

Enabling ECN for Datacenter Networks With RTT Variations

Junxue Zhang¹, Wei Bai, and Kai Chen²

Abstract—ECN has been widely employed in production datacenters to deliver high throughput low latency communications. Despite being successful, prior ECN-based transports have an important drawback: they adopt a fixed RTT value in calculating instantaneous ECN marking threshold while overlooking the RTT variations in practice. In this article, we reveal that the current practice of using a fixed high-percentile RTT for ECN threshold calculation can lead to persistent queue buildups, significantly increasing packet latency. On the other hand, directly adopting lower percentile RTTs results in throughput degradation. To handle the problem, we introduce ECN[‡], a simple yet effective solution to enable ECN for RTT variations. At its heart, ECN[‡] inherits the current instantaneous ECN marking (based on a high-percentile RTT) to achieve high throughput and burst tolerance, while further marking packets (conservatively) upon detecting long-term queue buildups to eliminate unnecessary queueing delay without degrading throughput. We implement ECN[‡] on a Barefoot Tofino switch and evaluate it through extensive testbed experiments and large-scale simulations. Our evaluation confirms that ECN[‡] can effectively reduce latency without hurting throughput. For example, compared to the current practice, ECN[‡] achieves up to 23.4% (31.2%) lower average (99th percentile) flow completion time (FCT) for short flows while delivering similar FCT for large flows under production workloads.

Index Terms—Datacenters, ECN, RTT variations, AQM

1 INTRODUCTION

DATACENTER applications desire high throughput (e.g., data mining and storage) and low latency (e.g., web search and memory cache) communications. To achieve this, Explicit Congestion Notification (ECN) has been widely used, and many ECN-based transports, such as DCTCP [1] and DCQCN [2], have been proposed and adopted by industry [3], [4], [5], [6].

ECN-based transports require both ECN-aware rate control at the end host and ECN marking at the switch. ECN marking is accomplished by an active queue management (AQM) policy. To handle transient bursts which are common and harmful in production environments [7], datacenter AQM solutions usually mark packets aggressively based on instantaneous congestion states. For example, DCTCP [1] modifies the original RED¹ [8] to leverage the instantaneous queue length to mark packets, thus achieving good burst tolerance.

1. We call this version of RED as DCTCP-RED in this article.

- Junxue Zhang and Kai Chen are with the iSING Lab, Hong Kong University of Science and Technology, Hong Kong. E-mail: {jzhangcs, kaichen}@cse.ust.hk.
- Wei Bai is with the Microsoft Research Lab, Redmond, WA 98052 USA. E-mail: baiwei0427@gmail.com.

Manuscript received 8 February 2022; revised 29 May 2022; accepted 5 September 2022. Date of publication 8 September 2022; date of current version 6 September 2023.

This work was supported in part by the Key-Area Research and Development Program of Guangdong Province under Grant 2021B0101400001, in part by Hong Kong RGC TRS under Grants T41-603/20-R, GRF-16215119, and GRF-16213621, and in part by NSFC under Grant 62062005.

(Corresponding author: Kai Chen.)

Recommended for acceptance by A. Louri.

Digital Object Identifier no. 10.1109/TCC.2022.3204988

For instantaneous ECN marking, the choice of the marking threshold is critical as it directly affects the tradeoff between throughput and latency [1], [9]. Despite being successful, prior ECN-based transports have shared an important drawback: they only adopt a fixed RTT² value in calculating the marking threshold without considering RTT variations.

However, large RTT variations are common in datacenters as different flows traverse different processing components, e.g., network stack, hypervisor and middlebox. For example, compared to intra-service flows, inter-service flows experience large extra processing delays from the layer 4 load balancer [10], [11]. Furthermore, the processing delay of a given component also varies depending on workloads. With a very simple testbed (Section 2.2), we show these factors can easily lead to $\sim 3\times$ RTT variations, which also implies even larger RTT variations in real production.

We point out, through testbed experiments, that the current practice of using a fixed high-percentile RTT for ECN threshold calculation can lead to persistent queue buildup, significantly increasing packet latency. For example, our experiment results show that using 90th percentile RTT (as suggested by [9]) can achieve high throughput, but the latency of short flows is increased by over 50% due to queueing delay (Section 2.3). On the other hand, directly adopting lower percentile RTTs does not solve the problem either, as it results in throughput degradation.

Motivated by this, we seek an ECN marking scheme that can deliver high throughput and low latency simultaneously, in the presence of high RTT variations. To this end, we present ECN[‡] (Section 3), a simple yet effective solution

2. In this article, we use RTT to denote the base RTT without the queueing delay in datacenters.

to achieve our goal. At its heart, ECN[#] marks packets based on both the instantaneous and persistent congestion states. On the one hand, ECN[#] inherits the advantage of current instantaneous ECN marking based on a high-percentile RTT to achieve high throughput and burst tolerance. On the other hand, to eliminate the persistent queueing caused by flows with small RTTs, ECN[#] further marks packets when observing long-term switch queue buildups. This is performed conservatively in order not to affect throughput (Section 3.2).

By nature, ECN[#] works with both queue length and sojourn time (i.e., the amount of time a packet spends in the queue) as congestion signals. In our implementation we adopt sojourn time in order to be compatible with packet scheduler with traffic dynamics [12], [13].

We implement an ECN[#] prototype with ~500 lines of P4 code and compile it to a Barefoot Tofino [14] switch using Barefoot Capilano SDE [15] (Section 4). Our implementation uses little switch resource, e.g., 5 32-bit register arrays and 2 64-bit register arrays. We address three practical challenges in the implementation: 1) To ensure the correctness, we use emulation to get a 32-bit microsecond-granularity system time; 2) To update switch states at line rate, we leverage match action tables to implement the complex control flow; 3) To implement the division function, we use approximation with math unit.

We build a small 10Gbps testbed with 8 servers connected to the above Tofino switch. Our experiments with realistic workloads [1], [16] show that, by eliminating persistent queue buildups, ECN[#] can deliver up to 23.4% lower average FCT and 31.2% lower 99th percentile FCT for short flows compared to DCTCP-RED (where the ECN marking threshold is based on a high-percentile RTT according to current practice) while still maintaining comparable throughput for large flows. Furthermore, we find that instantaneous ECN marking is integral to ECN[#] to tolerate bursts and reduce packet drops. For example, ECN[#] outperforms CoDel [13] (which marks packets only based on persistent queue buildups instead of instantaneous queueing) by more than 50.0% FCT reduction for short flows.

To complement small-scale testbed experiments, we further perform larger-scale simulations to deep-dive into ECN[#] (Section 5.3). Our simulation results further confirm the superior performance of ECN[#]. For example, compared to DCTCP-RED, ECN[#] can achieve up to 31.3% lower FCT for short flows. Furthermore, from the microscopic view of switch queues, we show ECN[#] can effectively eliminate queue buildups by keeping the switch queue length 95.6% lower than that of DCTCP-RED (from 182 packets to 8 packets). Finally, we show ECN[#] 1) is robust to parameter settings and arbitrary packet schedulers and 2) can further benefit QCDQN [2] and PIAS [17].

To make our work easy to reproduce, we have made our simulation code available at <https://github.com/snowzjx/ns3-ecn-sharp>.

2 BACKGROUND AND PROBLEMS

2.1 Instantaneous ECN Marking

ECN-based transports consist of two parts: ECN marking at the switch and ECN-aware rate control at the end host. ECN marking is accomplished by an active queue management

(AQM) policy. To handle transient bursts which are common in production datacenters [7], most datacenter AQM solutions [1], [2], [12] aggressively mark packets based on the instantaneous congestion states, e.g., instantaneous queue length.

For instantaneous ECN marking, the choice of the marking threshold is important as it directly affects the tradeoff between throughput and latency [1], [9]. Most ECN-based datacenter congestion control algorithms [1], [9], [18] set the two thresholds of RED to the same value, $K_{min} = K_{max} = K$. Given the low statistical multiplexing of large flows in datacenter environments [1], to fully utilize link bandwidth while delivering low latency, the ideal ECN marking threshold K is given as follows [1], [9], [19]:

$$K = \lambda \times C \times RTT \quad (1)$$

To calculate and configure the ideal ECN marking threshold at the switch, datacenter operators must get values for three parameters: λ , C and RTT . λ is a parameter determined by the congestion control algorithm at the end host. Different congestion control algorithms have different ECN reaction mechanisms, resulting in different λ . For example, ECN^{*} (regular ECN-enabled TCP) cuts the window by half in the presence of ECN marks, thus having $\lambda = 1$. DCTCP reduces the window in proportion to the fraction of ECN marked packets, thus having a much smaller λ (0.17 in theory [20]). C is the bottleneck link capacity, which can be easily obtained for threshold calculation as datacenter operators have full knowledge of the network.

RTT is the base round-trip time and is the focus of this paper. Here the base RTT does not include switch queueing delay. In current practice, people use a fixed RTT value for calculation, implicitly assuming that base RTTs in datacenters are relatively stable due to the small cable length (low propagation delay variations), e.g., 200-300 meters [21], and high link capacity (low transmission delay variations). However, we show this does not hold in reality (Section 2.2) and causes problems (Section 2.3).

2.2 RTT Variations in Datacenters

The base RTT consists of three parts: transmission delay, propagation delay, and processing delay. As discussed above, transmission delay and propagation delay are small inside datacenters. The transmission time of a 1.5KB packet on a 10Gbps link is only $\sim 1.2\mu s$. The propagation delay of a 1KM cable is only $3.3\mu s$. However, the processing delay can be as high as tens (or even hundreds) of μs , dominating the base RTT. Thus, the variation of the base RTT is actually caused by the variation of processing delay.

In fact, the processing delay in datacenters varies vastly as different flows traverse different processing components, e.g., network stack, hypervisor, and middlebox. The more processing components a flow traverse, the larger delay it experiences. For example, in our production datacenters, we use a layer-4 software load balancer (SLB) [10] to process and balance the inter-service traffic (e.g., traffic from compute to storage)³. Inter-service *inbound* traffic is first

3. A large portion of inter-service traffic stays within the same datacenter [10].

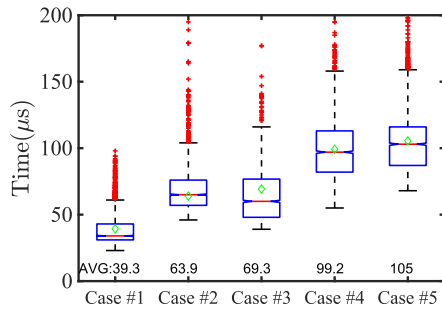


Fig. 1. [Testbed] RTT variations. Varying processing delay has enlarged RTT variations up to 2.68 times.

delivered to the SLB Multiplexer (a scalable set of dedicated servers), then forwarded to the application server, thus experiencing extra processing delay compared to intra-service traffic. Furthermore, the processing delay of a given component also varies depending on the loads.

We use testbed experiments to demonstrate the variations of processing delay. In our testbed, 3 hosts are connected to a Mellanox SN2100 switch via 100 Gbps links. Each server is equipped with a 100 Gigabit Ethernet Adapter. We enable DCTCP [1] at the end host. A host installed with APACHE HTTP SERVER works as the web server. We generate HTTP requests using APACHEBENCH on another host to fetch a 1-byte web page. On the client host, we install a modified TCP PROBE to record smoothed RTTs (SRTT) measured by Linux kernel. For each experiment, we generate 1000 requests and get ~ 3000 RTT samples. Note that a new request is sent when we receive the previous response. Since there is no congestion and hop difference, the RTT variation is caused by processing delay. We consider the following five cases and show corresponding results in Fig. 1 and Table 1.

