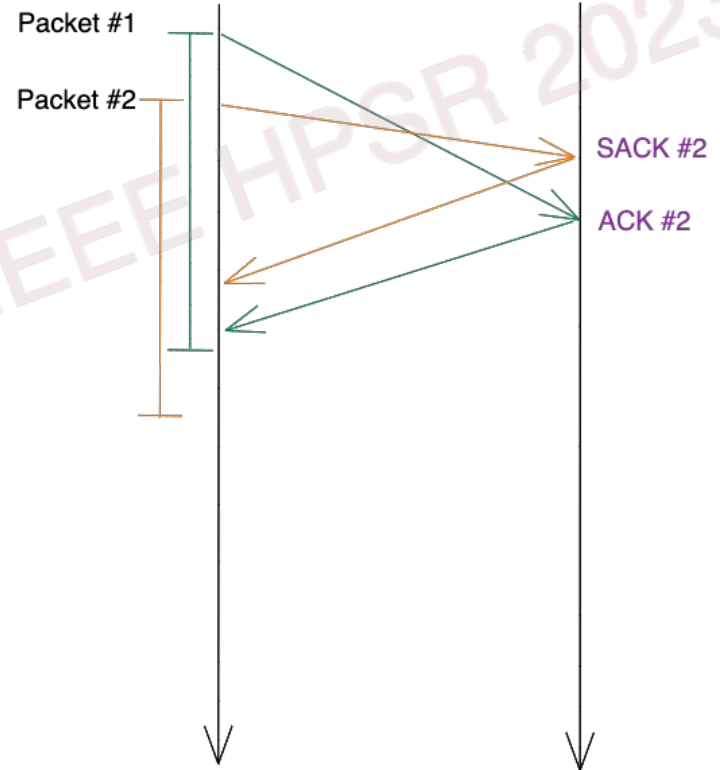
Arrival order is practically irrelevant

Tolerates reordering within time threshold!



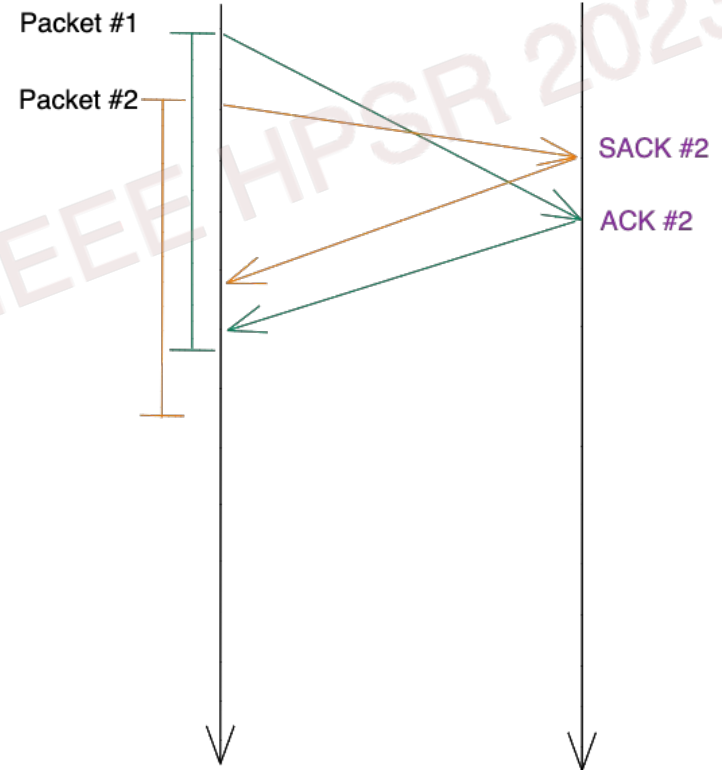
RACK Keeps a Dynamic Time Threshold

$$\text{Time Threshold} = \text{SRTT} + \text{reorder_wnd}$$

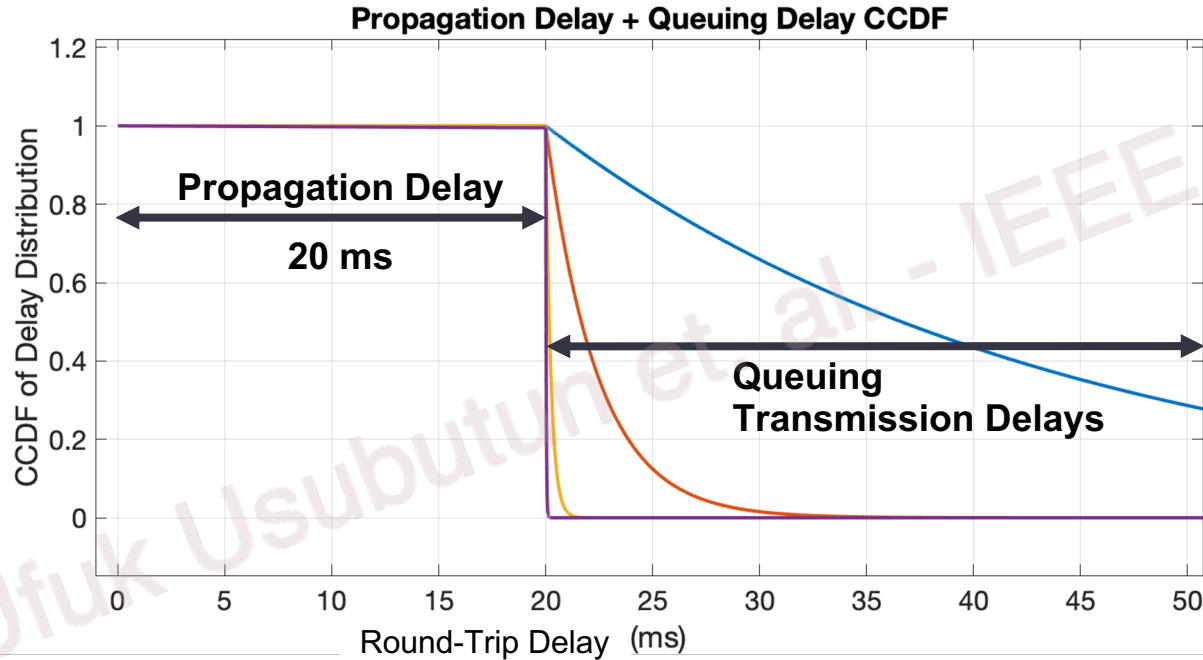


RACK Keeps a Dynamic Time Threshold

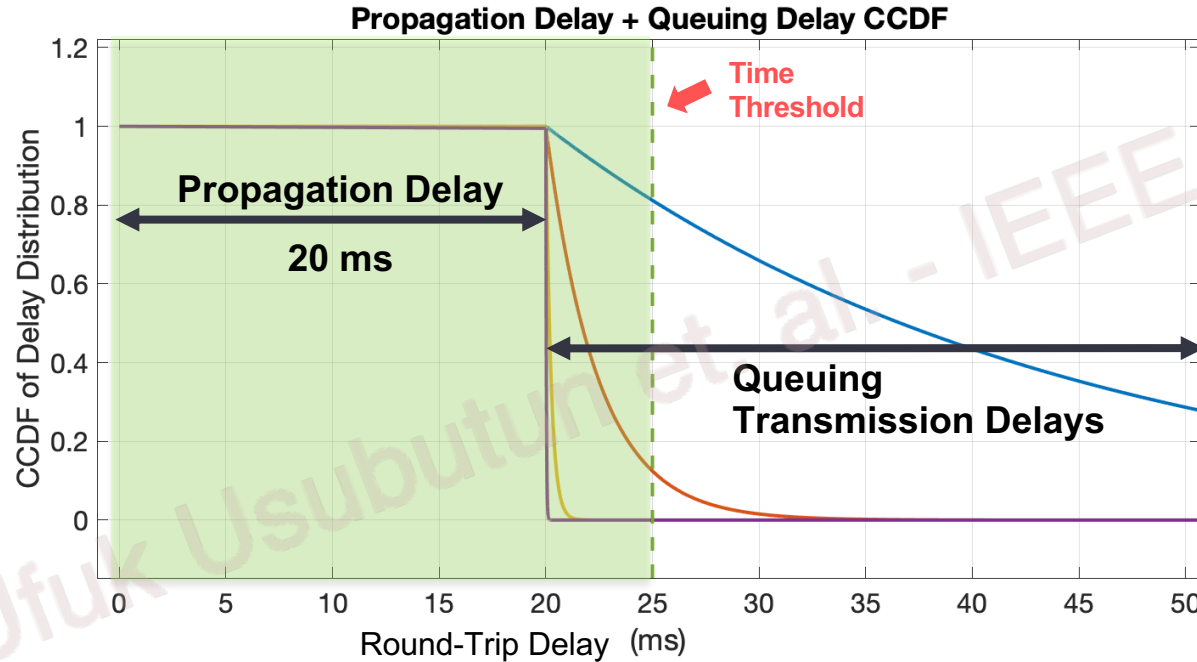
$$\text{Time Threshold} = \text{SRTT} + \text{reo_wnd}$$
$$(\text{min_RTT}/4 < \text{reo_wnd} < \text{SRTT})$$



