



Advanced Networked Systems SS24

Data Center Transport

Prof. Lin Wang, Ph.D. Computer Networks Group Paderborn University <u>https://en.cs.uni-paderborn.de/cn</u>



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?

Diverse applications and workloads

- Large variety in performance requirements

Traffic patterns

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



O PyTorch

Built out of commodity components: no expensive/customized hardware

Network

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



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



ApplicationReliableTCPIPLossyLossy

The transport layer in the network model:

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

Recall TCP AIMD





cwnd += 1 / cwnd

7

TCP incast problem

A data center application runs 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







Commodity off-the-shelf switches typically have shallow buffers. Do you know why?







How does TCP handle this?





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 alternaclient 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 millisecondsorders of magnitude greater than typical data fetch timessignificantly degrading overall throughput.

TCP incast

Packet drops due to the capacity overrun at shared commodify 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 (ACK, SACK)
- Various timeouts based optimizations

Whatever clever way you come up with can be over-optimized

- 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^{±1}, Albert Greenberg¹, David A. Maltz¹, Jitendra Padhye¹, Parveen Patel¹, Balaji Prabhakar[±], Sudipta Sengupta¹, Murari Sridharan¹

[†]Microsoft Research [‡]Stanford University {albert, dmaltz, padhye, parveenp, 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 using commodity, shallow biffered switchase. We find DCTCP do eral recent research proposals envision creating economical, easyto-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-

Explicit Congestion Notification (ECN)



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

- Part of the IP packet header (2 bits): capability and 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 (Random Early Detection, RED) that can probabilistically mark packets

Updated by: <u>4301, 6040, 8311</u>	PROPOSED STANDARD
Network Working Group	K. Ramakrishnan
Request for Comments: 3168	TeraOptic Networks
Updates: 2474, 2401, 793	S. Floyd
Obsoletes: 2481	ACIRI
Category: Standards Track	D. Black
	EPIL September 2001
	September 2001
The Addition of Explicit Congestion Notificati	on (ECN) to IP
Status of this Memo	
This document specifies an Internet standards tra	ck protocol for the
Internet community, and requests discussion and s	uggestions for
improvements. Please refer to the current edition	n of the "Internet
Official Protocol Standards" (STD 1) for the stan	dardization state
and status of this protocol. Distribution of thi	s memo is unlimited.
Copyright Notice	
Copyright (C) The Internet Society (2001). All R	lights Reserved.
Abstract	
This memo specifies the incorporation of ECN (Exp	licit Congestion
Notification) to TCP and IP, including ECN's use	of two bits in the
IP header.	



ECN bit in action

Assume that B is sending TCP segments 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?

The ECN bits location in TCP header



The TCP congestion window logic:

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

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 the ECE flag, until the sender ACKs back using CWR flag bit

DCTCP receiver

- Only mark ACKs corresponding to the ECN-marked packet

Sender's congestion control

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

DCTCP congestion window calculations

 $F = \frac{\text{\#ECN-marked ACKs}}{\text{\#Total ACKs}}$

 $\alpha \leftarrow (1 - g) \times \alpha + g \times F$

In every RTT, calculate the percentage of ECN-marked ACKs

Use a sliding window to estimate the average percentage of ECN-marked ACKs

Decide the congestion window based on the estimation

 $cwnd \leftarrow cwnd \times (1 - \alpha/2)$

DCTCP vs TCP example



DCTCP performance

At 1 Gbps, same bandwidth but lower switch queue occupancy At 10 Gbps, after certain K (queue occupancy parameter) threshold, the same bandwidth



Significantly lower queueing delay, similar throughput to TCP for large K

What about incast?

Query: 1 machine sending 1 MB/*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

Can we use a different congestion signal than the queue occupancy on the switch?

Recall BBR



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

1. INTRODUCTION

Datacenter transports aim to deliver low latency messaging together with high throughput. We show that simple packet delay, measured as round-tinj 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 queueing. Then we describe how TIMELY can adjust transmission rates using RTT gradients to keep packet latency low while delivering high handwidth. We implement our design 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 (>Gbps) and utilization at low latercy (<Cmsec), even though these aspects of performance are at odds. Consistently low latency matters because even a small fraction of late operations

However, getting precise RTT is challenging. Why?





