

X\_405082

# Advanced Computer Networks

## Data Center Transport

Lin Wang (lin.wang@vu.nl)

Period 2, Fall 2020

# Course outline

## Warm-up

- Fundamentals
- Forwarding and routing
- Network transport

## Data centers

- Data center networking
- **Data center transport** 🖱️

## Programmability

- Software defined networking
- Programmable forwarding

## Video

- Video streaming
- Video stream analytics

## Networking and ML

- Networking for ML
- ML for networking

## Mobile computing

- Wireless and mobile

# Learning objectives

What are the **new challenges** in data center transport?

What **design choices** do we have for data centers transport design?

# What is special about data center transport

## Network

- Extremely high speed (100+ Gbps)
- Extremely low latency (10-100s of us)

## Diverse applications and workloads

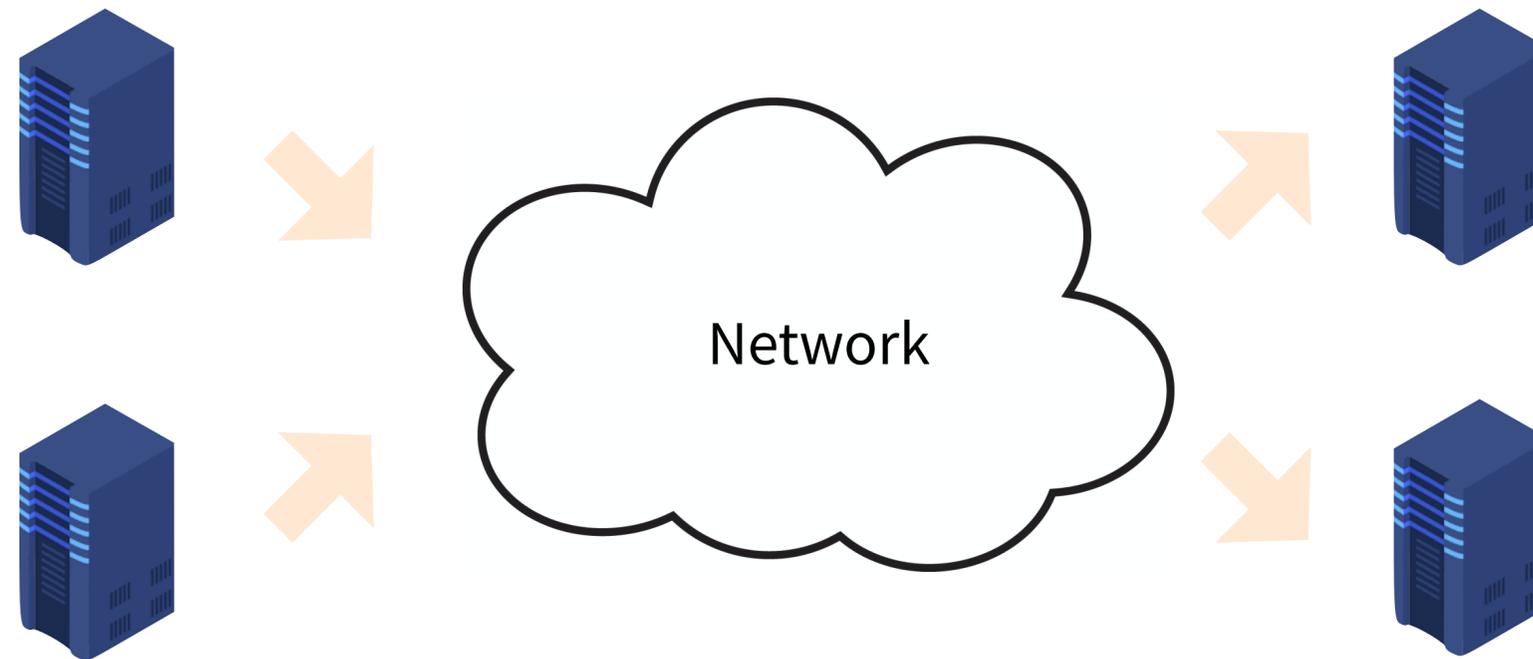
- Large variety in performance requirements

## Traffic patterns

- Large long-lived flows vs small short-lived flows
- Scatter-gather, broadcast, multicast

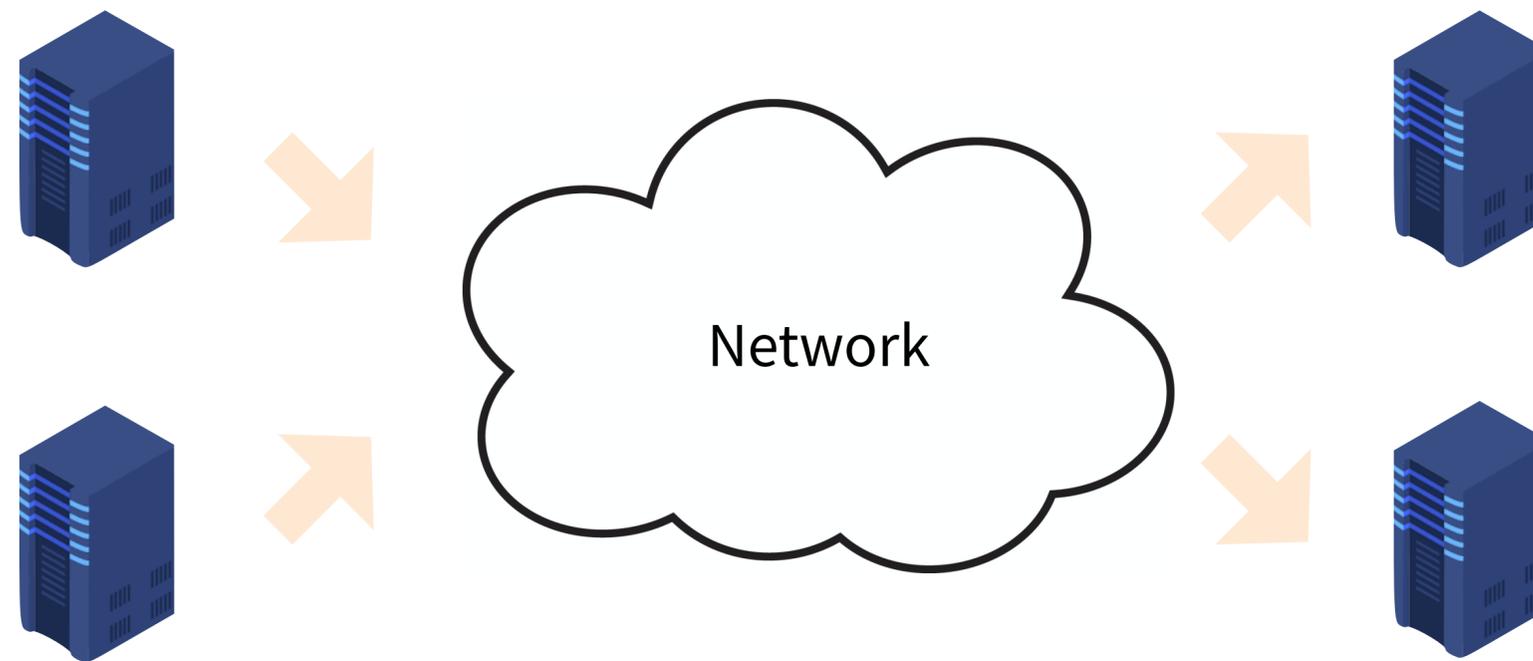
Built out of commodity components: no expensive/customized hardware

# Congestion control recall



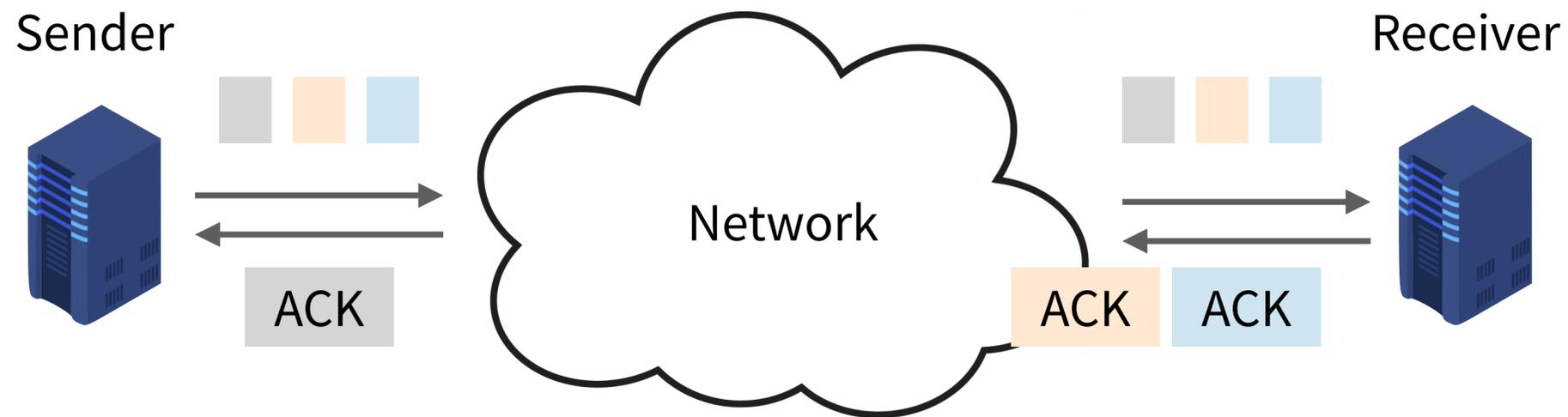
Do you still remember the goal of congestion control?

# Congestion control recall



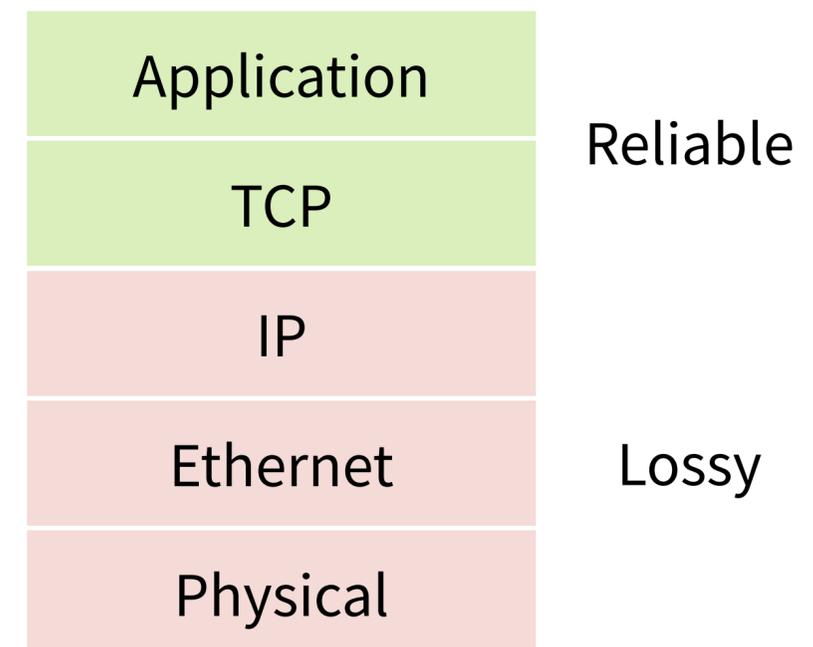
Congestion control aims to determine the **rate to send data** on a connection, such that (1) the sender does not overrun the network capability and (2) the network is efficiently utilized

# TCP



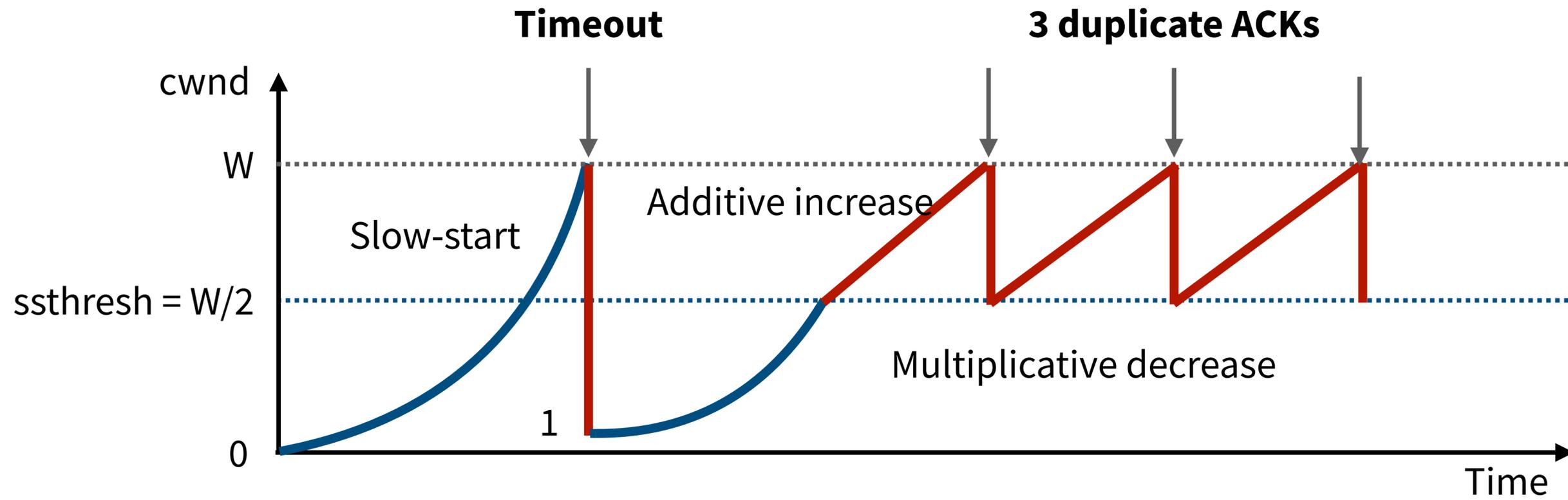
The **transport layer** in the network model:

- Reliable, in-order delivery using acknowledges
- Make sure not to overrun the receiver (receiving window, rwnd) and the network (congestion window, cwnd)
- What can be sent =  $\min(\text{rwnd}, \text{cwnd})$



# TCP AIMD

TCP Reno



What could possibly go wrong in a data center environment?