1. *Network Stack*. In this experiment, the requests from the client are directly sent to the server. Therefore, processing delay mainly comes from network stacks of the client and server. As shown in Table 1, the average RTT and standard deviation are $39.3 \mu s$ and $12.2 \mu s$, respectively. We think the variation is due to multiple factors, e.g., kernel scheduling, advanced NIC techniques such as TSO, GSO, GRO and TCP delayed ACK mechanism.

2. *Network Stack + SLB*. In this experiment, we add a layer 4 SLB between the client and the server. SLB is widely used to process inter-service traffic inside the datacenters [10]. Our SLB is another host installed with LINUX VIRTUAL SERVER (LVS)[22]. The requests are first sent to the SLB. Then SLB forwards them to the backend server. The responses are directly returned to the client, without traversing the SLB. With SLB, the average RTT and standard deviation are increased to $63.9 \mu s$ (1.62X) and $18.3 \mu s$, respectively. The processing latency imposed by the SLB is $\sim 24.6 \mu s$.

3. *Network Stack + Hypervisor*. In this experiment, we use a quad-core virtual machine (VM) as the web server. Therefore, compared to the first experiment, request and response packets experience extra processing latencies on the hypervisor. On the hypervisor, we only add a necessary bridge rule without installing other complex rules, which is much simpler than the configuration in production environments. The average RTT and standard deviation are $69.3 \mu s$ (1.76X) and $18.8 \mu s$, respectively. The extra latency caused by the hypervisor is $\sim 30 \mu s$.

TABLE 1
[Testbed] RTT Statistics

Combination of different processing components	Mean (μs)	STD (μs)	90 percentile (μs)	99 percentile (μs)
Networking Stack	39.3	12.2	59.0	79.0
Networking Stack + SLB	63.9	18.3	87.0	121.0
Networking Stack + Hypervisor	69.3	18.8	91.0	130.0
Networking Stack + SLB + Hypervisor	99.2	23.0	129.0	161.0
Networking Stack(high load) + SLB + Hypervisor	105.5	23.6	138.0	178.0

Different combinations of processing components result in up to $2.68\times$ RTT variations in a very simple testbed setting, implying that the actual RTT variations in production datacenters would be even larger.

4. *Network Stack + SLB + Hypervisor*. In this experiment, we combine the above two experiments. With the SLB and hypervisor, the average RTT is increased from $39.3 \mu s$ to $99.2 \mu s$ (2.52X), compared to the first one.

5. *Network Stack (High Load) + SLB + Hypervisor*. In this experiment, to emulate the high load, we use STRESS to run 4 workers spinning on `sqrt()` on the web server VM. As a result, the average RTT reaches $105.5 \mu s$, which is around 2.68X larger than that in the first experiment.

Summary. We use the above simple experiments to demonstrate the causes of RTT variations in datacenters. Network components, such as Network Stack, SLB, Hypervisor, and *etc.*, add varying and unpredictable variations to the RTT, especially when they are under different loads. Furthermore, different flows may traverse different numbers of components, making RTT variations more severe. Note that all of the experiments are conducted in very simple settings, e.g., only a single bridge rule in the hypervisor, and very low loads (except for the last one). We believe that the actual RTT variations in production environment would be much larger due to complex configurations, high system loads, and more processing components. For example, Rohan Gandhi et al. [23] showed that the SLB in Microsoft datacenter adds a median latency of $196 \mu s$ to packets while the 90th percentile can reach as high as $1 ms$. Moreover, the base RTT varies across both flows and time. As the last two experiments show, the base RTT of a given flow also varies when the load of a processing component changes.

2.3 Current Practice and Problems

Existing solutions [1], [9], [12], [19] only adopt a fixed RTT value for ECN marking threshold calculation while overlooking RTT variations, which can cause severe problems. For example, operators get RTT distributions using tools such as PingMesh [24] and in current practice they use a high percentile RTT (e.g., 90th percentile [9]) to derive the threshold. While this approach can achieve high throughput for flows with different RTTs in general, for flows with small RTTs, it results in persistent queue buildups, causing queueing delay. On the other hand, directly using low percentile RTTs (e.g., 50%) or average RTT can mitigate such queueing latency, but cause flows with high RTTs to suffer from throughput degradation.

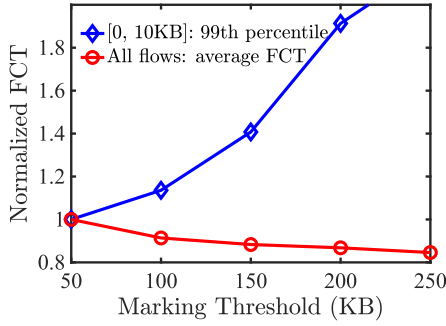


Fig. 2. [Testbed] Instantaneous marking cannot achieve high throughput and low latency simultaneously.

Since we only have three 100Gbps adapters, we build a new testbed to demonstrate this dilemma. We use 8 endhosts with 10Gbps Ethernet adapters, connecting to the Mellanox SN2100 switch. The cable speed is 10Gbps (with QSFP+ to SFP+ Adapter). We generate traffic according to the web search workload [1] (Fig. 7). The average load of the bottleneck link is 50%. To emulate RTT variations, we use netem [25]. netem is a Linux traffic control component which can emulate increased delay, packet loss, duplication, and *etc.* It allows adding delays in microseconds. We use netem to add extra delay at senders to emulate RTT variations.

Assume the minimum RTT and maximum RTT within the network are RTT_{min} and RTT_{max} respectively, we define RTT variation as RTT_{max}/RTT_{min} here.

Observation 1: Instantaneous ECN marking cannot achieve high throughput and low latency simultaneously. First, we let the RTT variation to be $3\times$ ($70\mu s$ to $210\mu s^4$) and choose different ECN marking thresholds. We change the marking threshold from 50KB to 250KB. The FCT statistics are shown in Fig. 2. The results have been normalized to the FCT achieved by the threshold 50KB. If we choose high percentile RTTs to derive the threshold, such as 90th percentile (250KB) according to the current practice, short flows suffer from 119.2% longer 99th percentile FCT ($581\mu s$ to $265\mu s$). On the contrary, setting threshold based on low percentile RTTs leads to throughput loss. For example, using average RTT ($\sim 100KB$) causes 8% throughput degradation ($3701\mu s$ to $3426\mu s$). Thus, none of the values in the range can simultaneously achieve high throughput and low latency.

Observation 2: Larger RTT variations enlarge performance loss. Second, we try to understand how the RTT variations affect the performance. We change the RTT variation from $2\times$ ($70\mu s$ to $140\mu s$) to $5\times$ ($70\mu s$ to $350\mu s$) and set the ECN marking threshold based on the average and 90th percentile RTTs accordingly. The FCT statistics are shown in Fig. 3, normalized to the result achieved by average RTT. When the variation is $2\times$, the threshold with the average RTT leads to 6.7% throughput loss ($3375\mu s$ to $3163\mu s$) compared to that with the 90th percentile. The gap increases to 29.8% ($4464\mu s$ to $3439\mu s$) when the RTT variation reaches $5\times$. Similarly, for 99th percentile FCT of short flows, the performance degradation caused by current practice has increased from 41.1% ($404\mu s$ to $287\mu s$) to 198.0% ($995\mu s$ to $334\mu s$) when the RTT variations rise from $2\times$ to $5\times$. Thus, larger RTT variation has enlarged the performance degradation.

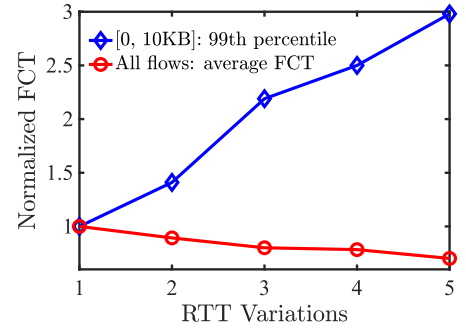


Fig. 3. [Testbed] Larger RTT variations cause more degradation to performance.

Problem Analysis. In this section, we try to analyze the root cause of the above-mentioned throughput loss or enlarged queueing delay caused by RTT variations. As shown in Fig. 4a, if there are no RTT variations and we choose a proper ECN threshold K , the buffer occupancy can achieve (1) high throughput, i.e., the buffer is always busy, and (2) low latency, i.e., there is no extra queue build-ups. However, if RTT variations exist, we will encounter the following problems:

- As shown in Fig. 4b, the ideal marking threshold for flows of large RTTs is K' . If we set the ECN threshold K based on the low percentile RTTs, $K < K'$. Consequently, if we use K instead of K' to mark ECN for these flows, these flows will suffer from throughput loss because a small K cause the queue to become idle from time to time.
- Similarly, as shown in Fig. 4c, the ideal marking threshold for flows of large RTTs is K' . If we set the threshold K based on the high percentile RTTs, $K > K'$. We can observe that these flows suffer from enlarged latency because a large K cause extra buffer occupancy, which does not contribute to the throughput but only increases the latency.

Summary. Instantaneous ECN marking cannot deliver both high throughput and low queueing delay under RTT variations. The performance degrades more as the RTT variation grows.

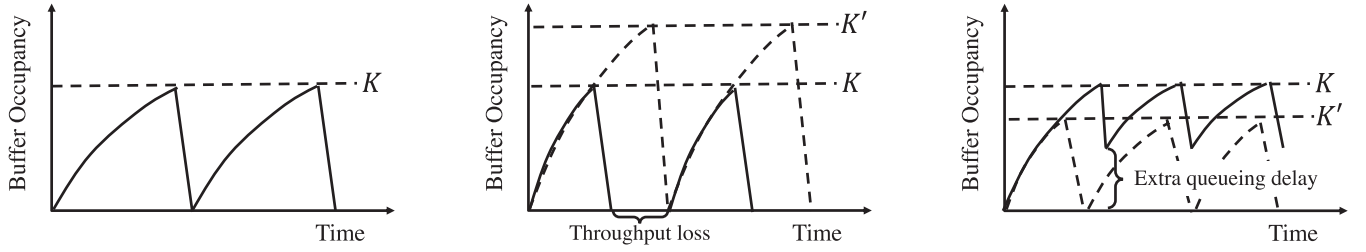
3 ECN[#]

3.1 Design Choice

Given the above dilemma, we seek to develop a solution to achieve high throughput and low queueing delay in data-centers even with high RTT variations. In the meanwhile, the solution should be resilient to traffic burstiness.

As shown in Section 2.3, the crux of existing instantaneous ECN marking is that the configured static threshold cannot adapt to the base RTTs of active flows. Given the emerging reconfigurable switching chips [14], a straightforward solution is a radical switch AQM that tracks the base RTTs of active flows (both spatially and temporally) and dynamically adjusts the instantaneous threshold on each individual switch correspondingly. For example, we can use the average base RTTs of all active flows to compute the threshold in Equation (1). However, it is challenging to measure real-time base RTTs. The reason is that, to exclude the noise of queueing delay, we first need to measure real-time

4. The emulated RTT denotes end-to-end base RTT.



(a) Buffer occupancy without RTT variations.

(b) Flows of large RTTs suffer from throughput loss.

(c) Flows of small RTTs suffer from extra queueing delay.

Fig. 4. RTT variations cause either throughput loss or enlarged queueing delay. For Fig. 4b, the dash line indicates the ideal buffer occupancy for flows of large RTTs while the solid line indicates the buffer occupancy if the ECN threshold is set according to low percentile RTTs. For Fig. 4c, the dash line indicates the ideal buffer occupancy for flows of small RTTs while the solid line indicates the buffer occupancy if the ECN threshold is set according to high percentile RTTs.

processing delay of each component in both directions and embed this information into packet headers. This is a huge burden for operators as they need to modify all network processing components in datacenters.

