

# DynaQ: Enabling Protocol-Independent Service Queue Isolation in Cloud Data Centers

Gyuyeong Kim<sup>✉</sup>, *Member, IEEE* and Wonjun Lee<sup>✉</sup>, *Fellow, IEEE*

**Abstract**—Switches in cloud data centers support multiple service queues per port to provide differentiated network performance among different traffic classes. To isolate service queues, recent solutions leverage the power of Explicit Congestion Notification (ECN). However, this causes a fundamental dependency on ECN-based transport protocols, making it hard to use generic transport protocols. To this end, we design DynaQ, a protocol-independent multi-queue management solution that enables service queue isolation with generic transport protocols. The key idea of DynaQ is to adjust the packet dropping threshold of service queues dynamically. Specifically, DynaQ allows a service queue to occupy free buffer space but prevents the queue from hurting other active queues. Our solution requires only a few additional clock cycles to implement on hardware. To evaluate DynaQ comprehensively, we conduct a series of testbed experiments and large-scale simulations. Our evaluation results show that, compared to alternative schemes, DynaQ is the only solution that achieves work-conserving weighted fair sharing and low latency without protocol dependency.

**Index Terms**—Data center networks, service queue isolation, buffer management

## 1 INTRODUCTION

DATA centers are shared by many cloud services having diverse network performance requirements. To provide differentiated performance, the operator groups services into different traffic classes, and maps the classes into service queues in a switch port [2], [3]. The switch enforces network policy across the queues through packet schedulers like strict priority queueing (SPQ) and weighted round-robin (WRR). Meanwhile, the port buffer is shared among service queues in a best-effort manner. In this regime, it is difficult to isolate service queues because an aggressive queue with many flows can monopolize the packet buffer regardless of the port buffer size. If a queue cannot occupy enough buffer space, it cannot achieve its fair share rate.

Recent solutions [2], [3], [4] leverage Explicit Congestion Notification (ECN), which maintains the maximum buffer occupancy around the ECN marking threshold  $K$ . Unfortunately, the existing solutions have a fundamental dependency on ECN-based transport protocols since ECN requires cooperation between end-hosts and switches. This protocol dependency is undesirable because the end-hosts cannot use generic transport protocols. The recent advance in transport protocols shows that non-ECN protocols can outperform ECN-based protocols using different congestion signals (e.g., In-band network telemetry [5], credit [6], and network delay [7], [8]). These works are motivated by the drawbacks of ECN like coarse-grained signaling and slow convergence time. Furthermore, in multi-tenant

environments like public clouds, it is hard to enforce a specific transport protocol to virtual machines (VMs) and bare-metal (BM) servers because tenants own the network stack.

In this context, we ask the following question: *how to isolate service queues in switch ports without dependency on transport protocols?* We may modify the ECN-based schemes to drop packets when the ECN marking conditions are met. However, we find that such a simple change is not enough. MQ-ECN [2] still does not support generic packet schedulers like SPQ, which is crucial to latency-sensitive applications. TCN [3] requires dropping dequeued packets, and this degrades throughput by causing idle time on the link. Due to per-queue ECN marking, PMSB [4] works like static per-queue buffer limit (PQL) [9]. However, PQL wastes bandwidth when few queues are active because a single queue cannot occupy the buffer larger than the assigned quota, which is generally less than the Bandwidth-Delay-Product (BDP) in shallow-buffered switches. Assigning a large buffer size does not help as well since the aggressive queue eventually exhausts the available buffer. It can also harm per-port fairness by occupying excessive buffer space.

In this paper, we present DynaQ, a protocol-independent multi-queue management scheme. DynaQ enables service queue isolation with generic transport protocols. The best-effort scheme and PQL provide us the following design guideline. First, to be work-conserving, a service queue should be able to occupy the buffer larger than or equal to the BDP. Second, to achieve weighted fair sharing, the switch should guarantee the buffer as much as the weighted BDP to a service queue. Third, to achieve the requirements simultaneously, the port buffer should be shared dynamically. Based on the guideline, DynaQ adjusts the packet dropping threshold of service queues dynamically so that a queue can occupy the buffer up to the port buffer size but does not take the buffer of unsatisfied active queues.

- The authors are with the Network and Security Research Lab., School of Cybersecurity, Korea University, Seoul 02841, South Korea. E-mail: {gykim08, wlee}@korea.ac.kr.

Manuscript received 8 Nov. 2020; revised 28 Aug. 2021; accepted 1 Sept. 2021. Date of publication 8 Sept. 2021; date of current version 8 Mar. 2023. (Corresponding author: Wonjun Lee.)

Recommended for acceptance by K. Chen.

Digital Object Identifier no. 10.1109/TCC.2021.3110276

DynaQ is simple and can be implemented on hardware inexpensively with up to 7 clock cycles. The overhead is relatively small because switch ASICs require at least hundreds of clock cycles to process a packet. For example, Broadcom Trident 3 ASIC consumes at least 800 clock cycles to process a packet. We also discuss how DynaQ can be implemented on a programmable switching chip like Barefoot Tofino [10]. We implement a software prototype of DynaQ as a Linux qdisc module to compare DynaQ with various solutions.

We build a small-scale testbed with 5 servers connected to a server-emulated switch supporting  $8 \times 1\text{GbE}$  ports with two Intel I350-T4 NICs. We conduct extensive experiments to validate the efficiency of DynaQ. Our experimental results show that DynaQ provides work-conserving weighted fair sharing regardless of the number of active queues, the number of competing flows, and different per-queue weights. Using multiple end-hosts with different transport protocols, we show that DynaQ works well with generic protocols as well. In addition, DynaQ outperforms compared schemes in the FCT for small and large flows. For example, DynaQ beats the best-effort by up to  $8.40\times$  in the 99th percentile FCT for small flows. Compared to TCN [3], DynaQ achieves the better average FCT for large flows by up to  $1.99 \times$ . To complement the small-scale testbed experiments, we also perform large-scale simulations using ns-2 [11]. With link capacities of 10Gbps and 100Gbps, we demonstrate that DynaQ can preserve work-conserving weighted fair sharing in high-speed data center networks.

In summary, we make the following contributions:

- We design DynaQ, a protocol-independent multi-queue management solution that enables service queue isolation with generic transport protocols.
- We discuss the hardware implementation of DynaQ and present a software prototype of DynaQ.
- We conduct extensive experiments and simulations to demonstrate that DynaQ ensures work-conserving weighted fair sharing and low latency simultaneously without protocol dependency.

## 2 BACKGROUND AND MOTIVATION

### 2.1 Impact of Buffers on Bandwidth Sharing

Modern cloud services have diverse network performance requirements. For example, web search requires low latency and data backup demands high throughput. To provide differentiated network performance based on network policy, the operator often leverages multiple service queues in a switch port. In multi-queue environments, the operator classifies services into different traffic classes. These classes are mapped into different service queues, and the queues are scheduled by weighted fair packet schedulers like WRR and deficit round robin (DRR) [2], [4], [12]. The operator also uses SPQ to prioritize latency-sensitive flows [3]. Shared SPQ queues have a higher priority and the other dedicated queues with fair schedulers have the lowest priority.

Unfortunately, network policy can be violated in spite of packet scheduling. To saturate the weighted bottleneck capacity, queue  $i$  requires the weighted BDP that

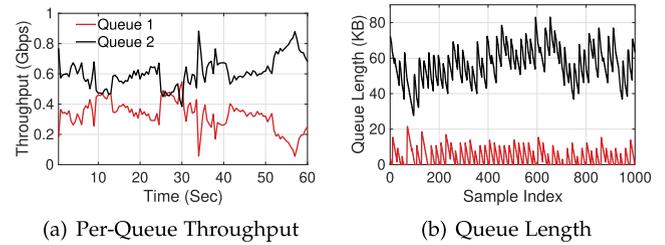


Fig. 1. [Experiment] Violated fair sharing by unfair buffer occupancy.

$$WBDP_i = C \times RTT \times \frac{w_i}{\sum w}, \quad (1)$$

where  $RTT$  is the base RTT,  $C$  and  $w_i$  denote the link capacity and the weight of queue  $i$ , respectively. However, because the port buffer is shared among service queues in a best-effort manner, a service queue may not obtain enough buffer space as much as the weighted BDP due to aggressive service queues. To show this, we conduct an experiment on a small-scale testbed consisting of 5 servers connected to a 1GbE server-emulated switch having an 85 KB of buffer per port. In the switch, we configure DRR with equal weight. We have one receiver and 4 senders. Three senders are mapped into service queue 2 while the other sender belongs to service queue 1. We start 8 flows from each of the senders at the same time and measure per-queue throughput for 60 seconds every 0.5 seconds. We also record every queue length evolution and obtain 1K sequential samples at a random time. Fig. 1 shows the results. Since the queues have the same weight, bandwidth should be shared equally. However, queue 1 cannot achieve its fair share rate because queue 1 cannot occupy the buffer larger than  $WBDP_1$  due to a large arrival rate of queue 2.