# TCP incast problem

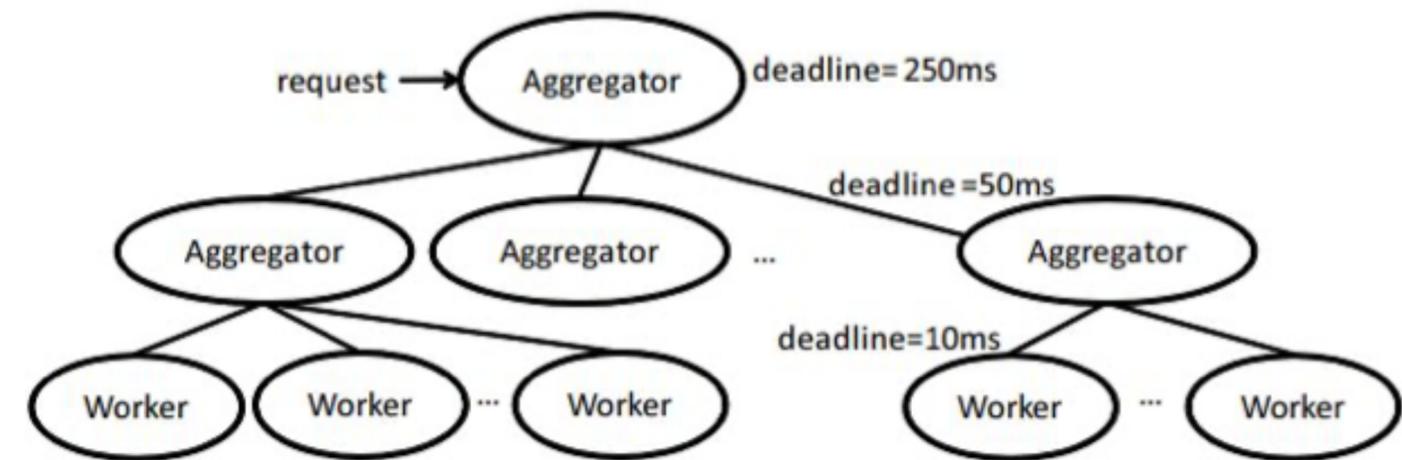
A data center application run on multiple servers

- Storage, cache, data processing (MapReduce)

They use a scatter-gather (or partition-aggregate) work pattern

- **[scatter]** A client sends a request to a bunch of servers for data
- **[gather]** All servers respond to the client

More broadly, a client-facing query might have to collect data from many servers

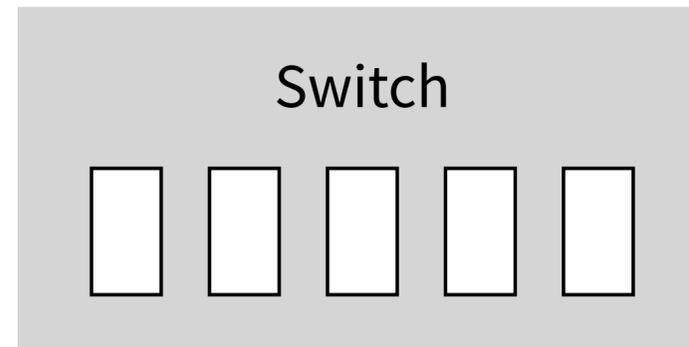


(DCTCP, ACM SIGCOMM 2010)

# From a switch point of view

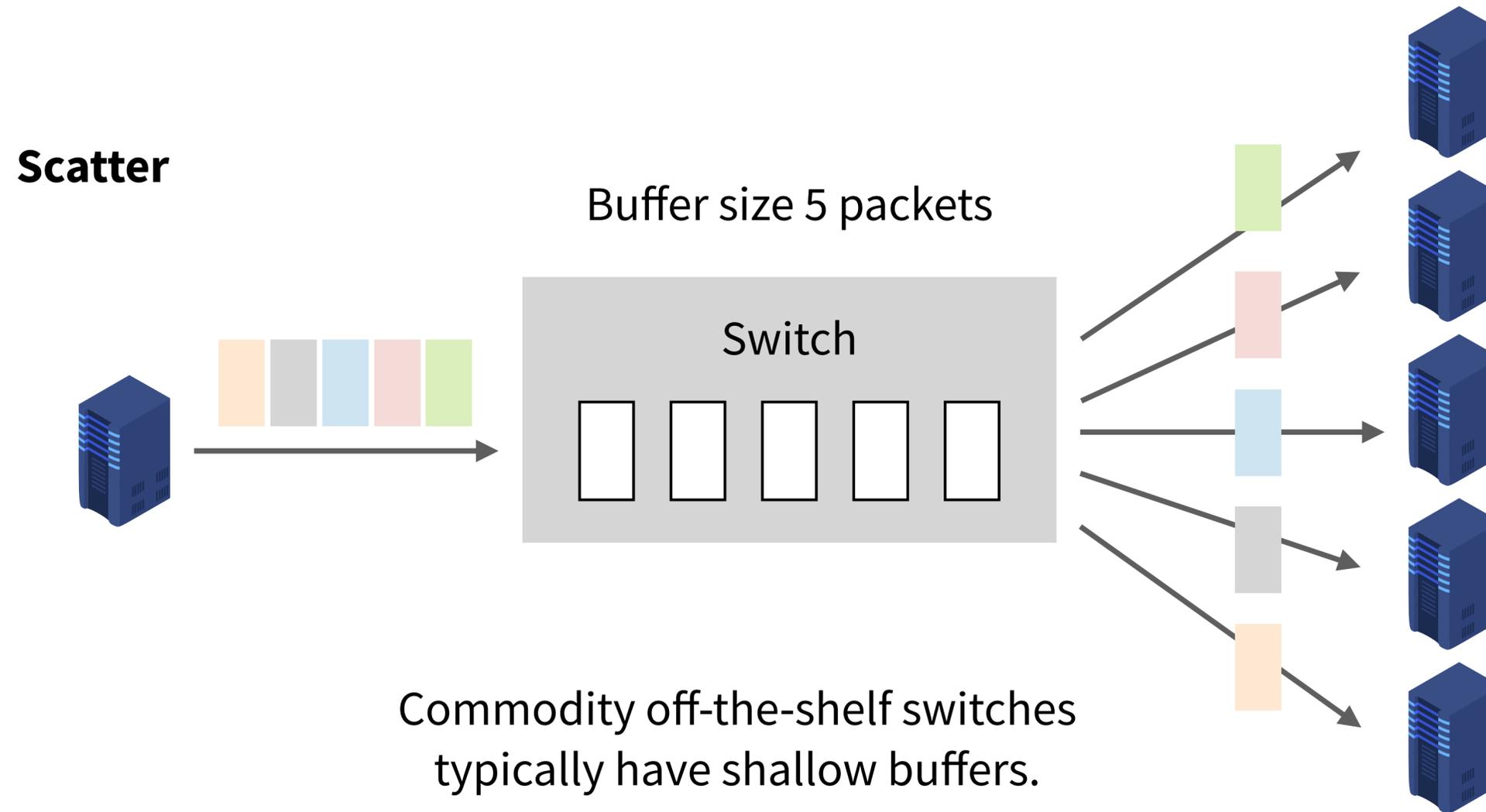


Buffer size 5 packets

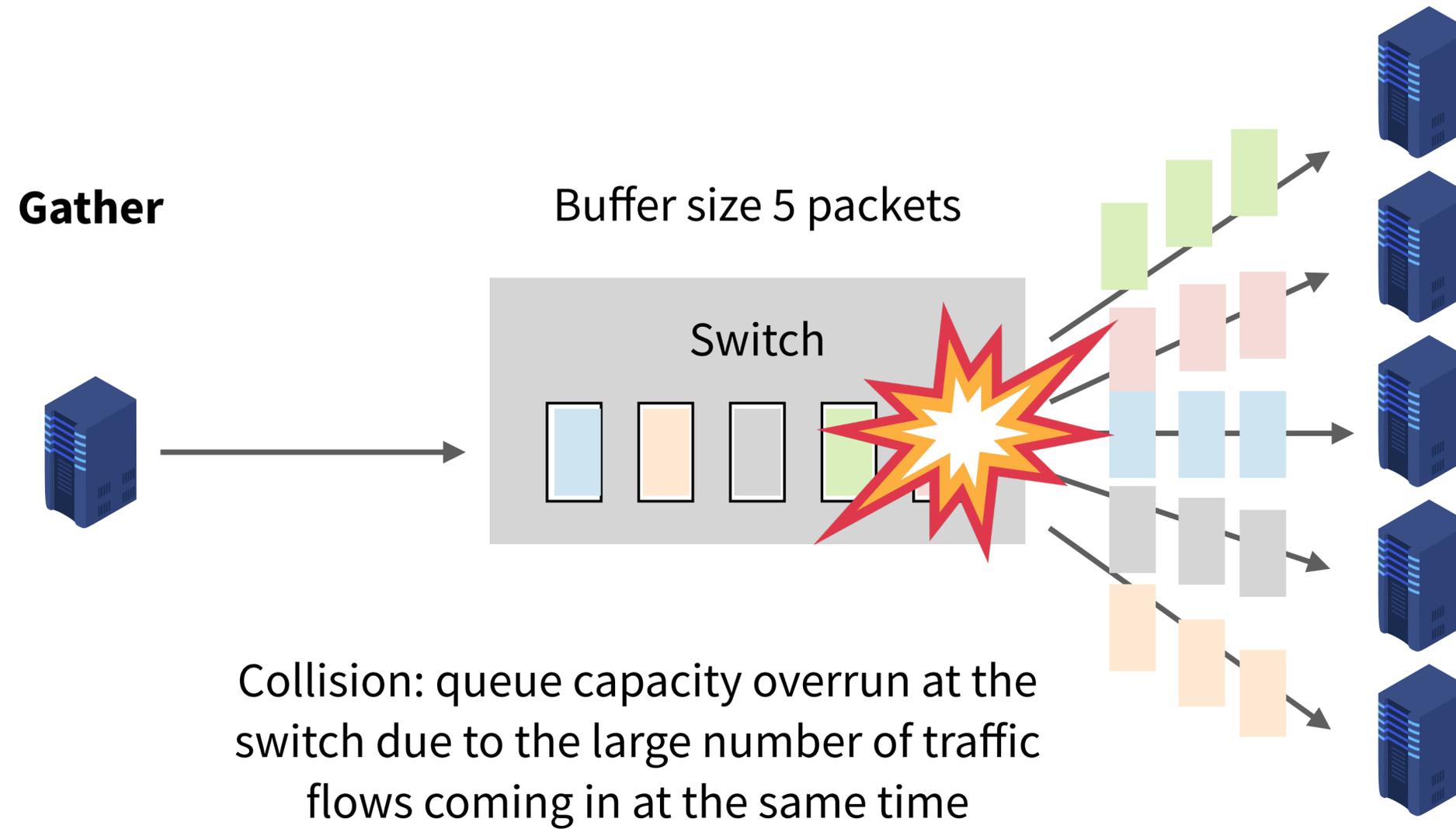


Commodity off-the-shelf switches typically have shallow buffers.