Therefore, we take one step back and seek a lightweight solution. Our lesson learned from Section 2.3 shows that the current practice of using a high percentile RTT to derive the instantaneous threshold is able to achieve high throughput and burst tolerance, but suffers from persistent queues when the base RTTs of active flows are small. Inspired by this, our design choice is to inherit the current practice for throughput and burst tolerance, and augment it with a scheme to detect and eliminate persistent queues. As a result, we introduce ECN[‡]: 1) inherits instantaneous ECN marking based on a high percentile RTT to handle bursts and maintain high throughput; and 2) further enables ECN marking upon persistent queue buildups to eliminate unnecessary queueing delay without incurring throughput loss.

3.2 ECN[‡] Mechanism

ECN[‡] is a new AQM solution. Following the above rationale, it marks packets based on both instantaneous and persistent congestion states. A packet is marked when either one of the following conditions decides to mark it:

- When there is a large instantaneous queue, ECN[‡] aggressively marks packets to avoid buffer overflow.
- When there is a persistent queue, ECN[‡] conservatively marks packets to reduce queueing delay.

By nature, ECN[‡] works with both queue length and sojourn time (i.e., the amount of time a packet spends in the queue) as congestion signals. According to Equation (1), the equivalent sojourn time based ECN marking threshold is:

$$T = \frac{K}{C} = \frac{\lambda \times C \times RTT}{C} = \lambda \times RTT \quad (2)$$

In our design and implementation, we adopt sojourn time over queue length in order to be compatible with packet scheduler with traffic dynamics [12], [13].

ECN Marking Based on Instantaneous Queue. To achieve burst tolerance, ECN[‡] employs instantaneous ECN marking. When a packet dequeues, ECN[‡] compares the sojourn time of the packet with the instantaneous marking threshold, $T = ins_target$, which can be calculated based on Equation (2) with

high percentile RTTs. If the sojourn time exceeds the threshold, the packet gets ECN marked.

ECN Marking Based on Persistent Queue Buildups. First, ECN[‡] tracks the minimal queueing over an interval to detect persistent queue buildups. Then, it performs conservative marking to eliminates those queues. Here, two parameters are involved: 1) *pst_interval*: the time interval used to observe the queue before deciding if there is persistent queueing; 2) *pst_target*: the persistent queue threshold used to compare against the packet sojourn time to control the long-term queueing. The algorithm mainly contains the following two parts as illustrated in Algorithm 1. Table 2 shows all parameters and variables.

Persistent Queue Buildups Detection. ECN[‡] measures whether the sojourn time has exceeded *pst_target* for at least one *pst_interval* to determine long-term queueing. ECN[‡] only uses one variable *first_above_time* to record the timestamp that sojourn time exceeds *pst_target* for the first time (*first_above_time* = now()). If sojourn time is lower than *pst_target*, which means the queue expires, ECN[‡] will reset *first_above_time* to 0. If ECN[‡] finds it has been over one *pst_interval* since *first_above_time* (*now()* > *first_above_time* + *pst_interval*), it confirms there is persistent queueing and triggers the conservative marking logic.

Conservative Marking. After detecting the persistent queue buildups, ECN[‡] conservatively marks one packet in each *pst_interval* in order not to adversely affect throughput. In the meanwhile, *pst_interval* is a time-varying interval that adapts to real workloads. If the sojourn times of packets continuously exceed the threshold, ECN[‡] reduces the interval of next

TABLE 2
Parameters and Variables Used in ECN[‡]

Configuration Parameters	
<i>pst_target</i>	persistent queueing target
<i>pst_interval</i>	interval to determine persistent queueing
<i>ins_target</i>	threshold for instantaneous marking
Variables	
<i>sojourn_time</i>	sojourn time of the current packet
<i>detected</i>	if persistent queueing is detected
<i>marking_state</i>	if ECN [‡] has marked packets
<i>marking_count</i>	number of marked packets
<i>marking_next</i>	the next marking time
<i>first_above_time</i>	when sojourn time starts to exceeds <i>pst_target</i>

marking to increase the marking probability ($\text{marking_next} += \text{pst_interval} / \text{marking_count}$), thus reducing the buffer occupancy.

Algorithm 1. ECN Marking Upon Persistent Congestion

```

input : Coming packet  $\text{pkt}$ 
output: Whether the  $\text{pkt}$  should be marked
1: Procedure ShouldPersistentMark ( $\text{pkt}$ )
2:   detected = IsPersistentQueueBuildups ( $\text{pkt}$ )
3:   if marking_state == true then
4:     if detected == false then
5:       marking_state = false
6:       return false
7:     else if now() > marking_next then
8:       marking_count ++
9:       marking_next +=  $\text{pst\_interval} / \text{marking\_count}$ 
10:      return true
11:    end
12:  else
13:    if detected == true then
14:      marking_state = true
15:      marking_count = 1
16:      marking_next = now() +  $\text{pst\_interval}$ 
17:      return true
18:    end
19:    return false
20:  end
21: Procedure IsPersistentQueueBuildups ( $\text{pkt}$ )
22:  if  $\text{pkt.sojourn\_time} < \text{pst\_target}$  then
23:    first_above_time = 0
24:    return false
25:  end
26:  if first_above_time == 0 then
27:    first_above_time = now()
28:    return false
29:  else if now() > first_above_time +  $\text{pst\_interval}$  then
30:    return true
31:  else
32:    return false
33:  end
  
```

3.3 Why ECN[#] Works?

We discuss why ECN[#] can achieve all the design goals.

High Throughput & Low Queueing Delay. Instead of directly using a low instantaneous ECN marking threshold to aggressively bring down the buffer occupancy at the cost of throughput loss, ECN[#] uses a more conservative mechanism. By tracking the minimal queueing over a time interval, ECN[#] only marks packets after it determines persistent queue buildups. These queue buildups are caused by providing excessive buffer space for flows with small RTTs, which do not contribute to the throughput but only increase the queueing delay. Eliminating these queues can benefit latency-sensitive short flows without throughput degradation. We carefully choose the interval to be around one worst case base RTT because TCP needs one RTT to react to ECN marking. Furthermore, the marking itself is conducted in a conservative way to mitigate throughput loss. Only after detecting persistent queue buildups, ECN[#] marks one packet and schedules next marking after an interval. Afterwards, ECN[#] checks again to see

whether the queue expires or not. If not, ECN[#] increases the marking frequency by reducing the interval. Through both persistent queueing detection and conservative ECN marking, ECN[#] can achieve both high throughput and low queueing delay simultaneously.

Burst Tolerance. ECN[#] leverages the instantaneous ECN marking to achieve good burst tolerance. As discussed above, the instantaneous marking scheme enables fast response to bursty traffic, such as incast. Once the queue exceeds the threshold, it aggressively marks the packets to tame the subsequent bursts from senders and bring down the buffer usage. This scheme ensures an upper bound of buffer occupancy to avoid packet loss. Furthermore, ECN[#]'s instantaneous marking threshold is derived based on a high percentile RTT, which provides sufficient buffer without causing throughput loss.

3.4 Parameter Setting

ECN[#] has 3 key parameters: ins_threshold , pst_target and pst_interval . We provide a rule-of-thumb to tune those parameters while leaving the optimality analysis as future work. We also conduct parameter sensitivity analysis in Section 5.4.2 to show that ECN[#] can achieve good performance by following the rule-of-thumb.

The instantaneous marking threshold, ins_threshold , can be calculated based on Equation (2). As discussed in Section 2.3, we use high percentile RTTs (e.g., 90th percentile) to derive a high instantaneous marking threshold to avoid throughput loss.

Both pst_interval and pst_target are critical for identifying and mitigating persistent queue buildups. pst_interval determines the time that ECN[#] spends tracking a queue before it confirms there is a persistent queue buildup. It has been recommended that the interval has to be around one RTT in order to accurately determine the persistent queueing [13]. Thus, we set pst_interval to be around the high percentile RTT to ensure that ECN[#] can precisely detect the persistent queue buildups caused by flows with smaller RTTs. However, in real environment where traffic is very bursty, we can set a smaller value for pst_interval to increase the marking rate to achieve fast reaction to those dynamics.

For pst_target , theoretically, we can set a very small value for it. The reason is that by observing the queue for an interval, we can detect persistent queue buildups, and a lower target will further bring down those unnecessary queues. However, in datacenters, the queue oscillation is large due to some network settings (MTU settings, etc.) or NIC off-loading (GSO, LSO, etc.) [1]. Thus, we also recommend setting it to be a more conservative value, which is larger than or equal to $\lambda \times \overline{RTT}$, where λ is determined by the transport protocol and \overline{RTT} is the average RTT. This setting can reduce the buffer occupancy to an acceptable value without risking throughput loss.

3.5 Discussion

ECN[#] versus CoDel. ECN[#] is partially inspired by CoDel [13], an AQM solution to address the bufferbloat problem in Internet. CoDel uses the minimal queueing over an interval to distinguish good queues (queues that contribute to throughput) and bad queues (queues that do not contribute

to throughput, such as the persistent queueing mentioned above). After detecting a bad queue, CoDel gradually reduces the marking interval to mitigate persistent queue buildups while not degrading throughput. However, such conservative marking makes CoDel react slowly to transient bursts, resulting in excessive packet losses. In contrast, ECN[‡] leverages the instantaneous ECN marking to keep good burst tolerance. Hence, ECN[‡] can greatly outperform CoDel in incast scenarios. We further analyze it through simulations in Section 5.4.1.

Sojourn Time versus Queue Length. As discussed above, in ECN[‡]'s design, we choose sojourn time instead of queue length as the congestion signal. CoDel is the first to use sojourn time as the congestion signal because CoDel is designed for the Internet, where lots of networking metrics, such as link bandwidth, queueing, are varying and unknown to the user. As a result, by using the sojourn time that directly reflect the time a packet stays in a queue, CoDel is much simpler to configure, making it a viable solution for Internet congestion control.

In contrast, in a datacenter, the networking metrics are relatively fixed and known to the datacenter operators. Thus, the sojourn time in datacenter networking is mainly used to solve the multiple queueing scenarios, as proposed in TCN [12]. In other words, in a single queue scenario, ECN[‡] is compatible with both sojourn time and queue length as its congestion signal. We have also implemented a version of ECN[‡] that uses queue length as congestion signal and achieves identical results as sojourn time in a single queue environment. Because the measurement of sojourn time requires attaching timestamp, which may not be compatible with some commodity switches, we can use ECN[‡] with queue length instead to enable ECN with RTT variations.

ECN[‡] versus TCN. TCN [12] combines the sojourn time and instantaneous ECN marking. However, under high RTT variations, it is challenging to decide a proper sojourn time threshold for TCN. Current practice with high percentile RTTs leads to persistent queue buildups. In contrast, ECN[‡] can conservatively mark packets based on persistent congestion state, thus delivering low latency.

Probabilistic Instantaneous Marking. Some transport protocols, e.g. DCQCN [2], need probabilistic instantaneous marking to ensure fairness. Two instantaneous marking thresholds, K_{min} and K_{max} are used, packets are marked with probabilities from 0 to 1 when queue length resides from K_{min} to K_{max} , instead of "cut-off" behaviors (mark all packets if queue length exceeds a single threshold). One possible solution for ECN[‡] to work with probabilistic marking is to change the original cut-off marking into probabilistic marking, and keep the marking based on persistent congestion unchanged since it is conducted in a probabilistic way. In this paper, by following the above-mentioned approach, we have provided a ECN[‡] + DCQCN prototype in our simulator and evaluated its performance in Section 5.