Round-Trip Delay CCDF of Packets from 4 TCP Flows



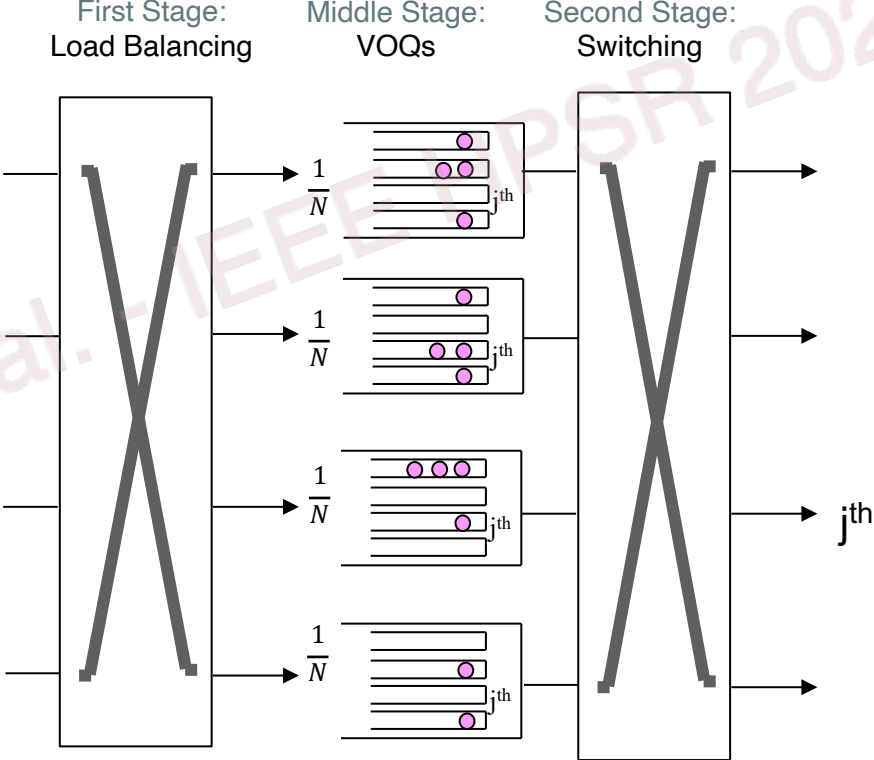
Round-Trip Delay CCDF of Packets from 4 TCP Flows



Low-Variation
fits within the
time window!

Load-Balanced Switches

A load balancer spreads all incoming packets uniformly

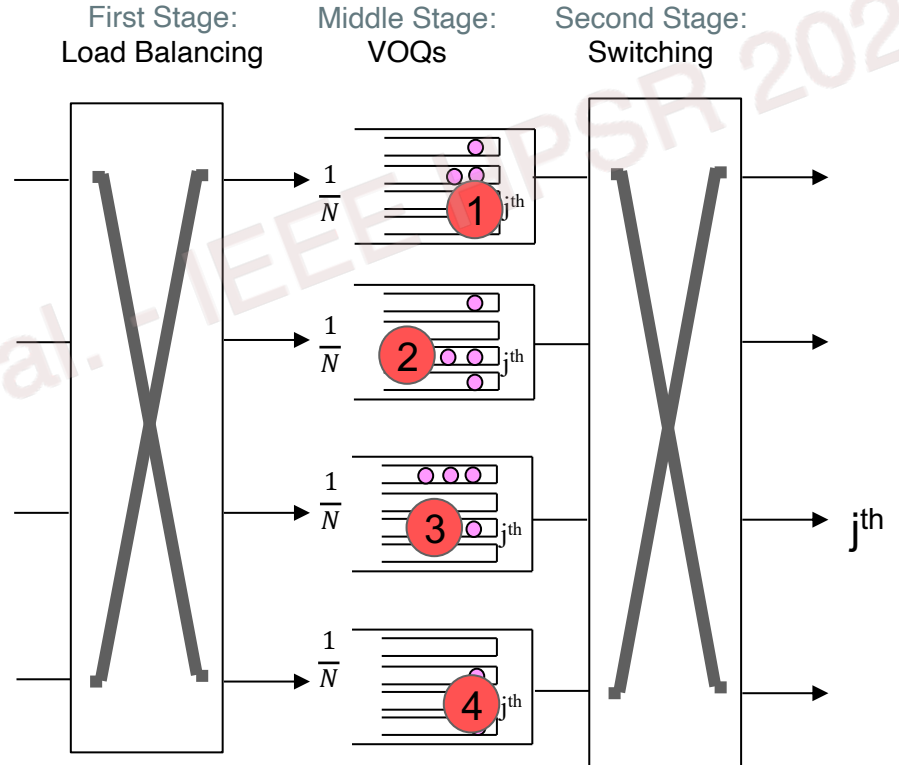


Ufuk Usubutun et. al. IEEE PSR 2023

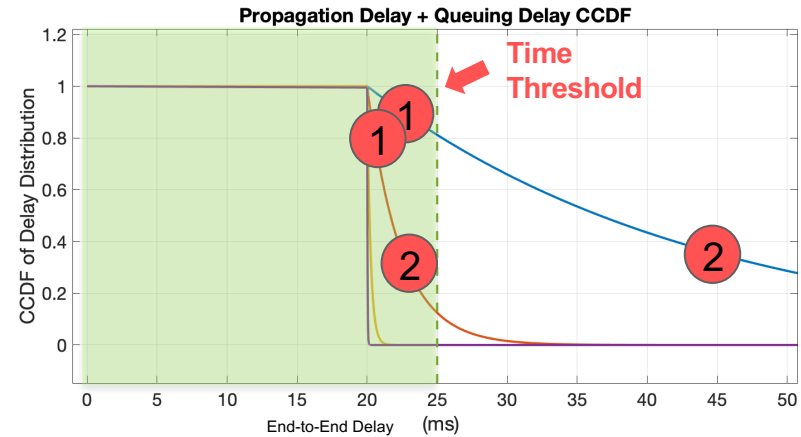
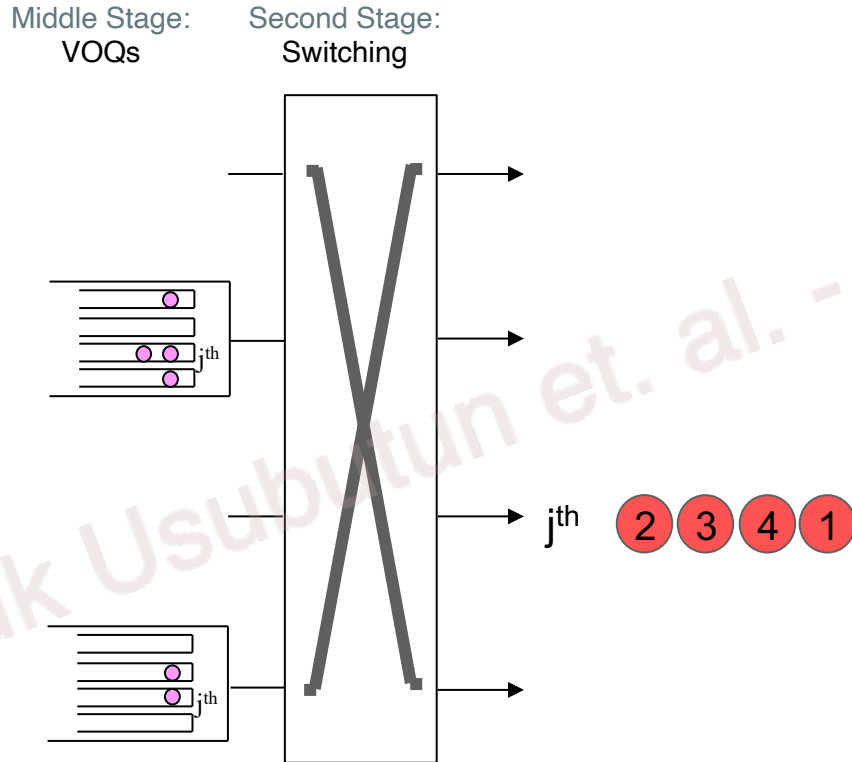
Load-Balanced Switches