### 2.2 Why Protocol Dependency Matters?

Existing solutions [2], [3], [4] leverage ECN to isolate service queues. The solutions have an underlying assumption: all end-hosts use ECN-based transport protocols. This is because ECN requires both ECN-enabled end-hosts and switches. The assumption implies that the existing solutions have a dependency on ECN-based transport protocols.

This is undesirable because the end-host network stack cannot adapt to the advance in transport protocols. We have witnessed the emergence of many non-ECN transport protocols, which have better performance than ECN-based transport protocols. For example, HPCC [5] leverages in-band network telemetry available in emerging switch ASICs because ECN does not know how to adjust sending rates exactly and causes a trade-off between latency and throughput with the ECN marking threshold [13]. ExpressPass [6] shows that credit packets are better congestion signals than ECN in terms of fairness and convergence time. There also exist delay-based protocols like DX [8] and TIMELY [7], which are motivated by that ECN cannot inform the extent of congestion quickly. We believe that transport protocols will continue to evolve and their congestion signals will not be limited to ECN. The assumption also does not hold in many multi-tenant environments like public clouds. VMs and BM servers are the major components in cloud data centers, and their network stack is owned by tenants, not

the network operator. Therefore, it is hard to enforce a specific transport protocol for all end-hosts.

### 2.3 Analysis of ECN-Based Schemes

One may wonder whether modifying the schemes to drop packets instead of ECN marking can solve the problem. However, we find that simple changes are not enough for the following reasons.

MQ-ECN [2] is the first work that addresses the service queue isolation problem caused by unfair buffer sharing. MQ-ECN [2] determines a ECN marking threshold of queue  $i$  as following that  $K_i = \min(\frac{quantum_i}{T_{round}}, C) \times RTT \times \lambda$  where  $quantum_i$  is the weight/quantum of queue  $i$  and  $\lambda$  is a coefficient for transport protocols. For example, theoretically,  $\lambda = 1$  for ECN\* [14] and  $\lambda = 0.17$  for DCTCP [15].  $T_{round}$  indicates the estimated total time to serve all queues once. The key drawback of MQ-ECN is that the solution relies on the concept of "round" with  $T_{round}$ . This means that MQ-ECN does not support packet schedulers like SPQ. Supporting SPQ is crucial to latency-sensitive services since it can accelerate small flows. Therefore, even we change the scheme to drop packets, the solution does not achieve our design goals since it fails to provide low latency.

TCN [3] uses the packet sojourn time as the threshold metric instead of queue length to support generic packet schedulers. Since the packet sojourn time can be calculated after passing through the queue, the solution performs dequeue marking. The standard ECN marking threshold is given by  $T = RTT \times \lambda$ . If we change TCN to drop the packet when the packet sojourn time exceeds  $T$ , we should drop the just dequeued packet. However, packet dropping at dequeue causes idle time on the link. This seriously degrades the effective throughput. In addition, dropping the buffered packet increases the FCT as much as the packet sojourn time in addition to the retransmission timeout (RTO).

PMSB [4] provides both generic packet schedulers and early congestion notification, which MQ-ECN and TCN do not support, respectively. PMSB only marks packets when per-port ECN marking and per-queue ECN marking conditions are met at the same time where the port ECN marking threshold is given by  $K = C \times RTT \times \lambda$ . The per-queue ECN marking threshold for queue  $i$  is given by  $K_i = \sum_w \frac{w_i}{w} C \times RTT \times \lambda$ . Since  $K_i \leq K$ , the dropping version of PMSB is similar to PQL, which is supported in some production switches.

PQL assigns a static buffer size to a service queue. Therefore, each queue can enjoy its fair share regardless of other queues. However, PQL is not work-conserving because the amount of buffer that a single queue can occupy is limited to the assigned quota. Therefore, when few queues are active, the link capacity can be underutilized since the aggregate buffer occupancy can be less than the BDP. One might argue that assigning a buffer of the BDP to all service queues can solve the problem. Unfortunately, the on-chip SRAM buffer is a scarce resource in switches [16]. Therefore, we do not have enough buffers to reserve a buffer size as much as the BDP for all service queues.

## 3 DYNAQ DESIGN

### 3.1 Design Goals

Our goal is to *design a multi-queue management solution that enables service queue isolation in switch ports over generic*

TABLE 1  
Used Notations

Notation	Description
$M$	Number of queues
$P$	Arriving packet
$p$	Queue index of packet $P$
$B$	Port buffer size
$w_i$	Weight of queue $i$
$T_i$	Packet dropping threshold of queue $i$
$q_i$	Queue length of queue $i$
$WBDP_i$	Weighted BDP of queue $i$
$S_i$	Satisfaction threshold of queue $i$
$T_i^{ex}$	Extra buffer size of queue $i$

transport protocols. We stipulate that a good solution should satisfy the following requirements *simultaneously*:

- *Protocol independence*: A solution should not be tied to a specific transport protocol.
- *Work conservation*: A solution should be able to utilize the whole link capacity regardless of any time.
- *Weighted fair sharing*: A solution should strictly preserve weighted fair sharing among service queues at any time regardless of traffic dynamics.
- *Low latency*: A solution should support arbitrary packet schedulers, especially SPQ, to minimize the FCT of small flows.
- *Practicality*: A solution should be inexpensive to implement on hardware.

### 3.2 Mechanisms

#### 3.2.1 Basic Idea

DynaQ is the first protocol-independent multi-queue management scheme that satisfies the above requirements simultaneously. We observe that the best-effort scheme and PQL provide the following design guideline that 1) to utilize the bottleneck link capacity fully, a service queue must be able to occupy the buffer larger than or equal to the BDP if there is free space in the port buffer.; 2) to share bandwidth fairly while respecting different weights of service queues, a service queue must be able to occupy buffer space larger than or equal to the weighted BDP regardless of other service queues.; 3) to guarantee the weighted fair share rate and sustain high link utilization at the same time, the switch should manage the port buffer dynamically among service queues. Without dynamic multi-queue management, we can meet only one of the two requirements at a time.

Following the above design guideline, DynaQ allows a single service queue to occupy free buffer space in the port but prevents the service queue from taking the buffer of unsatisfied active queues.<sup>1</sup> Table 1 summarizes mathematical notations used to describe our work. To realize the idea, DynaQ assigns packet dropping threshold  $T_i$  for each

1. We express that queue  $i$  is unsatisfied if the packet dropping threshold  $T_i$  of queue  $i$  is less than satisfaction threshold  $S_i$  defined in Eq. (4). Otherwise, queue  $i$  is satisfied.

service queue  $i$ , which means the total size of packets that can be buffered. The switch dynamically adjusts  $T_i$  every packet arrival. Since a service queue should be able to occupy the buffer up to the port buffer size, the switch does not simply drop the arriving packet  $P$  when the dropping threshold is exceeded. Instead, if the buffer occupancy of the service queue of packet  $P$  exceeds  $T$  with packet  $P$ , the switch increases the dropping threshold of the queue and decreases that of the victim queue, which is defined as the queue has the largest extra buffer size. However, to guarantee the weighted fair share rate, the switch drops packet  $P$  without threshold adjustment when the victim queue is an unsatisfied active queue.

---

**Algorithm 1.** Pseudocode of DynaQ
 

---

```

1: if  $q_p + size(P) > T_p$  then           ▷ Exceeds threshold?
2:    $v \leftarrow \arg \max_{i < M, i \neq p} T_i^{ex}$    ▷ Find victim queue
3:   if  $(T_v < size(P)) \vee (q_v > 0 \ \& \ T_v - size(P) < S_v)$  then
4:     Return Drop( $P$ )           ▷ Protect unsatisfied queues
5:   else                             ▷ It is okay to adjust dropping thresholds
6:      $T_v \leftarrow T_v - size(P)$            ▷ Decrease  $T$  of victim  $v$ 
7:      $T_p \leftarrow T_p + size(P)$          ▷ Increase  $T$  of queue  $p$ 
8:   end if
9: end if
  
```

---

### 3.2.2 Detailed Design

DynaQ operates before enqueueing decisions. As shown in Algorithm 1, the switch first compares dropping threshold  $T_p$  with the sum of queue length of queue  $p$  and the size of packet  $P$ . If enqueueing packet  $P$  does not make the queue length exceed  $T_p$ , nothing will be done. Otherwise, the switch begins to adjust packet dropping thresholds or to drop the packet.