Flow Scheduling. There are various flow scheduling mechanisms, such as PIAS [17], pFabric [26], QJUMP [27], etc., targeting at optimizing short flows as well. These works leverage switch queue(s) to enforce scheduling mechanisms for flows, such as shortest-job-first, jump-the-queue, etc. Since ECN[‡] targets at eliminating the queue buildups caused by RTT

variations, it can further benefit these flow scheduling mechanisms. For example, PIAS uses multiple priority queues to emulate a shortest-job-first scheduling without knowing the flow size and further adopts ECN to keep these queues short. However, the ECN configuration in PIAS also suffers from RTT variation, leading to degraded performance of short flows and compromising the overall performance gain. Therefore, applying ECN[‡] on these queues should further improve PIAS with RTT variations. We have provided a prototype of ECN[‡] + PIAS in our simulator and evaluated its performance in Section 5.

4 IMPLEMENTATION

We have implemented an ECN[‡] prototype in the egress pipeline of Barefoot Tofino switch [14]. ECN[‡] is essentially a *stateful* data plane algorithm as the per-packet processing involves updates of multiple switch states. Our implementation has around 500 lines of P4 [28] code and is compiled to Barefoot Tofino switch using Barefoot Capilano SDE 6.1.1 [15].

Resource Requirement. In our implementation, we use 7 match action tables in total. The number of table entries is less than 10 (most of match action tables use default actions, so no entry is needed explicitly). To store switch updates on all the 128 ports, we use 5 32-bit register arrays and 2 64-bit register arrays. The register memory consumption is only ~37KB. Each packet needs 124-bit metadata. Our ECN[‡] prototype only uses little switch resource, leaving more for other network functions.

We meet two main challenges in our implementation: 1) how to obtain the precise system time, and 2) how to update multiple switch states at line rate. In the following, we will describe our efforts to address them.

4.1 Emulate Precise System Time

As shown in Algorithm 1, ECN[‡] requires current system time to determine persistent queue buildups and mark packets. In the egress pipeline, each packet carries a metadata `egress_global_tstamp`, which has 64 bits and expresses current egress pipeline time in nanoseconds. However, we cannot directly use this 64-bit metadata. This is because we need to compare the current time with other states, e.g. `first_above_time`, using algorithm logical units (ALUs), but ALUs in Tofino only accept 32-bit input. Thus, we need an approximation method to derive a 32-bit system time. To solve this problem, there are several potential solutions:

- Use the lower 32 bits of `egress_global_tstamp`. The time is accurate enough but wrap around in every ~4 seconds (2^{32} nanoseconds). When the time wraps around, we will have a very small *marking_next* (line 16 of Algorithm 1), resulting in aggressive marks. A serious performance degradation in every 4 seconds is unacceptable.
- Use the higher 32 bits of `egress_global_tstamp`. The time is not accurate enough (around 4-second granularity) for low latency datacenters.
- Right shift `egress_global_tstamp` by 10 bits and then get lower 32 bits. The time is accurate enough (microsecond granularity) for datacenters and only

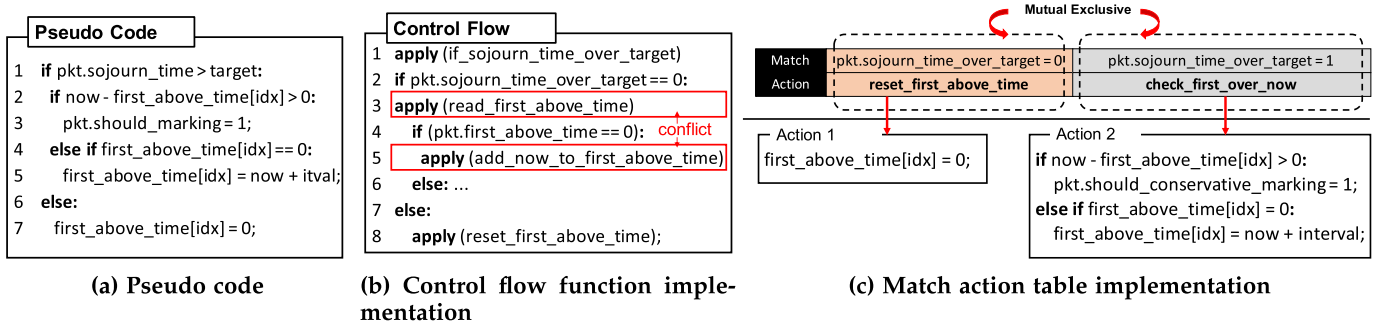


Fig. 5. Different control flow implementations. Control flow implementation that directly interprets the pseudo code into control flow cannot be compiled to Tofino as it causes multiple accesses to the same register. Hence, we use match action table to optimize the ECN[‡] control flow instead.

wraps around in every 4295 seconds (2^{32} microseconds), which is more than 1 hour. However, the operation `shift_right` in Barefoot Capilano SDE can only support 32-bit input data.

Faced with above problems, we decide to *emulate* a 32-bit microsecond-granularity system time. Our key idea is given as follows. We first use the lower 32 bits of `egress_global_timestamp`, then right shift the 32 bits by 10 bits to get a 22-bit microsecond-granularity time. For the remaining higher 10 bits, we increase it by 1 whenever we observe the lower 22 bits wrap around. Finally, we attach the emulated 32-bit microsecond-granularity time to the packet as a metadata. Algorithm 2 gives the pseudo code. In our implementation, we use two 32-bit registers to store values of the lower 22 bits and the higher 10 bits, respectively.

Algorithm 2. Emulate Precise System Time

Input: `pkt.egress_global_timestamp`
Output: `pkt.current_time`

```

1: tmp_timestamp = lower_32bits(pkt.egress_global_timestamp)
2: time_low = shift_right(tmp_timestamp, 10)
3: if time_low ≤ register_low then
4:   /* Wrap around */
5:   register_high = register_high + 1
6: end
7: register_low = time_low
8: pkt.current_time = register_high × 222 + register_low

```

4.2 Update Switch States at Line Rate

ECN[‡] keeps several switch states on the switch registers and updates them for every packet arrival. To achieve line rate, Tofino imposes a register update limit: a Tofino program can only access a register once (in a pipeline). Note that in Tofino, reading a register, comparing the register value with another value, and then updating the register correspondingly are also treated as one access. For example, the line 4 to 5 of Fig. 5a only has one access to `first_above_time[idx]`.

To implement the control flow of ECN[‡], we first directly use control flow function in P4 to apply different tables. A simple example is Fig. 5b. However, directly interpreting the logic of pseudo code into control flow function implementation causes problems, e.g., both tables `read_first_above_time` and `add_now_to_first_above_time` access the same register `first_above_time[idx]`. This

implementation cannot be compiled to the Tofino switch due to the hardware restrictions mentioned above. We notice that some reconfigurable hardware research work [29] also reports similar problem.

An alternative solution is to use `resubmit` primitive. The packet will be sent back to the beginning of the ingress pipeline with certain flags and regarded as a new packet. Then we can apply different tables based on the flags. This can guarantee the access to the register is mutually exclusive. However, resubmitting packets occupies more pipeline resources and largely degrades the throughput.

To this end, we use the match action table to optimize the control flow of ECN[‡]. Given a register, we first ensure that one register only has one table with actions accessing it. In this table, different actions reflect different control flow paths and they are naturally mutual exclusive. We further ensure the register is only accessed once within one action. Based on the two rules, we optimize the control flow by breaking it down into multiple match action tables. Fig. 5c is one example. Before applying a match action table, we first compute the result of condition and attach the result to the packet as a metadata. Then we use the metadata to match the desired action, where the register will only be accessed once, and update switch states correspondingly. In our evaluation (Section 5.2), we confirm ECN[‡] can achieve line-rate performance.

4.3 Implement Division Function

ECN[‡] needs to conduct a conservative marking by calculating $\text{pst_interval} / \text{marking_count}$ to decide the next marking time. Thus, the marking will become frequent only if there are lots of packets whose queueing time exceeds pst_target . However, current Tofino switch does not natively support division operation. To solve this challenge, we alternatively use approximation method with math unit provided by the switch.

Math unit provides a lookup table and provides 16 8-bit values that can be looked up, denoted as Table. In our paper, we generate the Table as follows. We then set the inverse table to `math_unit_lookup_table` property of the Tofino math unit.

$$\text{Table}_i = \begin{cases} 0 & , 0 \leq i \leq 7 \\ \lceil \frac{\text{pst_interval}}{i} \rceil & , 8 \leq i \leq 15 \end{cases} \quad (3)$$

Given a `marking_count`, the math unit will perform the computation as shown in Algorithm 3. The `table_offset` is

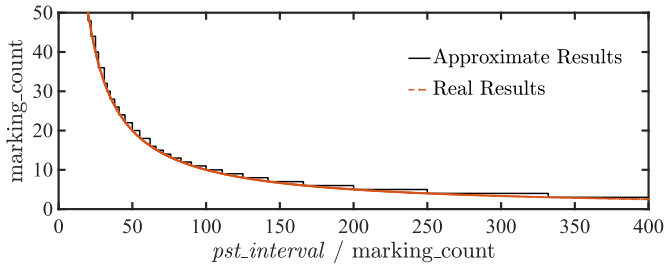


Fig. 6. Division function approximation can still preserve high accuracy and viable on current Tofino switches.

set to be 3 by the Tofino math unit. We will also set the scale at line 3 to the `math_unit_output_scale` property of the math unit which will be added at the exponent of the output (indicated by line 10) to scale the return results.

Algorithm 3. Approximate Division

Input : marking_count
Output: $pst_interval / marking_count$

- 1: $shift_val = \text{int}(\lfloor \log_2(x) \rfloor)$
- 2: $table_offset = 3$
- 3: $scale = \text{int}(\lceil \log_2(pst_interval / (8 * 127)) \rceil) + table_offset$
- 4: **if** $marking_count \geq 16$ **then**
- 5: $approx_val = marking_count \gg (shift_val - table_offset)$
- 6: **else**
- 7: $approx_val = marking_count \ll -(shift_val - table_offset)$
- 8: **end**
- 9: $ans = Table_{approx_val}$
- 10: **return** $ans * 2^{-(shift_val + scale)}$

We also compare our approximation method with real values to evaluate if our implementation can preserve high accuracy. Fig. 6 shows one concrete example where we set $pst_interval = 1000$. The results show that the values generated by our method overlaps with the real value. Furthermore, our experiment results in Section 5 also confirms that our approximation method is feasible and allows ECN[#] to work as expected.

5 EVALUATION

In this section, we conduct both testbed experiments and ns-3 simulations [30] to answer the following questions:

- *How does ECN[#] perform in practice?* Our testbed experiment results show that compared with current practice, ECN[#] achieves up to 23.4% smaller average FCT and 31.2% smaller 99th percentile FCT for small flows while providing the similar performance for large flows (Section 5.2). ECN[#]'s good burst control makes it largely outperform CoDel.
- *Does ECN[#] scale to large datacenters?* We evaluate the performance of ECN[#] using large scale ns-3 simulations (Section 5.3). Results show that ECN[#] achieves up to 36.9% smaller FCT for small flows compared to current practice.
- *How does ECN[#] control the queue length and react to bursty traffic?* We evaluate ECN[#] using simulations mixing both the background traffic and bursty query traffic (Section 5.4.1). We find that compared to current practice, ECN[#] keeps a much smaller

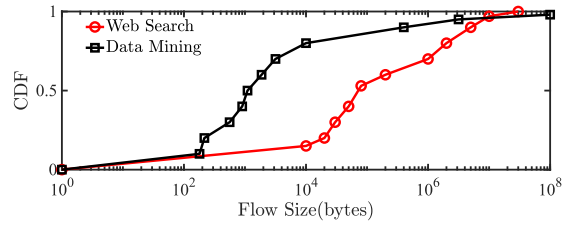


Fig. 7. Flow size distributions.

average queue occupancy (182 packets to 8 packets). By changing the fanout of the incast traffic, we further show that compared to CoDel, ECN[#] can support 1.75× more concurrent senders without TCP timeouts.