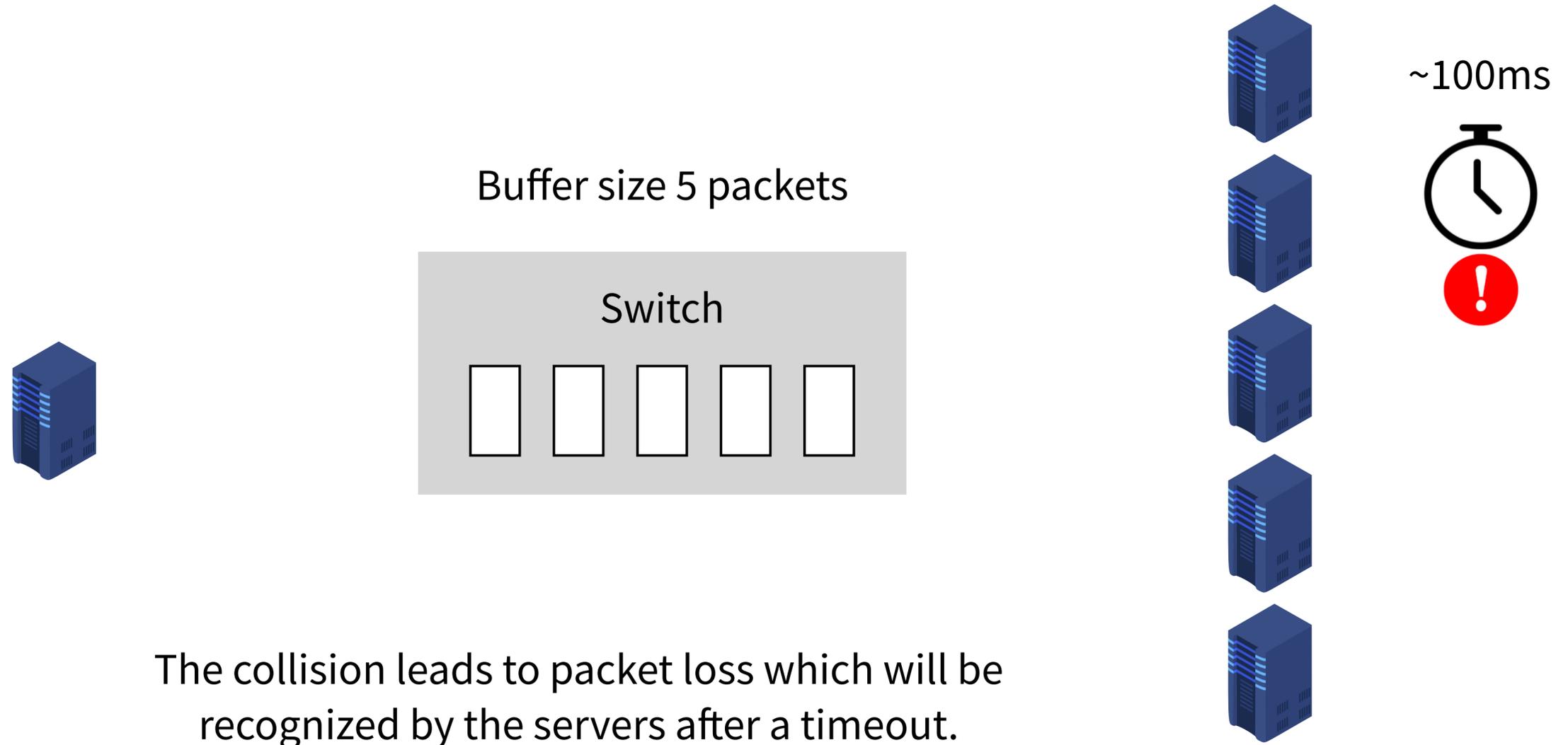
# From a switch point of view



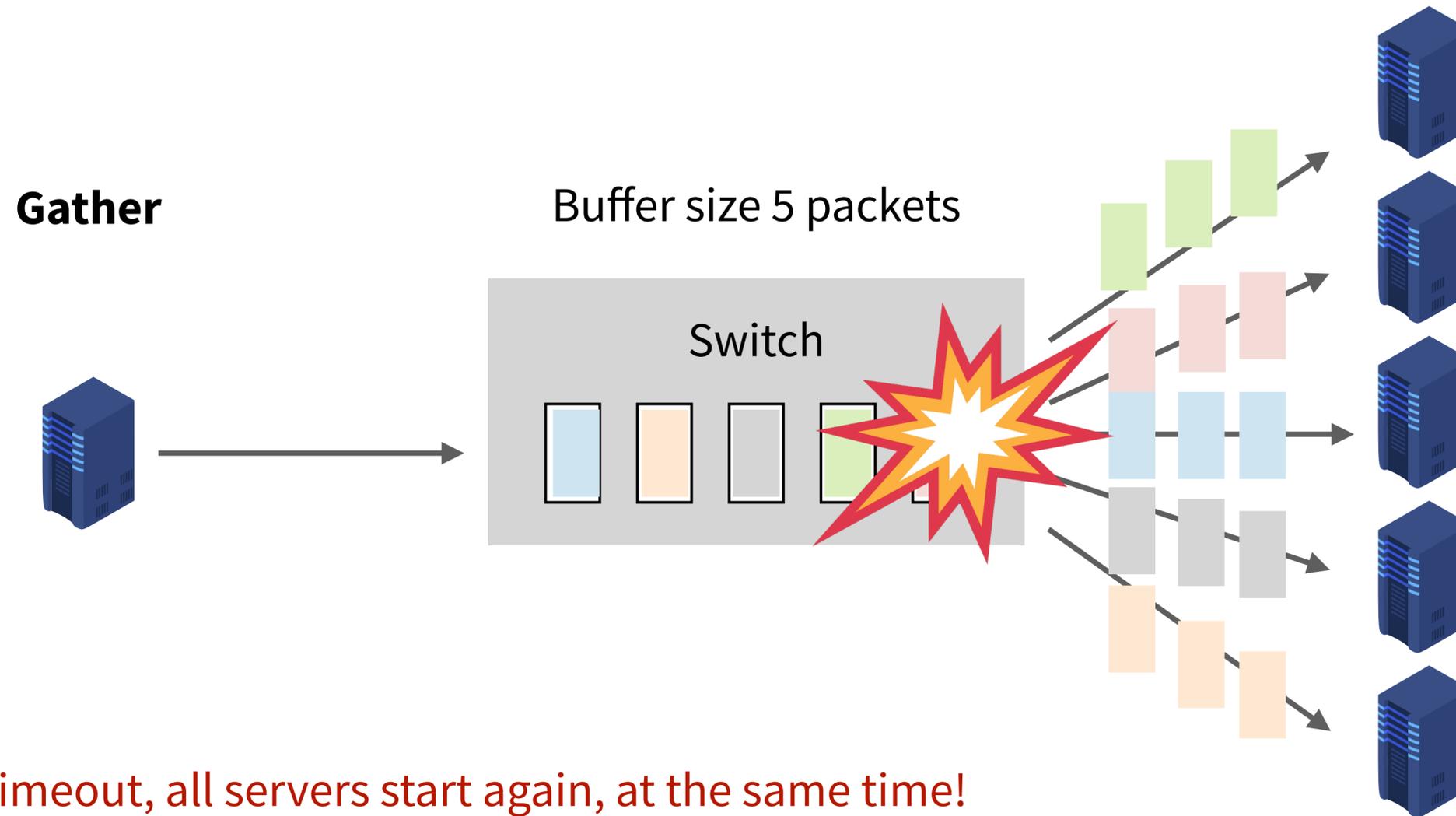
# From a switch point of view



# From a switch point of view

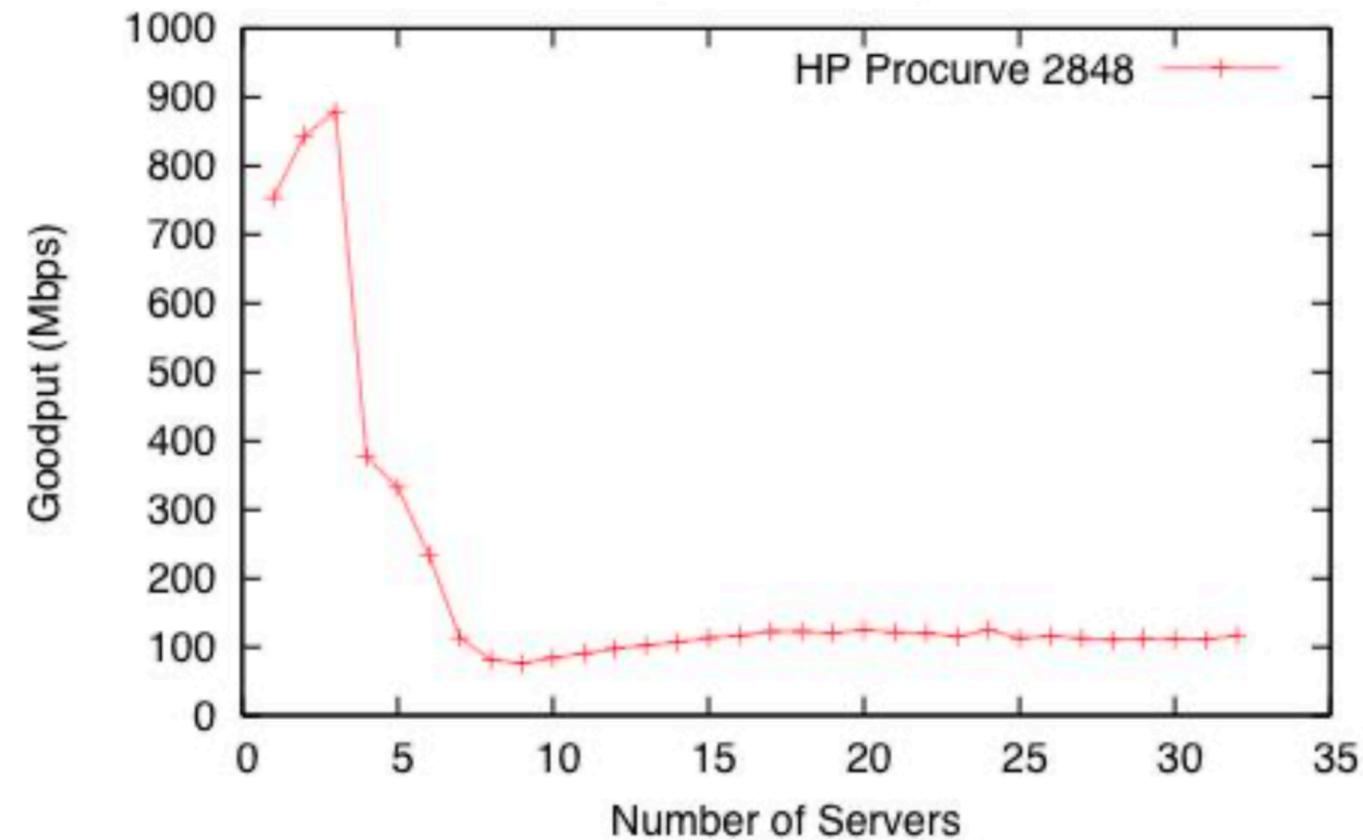


# From a switch point of view



After the timeout, all servers start again, at the same time!  
**TCP global synchronization** problem!

# TCP incast problem



## Measurement and Analysis of TCP Throughput Collapse in Cluster-based Storage Systems

Amar Phanishayee, Elie Krevat, Vijay Vasudevan,  
David G. Andersen, Gregory R. Ganger, Garth A. Gibson, Srinivasan Seshan

*Carnegie Mellon University*

### Abstract

Cluster-based and iSCSI-based storage systems rely on standard TCP/IP-over-Ethernet for client access to data. Unfortunately, when data is striped over multiple networked storage nodes, a client can experience a TCP throughput collapse that results in much lower read bandwidth than should be provided by the available network links. Conceptually, this problem arises because the client simultaneously reads fragments of a data block from multiple sources that together send enough data to overload the switch buffers on the client's link. This paper analyzes this *Incast* problem, explores its sensitivity to various system parameters, and examines the effectiveness of alterna-

client increases past the ability of an Ethernet switch to buffer packets. As we explore further in §2, the problem arises from a subtle interaction between limited Ethernet switch buffer sizes, the communication patterns common in cluster-based storage systems, and TCP's loss recovery mechanisms. Briefly put, data striping couples the behavior of multiple storage servers, so the system is limited by the request completion time of the *slowest* storage node [7]. Small Ethernet buffers are exhausted by a concurrent flood of traffic from many servers, which results in packet loss and one or more TCP timeouts. These timeouts impose a delay of hundreds of milliseconds—orders of magnitude greater than typical data fetch times—significantly degrading overall throughput.

USENIX FAST 2008

# TCP incast

## Packet drops due to the capacity overrun at shared commodity switches

- Can lead to TCP global synchronization and even more packet losses
- The link remains idle (hence, reduced capacity and poor performance)
- First discussed in Nagle et al., The Panasas ActiveScale Storage Cluster, SC 2004

## Some potential solutions

- Use lower timeouts: (1) can lead to spurious timeouts and retransmissions, (2) High operating system overhead
- Other variants of TCP (SACKS, Cubic): cannot avoid the basic phenomenon of TCP incast
- Larger switch buffer: helps to push the collapse point further, but is expensive and introduces higher packet delay

# Can we do better?

The basic challenge is that there are **only** limited number of things we can do once a packet is dropped

- Various acknowledgements schemes
- Various timeouts based optimizations

Whatever clever way you can come up with

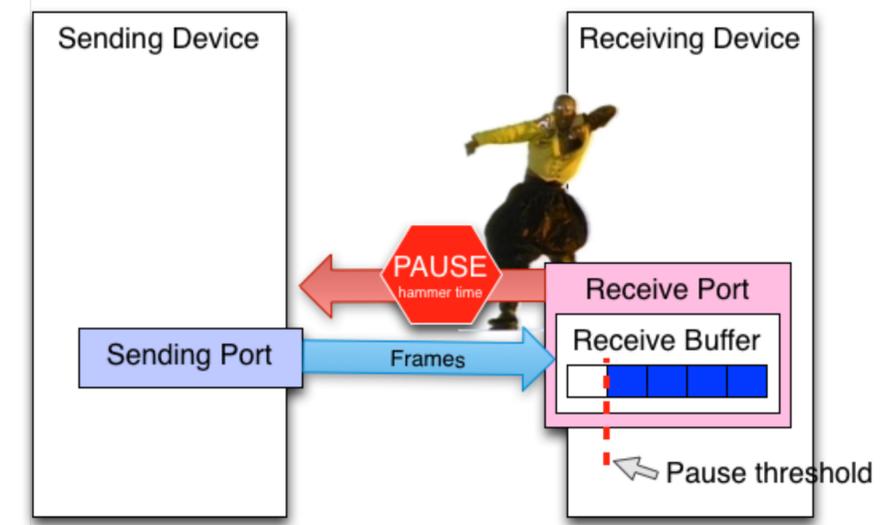
- Imagine deploying that with multiple workloads, flow patterns, and switches

Can we try to avoid packet drops in the first place? If so, how?

# Ethernet flow control

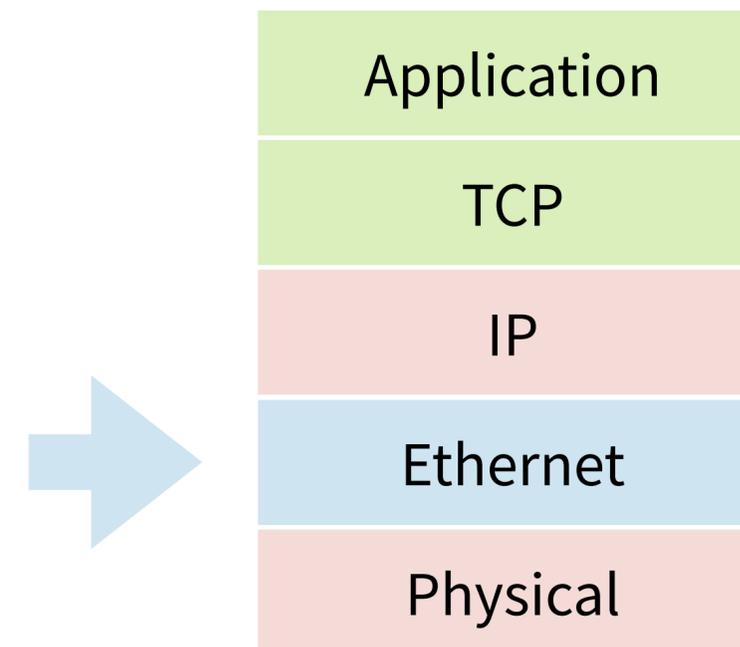
## Pause Frame (IEEE 802.3x)

- An overwhelmed Ethernet receiver/NIC can send a "PAUSE" Ethernet frame to the sender
- Upon receiving the PAUSE frame, the sender stops transmission for a certain duration of time



## Limitations

- Designed for end-host NIC (memory, queue) overruns, not switches
- Blocks all transmission at the Ethernet-level (port-level, not flow-level)



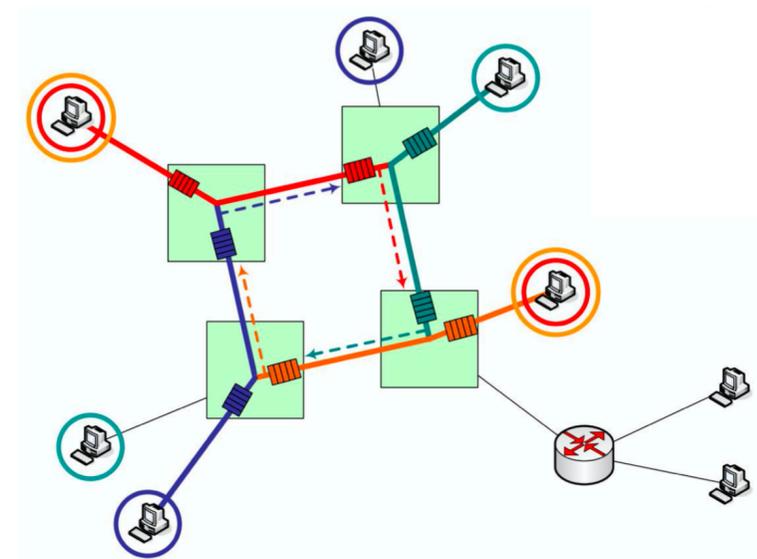
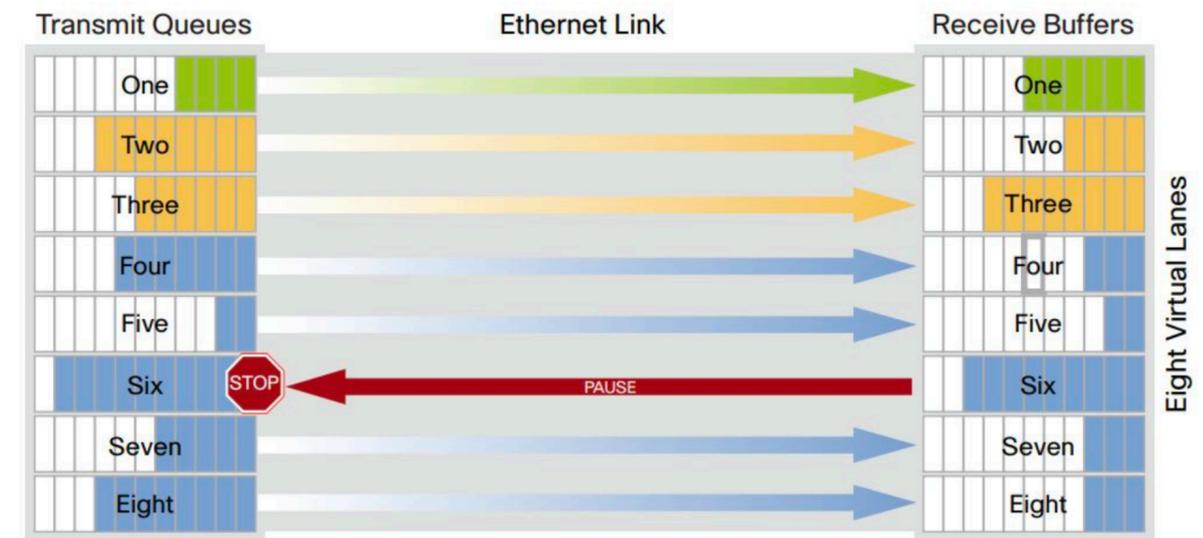
# Priority-based flow control

## PFC, IEEE 802.1Qbb

- Enhancement over PAUSE frames
- 8 virtual traffic lanes and one can be selectively stopped
- Timeout is configurable

## Limitations

- Only 8 lanes: think about the number of flows we may have
- Deadlocks in large networks
- Unfairness (victim flows)



# Data Center TCP (DCTCP)

TCP-alike congestion control protocol

Basic idea: pass information about switch queue buildup to senders

- From **where** to pass information?
- **How** to pass information?

At the sender, react to this information by slowing down the transmission

- By **how much**?
- **How frequent**?