*Packet Dropping Threshold.* DynaQ isolates service queues using dynamic packet dropping thresholds. Each queue  $i$  has its own dropping threshold  $T_i$ . The aggregate dropping threshold of all service queues is equal to the port buffer size. If the aggregate threshold exceeds the port buffer size, the buffer will be shared like the best-effort scheme. In contrast, if the aggregate threshold is less than the port buffer size, the buffer sharing policy becomes close to PQL. Therefore, when the switch is on, DynaQ initializes the packet dropping threshold of queue  $i$  that

$$T_i^{init} = B \times \frac{w_i}{\sum w}. \quad (2)$$

In addition, to ensure  $\sum_{i=1}^M T = B$ , DynaQ always decreases the dropping threshold of victim queue before increasing the dropping threshold of the queue of packet  $P$ .

*Victim Queue Selection.* Before adjusting dropping thresholds, the switch should find the victim queue  $v$ . One intuition is that the victim should be the queue who is expected to experience the minimal impact after decreasing its dropping threshold. Therefore, a natural way to select the victim is to find the queue with the largest threshold. However, this may not work when service queues are with different queue weights. For example, consider 3 service queues with weights of 1:2:3. The switch can choose queue 3 as the

victim queue although queue 3 has only a minimum required buffer size to enjoy the weighted fair share rate.

To respect different queue weights, DynaQ selects a service queue with the largest extra buffer size as the victim queue. The extra buffer size of queue  $i$  is defined as

$$T_i^{ex} = T_i - S_i. \quad (3)$$

The satisfaction threshold  $S_i$  specifies the minimum buffer size of queue  $i$  to achieve its weighted fair share rate regardless of other queues. Theoretically,  $WBDP_i$  is enough to saturate the weighted bottleneck link capacity  $C \frac{w_i}{\sum w}$ . However, we find that the switch does not preserve weighted fair sharing when  $S_i = WBDP_i$ . This is because  $T_i$  fluctuates over time, preventing queue  $i$  from enjoying its fair share rate stably. Therefore, we need to satisfy the inequality  $S_i > WBDP_i$  to make headroom to reduce the impact of the change of  $T_i$ . Thus, we simply use that

$$S_i = B \times \frac{w_i}{\sum w}. \quad (4)$$

Modern line-rate switches have enough buffer size to allocate a buffer size per port larger than BDP. Since  $B > BDP$ , it is obvious that  $S_i > WBDP_i$ .

*Victim Queue Search Without Loops.* Finding the victim queue can be done through a linear search using loops. However, modern switching ASICs prevent loop operations to guarantee a deterministic packet processing delay.

To deal with this constraint, DynaQ uses binary search to find the victim queue. We make `MaxIdx` function that returns the index of the larger queue after comparing the extra buffer size between the two service queues. For example, when the switch supports 4 service queues, we can find the index of victim queue by referring to the return value of `MaxIdx(MaxIdx(1, 2), MaxIdx(3, 4))`. This requires  $O(\log n)$  complexity bounded to the number of service queues that the switch supports. Modern switches typically support 4 or 8 service queues per port. Therefore, the complexity is fixed to  $O(2)$  or  $O(3)$  depending on the target switch architecture. DynaQ does not consume excessive clock cycles even with more than 8 queues. For example, when the number of queues is 32 [17], the complexity is only  $O(5)$ .

*Packet Dropping and Threshold Adjustment.* After finding the victim queue, the switch decides whether to drop the packet or adjust the dropping thresholds. The switch drops packet  $P$  if the dropping threshold of queue  $v$  is less than the size of packet  $P$  or queue  $v$  is an unsatisfied active queue. The former condition is to ensure  $T_i \geq 0, \forall i$ . The latter condition is to protect an unsatisfied active queue. If we allow a queue to take the buffer of unsatisfied active queues, aggressive queues with many flows can occupy the port buffer excessively. Meanwhile, the switch does not protect inactive queues from the aggressive queues to utilize free buffer space for high link utilization. When the above conditions are not met, DynaQ changes the dropping thresholds of queue  $v$  and queue  $p$  as much as the size of packet  $P$ . This finishes the operation of DynaQ. After this, the switch performs packet enqueueing decisions based on the port buffer occupancy.

### 3.2.3 Why It Works?

We show how DynaQ isolates service queues through a simple theoretical analysis. Consider an output port where  $M$  service queues share link capacity  $C$  and the buffer size is  $B$ . Each queue  $i$  has weight  $w_i$ . We suppose that queue  $i$  has traffic with  $\lambda_i(t)$ , the input rate of queue  $i$  at time  $t$ . For packet scheduling, we consider WRR as the underlying packet scheduler where packets in queue  $i$  are scheduled with the weight of  $w_i$  every round. Since we are interested in bandwidth sharing on the bottleneck, we consider cases when  $\sum_{i=1}^M \lambda_i(t) > C$ . The output rate of queue  $i$  at time  $t$  can be given by

$$\mu_i(t) = \min(\lambda_i(t), \alpha_i(t)), \quad (5)$$

where  $\alpha_i(t)$  represents the weighted fair share rate for queue  $i$  at time  $t$ , which is determined by the buffer sharing policy and the packet scheduling algorithm.

*Weighted Fair Sharing.* To ensure weighted fair sharing on the bottleneck, we must guarantee  $\alpha_i(t) \geq \frac{w_i}{\sum_{k=1}^M w_k} C$  when  $\lambda_i(t) \geq \frac{w_i}{\sum_{k=1}^M w_k} C$ . To achieve this, a buffer sharing solution must satisfy the following inequality that

$$q_i^{\max}(t) \geq WBDP_i, \quad (6)$$

where  $q_i^{\max}(t)$  indicates the maximum queue length of queue  $i$  at time  $t$  and  $WBDP_i$  is the weighted BDP of queue  $i$  in Eq. (1).

By adjusting packet dropping thresholds, DynaQ strictly ensures  $q_i^{\max}(t) = S_i$  regardless of the input rate of other queues where  $S_i$  is the satisfaction threshold in Eq. (4). Since  $S_i > WBDP_i$ , we can know that DynaQ satisfies the above inequality. Thus, DynaQ ensures weighted fair sharing between service queues on the bottleneck link. Note that the best-effort scheme cannot guarantee the satisfaction of the above inequality since it is possible  $q_i^{\max}(t) < WBDP_i$ .

*Work Conservation.* To achieve work conservation, we must guarantee  $\sum_{i=1}^M \mu_i(t) = C$  even if there exists queue  $j$  whose  $\lambda_j(t) < \frac{w_j}{\sum_{k=1}^M w_k} C$ . To preserve this, a buffer sharing solution must satisfy the following inequality that

$$\sum_{i=1}^M q_i(t) \geq C \times RIT. \quad (7)$$

Since DynaQ allows queue  $i$  to occupy the underutilized buffer space, we can say that DynaQ ensures the above inequality. Thus, DynaQ preserves work conservation. Note that PQL does not satisfy the inequality since it is possible  $\sum_{i=1}^M q_i(t) < C \times RIT$ .

### 3.2.4 Discussion

*ECN Support.* DynaQ should support ECN-based transport protocols since they are also generic transport protocols. Since there exist ECN-based solutions, we employ PMSB [4] rather than designing our own ECN-based mechanism. Specifically, when ECN is enabled in the switch, DynaQ does not adjust dropping thresholds but marks the packet when the port buffer occupancy exceeds the port ECN marking

threshold  $K = C \times RIT \times \lambda$  and the queue length of arriving packet queue exceeds the per-queue ECN marking threshold  $K_i = \sum_w \frac{w_i}{w} C \times RIT \times \lambda$  simultaneously.

*Port Buffer Size.* We have assumed that the port buffer size is constant so that the sum of dropping thresholds  $\sum T$  can be equal to the port buffer size  $B$ . However, the operator can change the port buffer size, breaking the equality between  $\sum T$  and  $B$ . This can be resolved by performing the initialization of the dropping thresholds via Equation (3) after adjusting the port buffer size.