- *How robust is ECN[#] to network settings?* Through targeted simulations, we show that ECN[#] is robust to parameter choices (Section 5.4.2) and the packet scheduler (Section 5.4.3).

5.1 Methodology

Transport Protocol. We use DCTCP [1] at the end host by default. In testbed experiments, we use the DCTCP in Linux kernel 4.9.0 [3]. We also implement DCTCP in ns-3 simulator. The parameters are set as suggested in [1]. All endhosts in our evaluation are running DCTCP.

Schemes Compared. We compare ECN[#] against the following three schemes:

- **DCTCP-RED:** [testbed, simulation] We use DCTCP-RED to refer to the modified RED in DCTCP paper [1], i.e., instantaneous ECN marking based on single threshold $K_{min} = K_{max} = K$. For testbed, we implement it on Barefoot Tofino switch. For simulation, we start from the ns-3's RED implementation and add instantaneous ECN marking. We set the threshold based on both the current practice, 90th percentile RTT (denoted as DCTCP-RED-Tail), and average RTT (denoted as DCTCP-RED-AVG). Please note, when only one queue is active, the performance of DCTCP-RED is identical to TCN.
- **CoDel:** [testbed, simulation] CoDel tracks minimal queueing over an interval to mark packets based on persistent congestion. For testbed, we implement CoDel on Barefoot Tofino to perform ECN marking.
- **TCN:** [simulation] TCN uses instantaneous sojourn time to adapt to packet schedulers. We implement TCN based on the software prototype provided by TCN paper [31].

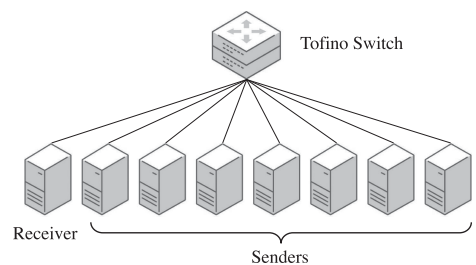


Fig. 8. Testbed setup.

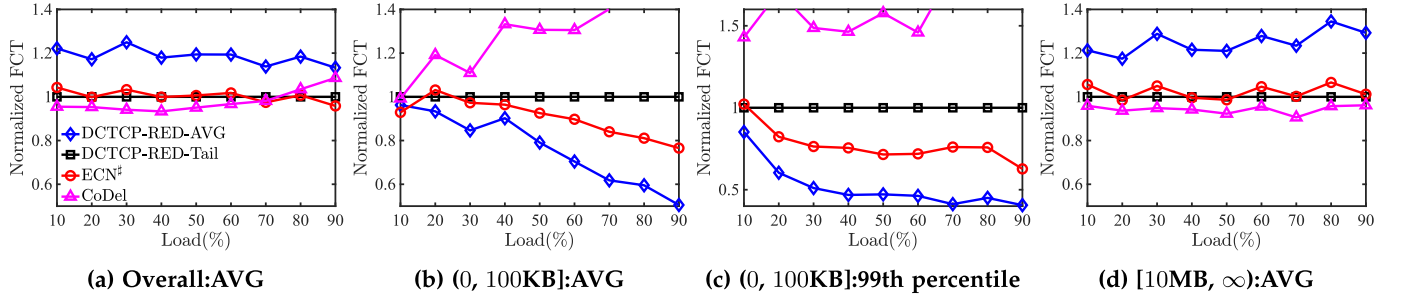


Fig. 9. [Testbed] FCT statistics with web search workload. All are normalized to the results achieved by DCTCP-RED with threshold derived from 90 percentile RTT.

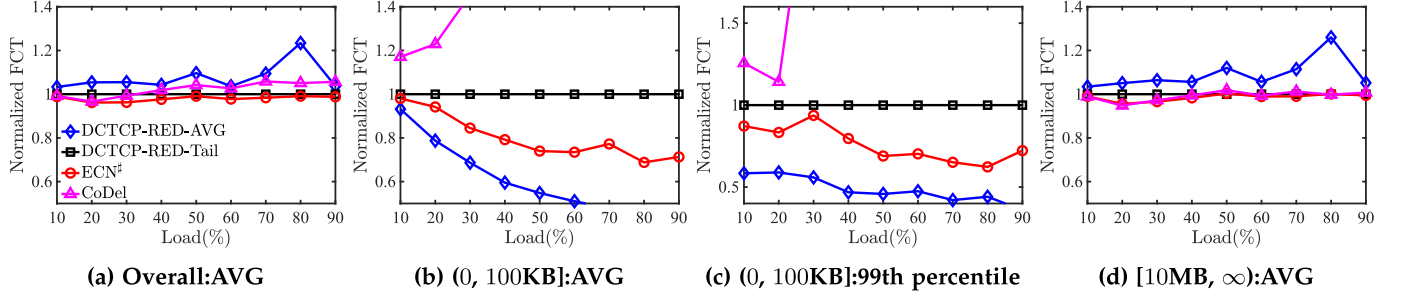


Fig. 10. [Testbed] FCT statistics with data mining workload. All are normalized to the results achieved by DCTCP-RED with threshold derived from 90 percentile RTT.

Instead of these schemes that we directly compare ECN# with, we also provide the following schemes that ECN# could work with:

- *ECN# + DCQCN*: [simulation] DCQCN [2] is a transport protocol, which uses probabilistic instantaneous marking to ensure fairness. As discussed in Section 3.5, we can combine ECN# with DCQCN by changing its instantaneous marking into probabilistic marking based on DCQCN's algorithm and leaving ECN marking upon persistent queueing unchanged. For DCQCN's probabilistic marking algorithm, we mainly refer to the simulator from DCQCN's paper.
- *ECN# + PIAS*: [simulation] PIAS [17] efficiently schedules flows without knowing flow sizes in advance. As discussed in Section 3.5, we can combine ECN# with PIAS to further benefit these flow scheduling mechanisms under the environment with high RTT variations.

Workloads. We generate traffic based on two realistic workloads in production: web search [1] and data mining [16]. The flow size distributions follow Fig. 7. Both workloads are heavy-tailed. In testbed experiments, we use an open source traffic generator [19], [32] to generate the benchmark traffic. Similar to previous work [17], [26], flows arrive according to a Poisson process to achieve the desired network utilization. We also use the same approach to generate traffic in simulations.

Metrics. We use the Flow Completion Time (FCT) as the primary metric. Besides the overall average FCT, we also breakdown FCT results across short flows ($< 100\text{KB}$) and large flows ($> 10\text{MB}$). We run experiments and simulations three times and report the average value.

5.2 Testbed Experiments

Testbed Setup. We use the testbed in Section 2.3 with 8 servers connected to a Barefoot Tofino switch with ECN#'s

implementation (Section 4). There are 7 senders and 1 receiver. The topology is also shown in Fig. 8 for better understanding.

We use netem to emulate a $3 \times$ RTT variations (from $70\mu\text{s}$ to $210\mu\text{s}$). The RTTs generated are based on the distribution in Fig. 1, which is long-tail distribution. For DCTCP-RED, we set the marking threshold of DCTCP-RED-AVG to 80KB and DCTCP-RED-Tail to 250KB . For CoDel, we set *interval* to be $200\mu\text{s}$ and *target* to be $85\mu\text{s}$. For ECN#, by following the rule-of-thumb (Section 3.4), we set the *ins_target* to be $200\mu\text{s}$, *pst_interval* to be $200\mu\text{s}$ and *pst_target* to be $85\mu\text{s}$.

5.2.1 Realistic Workloads

We evaluate ECN# using both web search and data mining workloads. The results are shown in Figs. 9 and 10. All results have been normalized to the result of current practice, DCTCP-RED-Tail. They are breakdown in terms of overall FCT (a), FCT of short flows (b, c) and FCT of large flows (d). All Figures show normalized FCT.

Short Flows. Compared to current practice, i.e., DCTCP-RED-Tail, ECN# can achieve up to 23.4% ($964\mu\text{s}$ to $738\mu\text{s}$ at 90% load) and 31.2% ($548\mu\text{s}$ to $377\mu\text{s}$ at 80% load) lower average FCT in web search and data mining workloads, respectively (Figs. 9b and 10b). At the 99th percentile, ECN# achieves up to 37.2% ($5242\mu\text{s}$ to $3287\mu\text{s}$ at 90% load) and 37.6% ($2161\mu\text{s}$ to $1347\mu\text{s}$ at 80% load) lower FCT with both workloads (Figs. 9c and 10c). These results show that, by well controlling the persistent queue buildups caused by RTT variations, ECN# can eliminate unnecessary queueing delay and effectively improve the performance of delay-sensitive short flows.

We observe that DCTCP-RED-AVG achieves the best performance for short flows, but it adversely hurts the FCT

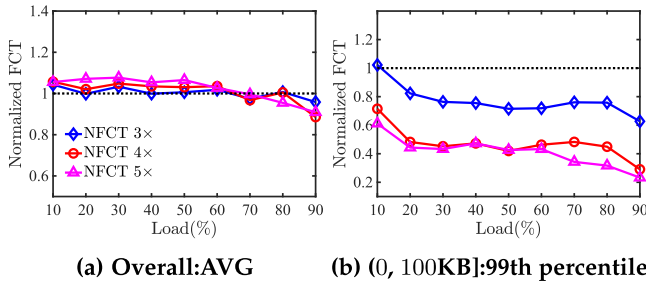


Fig. 11. [Testbed] FCT statistics with web search workload when RTT variation enlarges. NFCT $n\times$ denotes the normalized FCT of ECN[‡] to DCTCP-RED-Tail when RTT has $n\times$ variation. ECN[‡] achieves lower FCT for short flows and comparable throughput consistently.

of large flows, e.g., over 20% for the web search workload (Fig. 9d). This is because DCTCP-RED-AVG sets a lower ECN marking threshold, which strictly limits the queueing to a very low level, thus benefiting FCT of the short flows. In the meanwhile, such a low marking threshold throttles the throughput of large flows, affecting their FCT.

We also note that CoDel achieves very bad performance for short flows. The reason is that CoDel suffers from frequent TCP timeouts (1 TCP timeout adds > 1 milliseconds to FCT) because it does not react to the instantaneous queueing state, and thus experiences frequent packet loss under traffic bursts.

Large Flows. From Figs. 9d and 10d, we see ECN[‡] achieves similar results as DCTCP-RED-Tail and outperforms DCTCP-RED-AVG with both workloads. For web search workload, ECN[‡] achieves up to 25.6% (94411 μ s to 70192 μ s at 80% load) lower FCT compared to DCTCP-RED-AVG. For data mining workload, ECN[‡] achieves up to 20.5% (1049035 μ s to 833396 μ s at 80% load) lower FCT compared to DCTCP-RED-AVG. The reason is that by conservatively marking packets when observing persistent queue build-ups, ECN[‡] can bring down the buffer occupancy without hurting throughput for large flows.

Overall. In general, ECN[‡] achieves good overall performance among the four schemes. In data mining workload, ECN[‡] performs the best at all loads, whereas in web search workload its overall FCT degradation is within 4.2% (2240 μ s to 2148 μ s) compared to DCTCP-RED-Tail at 10% load (Fig. 9a). In contrast, DCTCP-RED-AVG achieves very bad performance. This further confirms that ECN[‡] can maintain the throughput of all flows.