3 2

0

RTT

RTT calculation challenges



RTT calculation challenges



RTT calculation challenges



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



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

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 \rightarrow rising RTT \rightarrow 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 ;		
$prev_rtt = new_rtt;$		
$rtt_diff = (1 - \alpha) \cdot rtt_diff + \alpha \cdot new_rtt_diff;$		
$\triangleright \alpha$: EWMA weight parameter		
normalized_gradient = rtt_diff / minRTT ;		
if $new_rtt < T_{low}$ then		
rate \leftarrow rate + δ ;		
$\triangleright \delta$: additive increment step		
_ return:		
if $new_rtt > T_{high}$ then		
rate \leftarrow rate $\cdot (1 - \beta \cdot (1 - \frac{T_{\text{high}}}{T_{\text{norm rft}}}));$		
$\triangleright \beta$: multiplicative decrement factor		
_ return;		
if normalized_gradient ≤ 0 then		
rate \leftarrow rate + N $\cdot \delta$;		
\triangleright N = 5 if gradient < 0 for five completion events		
(HAI mode); otherwise $N = 1$		
else		
\Box rate \leftarrow rate \cdot (1 - β \cdot normalized gradient)		



TIMELY performance: vs. DCTCP



TIMELY achieves much lower yet stable RTTs

TIMELY performance: incast



UR: Uniform Random

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

TIMELY issues

Gradient-based pacing

- **Complex** to tune the parameters

RTT measurement

- Only the fabric-related delay calculated with timestamps on NICs
- Does not differentiate between **fabric congestion** and **end-host congestion**

Extreme incast

- What happens if the number of flows is larger than the path BDP?
- Even one packet per flow would overrun the network

Swift

Swift: Delay is Simple and Effective for Congestion Control in the Datacenter

Gautam Kumar, Nandita Dukkipati, Keon Jang (MPI-SWS)^{*}, Hassan M. G. Wassel, Xian Wu, Behnam Montazeri, Yaogong Wang, Kevin Springborn, Christopher Alfeld, Michael Ryan, David Wetherall, and Amin Vahdat Google LLC

ABSTRACT

We report on experiences with Swift congestion control in Google datacenters. Swift targets an end-to-end delay by using AIMD control, with pacing under extreme congestion. With accurate RTT measurement and care in reasoning about delay targets, we find this design is a foundation for excellent performance when network distances are well-known. Importantly, its simplicity helps us to meet operational challenges. Delay is easy to decompose into fabric and host components to separate concerns, and effortless to deploy and maintain as a congestion signal while the datacenter evolves. In large-scale testbed experiments. Swift delivers a tail latency of <50 µs for short RPCs, with near-zero packet drops, while sustaining ~100Gbps throughput per server. This is a tail of $<3\times$ the minimal latency at a load close to 100%. In production use in many different clusters, Swift achieves consistently low tail completion times for short RPCs, while providing high throughput for long RPCs. It has loss rates that are at least 10× lower than a DCTCP protocol, and handles O(10k) incasts that sharply degrade with DCTCP.

CCS CONCEPTS

Networks → Transport protocols; Data center networks;

KEYWORDS

Congestion Control, Performance Isolation, Datacenter Transport

ACM Reference Format:

Gautam Kumar, Nandita Dukkipati, Keon Jang, Hassan M. G. Wassel, Xian Wu, Behnam Montazeri, Yaogong Wang, Kevin Springborn, Christopher Alfeld, Michael Ryan, David Wetherall, and Amin Vahdat. 2020. Swift: Delay is Simple and Effective for Congestion Control in the Datacenter. In Annual conference of the ACM Special Interest Group on Data Communication on the applications, technologies, architectures, and protocols for computer communication (SIGCOMM '20), August 10–14, 2020, Virtual Event, USA. ACM, New York, NY, USA, 15 pages. https://doi.org/10.1145/3387514.3406591

at 1M+ IOPS, else expensive servers sit idle while they wait for I/O [9]. Tight tail latency is also important because datacenter applications often use partition-aggregate communication patterns across many hosts [16]. For example, BigQuery [51], a query engine for Google Cloud, relies on a shuffle operation [11] with high IOPS per server [36]. Congestion control is thus a key enabler (or limiter) of system performance in the datacenter.

In this paper, we report on Swift congestion control that we use in Google datacenters. We found protocols such as DCTCP [1] inadequate because they commonly experience milliseconds of tail latency, especially at scale. Instead, Swift is an evolution of TIMELY [38] based on Google's production experience over the past five years. It is designed for excellent low-latency messaging performance at scale and to meet key operational needs: deploying and maintaining protocols while the datacenter is changing quickly due to technology trends; isolating the traffic of tenants in a shared fabric; efficient use of host CPU and NIC resources; and handling a range of traffic patterns including incast.

