Data Center TCP (DCTCP)



Presented by Zeyuan Hu



Measure and analyze production traffic Data Center TCP (DCTCP)



Measure and analyze production traffic



Measure and analyze production traffic

Case study: Microsoft Bing

- Measurements from 6000 server production cluster
- passively collects socket level logs, selected packet-level logs, and app-level logs describing latencies
- More than **I50TB** of compressed data over a month
- 99.91% of traffic in the Microsoft data center is TCP traffic
- Workloads
 - Partition/Aggregate [2KB 20KB] (Query) delay-sensitive

 - Large flows [IMB 50MB] (Data update) + throughput-sensitive



CS 395T – Data Centers @ UTCS Fall 2018



Part of Findings - Three Impairments





Incast





Queue Buildup





 The long, greedy TCP flows build up queues on their interfaces. Since buffer space is a shared resource (shallow buffered switches), the queue build up reduces the amount of buffer space available to absorb bursts of traffic from Partition/Aggregate traffic. — Packet loss & timeouts



Data Center Transport Requirements

- Low latency for short flows
- High burst tolerance (handle incast problem due to Partition/ Aggregation)
- High throughput for long flows
- Switch buffer occupancies need to be persistently low, while maintaining high throughput for the long flows



TCP in the Data Center

• TCP does not meet demands of apps.

• Incast

➢ Suffers from bursty pacl

- Builds up large queues:
 - > Adds significant latency.
 - ➤ Wastes precious buffer
- Operators work aroun
 - Ad-hoc, inefficient, ofter





DCTCP Algorithm



Congestion Control

• Sliding window: A flow control technique

- A generalization of stop-and-wait
- Allow up to W unack packets in flight at any time (W = window size)

• Congestion: Different sources compete for resources in the network

- Flows using up all link capacity $\Sigma r > C$
- Short flows compete with large flows on buffers

• Ways to perform congestion control:

- End-hosts (e.g.TCP congestion control)
- Network-based (e.g. ECN)



TCP Congestion Control

- TCP varies the number of outstanding packets in the network by varying the window size
 - Window size = min (Advertised Window, Congestion Window)
 - Congestion Window is denoted as "cwnd"
 - Packet dropped
 congestion
- How do we set cwnd? AIMD
 - Additive Increase, Multiplicative Decrease

If packet received OK:
$$W \leftarrow W + \frac{1}{W}$$

If a packet is dropped: $W \leftarrow \frac{W}{2}$





Review: The TCP/ECN Control Loop





Data Center Environment

- Low round trip times (less than 250 microseconds)
- Little Statistical Multiplexing
- Network is homogeneous
- A single administrative control
- Separate from external traffic



Two Categories of Congestion Control

- Delay-based protocols use increases in RTT measurements as a sign of growing queuing delay, and hence of congestion
 - Rely heavily on accurate RTT measurement susceptible to noise
- Active Queue Management (AQM) approaches use explicit feedback from congested switches
 - DCTCP
 - RED (randomly early marking): RED monitors the average queue size marks packets based on statistical probabilities. If the buffer is almost empty, then all incoming packets are accepted. As the queue grows, the probability for dropping an incoming packet grows too. When the buffer is full, the probability has reached 1 and all incoming packets are dropped. \longrightarrow average queue size is too slow for bursty flow



Balance Between Requirements



8 |



Two Key Ideas

- I. React in proportion to the **extent** of congestion, not its **presence**.
 - ✓ Reduces variance in sending rates, lowering queuing requirements.

| ECN Marks | ТСР | DCTCP |
|------------|--------------------------------|-------------------|
| 1011110111 | Cut window by <mark>50%</mark> | Cut window by 40% |
| 000000001 | Cut window by <mark>50%</mark> | Cut window by 5% |

2. Mark based on **instantaneous** queue length.

 \checkmark Fast feedback to better deal with bursts.

From Balaji Prabhakar, Stanford University



Data Center TCP Algorithm

Switch side:

Mark packets with Congestion Experienced (CE) code point when Queue Length > K.

Receiver side:

• Use state machine to decide whether to set ECN-Echo flag





CS 395T – Data Centers @ UTCS Fall 2Adapted from Balaji Prabhakar, Stanford University



Data Center TCP Algorithm

Sender side:

– Maintain running average of *fraction* of packets marked (α).

In each RTT (i.e., once every W of data):

$$F = \frac{\# of marked ACKs}{Total \# of ACKs} \qquad \alpha \leftarrow (1 - g)\alpha + gF \qquad 0 < g < 1$$

> Adaptive window decreases: $Cwnd \leftarrow (1 - \frac{\alpha}{2})Cwnd$

Adapted from Balaji Prabhakar, Stanford University



DCTCP in Action



Figure 1: Queue length measured on a Broadcom Triumph switch. Two long flows are launched from distinct 1Gbps ports to a common 1Gbps port. Switch has dynamic memory management enabled, allowing flows to a common receiver to dynamically grab up to 700KB of buffer.



Why it Works

I.High Burst Tolerance

✓ Large buffer headroom → bursts fit.

✓ Aggressive marking → sources react before packets are dropped.

2. Low Latency

✓ Small buffer occupancies \rightarrow low queuing delay.

3. High Throughput

 \checkmark ECN averaging \rightarrow smooth rate adjustments, cwind low variance.