A load balancer spreads all incoming packets uniformly

Packets from the same flow queue up at different VOQs



Lateness determines loss detection performance



Effect of Network Capacity Growth

Ufuk Usubutun et. al. IEEE HPSR 2023

Line rates increased from 100s of Mbps to 100s of Gbps

Then:



Now:



Both service rate μ and arrival rate λ are scaled up by a factor of $K = 1000$

Line rates increased from 100s of Mbps to 100s of Gbps

Then:



Now:



Both service rate μ and arrival rate λ are scaled up by a factor of $K = 1000$

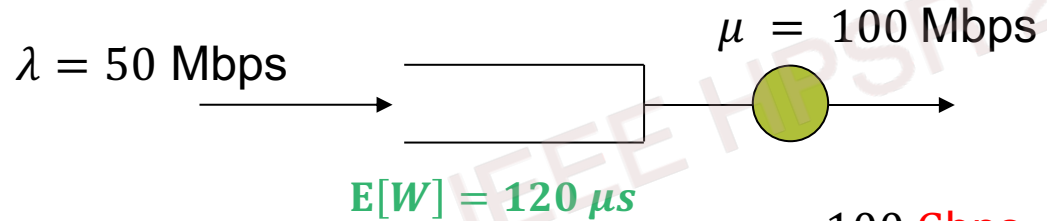
M/M/1 Queues:

Mean Queue Occupancy does not change

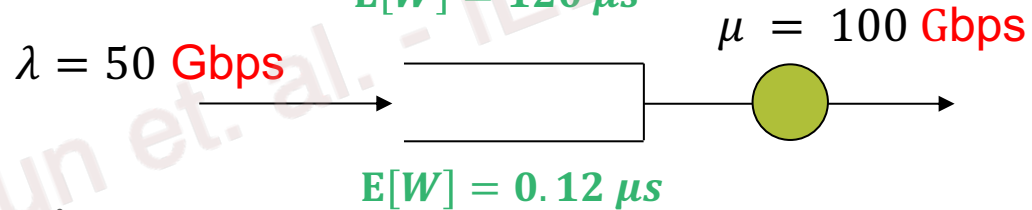
Mean Delay $E[W]$: scaled down by K

Line rates increased from 100s of Mbps to 100s of Gbps

Then:



Now:



Both service rate μ and arrival rate λ are scaled up by a factor of $K = 1000$

M/M/1 Queues:

Mean Queue Occupancy does not change

Mean Delay $E[W]$: scaled down by K

The delay distribution, too, shrinks with line rate increase

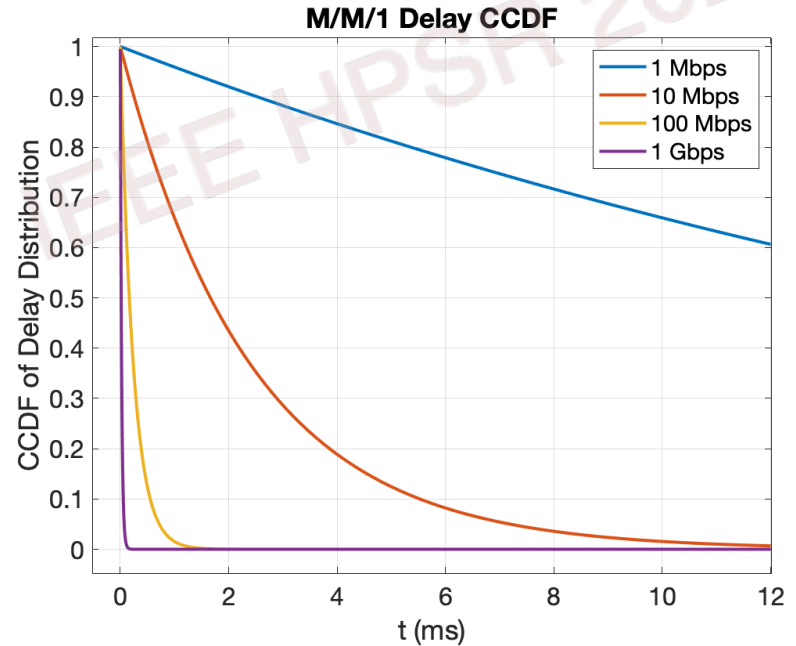
M/M/1 Queue

Tail Delay Probability:

$$P(W \geq \tau) = e^{-(\mu-\lambda)\tau} = e^{-K(\mu_0-\lambda_0)\tau}$$

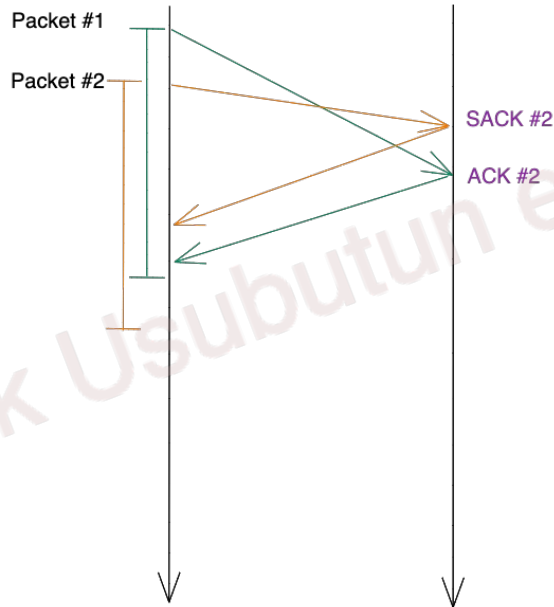
If both μ and λ are scaled by a factor of K

Tail probability is compressed by K

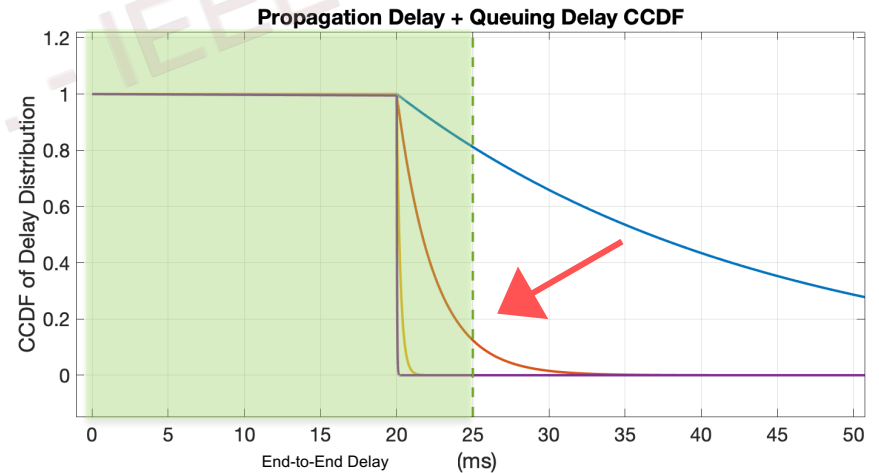


RECAP: Two Things Have Changed

- Time based loss detection, based on lateness of individual packets, is widely deployed.



- Increasing line rates lead to smaller delay and delay variation at network core switches.