## 4 IMPLEMENTATION

### 4.1 Hardware Implementation

DynaQ can be implemented on hardware inexpensively. Since we cannot program most switching chips, we first analyze the overhead in ASIC implementation. Next, we discuss the implementation on programmable switches.

*Processing Overhead in ASIC Implementation.* The processing overhead of DynaQ in ASIC implementation is relatively small since the switching ASIC consumes hundreds of clock cycles to process a packet. Consider typical hardware running at a clock frequency of 1 Ghz where 1 clock cycle is 1 ns. We also presume that the switch supports 8 service queues per port. Note that commodity switch ASICs support 4-8 service queues. With Algorithm 1, we can know that DynaQ requires up to 7 clock cycles. Broadcom Trident 3 offers a minimum per-packet processing delay of 800 ns. In this case, the overhead of DynaQ is 0.88%.

Let us show the detailed analysis. In the worst case, Lines 1-3 and Lines 6-7 in Algorithm 1 are performed. Line 1 consumes 1 clock cycle. Line 2 requires  $\log 8 = 3$  clock cycles. Line 3 consumes 2 clock cycles because  $(q_v > 0 \ \& \ T_v - size(P) < S_v)$  must be performed before  $\parallel$  operation with  $T_v < size(P)$ . Note that comparison operations like  $q_v > 0$  can be pipelined. Lines 6-7 require 1 clock cycle with pipelining because they have no read/write dependency.

*Implementation on Programmable Switches.* We analyze how DynaQ can be implemented on a programmable switch built with Barefoot Tofino [10]. With Tofino ASIC, we can program processing pipelines [18]. We target Tofino Native Architecture (TNA). TNA consists of 9 blocks whose 6 blocks are programmable and the other blocks only provide a fixed set of operations. DynaQ should be implemented in the Packet buffer and Replication Engine (PRE) where packet enqueueing and dequeueing decisions occur. However, the PRE is still a fixed-function block in TNA. Therefore, it is hard to implement packet buffering mechanisms directly due to the currently limited programmability. Instead, we should implement DynaQ in either ingress or egress pipeline indirectly. Since DynaQ operates before packet buffering, the ingress pipeline is the right place to implement DynaQ.

Most variables like packet dropping thresholds, queue weights, and satisfaction thresholds can be defined as user-defined metadata. In addition, we can manipulate them at runtime. MaxIdx function to find the victim queue also can be defined as a custom function. The function returns the index of queue with larger extra buffer size using comparators, which are supported in the current programmable switch architecture. The extra buffer size of each queue can be stored

in stateful registers as well. One challenge to implement DynaQ in TNA is to obtain the queue length information. Since the PRE is not programmable, it is hard to obtain the queue length of the arriving packet queue in the ingress pipeline. Note that this is not an issue in ASIC implementation. In TNA, the queue length metadata is included in egress intrinsic metadata, but it is prohibited to share metadata between the pipelines. Although TNA supports metadata bridging, we can only deliver ingress pipeline metadata to the egress pipeline. We note that new Tofino2 switch ASICs allow reading the queue length in the ingress pipeline directly, which greatly simplifies the implementation of DynaQ [19].

To deal with the issue, we may utilize multiple registers as follows. First, we define register arrays to count every packet arrival for each output port. Counting packets should be done after the port lookup stage since we need the index of the output port. Second, we define registers to store the current port buffer size in the ingress pipeline. This can be done via the switch control plane when the port buffer size changes because the control plane has switch configuration data and data plane statistics. Third, we add a stage in the egress pipeline that reports `deq_qdepth` metadata, the queue length at dequeue time, to the control plane. Upon the arrival of `deq_qdepth`, the control plane updates the counter registers in the ingress pipeline to reflect the correct queue length information. This may result in DynaQ operations based on inaccurate queue length information due to the time to update extern registers. However, with round-robin based schedulers, we believe that some inaccuracy is tolerable to isolate service queues. Note that it is also hard to implement ECN-based schemes in the ingress pipeline with the currently limited programmability since ECN also requires queue length information to mark packets. Owing to the simplicity of the mechanism, DynaQ only requires roughly 5 match-action stages and 5.68% of additional memory to be implemented with Tofino ASICs.

## 4.2 Software Implementation

To compare DynaQ with various existing works in a flexible environment, we implement a software prototype of DynaQ as a Linux qdisc module on a server-emulated switch. Our module is based on the software prototype of MQ-ECN [2]. Our module consists of two stages as follows.

*Enqueueing Stage.* When the packet from TCP/IP stack arrives at the qdisc layer, the module returns the index of the corresponding service queue by referring to the DSCP field in the IP header. Next, the switch checks whether the buffer is available to enqueue the arriving packet. Basically, the switch performs packet enqueueing decisions based on the port buffer occupancy or per-queue buffer occupancy relying on the switch configuration. When the switch uses DynaQ, the switch performs additional operations to adjust packet dropping thresholds before enqueueing. If the packet can be buffered, the switch enqueues the packet into the corresponding queue. When ECN is enabled, the switch performs ECN marking at the end of enqueueing stage.

*Dequeueing Stage.* The switch dequeues packets through work-conserving packet schedulers, which include SPQ, DRR, and WRR. The packet schedulers follow the data structure and mechanism of the current Linux qdisc

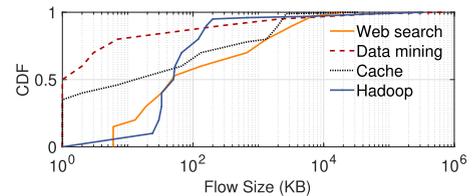


Fig. 2. Used workloads in dynamic flow experiments.

implementation. The dequeued packets go to NIC drivers and NIC hardware before it is transmitted to the wire. Our module uses a token-bucket rate limiter to shape outgoing traffic at 99.5% of the NIC capacity. This is to avoid excessive buffering in NIC drivers and NIC hardware, which can lead to inaccurate buffer occupancy in the qdisc.

## 5 PERFORMANCE EVALUATION

In this section, we evaluate the performance of DynaQ. We answer the following key questions:

- How does DynaQ perform in practice?
- Does DynaQ perform well in large-scale networks?
- How robust is DynaQ to network settings?

*Traffic Workloads.* We use four realistic workloads derived from production data centers, which are a web search workload [15], a data mining workload [20], a cache workload [21], and a hadoop workload [21]. As shown in Fig. 2, these workloads generate flows whose size distributions are heavy-tailed. For example, in the data mining workload, roughly 50% of flows are 1KB while 90% of bytes are from flows larger than 100 MB. Like MQ-ECN [2] and TCN [3], we also use the web search workload for all service queues in testbed experiments while using all the four workloads in simulations. This is because the web search workload is the most challenging workload due to its less skewed flow size distribution that generates multiple concurrent flows on the bottleneck. In addition, it is hard to make all workloads active during experiments because every workload results in different finish times.

*Compared Schemes.* We mainly compare DynaQ with the following schemes: BestEffort and PQL. BestEffort denotes the best-effort scheme that manages the buffer among service queues in a FIFO manner. PQL isolates service queues by reserving a static per-queue buffer size. When we consider low latency, we also compare DynaQ against the ECN-based solutions, TCN [3], and PMSB [4].

*Performance Metric.* Our primary performance metrics are throughput and FCTs. We also use the throughput share ratio and Jain's fairness index for a better understanding of throughput results. For the average FCT, we breakdown the FCT across different flow sizes to analyze the impact on small ( $\leq 100$  KB) and large flows ( $> 10$  MB). We also consider the 99th percentile FCT of small flows to evaluate tail latency. Due to space limitation, we omit the result of medium flows whose results are similar to overall flows. For a clear comparison, the FCT results are normalized by the values of DynaQ.

### 5.1 Testbed Experiments

*Testbed Setup.* We have built a small-scale testbed, which consists of 5 servers connected to a server-emulated switch. Authorized licensed use limited to: Korea University. Downloaded on May 18, 2023 at 05:27:30 UTC from IEEE Xplore. Restrictions apply.