## Data Center TCP (DCTCP)

Mohammad Alizadeh<sup>††</sup>, Albert Greenberg<sup>†</sup>, David A. Maltz<sup>†</sup>, Jitendra Padhye<sup>†</sup>,  
Parveen Patel<sup>†</sup>, Balaji Prabhakar<sup>†</sup>, Sudipta Sengupta<sup>†</sup>, Murari Sridharan<sup>†</sup>

<sup>†</sup>Microsoft Research   <sup>††</sup>Stanford University  
{albert, dmaltz, padhye, parveen, sudipta, muraris}@microsoft.com  
{alizade, balaji}@stanford.edu

### ABSTRACT

Cloud data centers host diverse applications, mixing workloads that require small predictable latency with others requiring large sustained throughput. In this environment, today's state-of-the-art TCP protocol falls short. We present measurements of a 6000 server production cluster and reveal impairments that lead to high application latencies, rooted in TCP's demands on the limited buffer space available in data center switches. For example, bandwidth hungry "background" flows build up queues at the switches, and thus impact the performance of latency sensitive "foreground" traffic.

To address these problems, we propose DCTCP, a TCP-like protocol for data center networks. DCTCP leverages Explicit Congestion Notification (ECN) in the network to provide multi-bit feedback to the end hosts. We evaluate DCTCP at 1 and 10Gbps speeds

eral recent research proposals envision creating economical, easy-to-manage data centers using novel architectures built atop these commodity switches [2, 12, 15].

Is this vision realistic? The answer depends in large part on how well the commodity switches handle the traffic of real data center applications. In this paper, we focus on soft real-time applications, supporting web search, retail, advertising, and recommendation systems that have driven much data center construction. These applications generate a diverse mix of short and long flows, and require three things from the data center network: low latency for short flows, high burst tolerance, and high utilization for long flows.

The first two requirements stem from the *Partition/Aggregate* (described in §2.1) workflow pattern that many of these applications use. The near real-time deadlines for end results translate into latency targets for the individual tasks in the workflow. These tar-

ACM SIGCOMM 2010

# Explicit Congestion Notification (ECN)

ECN is a standardized way of passing "the presence of congestion"

- Part of the IP packet header (2 bits): uses 1 bit for capability/ACK and 1 bit for congestion indication (yes/no)
- Supported by most commodity switches

Idea: For a queue size of  $N$ , when the queue occupancy goes beyond  $K$ , mark the passing packet's ECN bit as "yes"

- There are more sophisticated logics (Randomly Early Detection, RED) that can probabilistically mark packets

Updated by: [4301](#), [6040](#), [8311](#)      PROPOSED STANDARD  
 Network Working Group      [Errata Exist](#)  
 Request for Comments: 3168      K. Ramakrishnan  
 Updates: [2474](#), [2401](#), [793](#)      TeraOptic Networks  
 Obsoletes: [2481](#)      S. Floyd  
 Category: Standards Track      ACIRI  
    D. Black  
    EMC  
    September 2001

**The Addition of Explicit Congestion Notification (ECN) to IP**

Status of this Memo

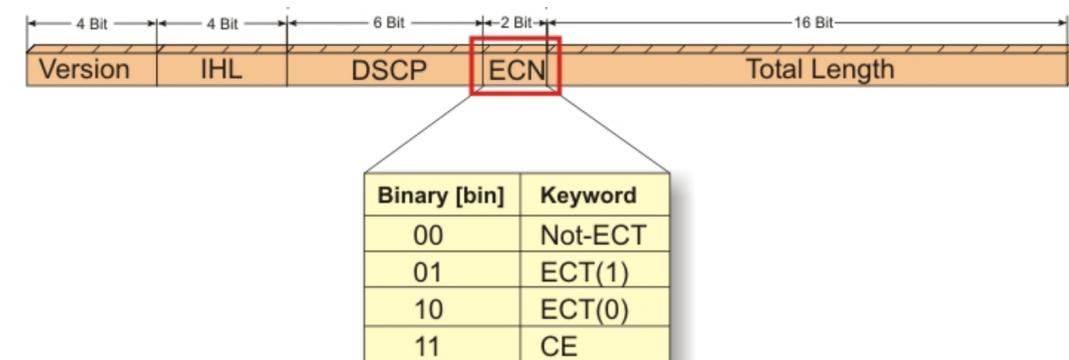
This document specifies an Internet standards track protocol for the Internet community, and requests discussion and suggestions for improvements. Please refer to the current edition of the "Internet Official Protocol Standards" (STD 1) for the standardization state and status of this protocol. Distribution of this memo is unlimited.

Copyright Notice

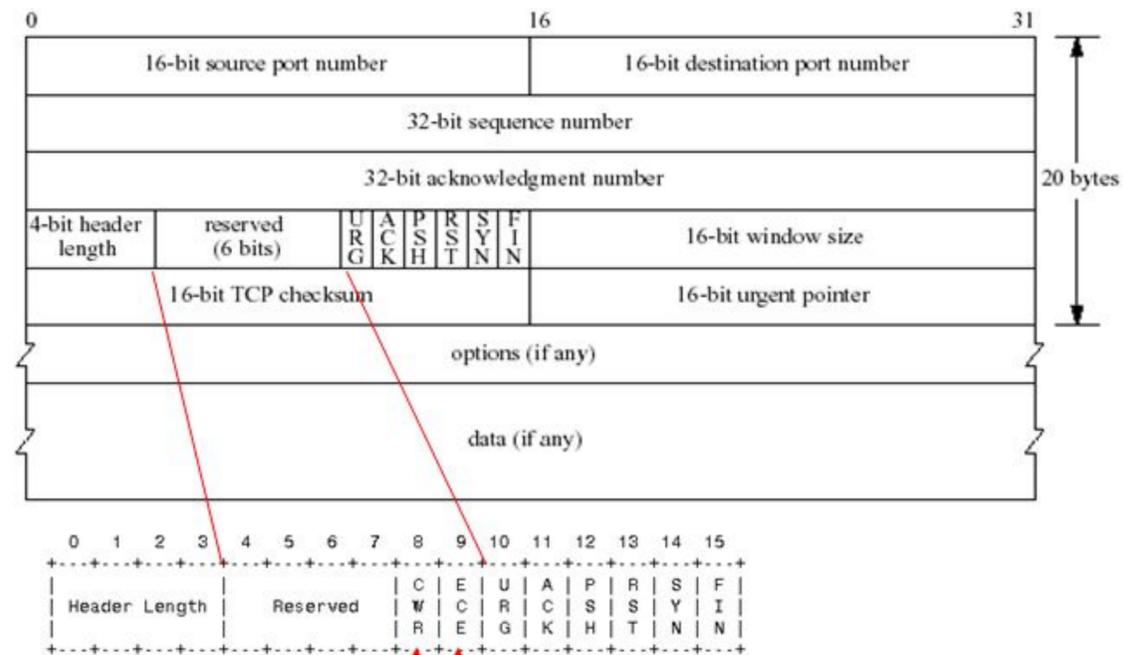
Copyright (C) The Internet Society (2001). All Rights Reserved.

Abstract

This memo specifies the incorporation of ECN (Explicit Congestion Notification) to TCP and IP, including ECN's use of two bits in the IP header.



# The ECN bits location in TCP header

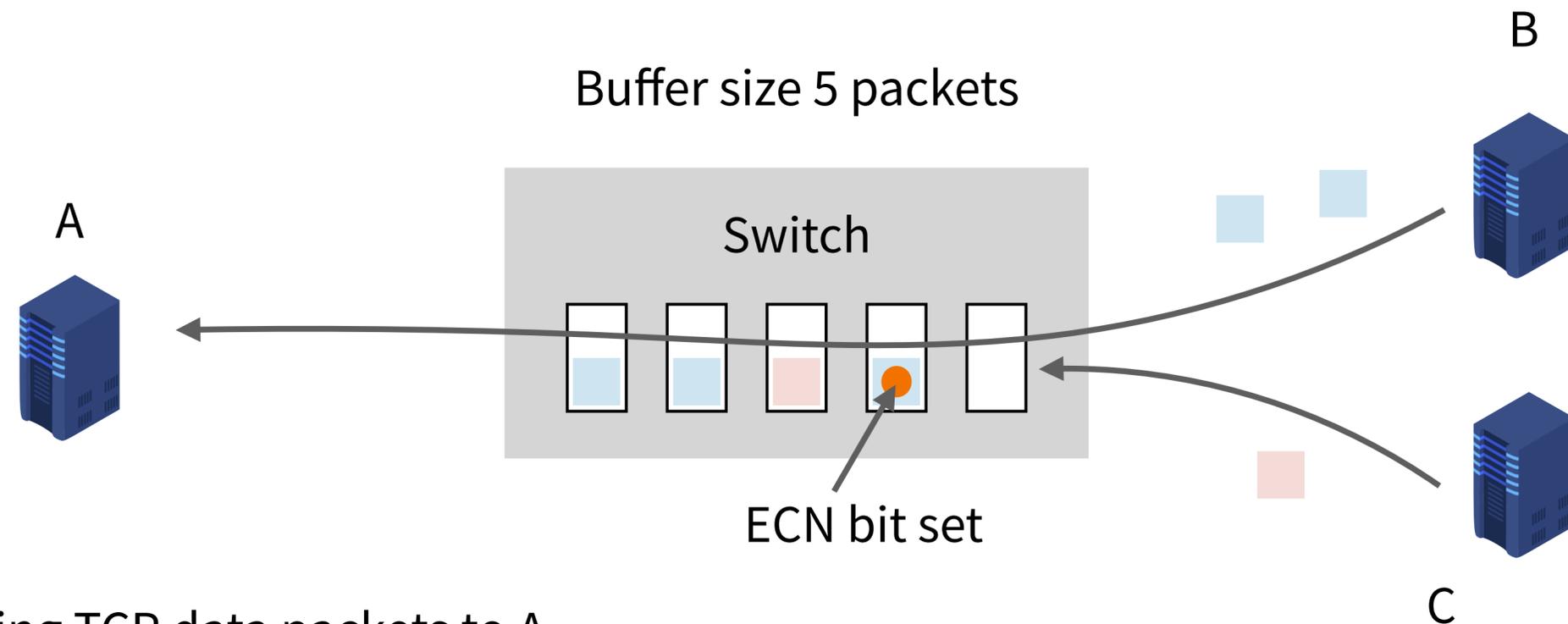


ECE flag - ECN-Echo flag  
 CWR flag - Congestion Window Reduced flag

The TCP congestion window logic:

- Additive increase:  $W \rightarrow W + 1$  per RTT
- Multiplicative decrease:  $W \rightarrow W / 2$ 
  - Packet loss
  - A packet received with ECN marked

# ECN bit in action



Assume that B is sending TCP data packets to A

At some point of time, C also starts to send packets, and the queue is getting full

The switch starts to mark packets with ECN bits

How does B get to know there was a congestion at the switch?

# DCTCP main idea

## Simple marking at the switch

- After threshold  $K$  start marking packets with ECN (instantaneous vs average marking)
- Uses instantaneous marking for fast notification

## Typical ECN receiver

- Mark ACKs with ECN, until the sender ACKs back (using CWR flag bit)

## DCTCP receiver

- Only mark ACKs corresponding to the ECN packet

## Sender's congestion control

- Estimate the packets that are marked with ECN in a running window

# DCTCP congestion window calculations

$$\text{each RTT : } F = \frac{\text{\#marked ECN ACKs}}{\text{\#Total Acks}}$$

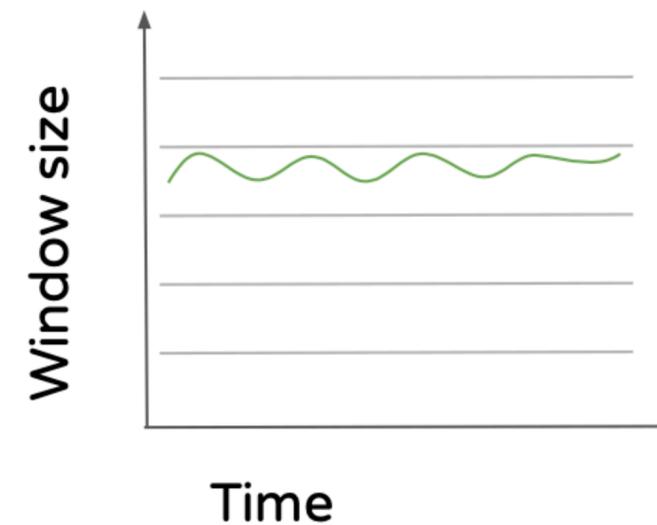
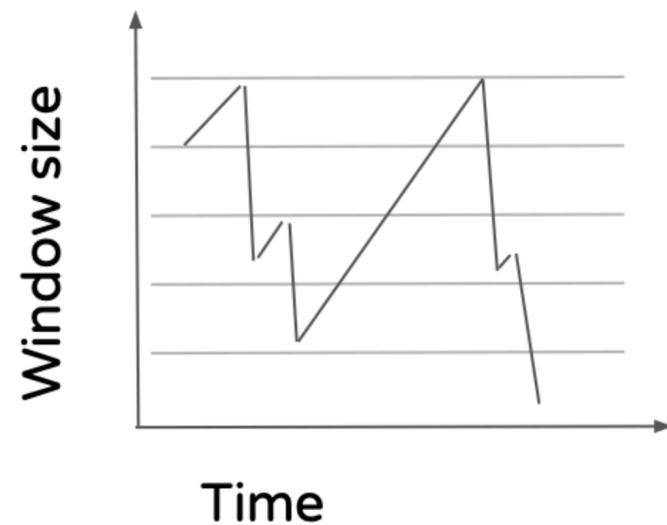
$$\alpha \leftarrow (1 - g) \alpha + g \times F \text{ (running estimate)}$$

$$\text{Congestion window (cnwd)} \leftarrow \text{cnwd} \times (1 - \alpha / 2)$$

Weightage for past vs present measurements  
 $0 < g < 1$

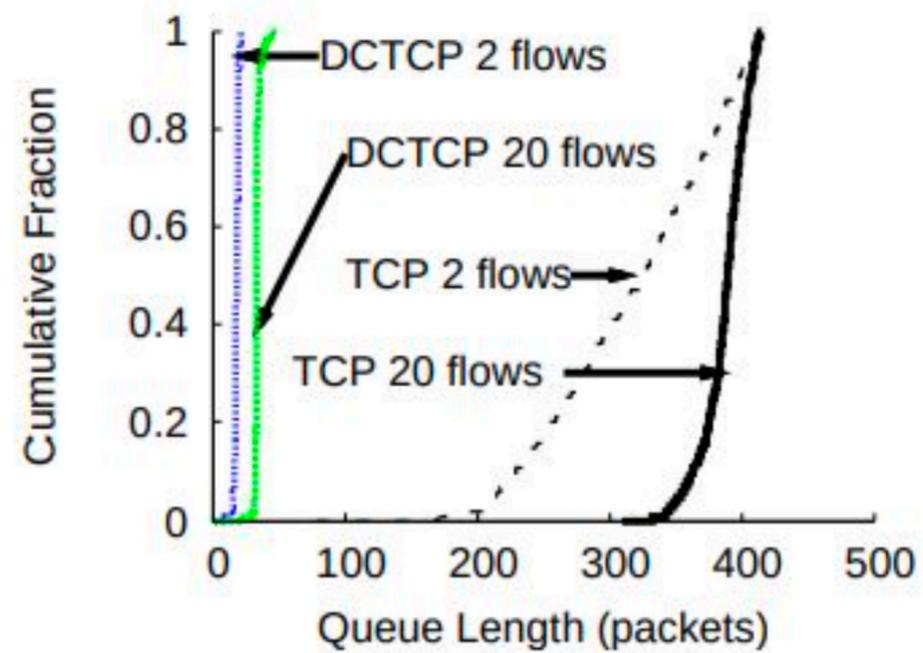
# DCTCP vs TCP example

ECN Marks on ACKs	TCP	DCTCP
0 1 1 0 0 0 1 0 0 1	Cut window by <b>50% (every time)</b>	Cut window by <b>40%</b>
0 0 0 0 0 0 0 0 0 1	Cut window by <b>50%</b>	Cut window by <b>10%</b>