Experimental Evaluation

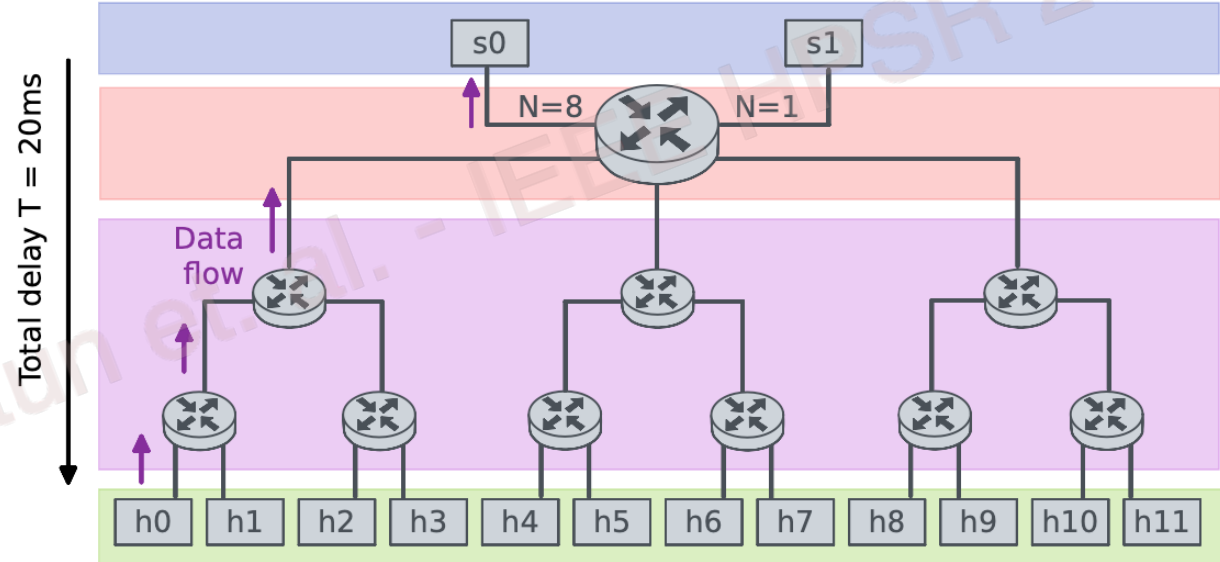
Ufuk Usubutun et. al. IEEE HPSR 2023



Testbed Experiments on Cloudlab

Goal:

Emulate a core network in a controlled environment



Ufuk Usubutun et al. - TELNWSR 2023

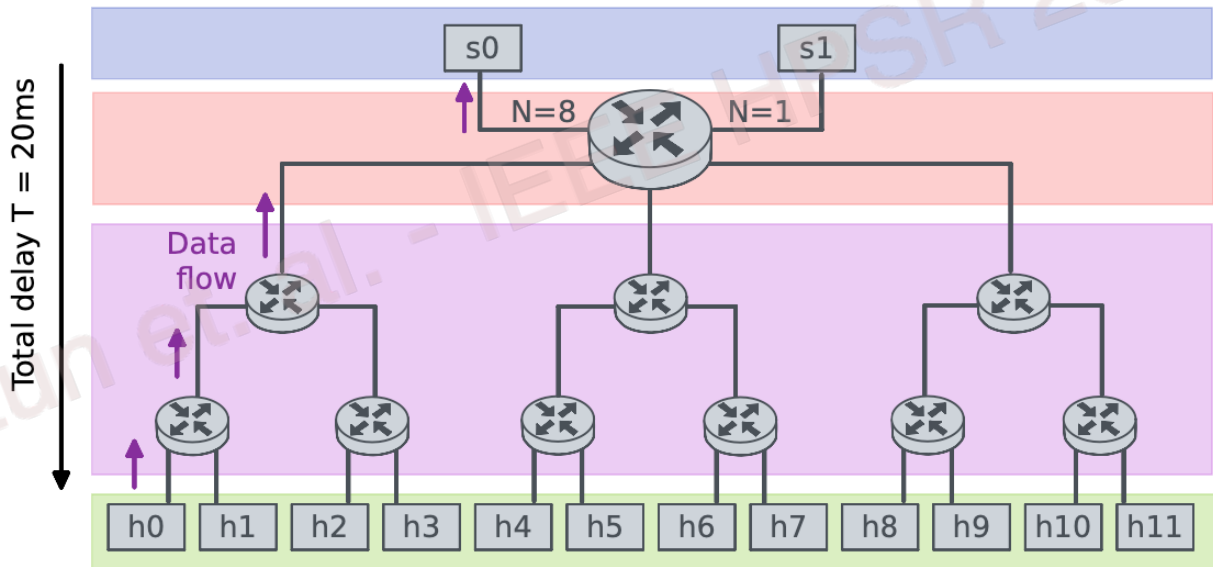


Testbed Experiments on Cloudlab

Goal:

Emulate a core network in a controlled environment

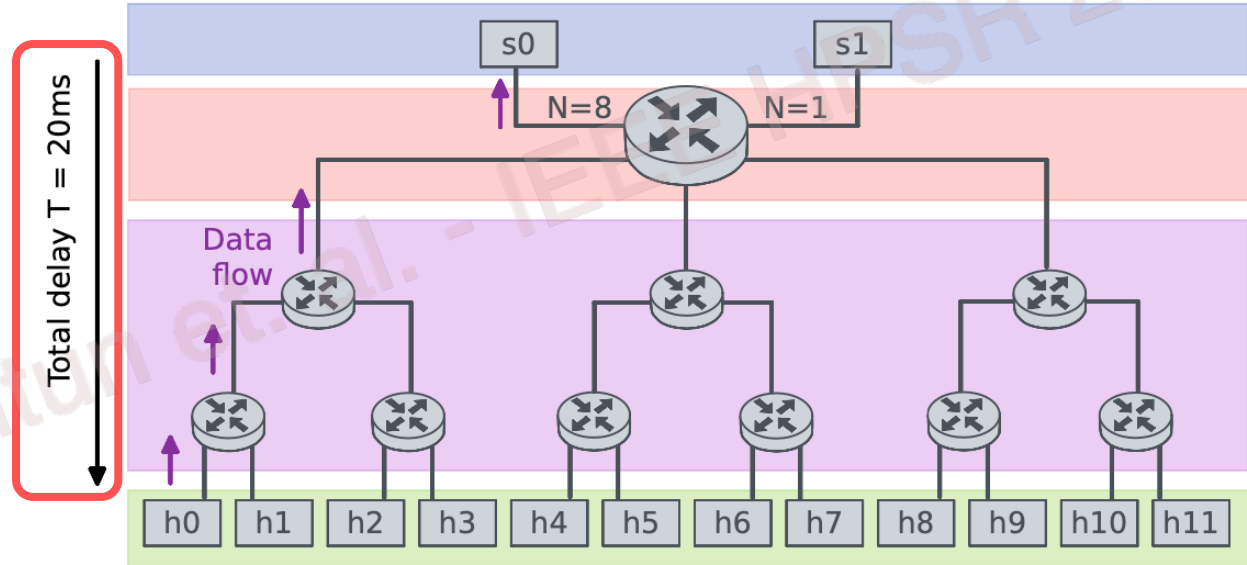
Thousands of TCP Cubic flows of different load sizes mixed through the topology