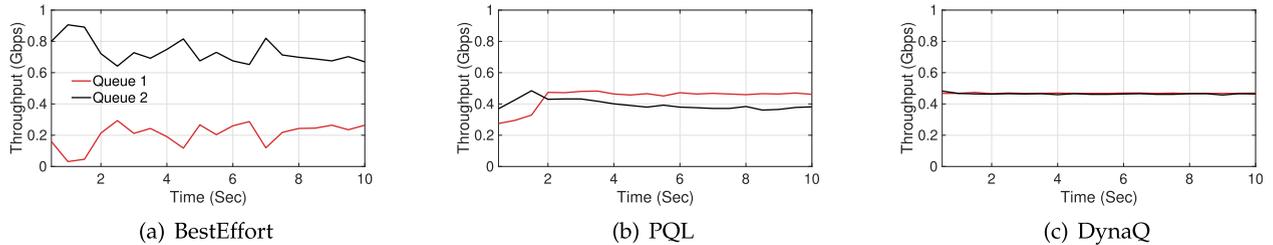


Fig. 3. [Experiment] Throughput convergence of two active DRR queues with equal weight.

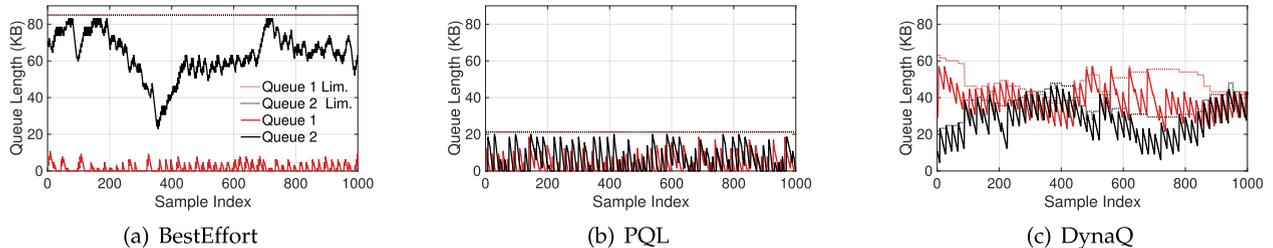


Fig. 4. [Experiment] Queue length evolution of two active DRR queues with equal weight.

The switch is equipped with two Intel I350-T4 v2 NICs where each NIC supports  $4 \times 1$ GbE ports. Each server is also equipped with an Intel gigabit NIC. All servers use Linux kernel 3.18.11. We use TCP for the non-ECN schemes and DCTCP [15] for the ECN-based schemes. We set TCP  $RTO_{min}$  to 10ms as suggested in many existing works [2], [3], [15], [22], [23]. The initial congestion window size is set to 10 packets as suggested in RFC6928. In end-hosts and the switch, we disable large send offload (LSO) and large receive offload (LRO) to reduce traffic burstiness and emulate switch hardware behaviors more correctly. To emulate a switch with Broadcom 56538 ASIC, the switch has an 85KB of the port buffer [24], which is completely shared by all service queues. PQL is the only exceptional scheme since it limits the per-queue buffer size. The base RTT is roughly  $500 \mu s$  and the corresponding BDP is 62.5 KB. We set ECN marking thresholds for DCTCP and TCN to 30 KB and  $240 \mu s$ , respectively. These are the best values experimentally found. Note that there is a theory-practice gap in determining ECN marking thresholds [3].

### 5.1.1 Static Flow Experiments

In static flow experiments, we focus on weighted fair sharing and work conservation.

*Convergence and Queue Evolution.* In this experiment, we show the basic results with two active queues between 4 DRR queues having an equal quantum of 1.5 KB. We use three servers where two servers are the senders for each service queue and the other one is the receiver. Using *iperf*, each sender starts flows to the receiver for 10 seconds. The sender of queue 2 generates 16 flows while the sender of queue 1 has only 2 flows. Ideally, the active queues should share bandwidth equally regardless of the number of competing flows and the inactive queues.

Fig. 3 shows the throughput of the active queues over time. It is easy to find that DynaQ is the only solution that shares bandwidth fairly. Throughput of the active queues in BestEffort does not converge, resulting in significant

unfairness. With PQL, the active queues share the bandwidth fairer than BestEffort but still results in considerable unfairness.

Fig. 4 reports the queue length evolution for the active queues. We measure per-queue buffer occupancy every enqueueing and dequeueing operations and obtain 1K sequential samples. The dotted lines indicate per-queue buffer size. The queue evolution samples explain the throughput in Fig. 3. In BestEffort, since queue 2 has more flows, queue 2 dominates the port buffer while queue 1 with smaller flows occupies a small buffer. In PQL, queue 2 can occupy buffer space more than that in BestEffort, but it is limited to a reserved size. Unlike the other schemes, DynaQ shares the buffer dynamically that dropping thresholds change over time. Thanks to this, each queue can occupy enough buffer regardless of the number of flows and active queues.

*Weighted Fair Sharing and Work Conservation.* In this experiment, we consider 4 DRR queues with equal quantum as same as the previous experiment. However, we now vary the number of active queues over time. Each queue has a different number of flows that the sender of queue  $i$  starts  $2^i$  flows to the receiver simultaneously. From 10 seconds, we change the number of active queues by stopping flows. At 10 seconds, the sender of queue 4 stops traffic. After 5 seconds, queue 3 becomes inactive. At 20 seconds, the sender of queue 2 no longer sends flows. The sender of queue 1 finishes at 25 seconds. We measure per-queue throughput every 0.5 seconds and obtain the aggregate throughput as well. Ideally, service queues should share bandwidth fairly regardless of the number of flows and the aggregate throughput always should be high regardless of the number of active queues.

Fig. 5 shows the results. We observe that BestEffort fails to achieve fair sharing. When all the four queues are active, queue 4 with 16 flows occupies the largest bandwidth share because the port buffer is dominated by the packets of queue 4. Due to this, queue 1 only obtains 0.14Gbps of average throughput for the first 10 seconds. Even when only

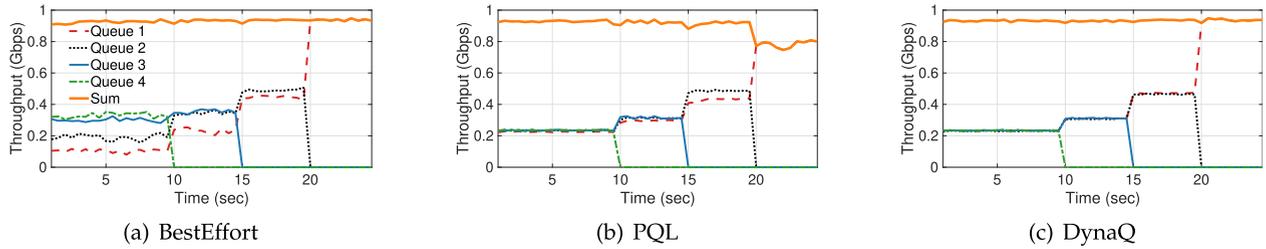


Fig. 5. [Experiment] Bandwidth sharing between 4 DRR queues with equal weight.

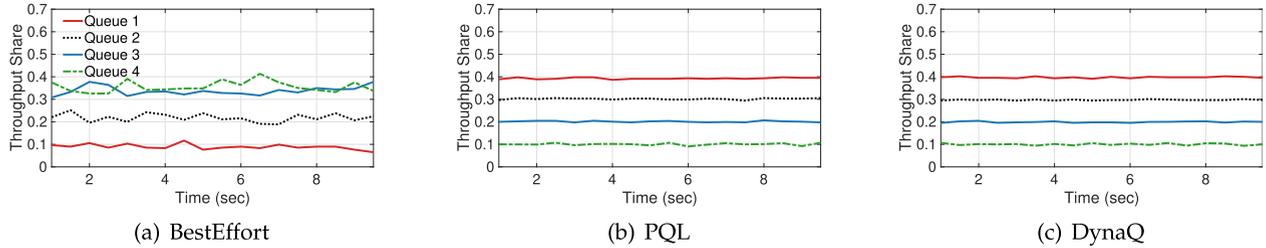


Fig. 6. [Experiment] Bandwidth sharing between 4 DRR queues with different queue weights of 4:3:2:1.

queues 1 and 2 are active at 15 seconds, the two queues do not share bandwidth fairly in spite of the scheduler because queue 2 has twice as many flows as queue 1.

PQL shows better results than BestEffort. When all queues are active, PQL preserves fair sharing. However, when the number of active queues decreases, fair sharing is violated. We can see the unfair bandwidth sharing between the two service queues at 15 ~ 20 seconds. More importantly, we can see that the aggregate throughput decreases as the queues become inactive. It is notable that the average aggregate throughput during 20 ~ 25 seconds is only 0.78 Gbps. This is because each service queue cannot utilize the remaining free buffer space. Unlike the compared schemes, DynaQ makes the service queues share bandwidth almost perfectly regardless of the number of flows. In addition, since DynaQ allows a queue to utilize free buffer space, the aggregate throughput does not decrease even when few queues are active.