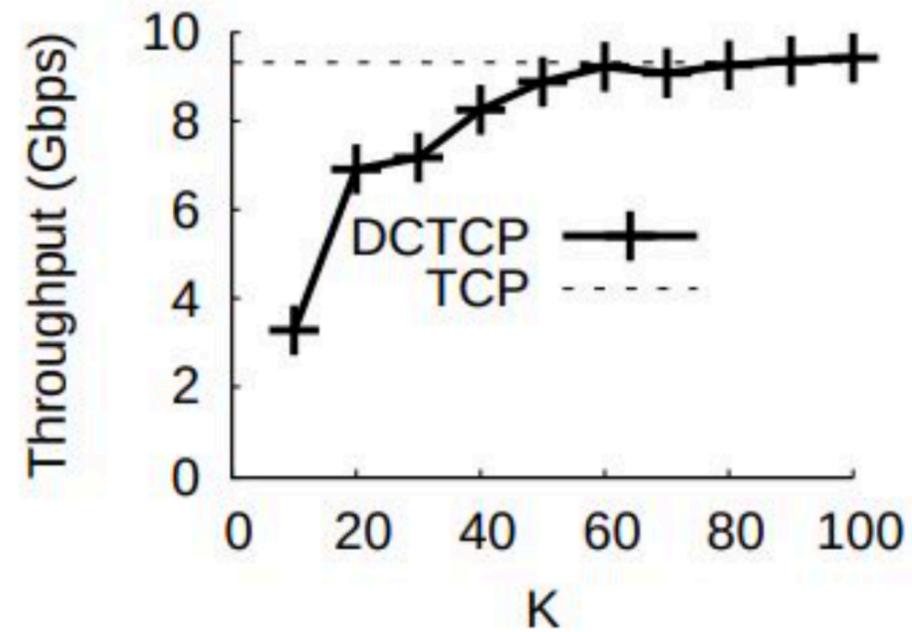


# DCTCP performance

1 Gbps, same bandwidth  
but lower switch occupancy

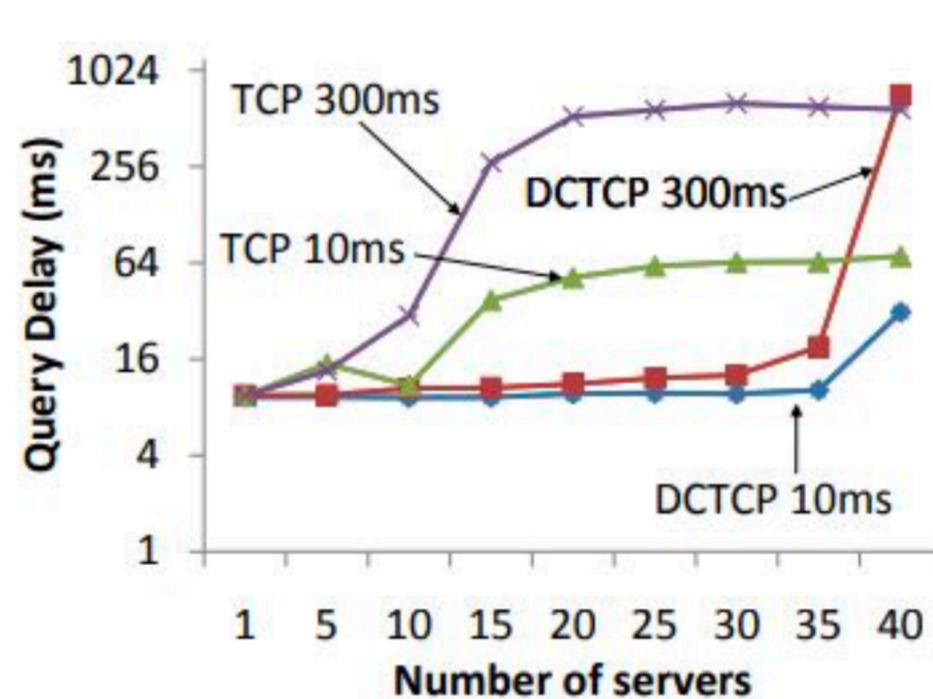


At 10 Gbps, after certain “K”  
threshold, the same bandwidth

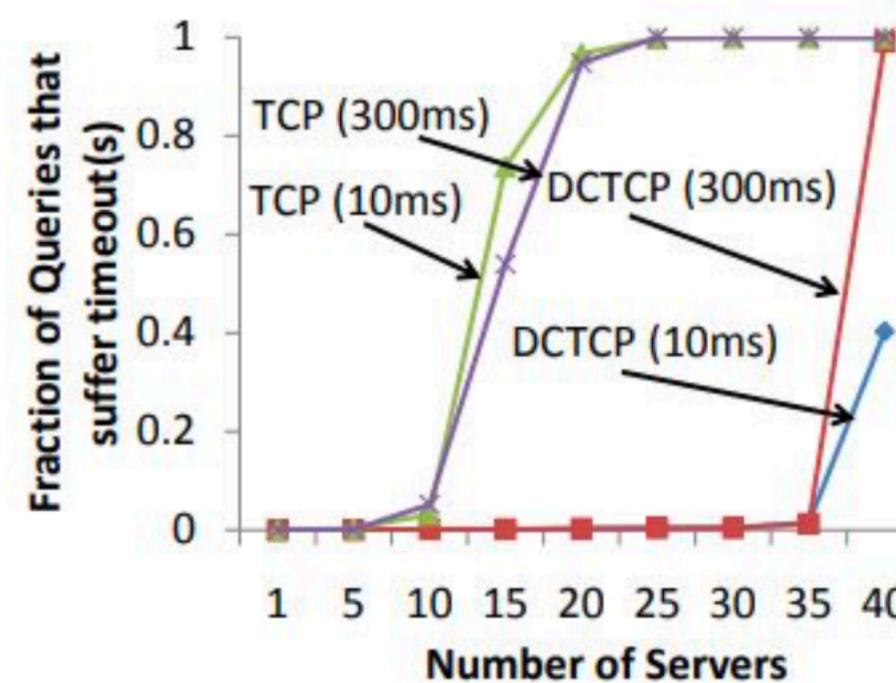


# What about incast?

Query: 1 machine sending 1MB/n data to n machines and waiting for the echo



Better performance than TCP up to a point where (#servers=35) where not even a single packet can pass from the switch

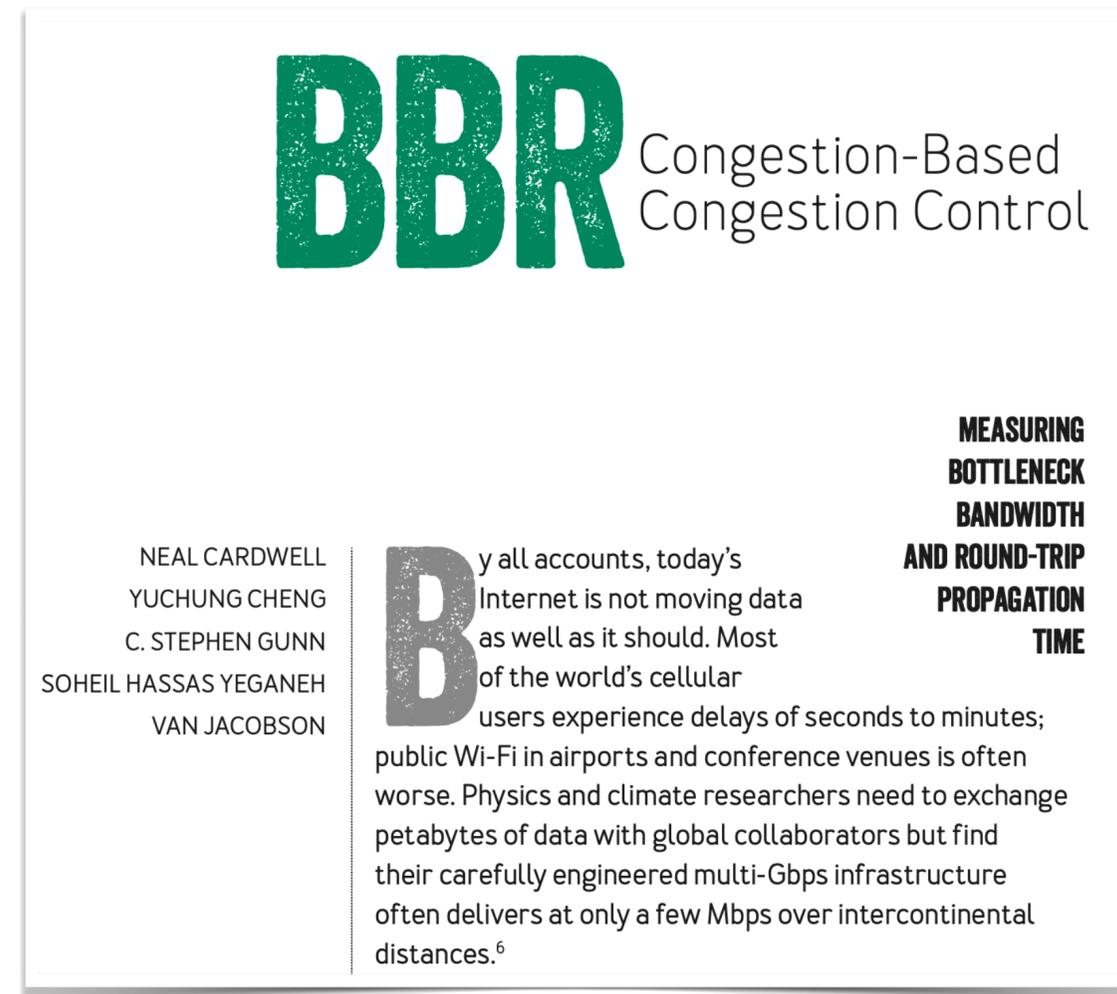


DCTCP has very low packet losses in comparison to TCP

DCTCP uses switch **queue occupancy** as congestion signal.

Can we use something else?

# Recall BBR



**BBR** Congestion-Based  
Congestion Control

NEAL CARDWELL  
YUCHUNG CHENG  
C. STEPHEN GUNN  
SOHEIL HASSAS YEGANEH  
VAN JACOBSON

**B**y all accounts, today's Internet is not moving data as well as it should. Most of the world's cellular users experience delays of seconds to minutes; public Wi-Fi in airports and conference venues is often worse. Physics and climate researchers need to exchange petabytes of data with global collaborators but find their carefully engineered multi-Gbps infrastructure often delivers at only a few Mbps over intercontinental distances.<sup>6</sup>

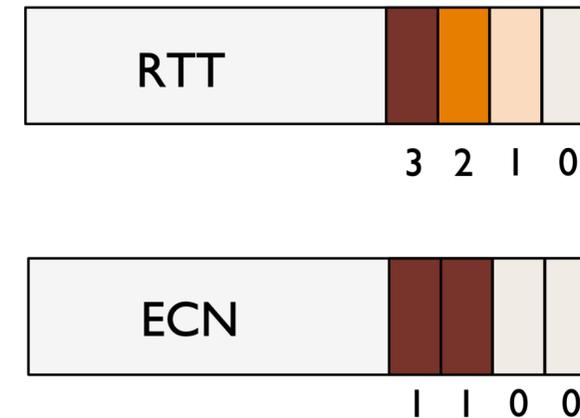
**MEASURING  
BOTTLENECK  
BANDWIDTH  
AND ROUND-TRIP  
PROPAGATION  
TIME**

Can we directly apply BBR on data center networks?

# TIMELY

Use Round Trip Time (RTT) as the indication of congestion signal

- RTT is a multi-bit signal indicating end-to-end congestion throughout the network — no explicit switch support required to do any marking
- RTT covers ECN signal completely, but not vice versa!



**TIMELY: RTT-based Congestion Control for the Datacenter**

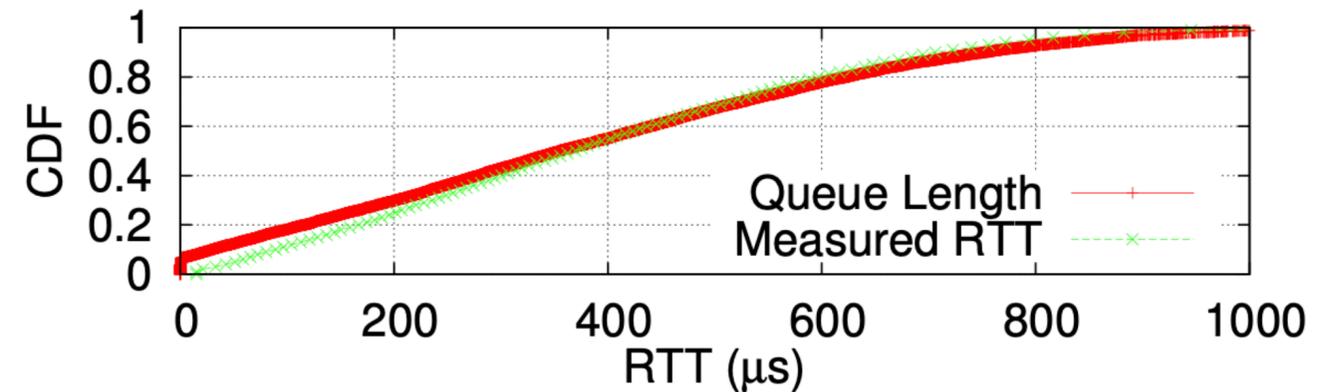
Radhika Mittal (UC Berkeley), Vinh The Lam, Nandita Dukkipati, Emily Blem, Hassan Wassel, Monia Ghobadi (Microsoft), Amin Vahdat, Yaogong Wang, David Wetherall, David Zats

Google, Inc.