5.2.2 Towards Larger RTT Variations

We increase the RTT variation from $3\times$ (70 μ s to 210 μ s) to $5\times$ (70 μ s to 350 μ s). We evaluate ECN[‡] using the more challenging web search workload, where the traffic is more bursty. We compare ECN[‡] with current practice, DCTCP-RED-Tail. The results are shown in Fig. 11. NFCT $n\times$ denotes the normalized FCT of ECN[‡] to DCTCP-RED-Tail when RTT has $n\times$ variation.

As shown in Fig. 11a, ECN[‡] generally achieves similar overall FCT as DCTCP-RED-Tail. When the RTT variation reaches $5\times$, the FCT degradation is within 7.6% (2854 μ s to 2650 μ s at 30% load). When the load is higher, ECN[‡] can achieve slightly better results than DCTCP-RED-Tail. The reason is that due to the conservative marking for persistent

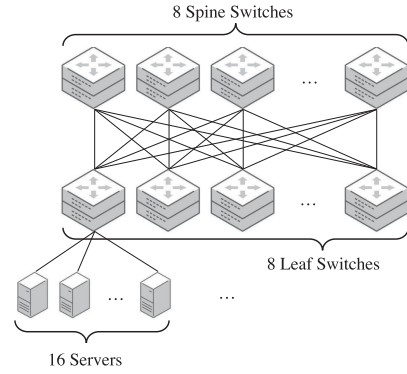


Fig. 12. Large-scale simulation setup.

congestion, ECN[‡] can consistently achieve comparable overall FCT no matter how RTT variation enlarges.

Furthermore, ECN[‡] greatly outperforms DCTCP-RED-Tail for short flows. As show in Fig. 11b, ECN[‡] can achieve up to 37.3% (5242 μ s to 3287 μ s at 90% load) lower FCT when the RTT variation is $3\times$. The advantage is enlarged to 71.2% and 73.4% when the variation is $4\times$ and $5\times$ respectively. This shows ECN[‡] can effectively reduce the FCT for short flows by controlling the persistent congestion. When the RTT variation is larger, ECN[‡] can bring more benefit to delay-sensitive short flows.

5.3 Large Scale Simulations

To complement the small-scale testbed, we further evaluate ECN[‡] on a larger-scale spine-leaf topology with production workloads (Section 5.3).

Setup. We simulate a 128-host leaf-spine topology with 8 spine and 8 leaf switches. Each leaf is connected to 16 servers via 10Gbps links. The spine and leaf switches are also connected via 10Gbps links. Fig. 12 shows the topology. We use ECMP for load balancing. The RTT has $3\times$ variations and varies from 80 μ s to 240 μ s. The average RTT here is $\sim 137\mu$ s and 90th percentile is $\sim 220\mu$ s. The distributions are similar to Fig. 1, which is a long-tail distribution.

In our simulation, the calculation of packets' sojourn time is implemented with packet tags in ns-3 [33]. When a packet enqueues, we add a packet tag with the enqueue timestamp to the packet. When the packet dequeues, we calculate the sojourn time by deducting the enqueue timestamp from the current timestamp.

We evaluate ECN[‡] using both web search and data mining workloads. The results are shown in Figs. 13 and 14. All results have been normalized to FCT achieved by DCTCP-RED-Tail.

Overall. With both workloads, ECN[‡] achieves similar results as DCTCP-RED-Tail, as shown in Figs. 13a and 14a. The results show that even in large-scale environment with high RTT variations, ECN[‡]'s marking mechanism does not hurt the overall networking performance.

Short Flows. ECN[‡] largely outperforms DCTCP-RED-Tail with both workloads. With web search workload, ECN[‡] achieves up to 23.4% better flow completion time (738 μ s to 964 μ s at 90% load), and achieves up to 31.3% better flow completion time with data mining workload (377 μ s to 548 μ s at 80% load). When considering 99th percentile short flows, the performance gain brought by ECN[‡] further

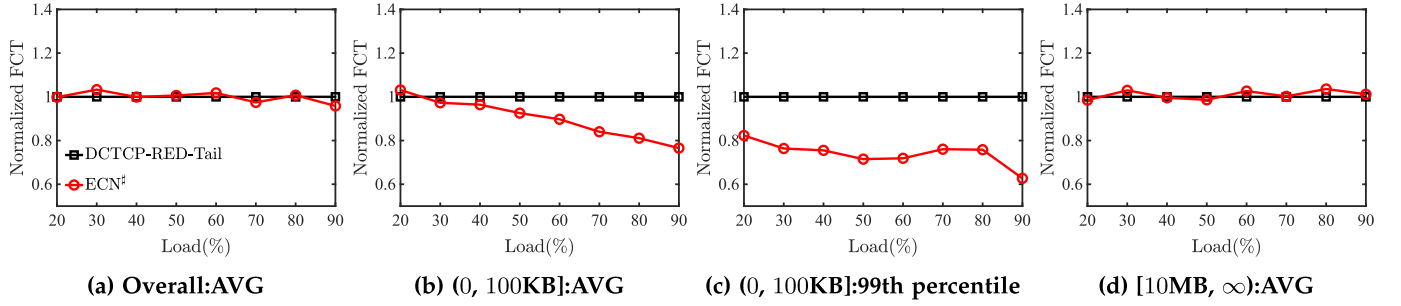


Fig. 13. [Simulations] FCT statistics with web search workload. All are normalized to the results achieved by DCTCP-RED with threshold derived from 90 percentile RTT.

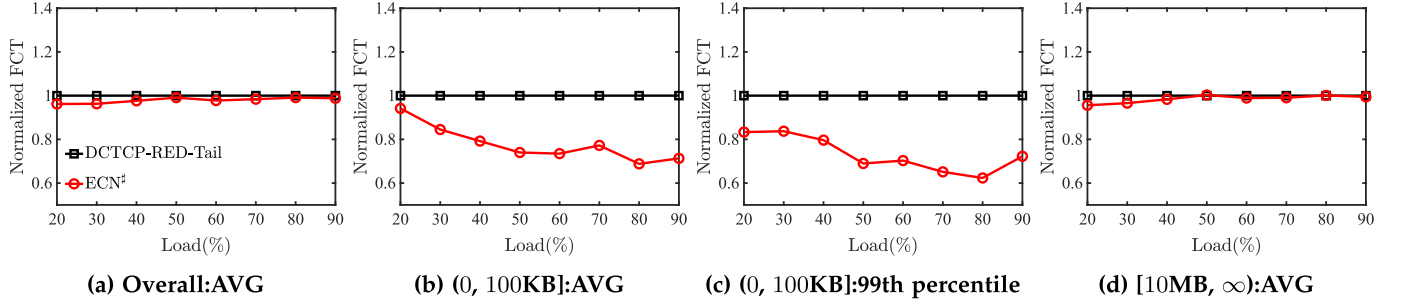


Fig. 14. [Simulations] FCT statistics with data mining workload. All are normalized to the results achieved by DCTCP-RED with threshold derived from 90 percentile RTT.

increases. ECN^{\sharp} outperforms DCTCP-RED-Tail by up to 37.3% ($3287\mu s$ to $5242\mu s$ at 90% load) with web search workload and 37.6% ($1347\mu s$ to $2161\mu s$ at 80% load) with data mining workload. The results are consistent with testbed results and show that ECN^{\sharp} can efficiently mitigate the queue buildup problems with high RTT variations and benefit the performance of short flows.

Large Flows. The performance of large flows are similar as the overall performance. Experiment results show that ECN^{\sharp} can well preserve the performance of throughput-sensitive large flows by using a threshold derived from a high RTT value.

5.4 ECN^{\sharp} Deep Dive

In this section, we first show a microscopic view of queues to explain the superior results achieved in both Sections 5.2 and 5.3. Then, we analyze ECN^{\sharp} 's parameter sensitivity. Finally, we compare ECN^{\sharp} with the latest datacenter AQM solution, TCN [12].

5.4.1 Microscopic View

Simulation Setup. As shown in Fig. 15, we use a simple 16 to 1 topology with 10Gbps links, 16 servers are senders and 1 receiver. Other settings are similar to Section 5.3. For CoDel, the *interval* is set to $240\mu s$ while *target* is set to $10\mu s$.

We send flows from 16 senders to the receiver. The size of both large and short flows are generated based on data mining workload. We also generate some query flows to emulate bursty incast traffic. The size of query flows follows a uniform distribution from 3KB to 60KB. We start N (100 in default) concurrent query flows at 4s. We further change the value of N to show how these bursty query flows affect the performance of all schemes.

Queue Occupancy. To better illustrate how different schemes manage their queues, we sample the queue length of the bottleneck link for 0.005 seconds and the results are shown in Fig. 16. Compared to other schemes, ECN^{\sharp} achieves two advantages:

- *Low persistent queueing:* ECN^{\sharp} keeps much lower queue occupancy (8 packets) compared to DCTCP-RED-Tail (182 packets). The low queue occupancy confirms ECN^{\sharp} 's effectiveness in handling persistent queue buildups caused by the mismatch between RTTs and marking threshold.
- *No packet drops with incast flows:* ECN^{\sharp} achieves comparable results in handling bursty query traffic (no packet drops) with DCTCP-RED-Tail, outperforming CoDel (drop 125 packets). The reason is that by leveraging the merit of instantaneous marking, ECN^{\sharp} can achieve good burst tolerance.

Impact of Bursty Query Flows. To have a deeper view on how those schemes handle bursty query flows, we change the concurrent query senders from 25 to 200 and trace the query completion time. The results are shown in Fig. 17. When there are 100 concurrent senders, CoDel begins to suffer from packet loss, seriously degrading

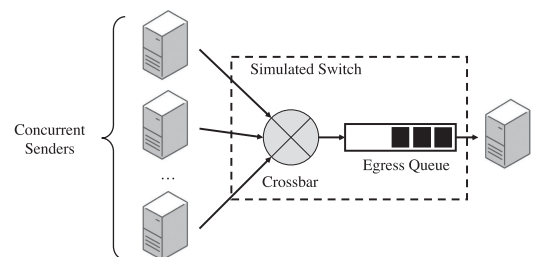


Fig. 15. Simulation setting for microscopic view.

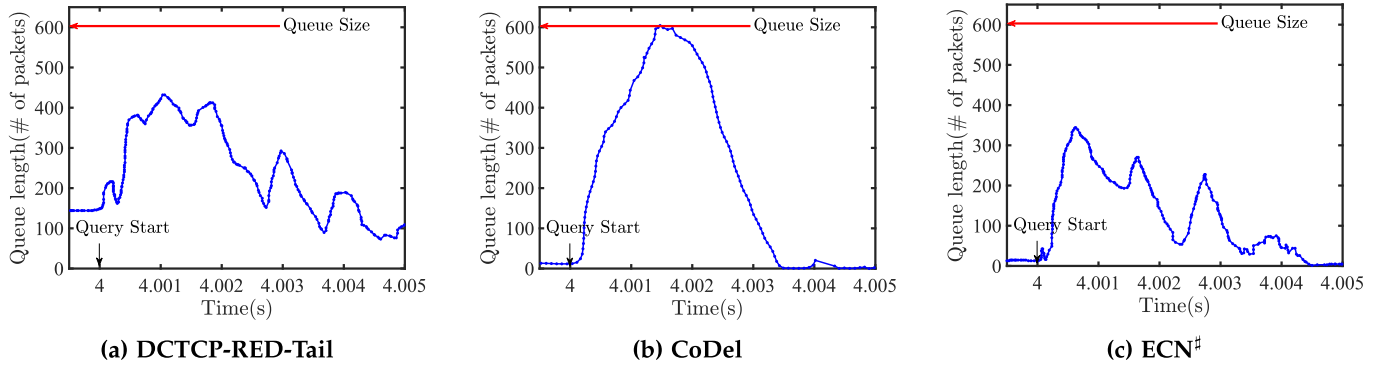


Fig. 16. [Simulations] Queue occupancy (with 100 concurrent query flows). ECN[#] achieves much lower queue occupancy (8 packets) compared to current practice, DCTCP-RED-Tail (182 packets). Meanwhile, ECN[#] achieves comparable performance in handling incast (no packets drop), large outperforms CoDel (drops 125 packets).

the performance. On the contrary, ECN[#] achieves comparable performance with DCTCP-RED-Tail, and begins to suffer from packet loss when the concurrent senders are increased to 175. In summary, due to its good burst tolerance, ECN[#] can support $1.75\times$ more concurrent senders than CoDel with bursty incast traffic. The reason is that ECN[#] leverages the instantaneous marking and has a good control over bursty traffic.