Swift is built on a foundation of hosts that independently adapt rates to a *target* end-to-end delay. We find that this design achieves high levels of performance when we accurately measure delay with NIC timestamps and carefully reason about targets, and has many other advantages. Delay corresponds well to the higher-level service-level objectives (SLOs) we seek to meet. It neatly decomposes into fabric and host portions to respond separately to different causes of congestion. In the datacenter, it is easy to adjust the delay target for different paths and competing flows. And using delay as a signal lets us deploy new generations of switches without concern for features or configuration because delay is always available, as with packet loss for classic TCP.

Compared to other work, Swift is notable for leveraging the simplicity and effectiveness of delay. Protocols such as DCTCP [1], PFC [49], DCQCN [59] and HPCC [34] use explicit feedback from

Swift designs

Simple target delay window control

Separating fabric and host congestion

Fractional congestion window to handle large-scale incast

Simple target delay window control

$\begin{array}{llllllllllllllllllllllllllllllllllll$	ces MD once every RTT		
4 On Receiving ACK			
5 retransmit_cnt $\leftarrow 0$			
$6 target_delay \leftarrow TargetDelay()$	⊳ See S3.5		
7 if <i>delay</i> < <i>target_delay</i> then	⊳ Additive Increase (AI)	Cumulative increase over	
s if $cwnd \ge 1$ then	if $cwnd \ge 1$ then	an RTT is equal to <i>ai</i>	
9 $cwnd \leftarrow cwnd + \frac{ai}{cwnd} \cdot num_acked$			
10 else			
11 $\left[cwnd \leftarrow cwnd + ai \cdot num_acket \right]$	d	Constrained to be max.	
12 else ⊳ Multij	plicative Decrease (MD)		
13 if can_decrease then		once per kri	
14 $cwnd \leftarrow \max(1 - \beta \cdot (\frac{delay - target_delay}{delay})),$			
$1 - max_mdf) \cdot cwnd$	5		
Proportional to the extent of			
	congestion		

Deciding target delay



Topology-based scaling:

- Fixed base delay plus a fixed per-hop delay
- Forward path hop count measured with IP TTL and reflected back with ACKs

Flow-based scaling:

- Target delay increases with the number of competing flows
- Average queue length grows as $O(\sqrt{N})$
- Adjust the target in proportion to $1/\sqrt{cwnd}$

Delay measurement in TIMELY



Fabric vs. host congestion



Timestamps measurement



Which part represents the remote queueing delay?

Timestamps measurement



Endpoint congestion control

Endpoint delay: $(t_4 - t_2) + (t_6 - t_5)$

Fabric delay: RTT - endpoint_delay



Follow the same target delay window control mechanism to decide *ecwnd*

Target delay for endpoints is decided based on Exponential Weighted Moving Average (EWMA) to remove noise

Actual cwnd = min(fcwnd, ecwnd)

fcwnd: fabric congestion window based on the fabric delay *ecwnd*: endpoint congestion window based on the endpoint delay

Handling large-scale incast



Even a **congestion window of one** is not able to prevent the incast congestion

Handling large-scale incast

Why not a problem in TIMELY?



[1] Ahmed Saeed, Nandita Dukkipati, Vytautas Valancius, Vinh The Lam, Carlo Contavalli, Amin Vahdat. Carousel: Scalable Traffic Shaping at End Hosts. ACM SIGCOMM 2017.

Pacing overhead in TIMELY



Pacing is done in chunks of 64KB in TIMELY, instead of MTU-size level in Swift

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
- Swift: simpler window control logic, handling endpoint congestion, and dealing with large-scale incase



Further reading material

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

Mohammad Alizadeh, Abdul Kabbani[†], Tom Edsall*, Balaji Prabhakar, Amin Vahdat^{†§}, and Masato Yasuda^{II} Stanford University [†]Google ^{*}Cisco Systems [§]U.C. San Diego [¶]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 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:

of-service capability to the Internet resulted in propos als 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 and, especially, at the tail. Recently, however, there has as file transfer or email, are not sensitive to the delive ery 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 near baseline fabric latency and high bandwidth utiliza- applications-such as web browsing and remote login-

Congestion Control for Large-Scale RDMA Deployments

Yibo Zhu^{1,3} Haggai Eran² Daniel Firestone¹ Chuanxiong Guo¹ Marina Lipshteyn¹ Yehonatan Liron² Jitendra Padhye¹ Shachar Raindel² Mohamad Haj Yahia² Ming Zhang¹