**ABSTRACT**  
Datacenter transports aim to deliver low latency messaging together with high throughput. We show that simple packet delay, measured as round-trip times at hosts, is an effective congestion signal without the need for switch feedback. First, we show that advances in NIC hardware have made RTT measurement possible with microsecond accuracy, and that these RTTs are sufficient to estimate switch queuing. Then we describe how TIMELY can adjust transmission rates using RTT gradients to keep packet latency low while delivering high bandwidth. We implement our design

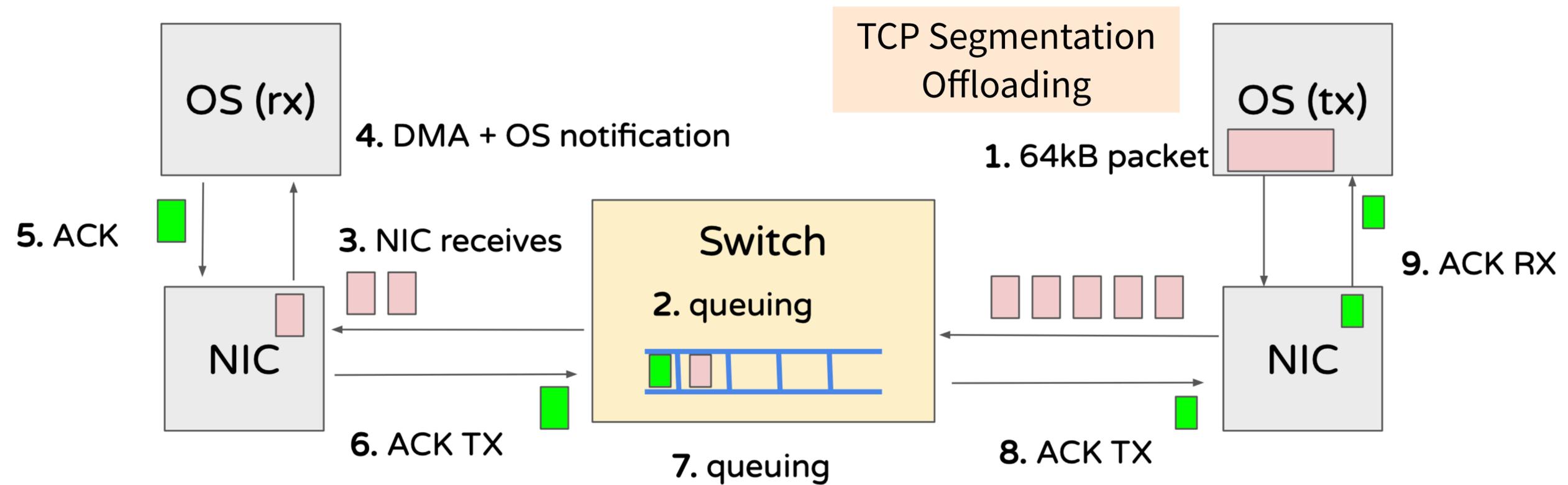
**1. INTRODUCTION**  
Datacenter networks run tightly-coupled computing tasks that must be responsive to users, e.g., thousands of backend computers may exchange information to serve a user request, and all of the transfers must complete quickly enough to let the complete response to be satisfied within 100 ms [24]. To meet these requirements, datacenter transports must simultaneously deliver high bandwidth ( $\gg$ Gbps) and utilization at low latency ( $\ll$ msec), even though these aspects of performance are at odds. Consistently low latency matters because even a small fraction of late operations

ACM SIGCOMM 2015



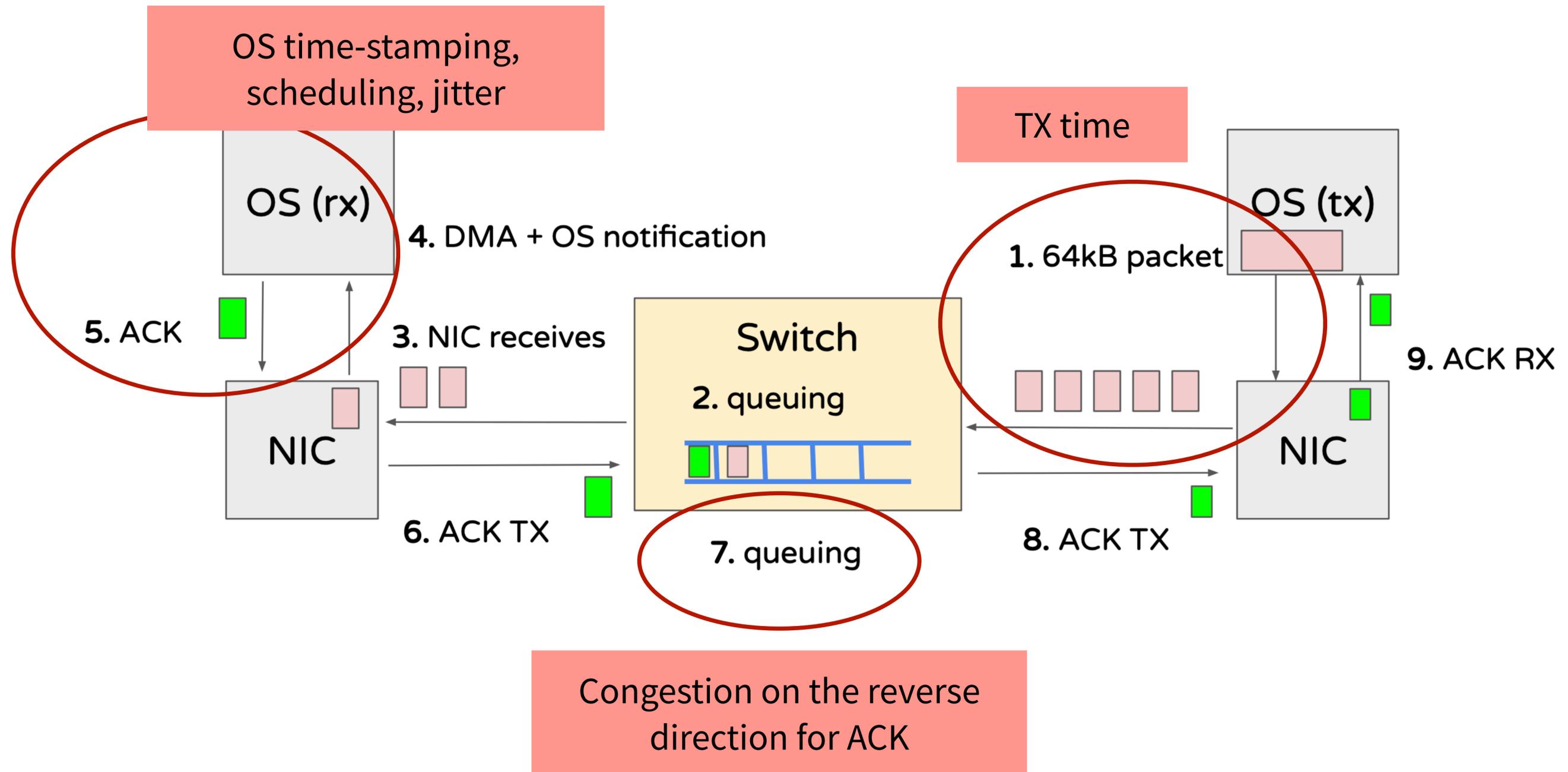
However, getting precise RTT is challenging. Why?

# RTT calculation challenges



Many steps are involved in the packet processing in one RTT

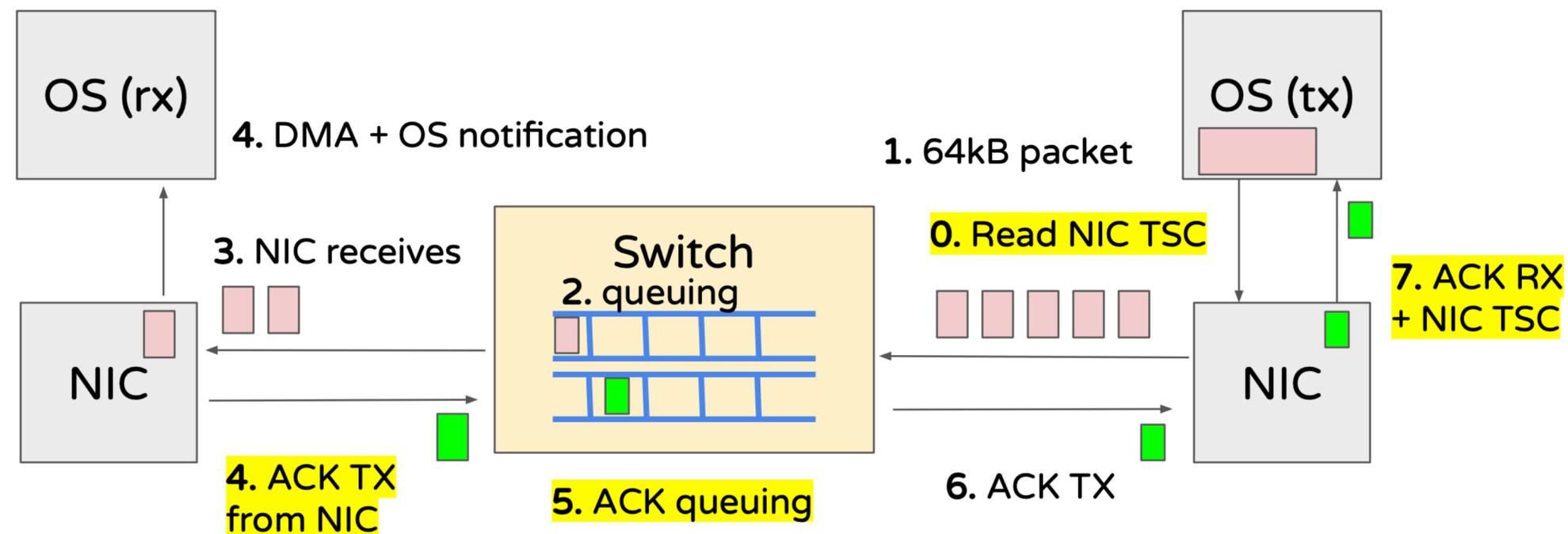
# RTT calculation challenges



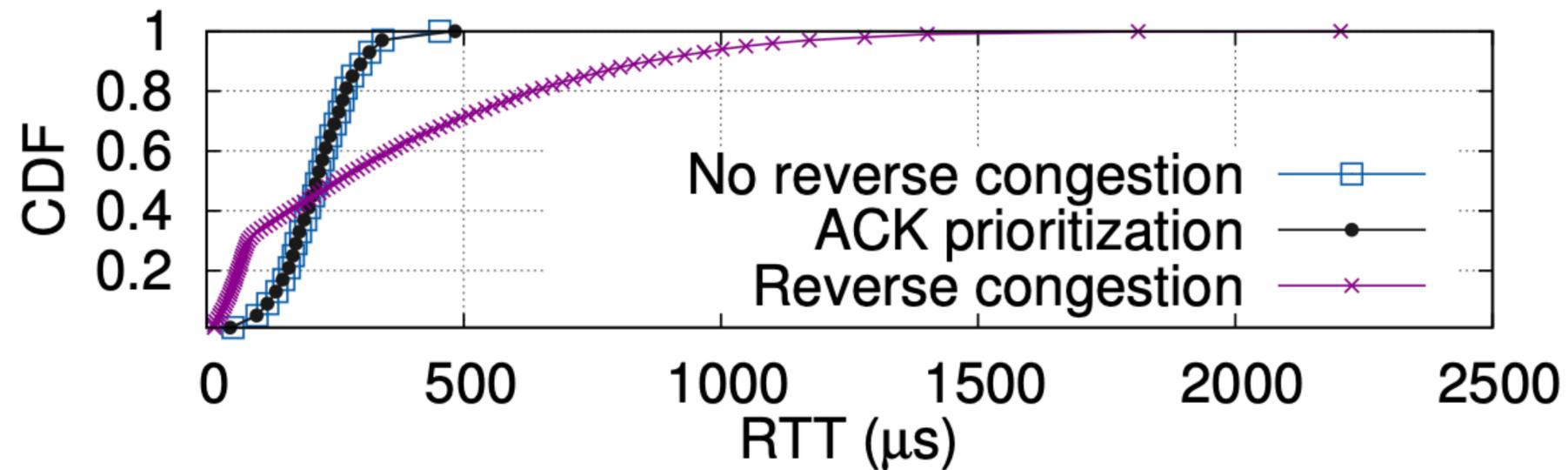
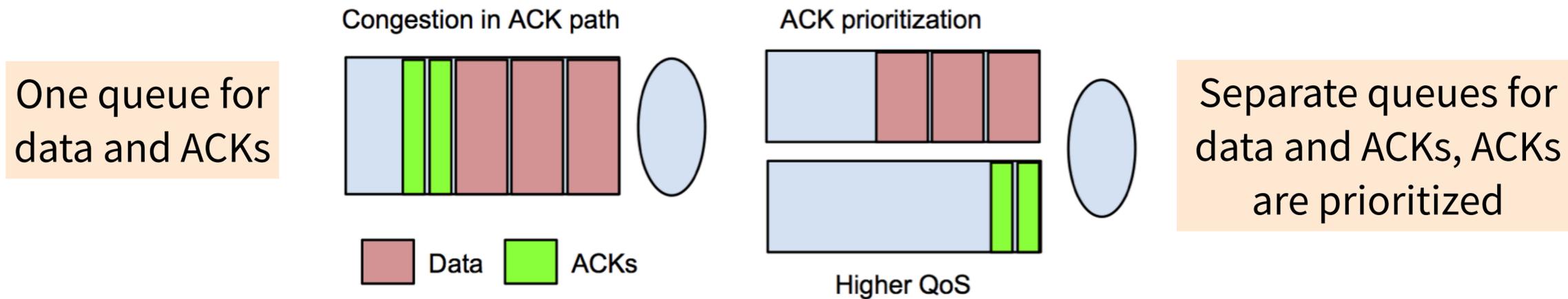
# RTT calculation: support from NICs

TIMELY assumes that

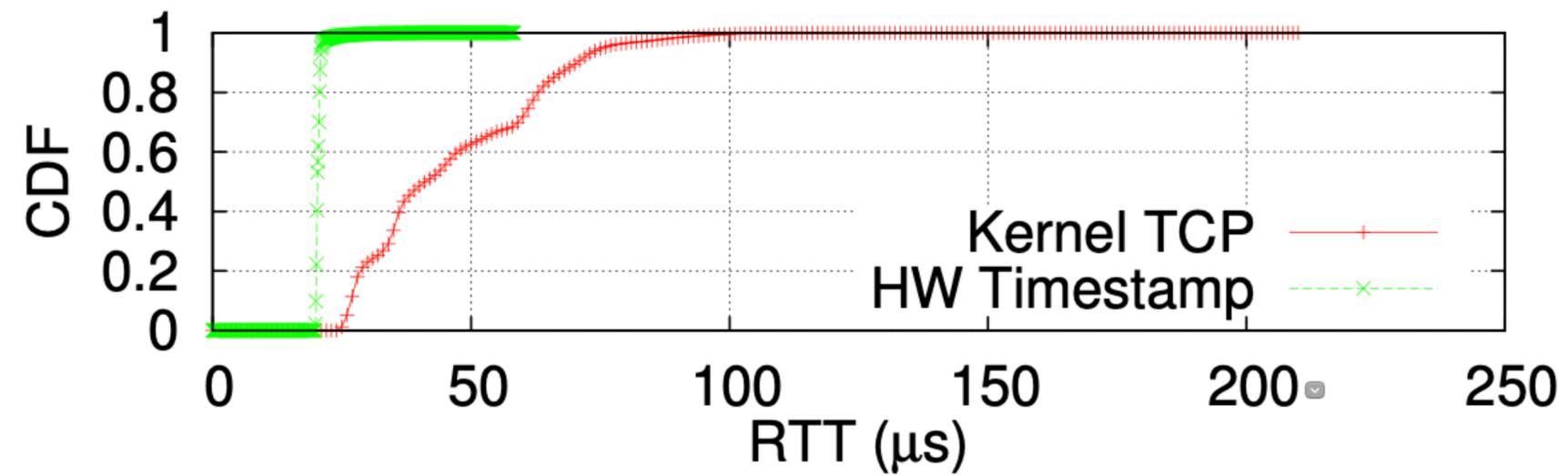
- The TX NIC can generate completion timestamps so that OS knows when a transmission is finished
- The RX NIC can generate ACKs in hardware without OS involvement
- At switches ACKs go through a high-priority separate queue