*Impact of Queue Weights.* In this experiment, we use the same scenario in the previous experiment. The difference is that we configure different quantum for the DRR queues. Our default quantum is 1.5 KB of MTU. We set the weights of the queues to  $\{4, 3, 2, 1\}$ , which result in  $\{6, 4.5, 3, 1.5\}$  KB of quantum. We measure the per-queue throughput every 0.5 seconds for 10 seconds. Ideally, the queues should share the bandwidth by respecting their assigned weights.

Fig. 6 shows the throughput share of each service queue. The throughput share is defined as  $R_i(t)/\sum R(t)$  where  $R_i(t)$  denotes the throughput of queue  $i$  at time  $t$ . The result with BestEffort shows that BestEffort does not preserve weighted fair sharing. In spite of different quantum, BestEffort allows a queue with many flows to occupy more throughput share. For example, the average throughput share of queue 4 for 10 seconds is 0.35 although its desirable share is 0.1. PQL achieves weighted fair sharing in this experiment. This is not surprising because PQL assigns a static buffer size to each service queue. However, as we have shown in the previous experiment, PQL loses throughput as the number of active queues decreases although it is

omitted in this experiment. The result demonstrates that DynaQ achieves weighted fair sharing by respecting the assigned weights regardless of the number of competing flows.

*Impact of Transport Protocols.* We consider a scenario where senders use different transport protocols. We have tried to use emerging protocols, but it is hard to obtain their codes. Instead, we use TCP and CUBIC. Unlike the previous experiments where all senders use TCP, we now make the senders of queues 3 and 4 use CUBIC. Except for transport protocols, we use the same scenario in the equal sharing experiment. Fig. 7 shows the results. We can see that DynaQ achieves fair sharing regardless of employed transport protocols. One notable point is that, at 10 seconds and 15 seconds, we can see that the aggregate throughput decreases slightly for a moment when queue 4 and queue 3 become inactive. This is because of the time for the ramp-up of the other queues, not due to our buffer sharing policy. Delay-based protocols and INT-based protocols also can work well with DynaQ since the adjustment of dropping threshold only changes the buffer occupancy that each packet sees, not the sending rate of senders directly. The protocols simply adjust the sending rate with updated congestion signals, which are affected by the buffer occupancy.

### 5.1.2 Dynamic Flow Experiments

*Methodology.* At the switch, we configure SPQ/DRR where one queue has a higher priority than the other four DRR

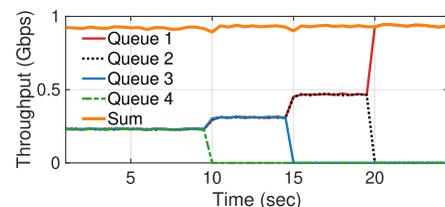


Fig. 7. [Experiment] Throughput with 2 TCP and 2 CUBIC senders.

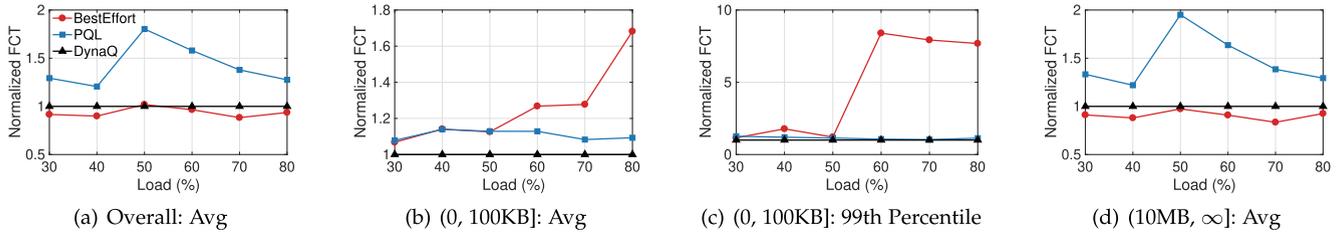


Fig. 8. [Experiment] FCT comparison against non-ECN schemes with SPQ (1 queue)/DRR (4 queues).

queues. Packets in the DRR queues can be dequeued only when the SPQ queue is empty. This is a common switch configuration to accelerate latency-sensitive small flows [3]. We use 1.5 KB of an equal quantum for all DRR queues. Note that SPQ does not require threshold configurations of DynaQ because it is to dequeue packets of queues with priorities, not to ensure fairness among the queues.

We use a client/server application [2] to generate traffic with the web search workload. We have 4 servers and 1 client. The client initially opens 5 persistent TCP connections to each server. The client generates requests to the servers through available connections. When there is no available connection, the client creates a new connection. The inter-arrival time of generated requests follows a Poisson process. When a request arrives, each of the servers responds with the requested data size. The server application sets DSCP values for outgoing packets using `setsockopt` and a flow is mapped to one of the service queues randomly. We also employ a two-level PIAS [25] to classify small flows with 100KB of a priority demotion threshold. Therefore, the first 100K bytes are buffered into the SPQ queue and the remaining bytes are enqueued into the DRR queues in the low priority. Overall, we generate 10K flows by varying the traffic load from 30% to 80%.

*Comparison With Non-ECN Schemes.* Fig. 8 reports the results that compare DynaQ with non-ECN schemes. The average FCTs of overall and large flows are similar since most of the bytes in the benchmark traffic come from large flows. Compared to PQL, DynaQ achieves the better average FCT for large flows by up to  $1.95\times$ . Similarly, DynaQ outperforms PQL in the average FCT for overall flows by up to  $1.80\times$ . This is because a service queue in PQL can occupy a limited buffer space. The results against BestEffort are mixed whose gaps are within  $0.90\times \sim 1.02\times$  and  $0.83\times \sim 0.97\times$  for the average FCT of overall flows and large flows, respectively. The reason why BestEffort outperforms DynaQ is that large flows block small flows. Not surprisingly, we can see that DynaQ outperforms BestEffort in the FCT of small flows. Although BestEffort outperforms DynaQ in the FCT of large flows, BestEffort cannot preserve weighted fair sharing as we have shown.

For the average FCT of small flows, DynaQ achieves the best performance between the schemes. DynaQ beats BestEffort by  $1.26\times$  on average across the traffic loads. The performance gap increases as the traffic load increases. Compared to PQL, DynaQ has the better performance within  $1.08\times \sim 1.14\times$ . From the 99th percentile FCT result, we find that BestEffort shows a significantly worse performance than DynaQ when traffic load increases. For example, the FCT gap in 60% of load is  $8.40\times$ . Unlike BestEffort, PQL shows relatively stable performance. However, DynaQ still outperforms PQL by  $1.14\times$  on average across the loads.

Authorized licensed use limited to: Korea University. Downloaded on May 18, 2023 at 05:27:30 UTC from IEEE Xplore. Restrictions apply.

*Comparison With ECN-Based Schemes.* We now discuss the FCT results of DynaQ compared to the ECN-based schemes. Fig. 9 shows the results. Like the results in Fig. 8, the results for overall and large flows are similar. DynaQ shows the mixed performance against the ECN-based solutions but generally outperforms the comparisons. TCN shows similar average FCTs for overall and large flows when traffic loads are within 30% ~ 40%. However, the maximum gap is only  $0.95\times$  for overall flows when the load is 30%. DynaQ outperforms TCN for the rest traffic loads whose ranges of performance gaps are within  $1.28\times \sim 1.85\times$  and  $1.29\times \sim 1.99\times$  for overall and large flows, respectively. PMSB also has a similar performance to TCN. Per-Queue ECN shows the worst performance between the schemes. When we consider the results for small flows, DynaQ beats the other ECN-based schemes in both average and 99th percentile FCTs. For example, DynaQ is better than PMSB by up to  $1.29\times$  in the average FCT. The ECN-based schemes show the worse performance at lower traffic loads than high traffic loads. DynaQ outperforms PMSB and Per-Queue ECN by  $12.23\times$  and  $12.63\times$  at 30% of traffic load, respectively.