Testbed Experiments on Cloudlab

Fixed base delay applied to reverse direction



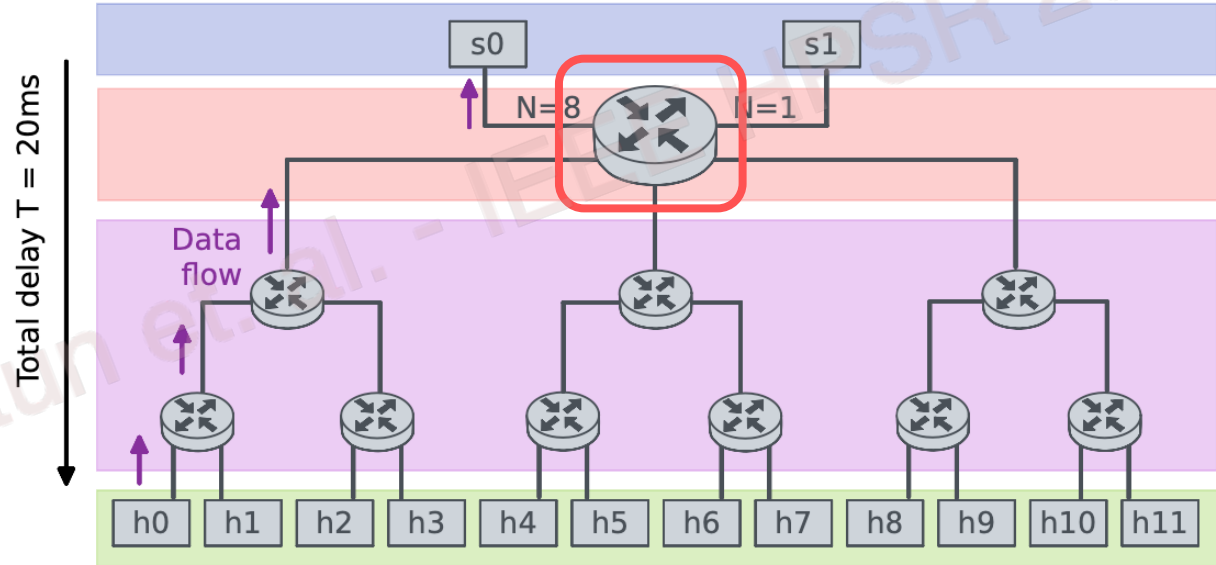
Ufuk Usubui et al. - TELNWSR 2023



Testbed Experiments on Cloudlab

Fixed base delay
applied to reverse
direction

Main switch was the
bottleneck



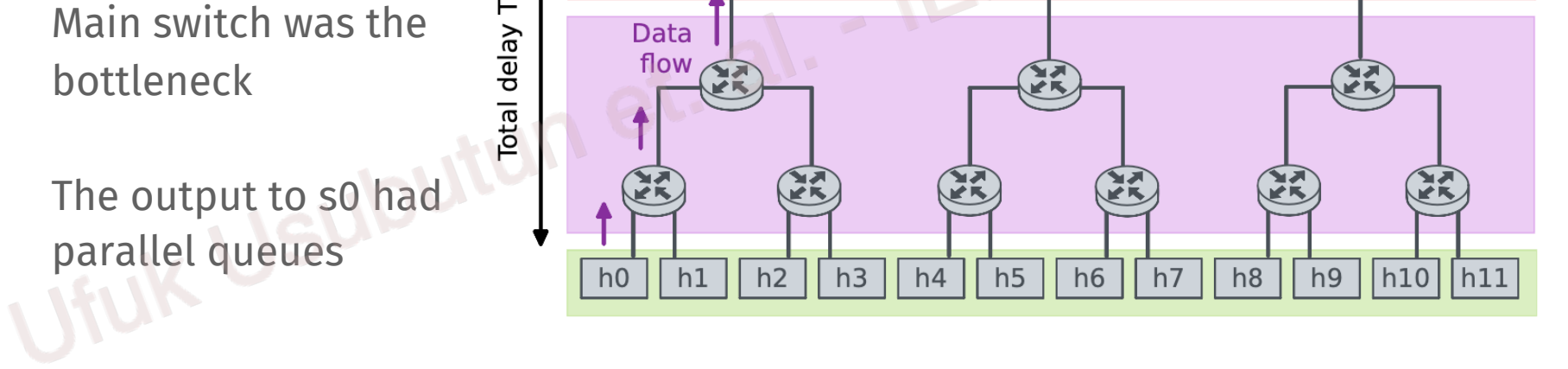
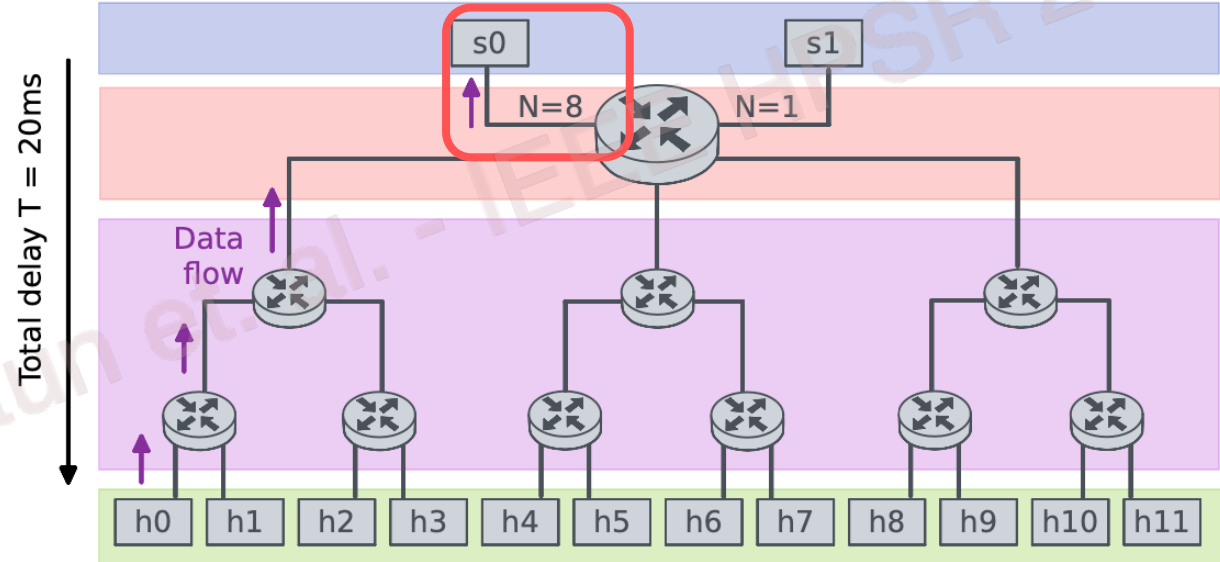


Testbed Experiments on Cloudlab

Fixed base delay applied to reverse direction

Main switch was the bottleneck

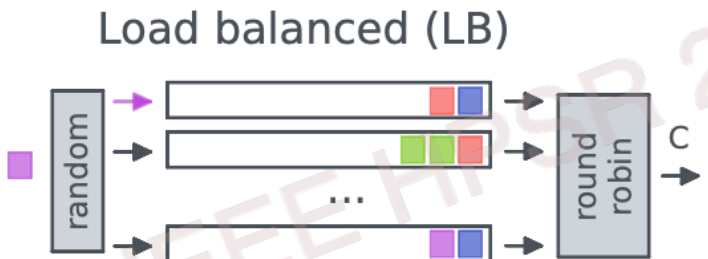
The output to s0 had parallel queues



The experiment interface emulated a load balancer

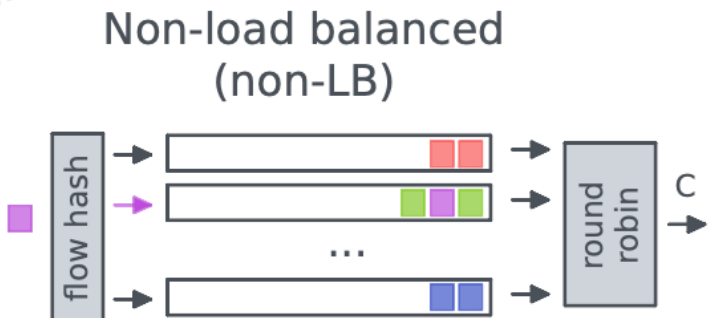
LB Configuration:

Probabilistic placement,
produces reordering



non-LB Configuration:

Hashed placement,
No reordering

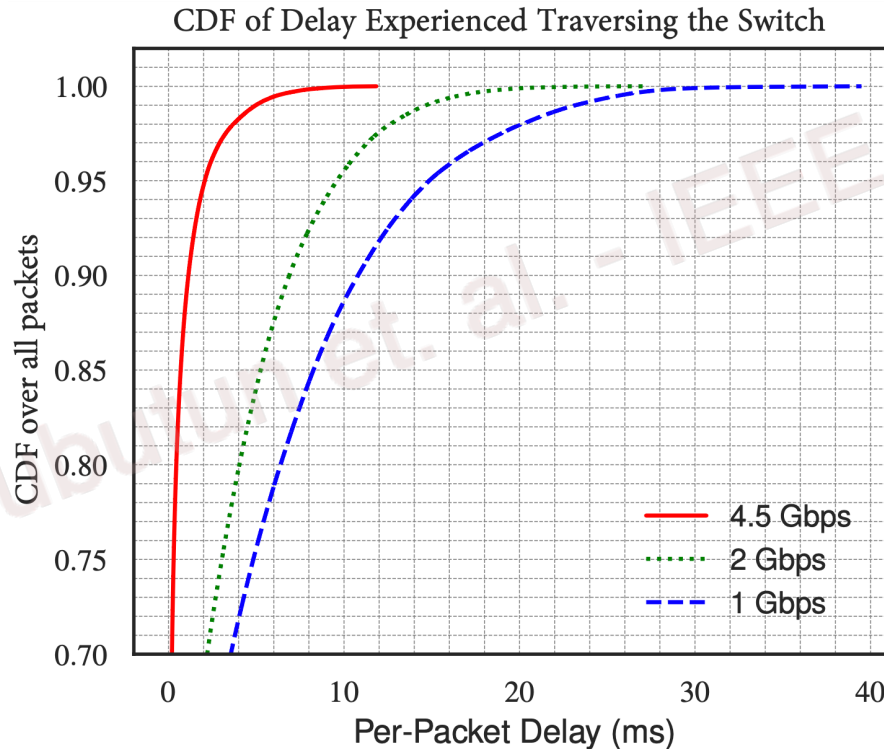


Both served in round-robin order.