# Separate ACK queuing to solve reverse congestion

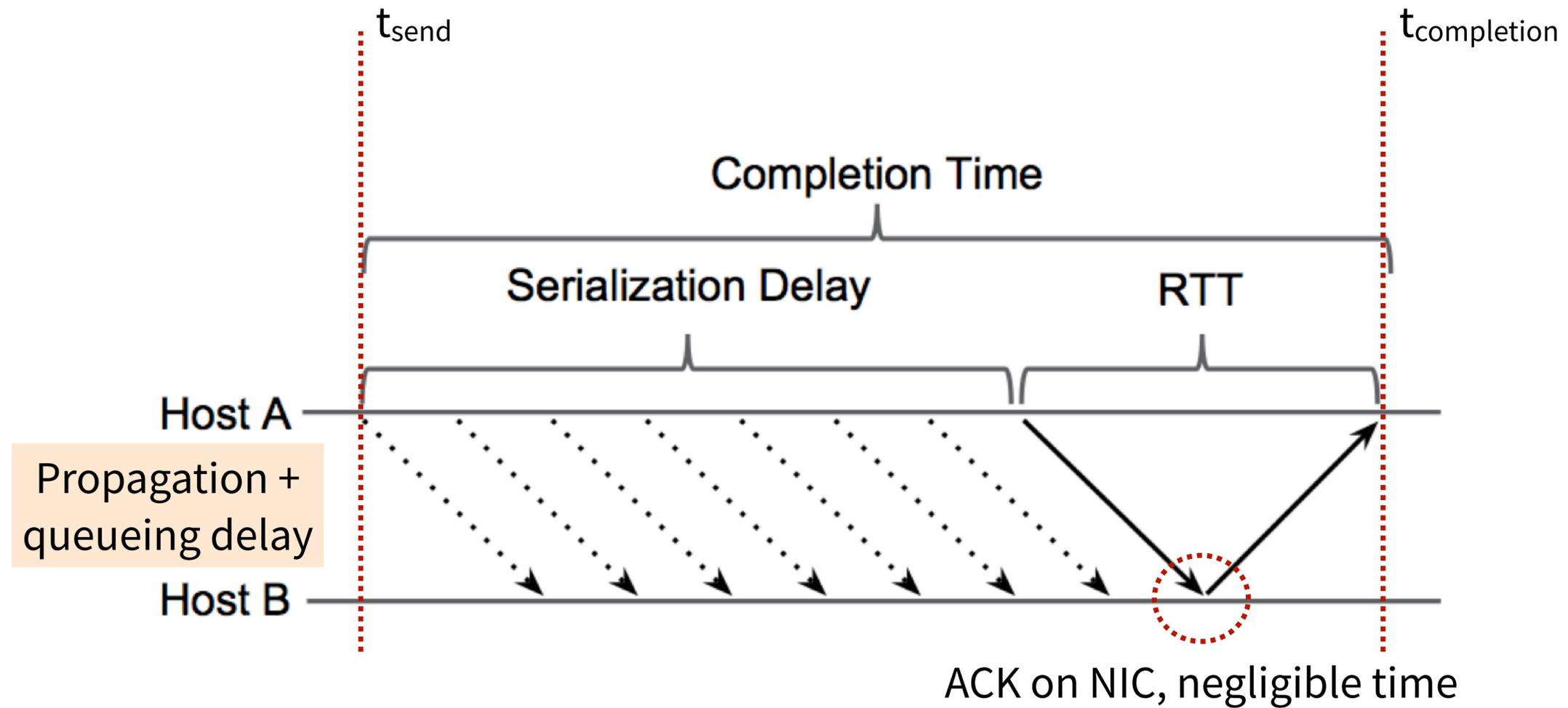


# Can we measure RTTs precisely?



Yes, the random variance is much smaller than the kernel TCP measurements

# RTT calculation



$$\text{RTT} = t_{\text{completion}} - t_{\text{send}} - \text{seg.size} / \text{NIC.linerate}$$

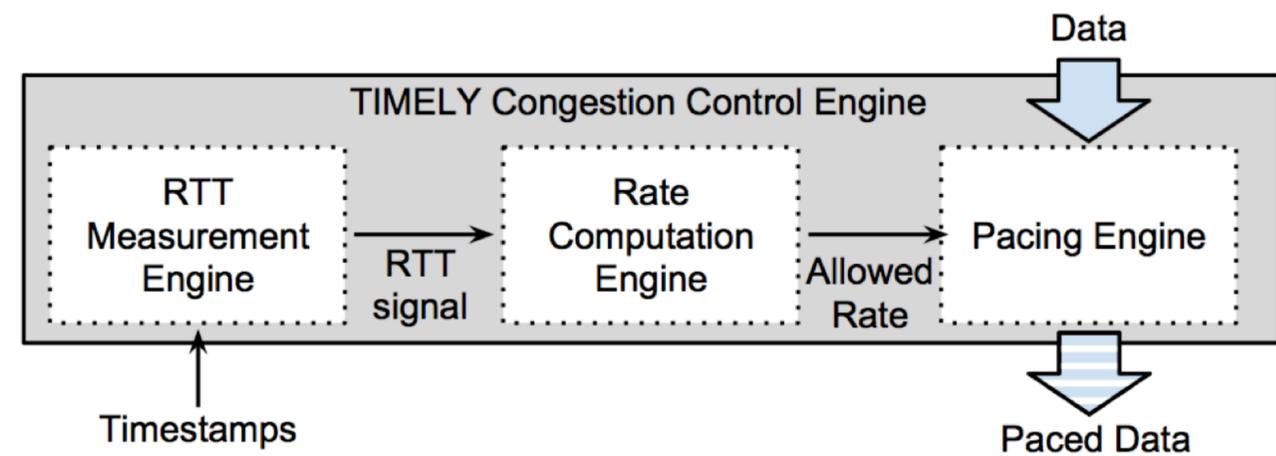
# TIMELY

## Independent of the transport used

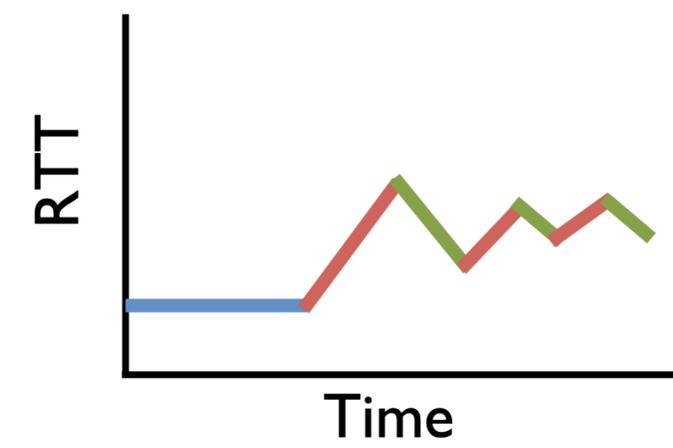
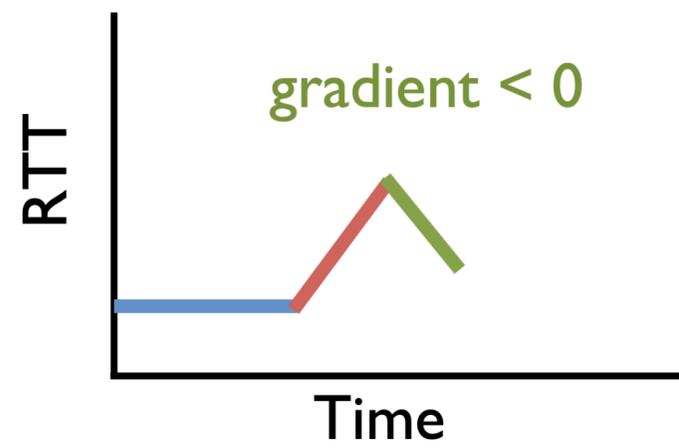
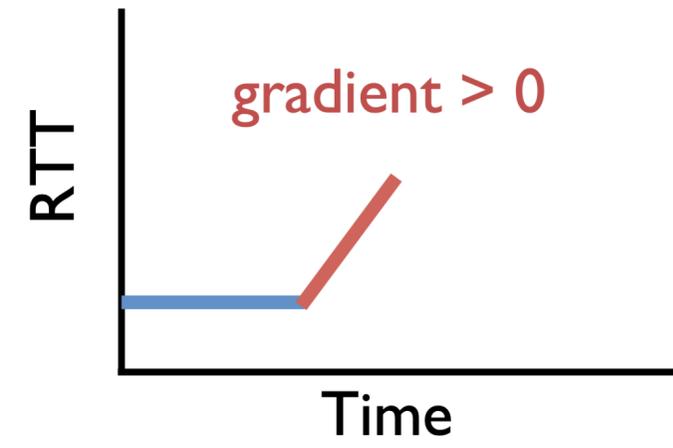
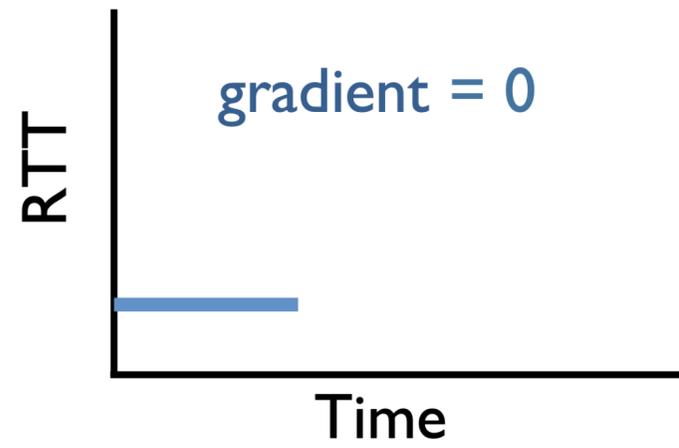
- Assumes an ACK-based protocol (TCP)
- Receivers must generate ACKs for incoming data

## Key concept

- Absolute RTTs are not used, only the **gradient** of the RTTs
- Positive gradient → rising RTT → queue buildup
- Negative gradient → decreasing RTT → queue depletion



# RTT gradient



# TIMELY pacing engine

**Algorithm 1:** TIMELY congestion control.

**Data:** new\_rtt

**Result:** Enforced rate

new\_rtt\_diff = new\_rtt - prev\_rtt ;

prev\_rtt = new\_rtt ;

rtt\_diff =  $(1 - \alpha) \cdot \text{rtt\_diff} + \alpha \cdot \text{new\_rtt\_diff}$  ;

▷  $\alpha$ : EWMA weight parameter

normalized\_gradient = rtt\_diff / minRTT ;

**if** new\_rtt <  $T_{\text{low}}$  **then**

rate  $\leftarrow$  rate +  $\delta$  ;

▷  $\delta$ : additive increment step

return;

**if** new\_rtt >  $T_{\text{high}}$  **then**

rate  $\leftarrow$  rate  $\cdot$   $\left( 1 - \beta \cdot \left( 1 - \frac{T_{\text{high}}}{\text{new\_rtt}} \right) \right)$  ;

▷  $\beta$ : multiplicative decrement factor

return;

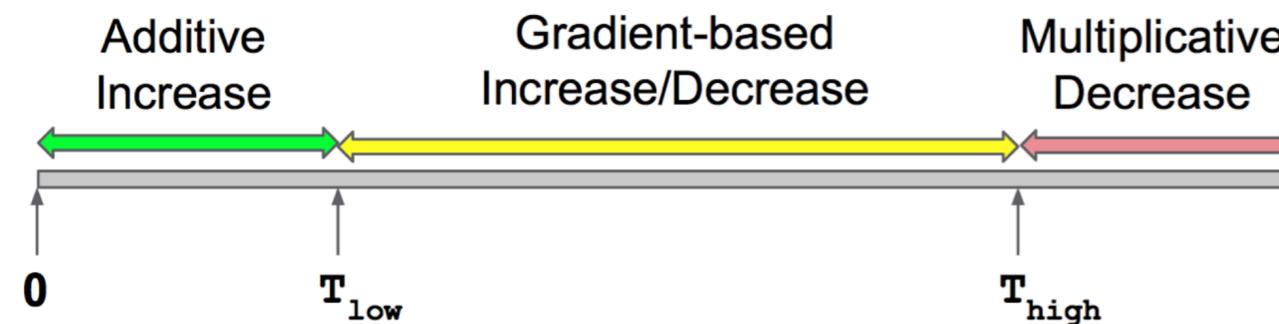
**if** normalized\_gradient  $\leq 0$  **then**

rate  $\leftarrow$  rate +  $N \cdot \delta$  ;

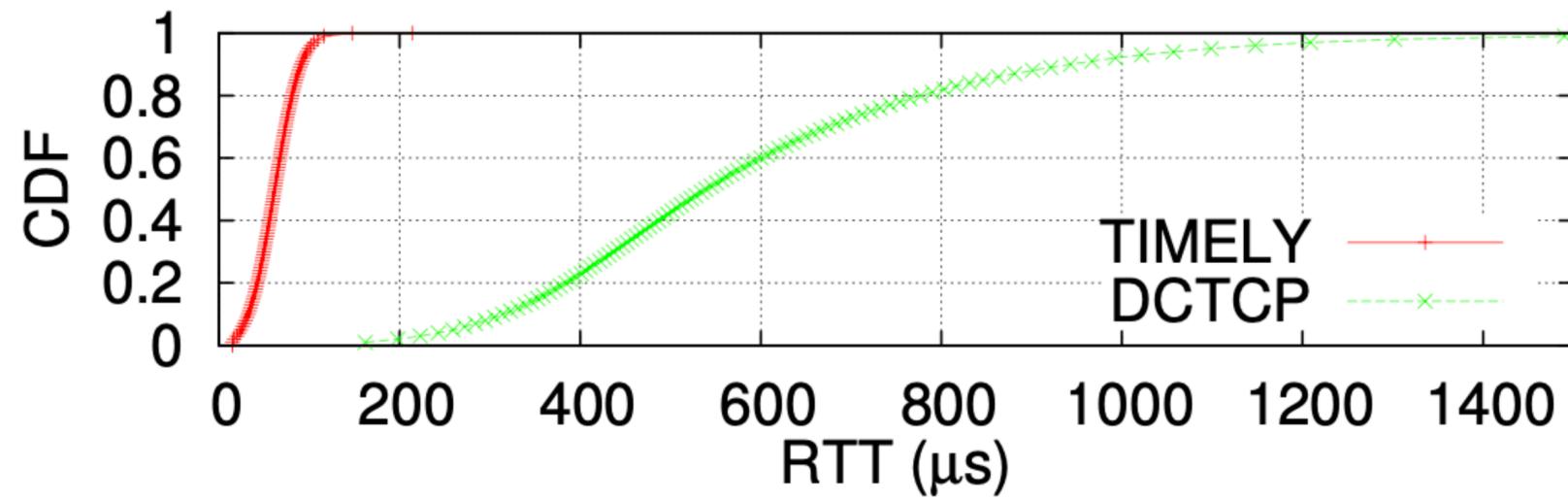
▷  $N = 5$  if gradient < 0 for five completion events (HAI mode); otherwise  $N = 1$

**else**

rate  $\leftarrow$  rate  $\cdot$   $(1 - \beta \cdot \text{normalized\_gradient})$