*Impact of Packet Schedulers.* For deep dive, we also conduct an experiment with different scheduling configurations. The methodology of this experiment is the same as the previous experiments except that we have only four DRR queues with equal quantum. This configuration deserves consideration in practice because the switch may not configure a SPQ queue for some reason. Specifically, we can consider some scenarios when many traffic classes can be required or the switch supports only a few service queues.

We plot the results at a moderate load of 60% in Fig. 10. We observe that DynaQ outperforms BestEffort and PQL in the average FCTs for overall, small, and large flows as well as the 99th percentile FCT for small flows. Compared to the results in Fig. 8 where SPQ exists, we can find the following differences. First, for the FCT of small flows, the gap between DynaQ and BestEffort is decreased because small flows are not prioritized anymore. Second, DynaQ now beats BestEffort in the average FCT for large flows. This is because, in this configuration, DynaQ does not penalize large flows to expedite small flows. We do not observe any negative results, and this suggests that DynaQ is robust to different packet scheduling configurations.

## 5.2 Large-Scale Simulations

We conduct ns-2 simulations with large-scale environments.

### 5.2.1 Static Flow Simulations

*Methodology.* Like our testbed, we build a star topology to emulate a compute rack. We consider two high-speed links:

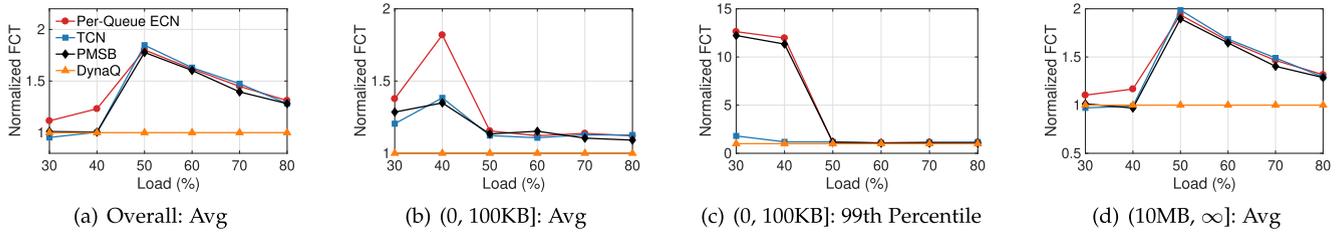


Fig. 9. [Experiment] FCT comparison against ECN-based schemes with SPQ (1 queue)/DRR (4 queues).

10 Gbps and 100 Gbps. The base RTTs are  $84\mu s$  and  $40\mu s$  for each of the links. In the switch, we configure WRR with equal weight for packet scheduling. We have enabled the Jumbo frame for 100Gbps links. We consider Broadcom Trident+ and Trident 3 ASICs for 10 Gbps and 100 Gbps links, respectively. Therefore, each port has 192 KB and 1 MB of buffers, which are completely shared by service queues, except in PQL. We use TCP for the transport protocol and set  $RTO_{min}$  to 5ms, the lowest stable value in `jitfy` timer [23].

**Impact of Link Capacity.** We first perform a simulation with 10 Gbps links. We have 8 services, which are mapped to each of 8 service queues. There exist  $2 \times i$  senders for queue  $i$ . Every sender of service queues starts a flow at the beginning simultaneously. From 200 ms, the senders of queues 2 ~ 8 stop their transmissions every 50 ms in order. For example, the senders of queue 3 finish flows at 250 ms, and queue 1 is the only active queue after 500 ms. We measure per-queue throughput every 10 ms. Using the measured data, we calculate Jain’s fairness index between active queues and the aggregate throughput. The aggregate throughput is to inspect whether the link is fully utilized at any time. If a solution preserves weighted fair sharing and work conservation, the two metrics always should be high.

Fig. 11 shows the bandwidth sharing results with 10 Gbps links. Not surprisingly, DynaQ and PQL achieve the near-optimal fairness index while BestEffort causes fluctuations. For example, the fairness index plunges to 0.67 at 410 ms. This is because BestEffort cannot handle service queues with many flows. Fig. 11b shows that DynaQ is the only solution that maintains high link utilization between the schemes. When queue 8 finishes at 500 ms, PQL causes a huge throughput collapse. PQL maintains the aggregate

throughput around 8.5 Gbps after 500 ms. This is because queue 1, the only active service queue, cannot occupy the buffer as much as the BDP.

We perform the same simulation with 100 Gbps and Fig. 12 reports the results. We find that DynaQ can achieve work conservation and preserve weighted fair sharing with a high link capacity. The overall tendency is very similar to the 10 Gbps results. BestEffort cannot provide per-queue fairness and PQL loses a significant amount of throughput when service queue 1 is the only active queue. In addition, DynaQ does not lose throughput much at 500 ms. Unlike DynaQ, BestEffort causes 9.2 Gbps of throughput loss.

**Impact of Traffic Dynamics.** We now inspect how DynaQ is robust to traffic dynamics. We conduct a simulation with the almost same scenario as the previous simulation with 100 Gbps. The only different setting is the number of senders per service queue. We consider a very extreme scenario where service queue  $i$  has  $2^{(3+i)}$  senders generating a single flow. For example, queue 8 has 2048 flows. Fig. 13 plots the results. We observe that DynaQ is robust to the extreme traffic scenario. BestEffort shows the worst performance in the fairness index that it achieves only 0.24 of fairness index for the first 200 ms. The scheme also loses throughput at 300 ms for a moment. PQL still fails to achieve work conservation. The aggregate throughput stays below 94.5 Gbps from 500 ms.

5.2.2 Dynamic Flow Simulations

**Methodology.** We build a non-blocking leaf-spine topology, a widely used data center network topology design. Our topology has 12 leaf (ToR) switches and 12 spine (Core)

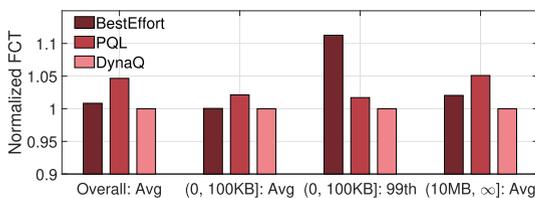


Fig. 10. [Experiment] FCT comparison with DRR (4 queues).

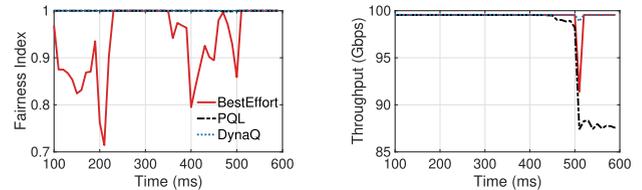


Fig. 12. [Simulation] Bandwidth sharing on 100Gbps links.

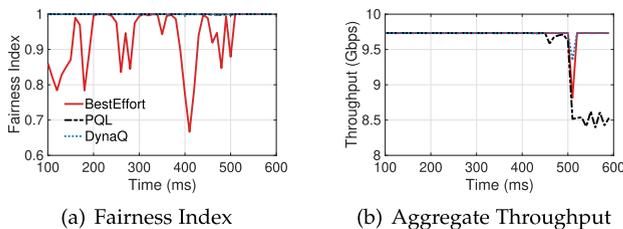


Fig. 11. [Simulation] Bandwidth sharing on 10Gbps links.

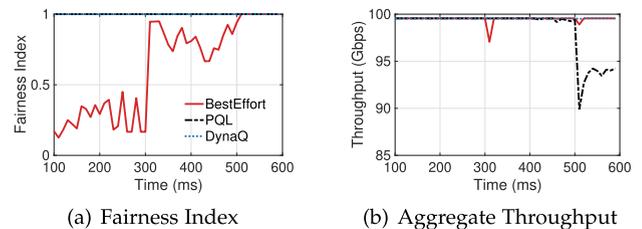


Fig. 13. [Simulation] Bandwidth sharing with many flows.

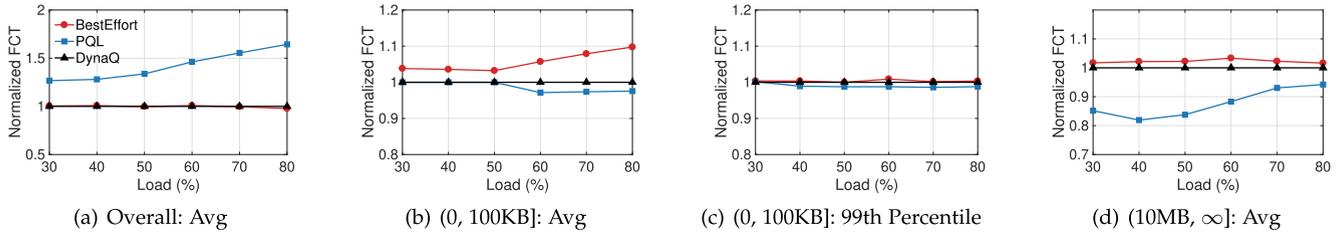


Fig. 14. [Simulation] FCT comparison against non-ECN schemes with SPQ (1 queue)/DRR (7 queues).