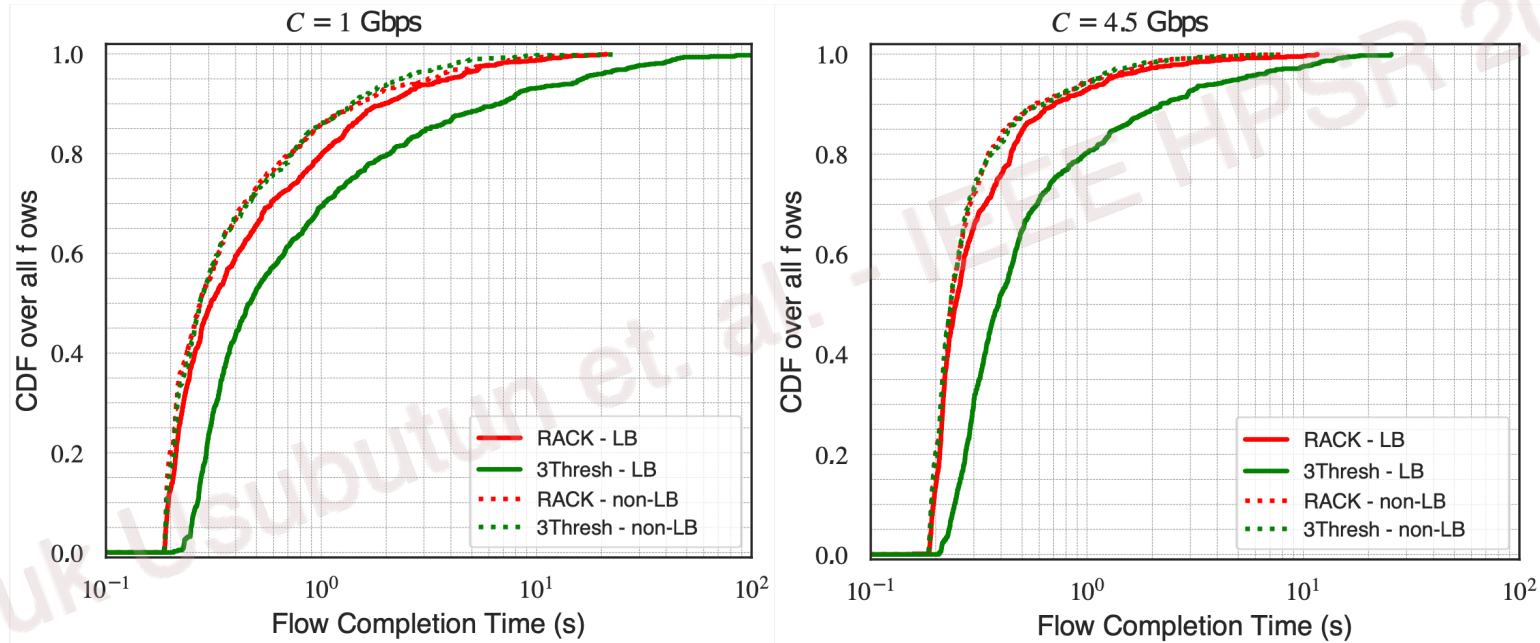
Results

Ufuk Usubutun et. al. - IEEE HPSR 2023

Higher line rates result in narrower delay distributions through the reordering switch

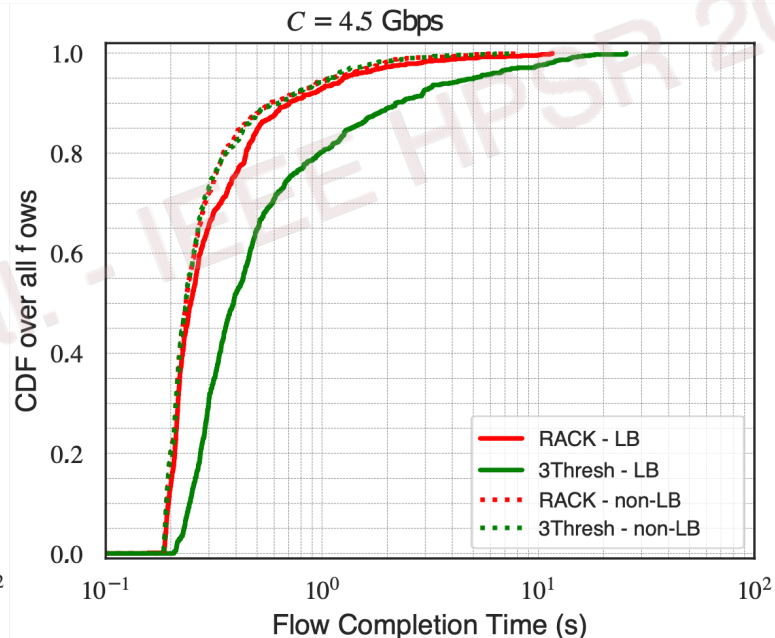
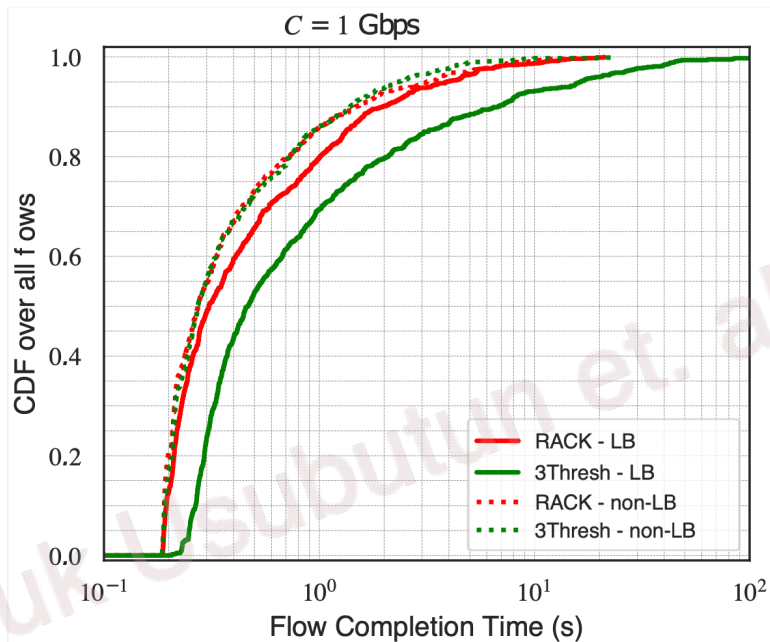


All algorithms perform the same without reordering (dashed lines)