5.4.2 Parameter Sensitivity

We next show how robust ECN[#] is to different parameter settings. We set both *pst_target* and *pst_interval* based on the rule-of-thumb proposed in Section 3.4. We conduct analysis on both workloads. We use the setup in Section 5.3 and the results are shown in Fig. 18. We have the following two observations from the results:

- Different parameter settings do not have a large impact on the performance of ECN[#] considering the overall FCT. The difference is within 1% for web search workload and 0.2% for data mining workload.
- After reducing the value of *pst_interval*, ECN[#] behaves differently on the two workloads. For data mining workloads, ECN[#] performs slightly worse (0.2% longer FCT). However, for web search workload, ECN[#] performs 1% better. The reason is that web search workload is more bursty, and a more aggressive marking mechanism is better. We can reduce the *pst_interval* to increase the marking frequency.

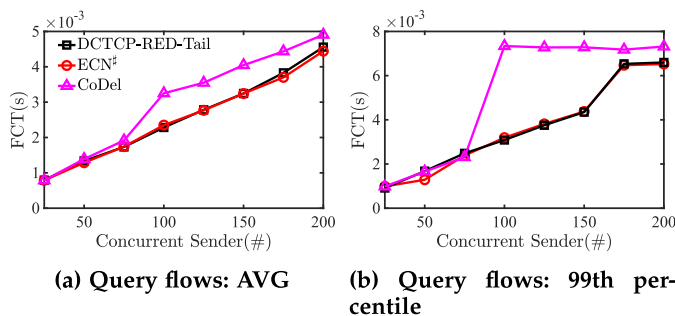


Fig. 17. [Simulations] Statistics of query flow completion time. ECN[#] can achieve comparable performance as DCTCP-RED-Tail and largely outperforms CoDel.

5.4.3 Packet Scheduler

In this section, we compare ECN[#] with TCN [12], the latest AQM solution to show how ECN[#] works with arbitrary packet scheduler. The experiment setting is shown in Fig. 19. We mainly focus on the following questions.

- Can ECN[#] preserve the packet scheduling policy?
- Can ECN[#] mitigate the persistent queueing exacerbated by RTT variations even with packet scheduler?

To evaluate this, we configure the switch with Deficit Weighted Round Robin (DWRR) with 3 queues/services. The weights among the 3 queues/services are 2 : 1 : 1. We first start a long-lived TCP flow from sender 1 and classify this flow into queue/service 1, then from sender 2 and classify into queue/service 2, finally from sender 3 into queue/service 3. We also randomly start short flows (the size is from 3KB to 60KB) from the rest of senders to probe the queue occupancies. We set the TCN marking threshold to $150\mu s$ to avoid throughput loss. Fig. 19 shows our settings.

Fig. 20a shows the goodput of flows achieved by ECN[#]. We observe at the beginning, only queue 1 is active, flow 1 achieves around ~ 9.6 Gbps goodput (Goodput is slightly smaller than throughput due to packet header overhead). After flow 2 starts, queue 2 becomes active and flow 1 achieves ~ 6.42 Gbps goodput while flow 2 achieves ~ 3.18 Gbps goodput. Finally, when three queues all become active, flow 1, flow 2 and flow 3 achieve around 4.82 Gbps, 2.40 Gbps and 2.40 Gbps goodput respectively, which strictly preserves the packet scheduling policy.

We also measure the FCT of short flows among all queues and the results are shown in Fig. 20b. Compared to

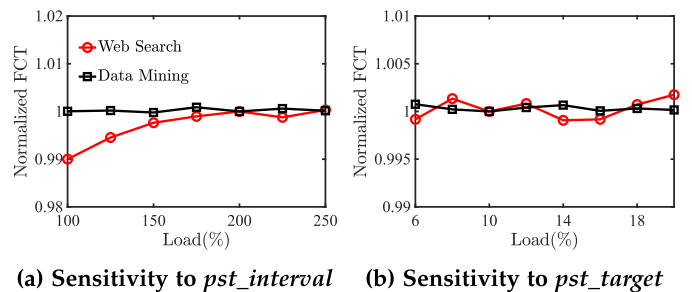


Fig. 18. [Simulations] Parameter Sensitivity. ECN[#] is robust to parameter choices with only $<1\%$ variations on overall FCT with different parameter settings.

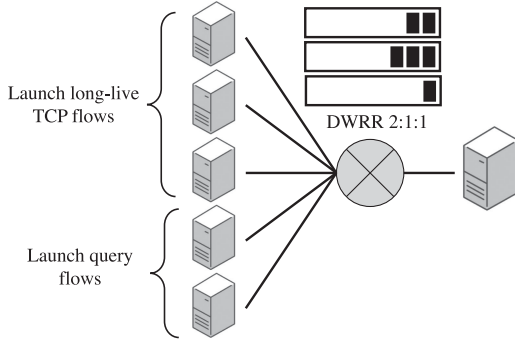


Fig. 19. Simulation setting for packet scheduler.

TCN, ECN[#] achieves 19.6% better average FCT ($2913\mu s$ to $2341\mu s$) for short flows.

In summary, ECN[#] can strictly preserve the packet scheduling policy with multiple queues/classes. Furthermore, compared to TCN, ECN[#] achieves better performance for short flows because ECN[#] eliminates persistent queue buildups exacerbated by the RTT variations.

5.4.4 ECN[#] + DCQCN

In this section, we will demonstrate how ECN[#] can benefit DCQCN, a transport protocol that uses probabilistic instantaneous queueing. The experiment setting is identical as our large-scale simulation experiment (Section 5.3). Fig. 21 illustrates the results. The parameters of DCQCN is derived based on high percentile RTTs. In general, as shown in Fig. 21a, we can observe that DCQCN and DCQCN+ECN[#] does not impact the overall FCT of all flows. It further confirms that ECN[#] does not hurt networking throughput. Moreover, as shown in Fig. 21b, we can observe that compare to original DCQCN, DCQCN+ECN[#] can further improve the FCT of short flows by up to 14.9%. The results indicate that even for transport protocols that leverage probabilistic instantaneous marking, ECN[#] can benefit them by eliminating the queue buildups caused by RTT variations, improving the performance of short flows.

5.4.5 ECN[#] + PIAS

We finally show how ECN[#] can benefit PIAS, a flow scheduling mechanism. We use the experiment settings in Section 5.3 and the results are shown in Fig. 22. Similarly, as shown in Fig. 22a, ECN[#] does not impact the overall FCT of all flows. Furthermore, while PIAS can already optimize

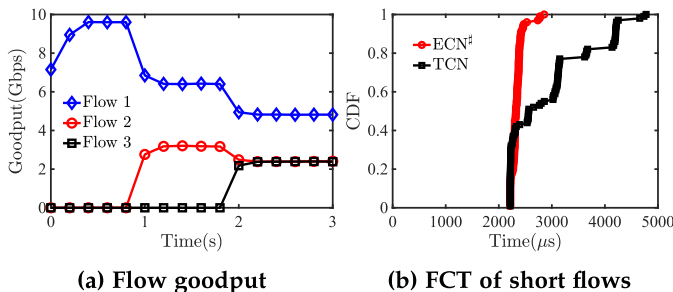


Fig. 20. [Simulations] ECN[#] with packet schedulers. ECN[#] can preserve the network scheduling policy while achieve better performance for short flows.

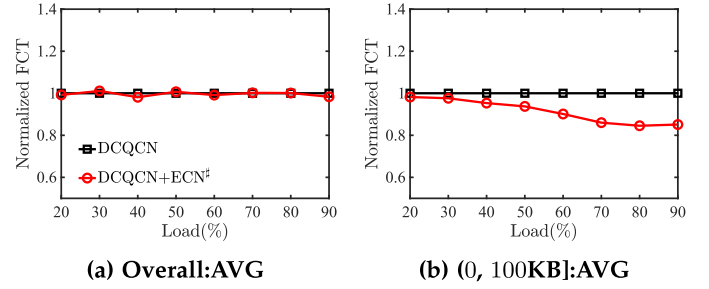


Fig. 21. [Simulations] FCT statistics with web search workload. All are normalized to the results achieved by DCQCN.

FCT for short flows by granting these flows higher priorities in the network, we observe that ECN[#] can further optimize FCT for short flows by up to 10.1% as shown in Fig. 22b. We admit that for PIAS, ECN[#] does not provide as large improvement as that for DCTCP/DCQCN/Code/TCN since PIAS has already provided superior optimization for short flows. However, we can see that ECN[#] can still provide additional improvements for short flows by eliminating unnecessary queue buildups. On a more general note, we believe that ECN[#] can also benefit other flow scheduling mechanisms, such as pFabric [26], QJUMP [27], *etc.* when RTT variation exist.

6 RELATED WORK

ECN-Based Transports in Datacenters. The ECN-related literature on datacenter networks is vast. Alizadeh et al. [1] identified the transport requirements in production datacenters and proposed DCTCP to address the challenges. To achieve good burst tolerance, DCTCP uses a simplified version of ECN/RED [8] to mark packets at the switch.

ECN^{*} [9] uses the standard TCP congestion control algorithm at the end host. D² TCP [18] and L² DCT [34] modified the DCTCP congestion control algorithm to meet flow deadlines. Shan et al. proposed CEDM [35] to accurately mark ECN to reduce throughput loss by reducing queue oscillations. These solutions apply the same AQM marking as DCTCP. Thus they all suffer from increased queueing delay with the current practice under RTT variations (Section 2.3). ECN[#] can benefit all of them. HULL [36] delivers near zero queueing delay by sacrificing a little throughput. It uses phantom queues instead of RED to mark packets at the switch. DCQCN [2] is a rate-based congestion control to enable RDMA in datacenters. At the switch, it requires RED-like probabilistic marking to ensure convergence. ECN[#] can be extended to work with it. MBECN [37] uses a queue-occupa-

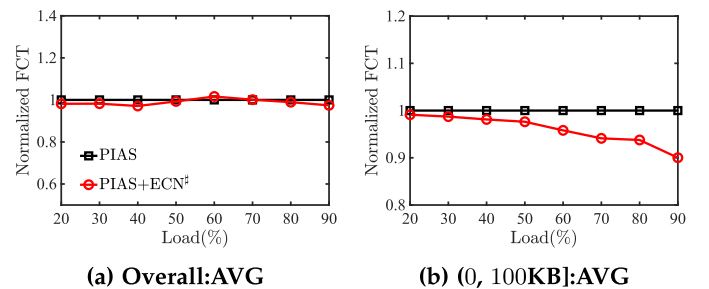


Fig. 22. [Simulations] FCT statistics with web search workload. All are normalized to the results achieved by PIAS.

tion-based dynamically adjusting algorithm to solve microburst problem, which can also be used to solve RTT variation problems but not feasible to implement on switches.