¹Microsoft ³Mellanox ³U. C. Santa Barbara

ABSTRACT brutal economics of cloud services business dictates that CPU Modern datacenter applications demand high throughput (40Gbps) and ultra-low latency (< 10 µs per hop) from the network, with low CPU overhead. Standard TCP/IP stacks usage that cannot be monetized should be minimized: a core usage that cannot be monetized should be minimized a core spent on supporting slight TC throughput is a core that can-net be sold as a VM. Other applications such an distributed memory caches [10, 30] and large-scale machine idensing demand ultra-low latency (sees than 10 µe per hop) mes-age transition. Toulitional TCDPH zooks have for higher largeny [10]. Memory in Memory in Amountain to passed have for higher memory in Memory in Amountain to passed have for how low cannot meet these requirements, but Remote Direct Merr ory Access (RDMA) can. On IP-routed datacenter networks, RDMA is deployed using RoCEv2 protocol, which relies on BDMA is displayed using BeGT-2 protocol, which relies on Philoticy-based Bey Control (PC) to create desplication special and the philotic spectra of the spectra mass due to protein like based of early the biologing at us-most of the protein like based of early the biologing at us-locy optimize DC/Cy performance, we have the full full of coprimize DC/Cy performance, which is a full model, and other protocol guarantees the spectra of the spectra and other protocol guarantees. Using a shift colors market hand the protocol guarantees that the full full pro-teement of the Merican NRL, and is bine to chemed in this demonstration Merican NRL, and is bine to chemed in the technology in Microsoff's datacenters to provide ultra-low latency and high throughput to applications, with very low CPU overhead. With RDMA, network interface cards (NICs) transfer data in and out of pre-registered memory buffers at both end hosts. The networking protocol is implemented entirely on the NICs, bypassing the host networking stack. The bypass significantly reduces CPU overhead and overall to tency. To simplify design and implementation, the protocol ssemes a lossless networking fabric. While the HPC community has long used RDMA in specialplemented in Mellanox NICs, and is being deployed in Microsoft's datacenters.

purpose clusters [11, 24, 26, 32, 38], deploying RDMA on a

pFabric: Minimal Near-Optimal Datacenter Transport

Mohammad Alizadeh11, Shuang Yang1, Milad Sharif1, Sachin Katti1, Nick McKeown[†], Balaji Prabhakar[†], and Scott Shenker^s

¹Stanford University ¹Insieme Networks ⁵U.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 comple tion 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 sim ple 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

Homa: A Receiver-Driven Low-Latency **Transport Protocol Using Network Priorities**

Behnam Montazeri, Yilong Li, Mohammad Alizadeh[†], and John Ousterhout Stanford University, †MIT

ABSTRACT

Homa is a new transport protocol for datacenter networks. It provides exceptionally low latency, especially for workloads with a high volume of very short messages, and it also supports large messages and high network utilization. Homa uses in-network priority queues to ensure low latency for short messages; priority allocation is managed dynamically by each receiver and integrated with a receiver-driven flow control mechanism. Homa also uses controlled overcommitment of receiver downlinks to ensure efficient bandwidth utilization at high load. Our implenentation of Homa delivers 99th percentile round-trip times less than 15 us for short messages on a 10 Gbps network running at 80% load. These latencies are almost 100x lower than the best published measurements of an implementation. In simulations Homa's latency is roughly equal to pFabric and significantly better than pHost, PIAS, and NDP for almost all message sizes and workloads. Homa can also sustain higher network loads than pFabric, pHost, or PIAS.

these conditions, so the latency they provide for short messages is far higher than the hardware potential, particularly under high network loads.

Recent years have seen numerous proposals for better transport protocols, including improvements to TCP [2, 3, 31] and a variety of new protocols [4, 6, 14, 15, 17, 25, 32]. However, none of these designs considers today's small message sizes; they are based on heavy-tailed workloads where 100 Kbyte messages are considered "small," and latencies are often measured in milliseconds, not microseconds. As a result, there is still no practical solution that provides near-hardware latencies for short nessages under high network loads. For example, we know of no existing implementation with tail latencies of 100 µs or less at high network load (within 20x of the hardware potential). Homa is a new transport protocol designed for small messages in low-latency datacenter environments. Our implementation of Homa achieves 99th percentile round trip latencies less than 15 µs for small messages at 80% network load with 10 Gbps link speeds, and it does this even in the presence of competing

HUIT (ECN-based)

DCOCN (Explicit feedback)

(Packet scheduling)

Homa (Credit-based)

Next week: software defined networking

How do we manage a complex network?

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