From Balaji Prabhakar, Stanford University



Analysis

• Assumptions:

- N infinitely long-lived flows with identical round-trip times RTT, sharing a single bottleneck link of capacity C.
- N flows are synchronized (i.e. their "sawtooth" window dynamics are in-phase).





Analysis

• Goals:

- Mathematically characterize the "sawtooth" by computing:
 - The maximum queue size Q_{max}
 - The amplitude of queue oscillations A
 - The period of oscillations $T_{m{C}}$





Figure 11: Window size of a single DCTCP sender, and the queue size process.



Compute A

- By assumption, A = ND, where D is the amplitude of oscillation in window size of a single flow
- $D = (W^* + 1) (W^* + 1)(1 \alpha/2).$ (7)
- \bullet W* is the window size at which the queue size reaches K
- $W^* = (C \times RTT + K)/N$
- Thus, to calculate A, we need to compute α
 - Key observation: the queue size exceeds K for exactly one RTT in each period of the "sawtooth", before the sources receive ECN marks and reduce their window sizes accordingly.
 - We can compute α (the fraction of marked packets) by :
 - the number of packets sent during the last RTT of the period / the total number of packets sent during a full period of the sawtooth



Compute A

$$S(W_1, W_2) = (W_2^2 - W_1^2)/2.$$
(4)

$$\alpha = S(W^*, W^* + 1) / S((W^* + 1)(1 - \alpha/2), W^* + 1).$$
 (5)

$$\alpha^{2}(1-\alpha/4) = (2W^{*}+1)/(W^{*}+1)^{2} \approx 2/W^{*},$$
 (6)

 $lpha pprox \sqrt{2/W^*}$



Compute A

$$A = ND = N(W^* + 1)\alpha/2 \approx \frac{N}{2}\sqrt{2W^*}$$
$$= \frac{1}{2}\sqrt{2N(C \times RTT + K)},$$
(8)



Compute T_C

$$T_C = D = \frac{1}{2}\sqrt{2(C \times RTT + K)/N} \quad \text{(in RTTs).} \tag{9}$$



Compute Q_{max}

(3)

$$Q(t) = NW(t) - C \times RTT,$$

• Queue size at time t = Arrival rate at t – Departure rate at t

$$Q_{max} = N(W^* + 1) - C \times RTT = K + N.$$
 (10)



How good is the Analysis?



Figure 12: Comparison between the queue size process predicted by the analysis with NS-2 simulations. The DCTCP parameters are set to K = 40 packets, and g = 1/16.



- Implemented in Windows stack.
- Real hardware, IGbps and IOGbps experiments
 - 90 server testbed
 - Broadcom Triumph 48 I G ports 4MB shared memory
 - Cisco Cat4948
 - 48 IG ports 16MB shared memory
 - Broadcom Scorpion 24 10G ports 4MB shared memory
- Numerous benchmarks
 - -Throughput and Queue Length
 - Multi-hop
 - Queue Buildup
 - Buffer Pressure

- Fairness and Convergence
- Incast
- Static vs Dynamic Buffer Mgmt

From Balaji Prabhakar, Stanford University



Microbenchmarks: Incast



Figure 18: DCTCP performs better than TCP and then converges at 35 senders (log scale on Y axis; 90% confidence intervals for the means are too small to be visible).





Figure 19: Many-to-one: With dynamic buffering, DCTCP does not suffer problems at even high load.







| | Without background traffic | With background traffic |
|-------|----------------------------|-------------------------|
| TCP | 9.87ms | 46.94ms |
| DCTCP | 9.17ms | 9.09ms |

Table 2: 95th percentile of query completion time. DCTCP prevents background traffic from affecting performance of query traffic. $RTO_{min} = 10ms$, K = 20.





From Balaji Prabhakar, Stanford University

37





✓ Low latency for short flows.

CS 395T – Data Centers @ UTCS Fall 2018

From Balaji Prabhakar, Stanford University

38





From Balaji Prabhakar, Stanford University

CS 395T – Data Centers @ UTCS Fall 2018







- 99.91% of data center traffic is TCP
- Throughput-sensitive large flows, delay sensitive short flows and bursty query traffic co-exists in a data center network
- Switch buffer occupancies need to be persistently low, while maintaining high throughput for the long flows
- Data Center Transport Requirements
 - Low latency for short flows
 - High burst tolerance
 - High utilization for long flows



Conclusions: DCTCP

• DCTCP satisfies all our requirements for Data Center packet transport.

- ✓ Handles bursts well
- \checkmark Keeps queuing delays low
- \checkmark Achieves high throughput

• Features:

- \checkmark Very simple change to TCP and a single switch parameter K.
- ✓ Based on ECN mechanisms already available in commodity switch.



- MMU (Memory Management Unit): perform translation of virtual memory address to physical address; effectively perform virtual memory management
- RTO (Recovery Time Objective): duration of time before packet retransmission (i.e. TCP minimum retransmission timeout)
- Statistical multiplexing: the idea of taking a single resource and sharing it across multiple users in a probabilistic or statistical
 - It's statistical in that each user receives a statistical share of the resource based on how much others are using it
 - For example, if your friend is reading, you can use all of the link. If both of you are loading a page, you receive half of the link capacity
 - Packet switching