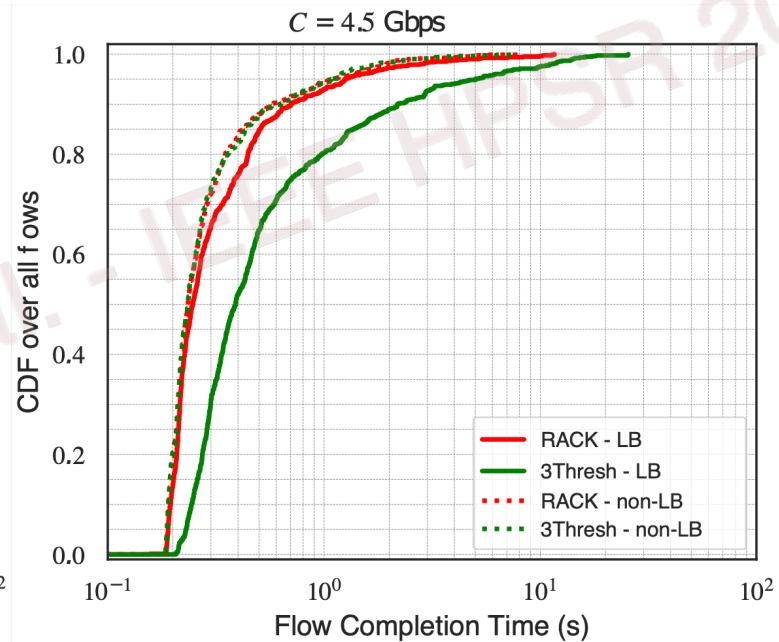
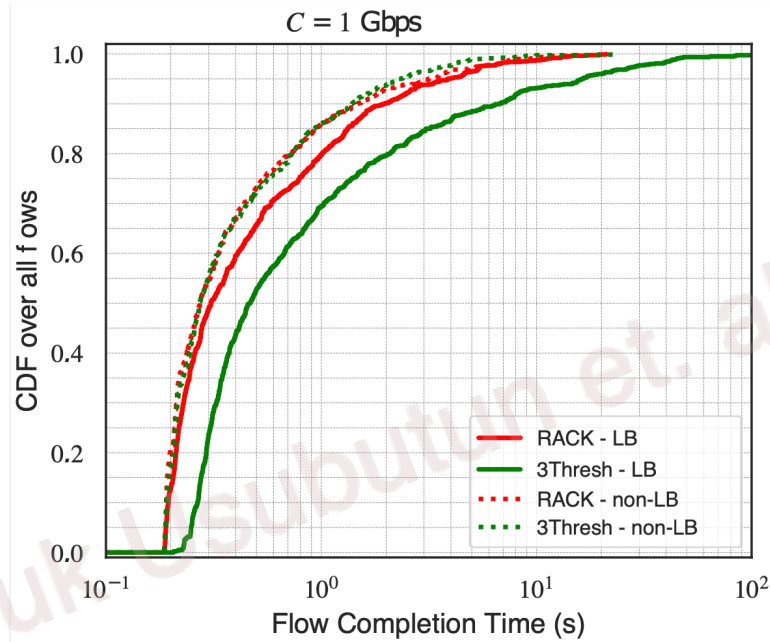
Approx $(51 \pm 4)\%$ utilization at each scenario

Triple Duplicate ACK does poorly under reordering



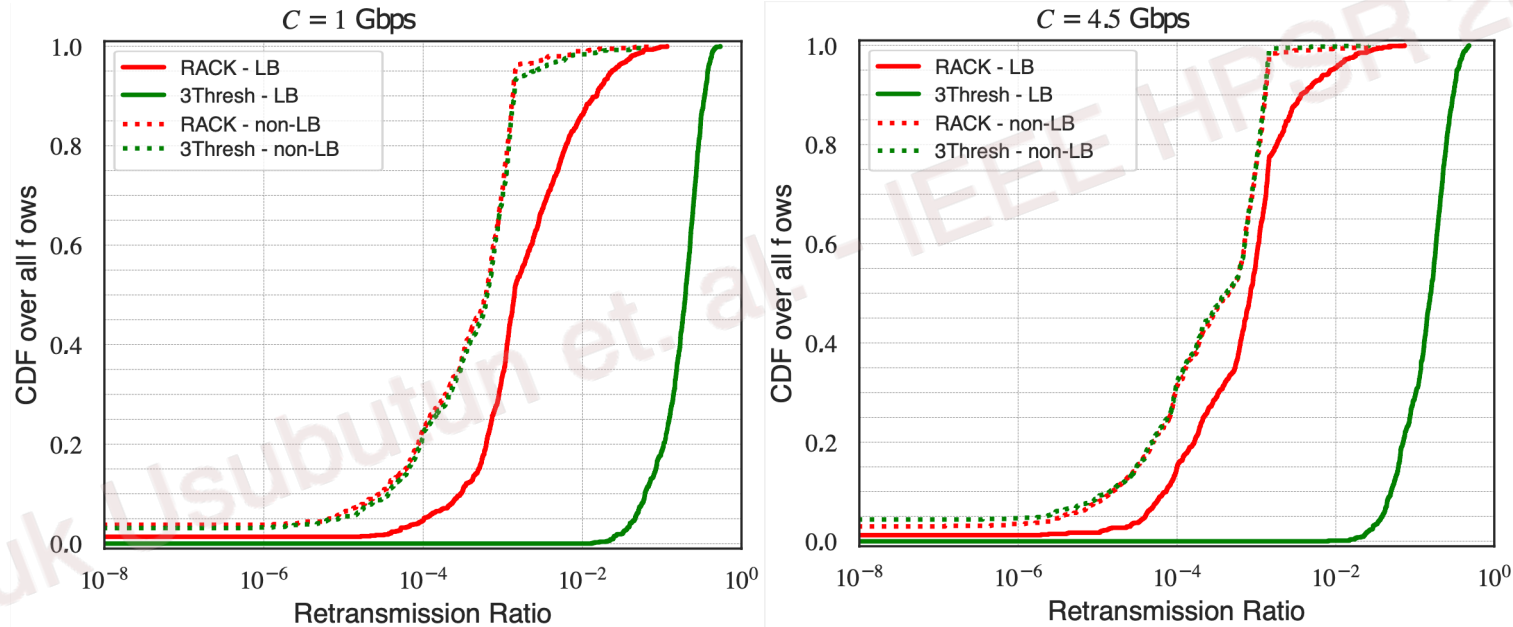
Approx (51 ± 4)% utilization at each scenario

Higher line rates lead to better performance for RACK (time based)



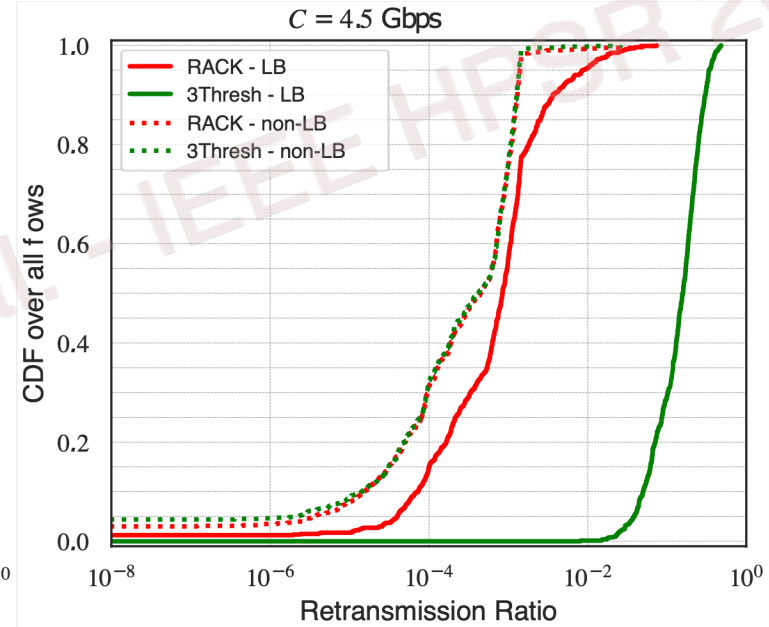
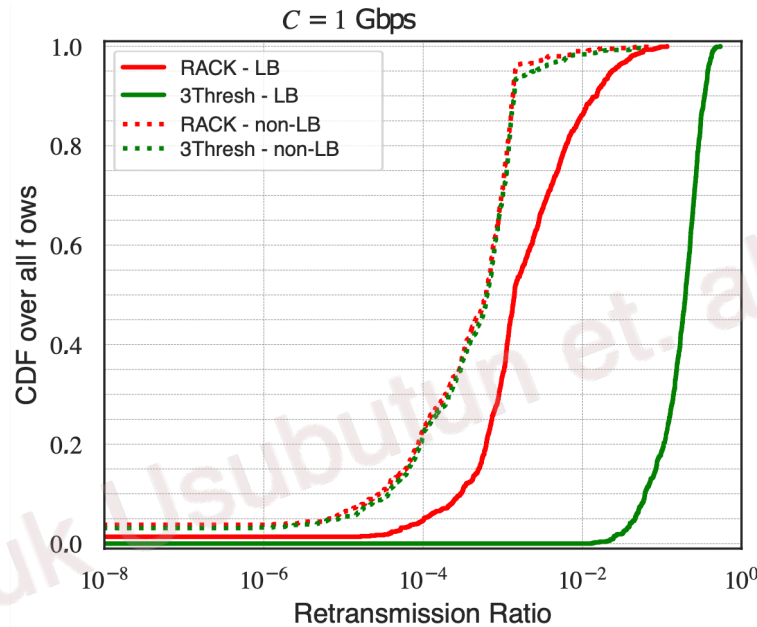
Approx $(51 \pm 4)\%$ utilization at each scenario

Triple Duplicate ACK retransmits significantly



Approx $(51 \pm 4)\%$ utilization at each scenario

High line rates results in less retransmissions for RACK (time based)



Approx $(51 \pm 4)\%$ utilization at each scenario

Conclusion

Ufuk Usubutun et. al. - IEEE HPSR 2023

Conclusions & Future Work

- Traditional wisdom on in-sequence delivery requirements for switches should be revisited
- This result also has implications for Data Center Networks and multi Radio Access Technology wireless systems.
- Similar implications for UDP based QUIC.
- The case of TCP BBR should be investigated

GitHub Repo for Artifacts →



THANK YOU! QUESTIONS?