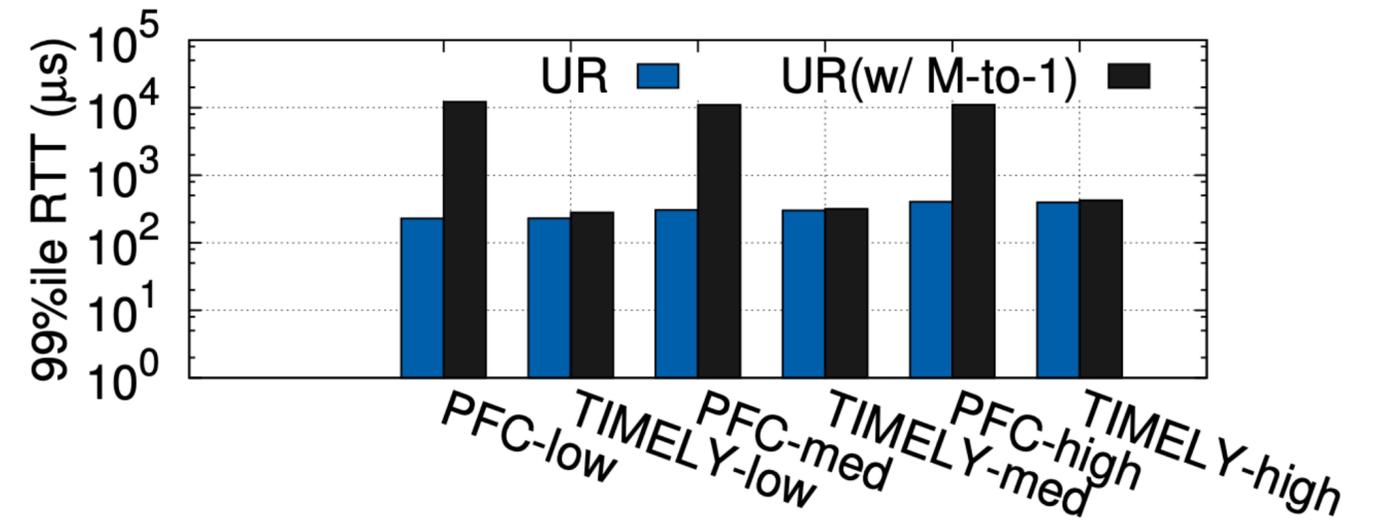
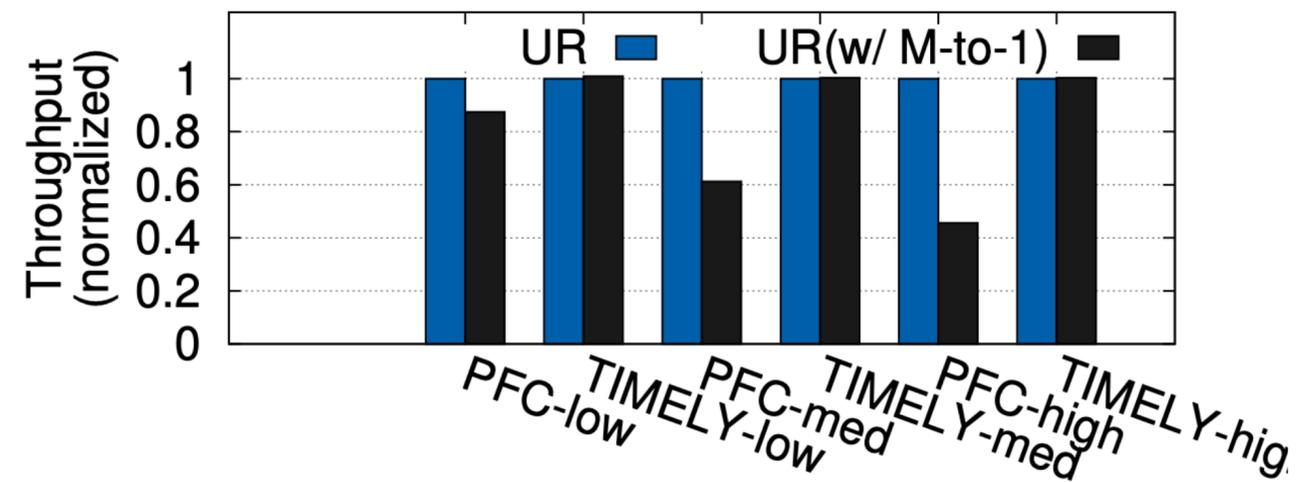


# TIMELY performance: vs. DCTCP



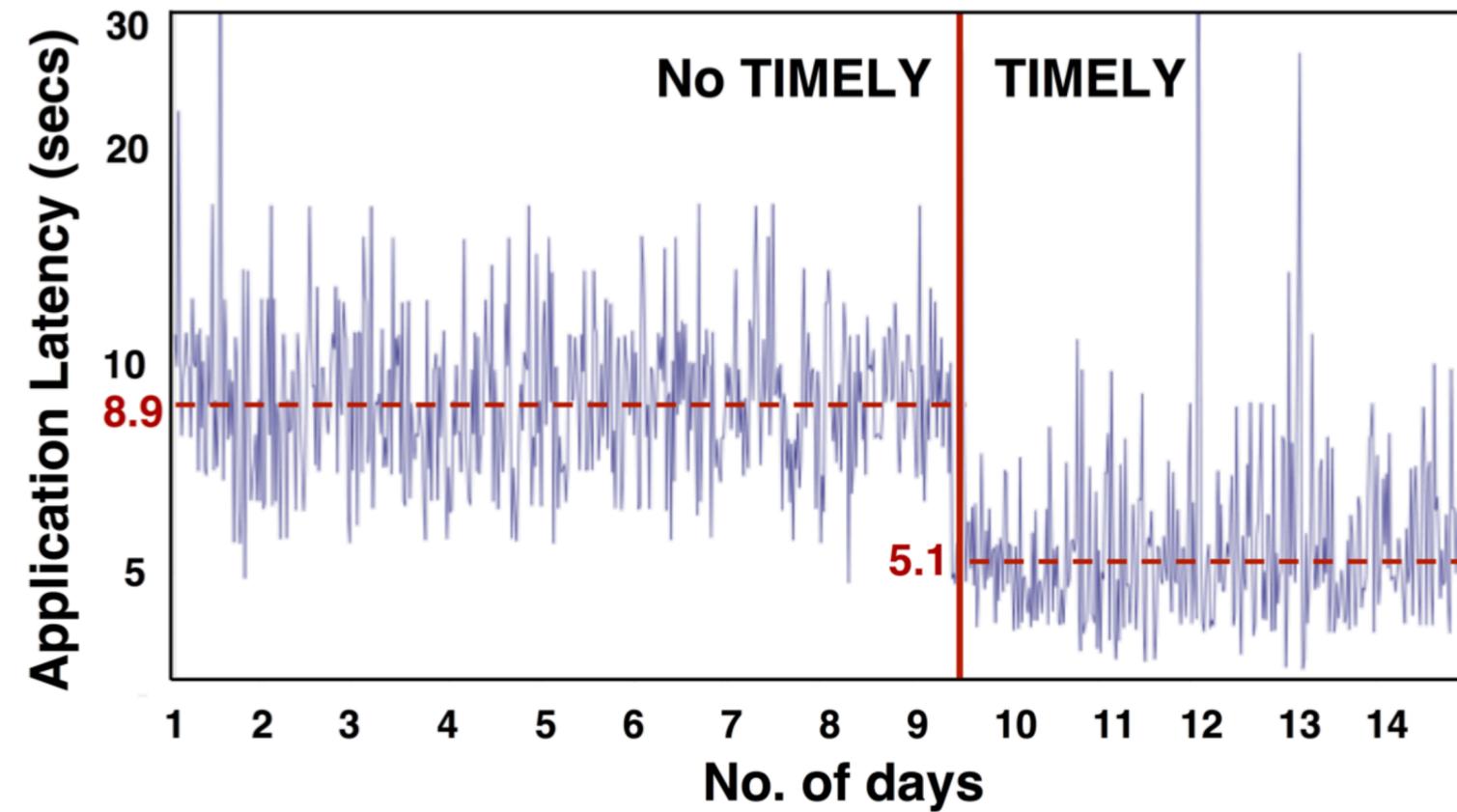
TIMELY achieves much lower yet stable RTTs

# TIMELY performance: incast



TIMELY throughput and latency stay the same with and without incast traffic

# TIMELY performance: application level



A (unknown) RPC latency of a data center storage benchmark

# Summary

## Congestion control challenges in data centers

- Network has low latency and high throughput
- Applications are diverse with different requirements
- Incast congestion

## Transport in data centers

- PFC has limited capability
- DCTCP: ECN-based congestion control
- TIMELY: RTT-based congestion control

# Other reading material

## Less is More: Trading a little Bandwidth for Ultra-Low Latency in the Data Center

Mohammad Alizadeh, Abdul Kabbani<sup>†</sup>, Tom Edsall<sup>\*</sup>, Balaji Prabhakar, Amin Vahdat<sup>‡§</sup>, and Masato Yasuda<sup>¶</sup>

Stanford University <sup>†</sup>Google <sup>\*</sup>Cisco Systems <sup>‡</sup>U.C. San Diego <sup>§</sup>NEC Corporation, Japan

### Abstract

Traditional measures of network goodness—goodput, quality of service, fairness—are expressed in terms of bandwidth. Network latency has rarely been a primary concern because delivering the highest level of bandwidth essentially entails driving up latency—at the mean and, especially, at the tail. Recently, however, there has been renewed interest in latency as a primary metric for mainstream applications. In this paper, we present the HULL (High-bandwidth Ultra-Low Latency) architecture to balance two seemingly contradictory goals: near baseline fabric latency and high bandwidth utilization.

of-service capability to the Internet resulted in proposals such as RSVP [55], IntServ [10] and DiffServ [35], which again focussed on bandwidth provisioning.

This focus on bandwidth efficiency has been well justified as most Internet applications typically fall into two categories. Throughput-oriented applications, such as file transfer or email, are not sensitive to the delivery times of individual packets. Even the overall completion times of individual operations can vary by multiple integer factors in the interests of increasing overall network throughput. On the other hand, latency-sensitive applications—such as web browsing and remote login—

HULL, USENIX NSDI 2012

## Deadline-Aware Datacenter TCP (D<sup>2</sup>TCP)

Balajee Vamanan  
Purdue University  
bvamanan@ecn.purdue.edu

Jahangir Hasan  
Google Inc.  
jahangir@google.com

T. N. Vijaykumar  
Purdue University  
vijay@ecn.purdue.edu

### ABSTRACT

An important class of datacenter applications, called Online Data-Intensive (OLDI) applications, includes Web search, online retail, and advertisement. To achieve good user experience, OLDI applications operate under soft-real-time constraints (e.g., 300 ms latency) which imply deadlines for network communication within the applications. Further, OLDI applications typically employ tree-based algorithms which, in the common case, result in bursts of children-to-parent traffic with tight deadlines. Recent work on datacenter network protocols is either deadline-agnostic (DCTCP) or is deadline-aware (D<sup>2</sup>) but suffers under bursts due to race conditions. Further, D<sup>2</sup> has the practical drawbacks of requiring changes to the switch hardware and not being able to coexist with legacy TCP.

We propose Deadline-Aware Datacenter TCP (D<sup>2</sup>TCP), a novel transport protocol, which handles bursts, is deadline-aware, and is compatible with legacy TCP.

The applications employ tree-based, divide-and-conquer algorithms where every query operates on data spanning thousands of servers [12].

An OLDI query's overall time budget gets divided among the nodes of the algorithm's tree (e.g., 40 ms for a parent-to-leaf RPC). To avoid missing its deadline, a parent node sends out an incomplete response without waiting for slow children that have missed their deadlines. Because such incomplete responses adversely affect response quality, and hence revenue, achieving fewer missed deadlines is important. While a node's response time includes both computational and network latencies, we focus on reducing the network delay. We note that in addition to achieving fewer missed deadlines, a network protocol that allows tighter network budgets is invaluable as it allows more time for computation, thus producing higher-quality responses.

A key reason for increased network delay is that all the children of a parent node face the same deadline and are likely to

D2TCP, ACM SIGCOMM 2012

## pFabric: Minimal Near-Optimal Datacenter Transport

Mohammad Alizadeh<sup>†</sup>, Shuang Yang<sup>†</sup>, Milad Sharif<sup>†</sup>, Sachin Katti<sup>†</sup>, Nick McKeown<sup>†</sup>, Balaji Prabhakar<sup>†</sup>, and Scott Shenker<sup>‡</sup>

<sup>†</sup>Stanford University <sup>‡</sup>Insieme Networks <sup>§</sup>U.C. Berkeley / ICSI  
{alizade, shyang, msharif, skatti, nickm, balaji}@stanford.edu shenker@icsi.berkeley.edu

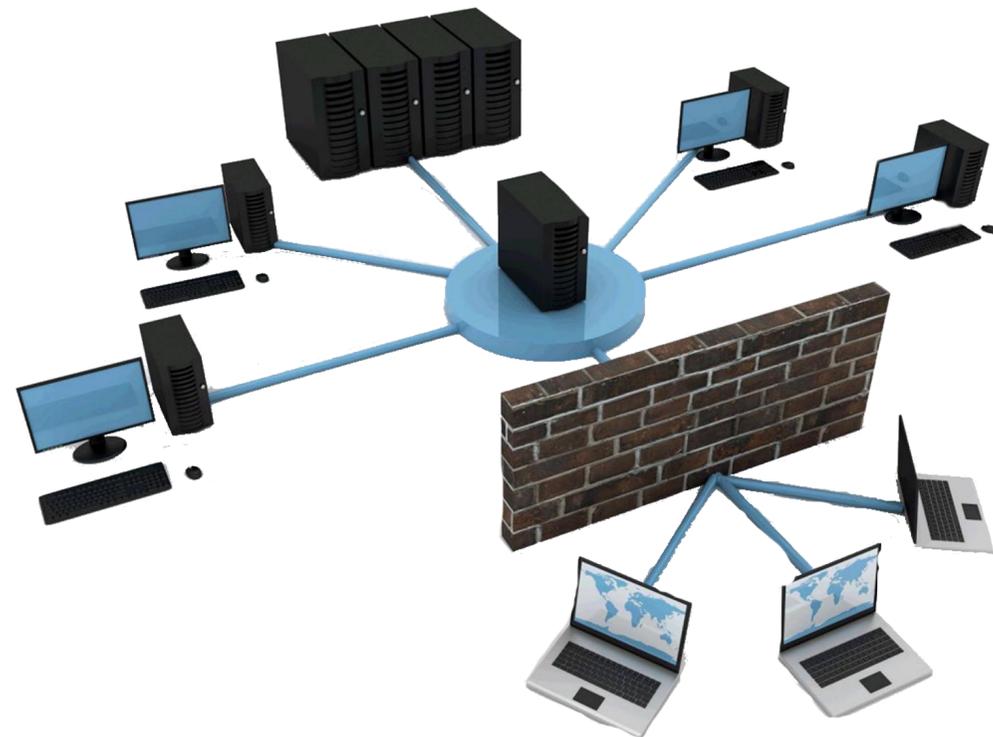
### ABSTRACT

In this paper we present pFabric, a minimalistic datacenter transport design that provides near theoretically optimal flow completion times even at the 99th percentile for short flows, while still minimizing average flow completion time for long flows. Moreover, pFabric delivers this performance with a very simple design that is based on a key conceptual insight: datacenter transport should decouple flow scheduling from rate control. For flow scheduling, packets carry a single priority number set independently by each flow; switches have very small buffers and implement a very simple priority-based scheduling/dropping mechanism. Rate control is also correspondingly simpler; flows start at line rate and throttle

Motivated by this observation, recent research has proposed new datacenter transport designs that, broadly speaking, use rate control to reduce FCT for short flows. One line of work [3, 4] improves FCT by keeping queues near empty through a variety of mechanisms (adaptive congestion control, ECN-based feedback, pacing, etc) so that latency-sensitive flows see small buffers and consequently small latencies. These *implicit* techniques generally improve FCT for short flows but they can never precisely determine the right flow rates to optimally schedule flows. A second line of work [21, 14] *explicitly* computes and assigns rates from the network to each flow in order to schedule the flows based on their sizes or deadlines. This approach can potentially provide very good performance, but it is rather complex and challenging to implement in

pFabric, ACM SIGCOMM 2013

# Next week: software defined networking



How do we manage the complex networks?

- Remember all the protocols
- Remember the configurations with every protocol
- Diagnose problems with networking tools like ping, traceroute, tcpdump?

[Home](#) > [News](#)

## Misconfigured CenturyLink database caused global internet outage

By [Sead Fadilpašić](#) 20 days ago

Global internet traffic dropped by 3.5 percent.