MQ-ECN [19] first pointed out the drawbacks of existing ECN/RED implementations in packet scheduling context. To adapt to the varying queue capacity caused by packet schedulers, TCN [12] proposed to use instantaneous sojourn time to mark packets. ECN[#] inherits the merit of TCN, but further tracks the persistent congestion state to reduce long-term queue buildups.

Buffer Sizing and AQM in Internet. Appenzeller et al. [38] showed that the minimum buffer size should be $C \times RTT/\sqrt{n}$ when there are a large number of n concurrent large flows. However, datacenter networks typically have a small number of concurrent large flows [1]. Therefore, people still applied the old rule-of-thumb $C \times RTT$ for buffer sizing for regular TCP.

Many AQM solutions have been proposed to handle bufferbloat in Internet. CoDel [13] tracks the minimum sojourn time over a varying time interval to detect bad queues. PI [39] and PIE [40] use a proportional-integral controller to keep the queueing to a constant target. Relative to them, ECN[#] further applies the aggressive instantaneous ECN marking to handle transient bursts.

New Design Over Programmable Switches. With the emergence of programmable switches, such as Barefoot Tofino. Many clean designs, such as NDP [41] and PERC [42] have been proposed to deliver ultra low latency communication. However, they are not compatible with state-of-the-art ECN-based solutions in datacenters. On the contrary, ECN[#] is fully compatible with ECN-based solutions and utilizes the networking programmability to mitigate performance degradation caused by RTT variations. Furthermore, AIFO [43] enables programmable packet scheduling with a single queue, which we believe could combine with ECN[#] in the future works.

7 CONCLUSION

In this paper, we have presented ECN[#] a simple yet effective solution that enables ECN in datacenters with high RTT variations. ECN[#] marks packets based on both instantaneous and persistent congestion states. It inherits the merit of current instantaneous ECN marking based on a high-percentile RTT to achieve high throughput and burst tolerance, and further marks packets when observing long-term switch queue buildups to eliminate the persistent queueing caused by flows with smaller RTTs. We have implemented ECN[#] on a Barefoot Tofino switch and evaluated it through both testbed experiments and ns-3 simulations. Our evaluation shows that ECN[#] is a viable solution that achieves all our design goals.

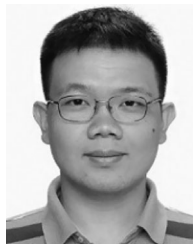
REFERENCES

- [1] M. Alizadeh et al., "Data center TCP (DCTCP)," in *Proc. ACM SIGCOMM Conf.*, 2010, pp. 63–74.
- [2] Y. Zhu et al., "Congestion control for large-scale RDMA deployments," in *Proc. ACM Conf. Special Int. Group Data Commun.*, 2015, pp. 523–536.
- [3] DCTCP in Linux kernel 3.18, 2014. [Online]. Available: https://kernelnewbies.org/Linux_3.18
- [4] DCTCP in Windows Server 2012, 2012. [Online]. Available: [https://technet.microsoft.com/en-us/library/hh997028\(v=ws.11\).aspx](https://technet.microsoft.com/en-us/library/hh997028(v=ws.11).aspx)
- [5] A. Singh et al., "Jupiter rising: A decade of Clos topologies and centralized control in Google's datacenter network," in *Proc. ACM Conf. Special Int. Group Data Commun.*, 2015, pp. 183–197.
- [6] G. Judd, "Attaining the promise and avoiding the pitfalls of TCP in the datacenter," in *Proc. USENIX Conf. Netw. Syst. Des. Implementation*, 2015, pp. 145–157.
- [7] Q. Zhang, V. Liu, H. Zeng, and A. Krishnamurthy, "High-resolution measurement of data center microbursts," in *Proc. Internet Meas. Conf.*, 2017, pp. 78–85.
- [8] S. Floyd and V. Jacobson, "Random early detection gateways for congestion avoidance," *IEEE/ACM Trans. Netw.*, vol. 1, no. 4, pp. 397–413, Aug. 1993.
- [9] H. Wu, J. Ju, G. Lu, C. Guo, Y. Xiong, and Y. Zhang, "Tuning ECN for data center networks," in *Proc. Int. Conf. Emerg. Netw. Experiments Technol.*, 2012, pp. 25–36.
- [10] P. Patel et al., "Ananta: Cloud scale load balancing," in *Proc. ACM SIGCOMM Conf.*, 2013, pp. 207–218.
- [11] D. E. Eisenbud et al., "Maglev: A fast and reliable software network load balancer," in *Proc. USENIX Conf. Netw. Syst. Des. Implementation*, 2016, pp. 523–535.
- [12] W. Bai, K. Chen, L. Chen, C. Kim, and H. Wu, "Enabling ECN over generic packet scheduling," in *Proc. Int. Conf. Emerg. Netw. Experiments Technol.*, 2016, pp. 191–204.
- [13] K. Nichols and V. Jacobson, "Controlling queue delay," *ACM Queue*, vol. 10, no. 5, pp. 20:20–20:34, 2012.
- [14] Barefoot Tofino, 2016. [Online]. Available: <https://www.barefootnetworks.com/technology/>
- [15] Barefoot Capilano SDE, 2018. [Online]. Available: <https://www.barefootnetworks.com/products/brief-capilano/>
- [16] A. Greenberg et al., "V12: A scalable and flexible data center network," in *Proc. ACM SIGCOMM Conf. Data Commun.*, 2009, pp. 51–62.
- [17] W. Bai, L. Chen, K. Chen, D. Han, C. Tian, and H. Wang, "Information-agnostic flow scheduling for commodity data centers," in *Proc. USENIX Conf. Netw. Syst. Des. Implementation*, 2015, pp. 455–468.
- [18] B. Vamanan, J. Hasan, and T. Vijaykumar, "Deadline-aware data-center TCP (D2TCP)," in *Proc. ACM SIGCOMM Conf. Appl. Technol. Archit. Protoc. Comput. Commun.*, 2012, pp. 115–126.
- [19] W. Bai, L. Chen, K. Chen, and H. Wu, "Enabling ECN in multi-service multi-queue data centers," in *Proc. USENIX Conf. Netw. Syst. Des. Implementation*, 2016, pp. 537–549.
- [20] M. Alizadeh, A. Javanmard, and B. Prabhakar, "Analysis of DCTCP: Stability, convergence, and fairness," in *Proc. ACM SIGMETRICS Joint Int. Conf. Meas. Model. Comput. Syst.*, 2011, pp. 73–84.
- [21] C. Guo et al., "RDMA over commodity ethernet at scale," in *Proc. ACM SIGCOMM Conf.*, 2016, pp. 202–215.
- [22] Linux Virtual Server, 2019. [Online]. Available: <http://www.linuxvserver.org>
- [23] R. Gandhi et al., "Duet: Cloud scale load balancing with hardware and software," in *Proc. ACM Conf. SIGCOMM*, 2014, pp. 27–38.
- [24] C. Guo et al., "Pingmesh: A large-scale system for data center network latency measurement and analysis," in *Proc. ACM Conf. Special Int. Group Data Commun.*, 2015, pp. 139–152.
- [25] NetEm in Linux, 2016. [Online]. Available: <https://wiki.linuxfoundation.org/networking/netem>
- [26] M. Alizadeh et al., "pFabric: Minimal near-optimal datacenter transport," in *Proc. ACM SIGCOMM Conf.*, 2013, pp. 435–446.
- [27] M. P. Grosvenor et al., "Queues {Don't} matter when you can {JUMP} them!," in *Proc. USENIX Conf. Netw. Syst. Des. Implementation*, 2015, pp. 1–14.
- [28] P4 Language, 2018. [Online]. Available: <https://p4.org/>
- [29] B. Li et al., "ClickNP: Highly flexible and high performance network processing with reconfigurable hardware," in *Proc. ACM SIGCOMM Conf.*, 2016, pp. 1–14.
- [30] The Network Simulator ns-3, 2015. [Online]. Available: <https://www.nsnam.org>
- [31] TCN Prototype, 2016. [Online]. Available: <https://github.com/HKUST-SING/TCN-Software>
- [32] Traffic Generator, 2016. [Online]. Available: <https://github.com/HKUST-SING/TrafficGenerator>
- [33] The Network Simulator ns-3. Packets, 2015. [Online]. Available: <https://www.nsnam.org/docs/release/3.10/manual/html/packets.html#tags-implementation>

- [34] A. Munir et al., "Minimizing flow completion times in data centers," in *Proc. IEEE Int. Conf. Comput. Commun.*, 2013, pp. 2157–2165.
- [35] D. Shan and F. Ren, "Improving ECN marking scheme with micro-burst traffic in data center networks," in *Proc. IEEE Int. Conf. Comput. Commun.*, 2017, pp. 1–9.
- [36] M. Alizadeh, A. Kabbani, T. Edsall, B. Prabhakar, A. Vahdat, and M. Yasuda, "Less is more: Trading a little bandwidth for ultra-low latency in the data center," in *Proc. USENIX Conf. Netw. Syst. Des. Implementation*, 2012, Art. no. 19.
- [37] K. Kang et al., "MBECN: Enabling ECN with micro-burst traffic in multi-queue data center," in *Proc. IEEE Int. Conf. Cluster Comput.*, 2019, pp. 1–12.
- [38] G. Appenzeller, I. Keslassy, and N. McKeown, "Sizing router buffers," in *Proc. Conf. Appl. Technol. Archit. Protoc. Comput. Commun.*, 2004, pp. 281–292.
- [39] C. V. Hollot, V. Misra, D. Towsley, and W.-B. Gong, "A control theoretic analysis of RED," in *Proc. IEEE Int. Conf. Comput. Commun. 20th Annu. Joint Conf. IEEE Comput. Commun. Soc.*, 2001, pp. 1510–1519.
- [40] R. Pan et al., "PIE: A lightweight control scheme to address the bufferbloat problem," in *Proc. IEEE 14th Int. Conf. High Perform. Switching Routing*, 2013, pp. 148–155.
- [41] M. Handley et al., "Re-architecting datacenter networks and stacks for low latency and high performance," in *Proc. ACM Conf. Special Int. Group Data Commun.*, 2017, pp. 29–42.
- [42] L. Jose, L. Yan, M. Alizadeh, G. Varghese, N. McKeown, and S. Katti, "High speed networks need proactive congestion control," in *Proc. 14th ACM Workshop Hot Topics Netw.*, 2015, Art. no. 14.
- [43] Z. Yu et al., "Programmable packet scheduling with a single queue," in *Proc. ACM SIGCOMM Conf.*, 2021, pp. 179–193.



Junxue Zhang is currently working toward the PhD degree in computer science with the Hong Kong University of Science and Technology. His current research interests include data center networking, machine learning systems, and privacy-preserving computing. His research work has been published in many top conferences such as SIGCOMM, CoNEXT.



Wei Bai received the PhD degree in computer science from the Hong Kong University of Science and Technology. He is a senior researcher with Microsoft Research Redmond. He is broadly interested in computer networking with a special focus on data center networking. His research work has been published in many top conferences and journals, such as SIGCOMM, NSDI, and *IEEE/ACM Transactions on Networking*. He also has rich experience in developing and operating production cloud networks. Currently, he is mainly focusing on network infrastructure to support cloud-scale RDMA deployments.



Kai Chen received the PhD degree in computer science from Northwestern University, Evanston, Illinois, in 2012. He is an associate professor with the Department of Computer Science and Engineering, Hong Kong University of Science and Technology, Hong Kong. His current research interests include data center networking, machine learning systems, and privacy-preserving computing.

► For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/csdl.