Ufuk Usubütün
usubutun@nyu.edu

Ufuk Usubütün et. al. - IEEE TNSR 2023

Backup Slides



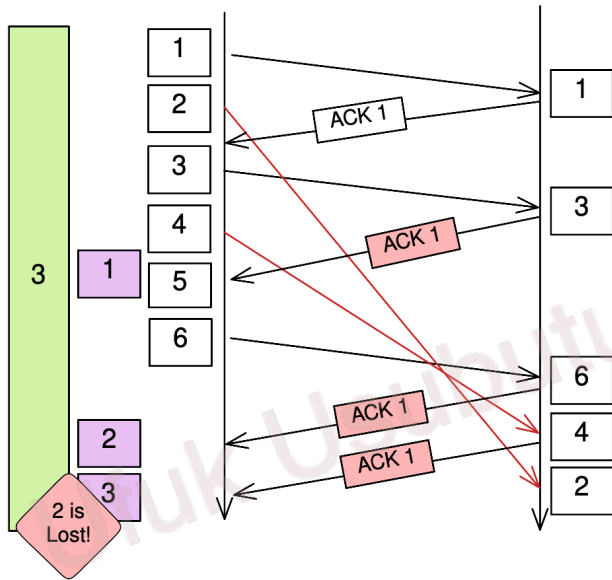
NYU

TANDON SCHOOL
OF ENGINEERING

Ufuk Usubutun et. al. - IEEE HPSR 2023

A look into classical loss detection

Triple Duplicate ACK



Evolved from 3 duplicate ACK

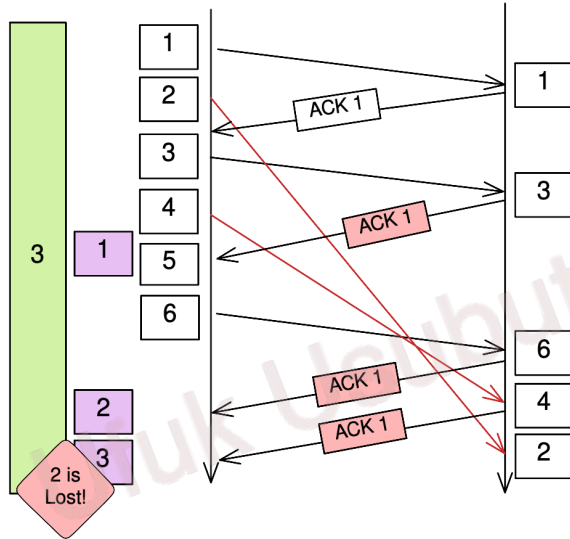
Selective ACK (SACK) and Duplicate SACK (DSACK) emerges.

- Expansion of the same ACK counting idea
adaptive dupthresh
- Time based approach
RACK

SACKs allow better knowledge of reception

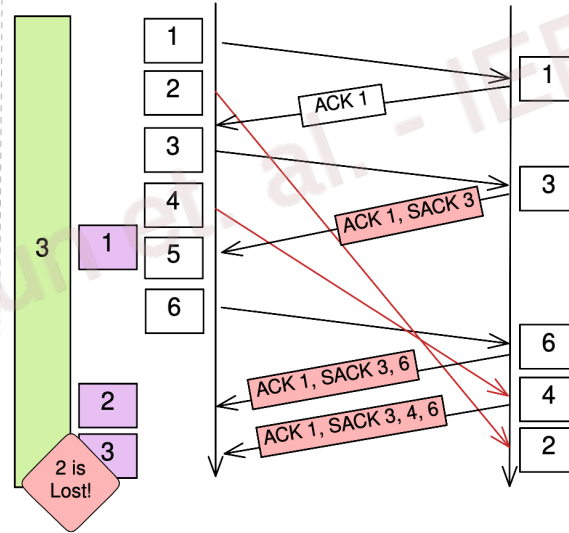
Triple Duplicate ACK

Triple Duplicate ACK



dupthresh (**3Threshold**)

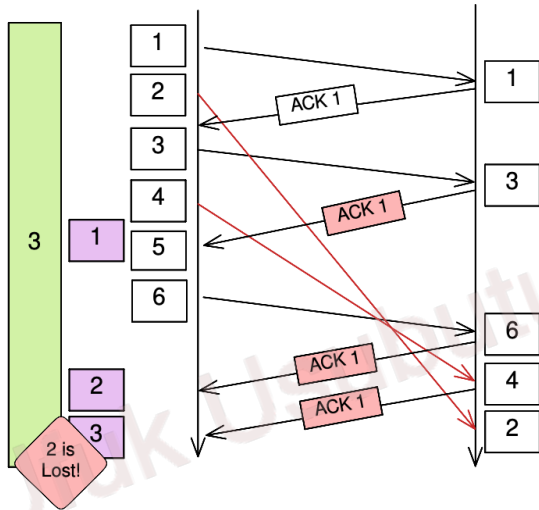
dupthresh with fixed threshold = 3



Adaptive mechanisms emerged

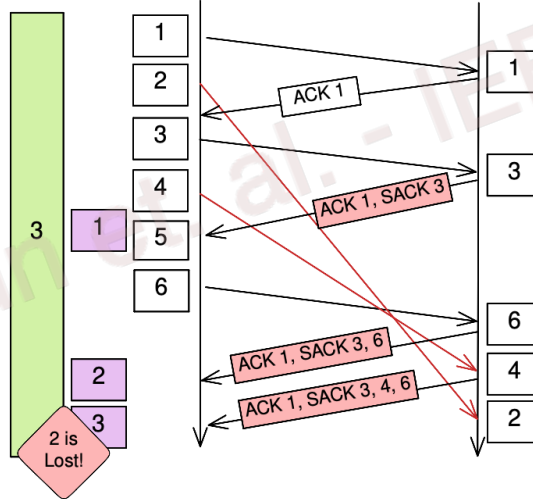
Triple Duplicate ACK

Triple Duplicate ACK



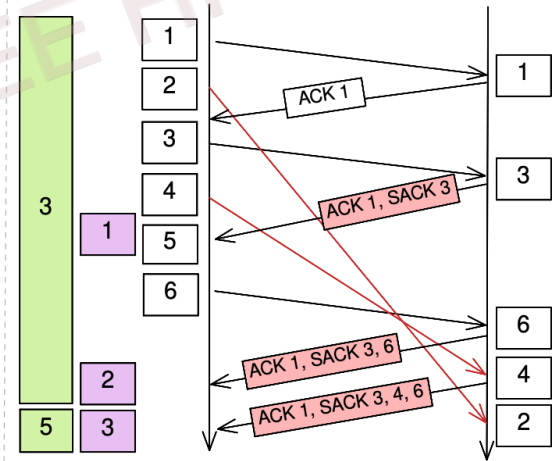
dupthresh (**3Thresh**)

dupthresh with fixed threshold = 3



adaptive dupthresh (**adapThresh**)

dupthresh with adaptive threshold



adapThresh gains resilience to reordering after an adaptation episode.

How to generate flows: Set up F flow generators at each node

At each generator:

Sample a flow size from a WAN TCP Traffic study (90% mice, 10% elephant)

At each generator:

Random wait between start of flows,

Always greater than flow completion time. Wait time scaled wrt capacity C .



Aim: Keep the total bytes transferred independent of algorithm.

Experimenting within a closed loop

We tested **different loss detection algorithms** at different line rates.

This results in a closed loop

We achieved Approx. $(51 \pm 4)\%$ utilization at each scenario