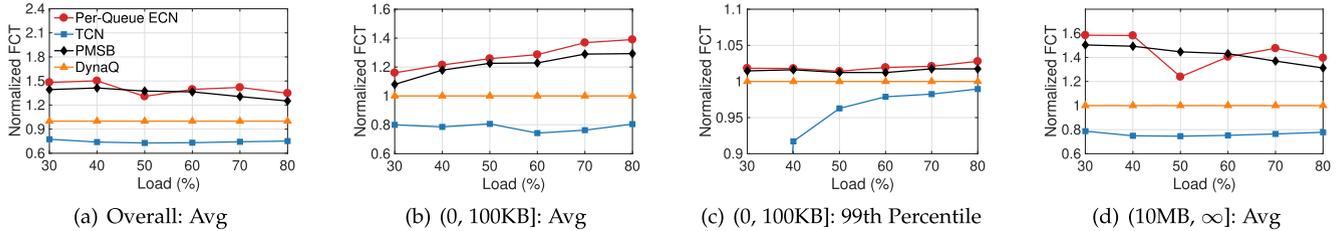


Fig. 15. [Simulation] FCT comparison against ECN-based schemes with SPQ (1 queue)/DRR (7 queues).

switches. Each leaf switch has  $12 \times 10$  Gbps downlinks and  $12 \times 10$  Gbps uplinks. The base RTT across spine switches is  $85.2 \mu s$ . We use ECMP as the load balancing scheme [26]. We use TCP and  $RTO_{min}$  is 5ms [23]. Each switch port has a 192KB buffer [24] shared by service queues except in PQL.

We configure SPQ/DRR with 8 service queues where one queue is the shared high priority queue and the other 7 queues are dedicated DRR queues with equal quantum. Like the testbed experiments, we also employ a two-level PIAS [25] whose priority demotion threshold is 100 KB to classify small flows from large flows. We evenly classify the  $144 \times 143$  communication pairs into 7 services and each service has its own service queue. Different services use different traffic distributions in Fig. 2. We generate 10K flows by varying traffic load from 30% to 80%.

*Comparison With Non-ECN Schemes.* Fig. 14 compares DynaQ against non-ECN schemes. In the average FCT for overall flows, we observe that DynaQ has mixed results compared to BestEffort and outperforms PQL. This is similar to the results in the testbed experiments. The gaps with BestEffort are within  $0.98 \times \sim 1.01 \times$ . For the average FCT of small flows, we can see the mixed results. DynaQ beats BestEffort across the loads but underperforms PQL from 60% to 80%. However, the gaps are only within  $0.98 \times \sim 1.02 \times$ . For the 99th percentile FCT of small flows, we can see that DynaQ achieves similar performance to the compared schemes. PQL slightly outperforms DynaQ that the maximum gap is  $0.98 \times$ . In the testbed experiments with 1Gbps links, BestEffort results in performance degradation when loads are high. However, with 10Gbps links, DynaQ beats BestEffort only by up to  $1.01 \times$  due to the increased link capacity. For the average FCT of large flows, DynaQ is better than BestEffort but is worse than PQL. The maximum gap is  $0.82 \times$  at 40% of traffic load. We suspect that this is because PQL does not share buffer space between the queues in the leaf and spine switches.

*Comparison With ECN-Based Schemes.* Fig. 15 shows the FCT results with ECN-based schemes. It is easy to find that DynaQ outperforms Per-Queue ECN and PMSB across all the flow sizes, and this is similar to the results in the testbed experiments. One difference from the testbed experiments is that TCN outperforms DynaQ. For the average FCT of

overall flows, DynaQ underperforms TCN by up to  $0.73 \times$ . The results for the average FCT of small and large flows are also similar. Specifically, DynaQ is worse than TCN by up to  $0.77 \times$  and  $0.75 \times$  in the average FCT of small and large flows, respectively. For the 99th percentile FCT of small flows, the gap between DynaQ and TCN decreases as traffic load increases. We omit the result at 30% of  $0.10 \times$  because it harms the visibility of the figure. This is because TCN can provide low tail latency, especially in a low traffic load. Despite its superior performance in the FCT, TCN does not support generic transport protocols.

## 6 RELATED WORK

We briefly review existing works addressing buffer sharing in switches. LossPass [13], EDT [27], CEDM [28], and DTS [29] propose buffer sharing mechanisms to absorb microburst traffic. DTS reserves a dedicated buffer for each service queue and utilizes shared buffers to absorb microbursts. In terms of service queue isolation, DTS is similar to PQL because it statically reserves a buffer size, leading to violating work conservation when few queues active. Aelous [30] and BCC [24] concern buffer sharing with ECN. These solutions are unaware of service queue isolation.

There exist per-port buffer sharing solutions like the dynamic threshold algorithm [31]. In addition, commodity switches allow a port to occupy many buffers [32]. However, even we allocate a large buffer to a port, bandwidth cannot be shared fairly since aggressive queues eventually fill up the buffer. It also harms per-port fairness by taking excessive buffers that can be assigned to the other ports. Meanwhile, there exist deep buffered switches with an external large DRAM buffer. Unfortunately, the switches have low switching throughput and insufficient processing delay to satisfy the requirements of user-facing applications.

## 7 CONCLUSION

This paper proposed DynaQ, a protocol-independent multi-queue management solution that can isolate service queues with generic transport protocols through dynamic packet

dropping thresholds. We have discussed how DynaQ can be implemented on hardware. To compare DynaQ with existing solutions, we have implemented a software prototype of DynaQ as a Linux qdisc module. Our extensive experiment and simulation results demonstrated that DynaQ ensures work conservation, weighted fair sharing, and low latency without protocol dependency.

## ACKNOWLEDGMENTS

This work was supported in part by the National Research Foundation of Korea (NRF) Grant funded by the Ministry of Science and ICT under Grant 2019R1A2C2088812. A preliminary version of this article appeared in the Proceedings of the 40th IEEE International Conference on Distributed Computing Systems (IEEE ICDCS 2020) [1].

## REFERENCES

- G. Kim and W. Lee, "Protocol-independent service queue isolation for multi-queue data centers," in *Proc. IEEE 40th Int. Conf. Distrib. Comput. Syst.*, 2020, pp. 355–365.
- W. Bai, L. Chen, K. Chen, and H. Wu, "Enabling ECN in multi-service multi-queue data centers," in *Proc. 13th USENIX Conf. Netw. Syst. Des. Implementation*, 2016, pp. 537–549.
- W. Bai, K. Chen, L. Chen, C. Kim, and H. Wu, "Enabling ECN over generic packet scheduling," in *Proc. 12th Int. Conf. Emerg. Netw. Experiments Technol.*, 2016, pp. 191–204.
- Y. Pan *et al.*, "Support ECN in multi-queue datacenter networks via per-port marking with selective blindness," in *Proc. IEEE 38th Int. Conf. Distrib. Comput. Syst.*, 2018, pp. 33–42.
- Y. Li *et al.*, "HPCC: High precision congestion control," in *Proc. ACM Special Interest Group Data Commun.*, 2019, pp. 44–58.
- I. Cho, K. Jang, and D. Han, "Credit-scheduled delay-bounded congestion control for datacenters," in *Proc. ACM Special Interest Group Data Commun.*, 2017, pp. 239–252.
- R. Mittal *et al.*, "TIMELY: RTT-based congestion control for the datacenter," in *Proc. ACM Conf. Special Interest Group Data Commun.*, 2015, pp. 537–550.
- C. Lee, C. Park, K. Jang, S. Moon, and D. Han, "Accurate latency-based congestion feedback for datacenters," in *Proc. USENIX Conf. USENIX Annu. Techn. Conf.*, 2015, pp. 403–415.
- Mellanox MLNX-OS user manual for ethernet, Sunnyvale, CA, USA, Mellanox Technologies, White Paper, 2017. [Online]. Available: [https://www.mellanox.com/related-docs/prod\\_management\\_software/MLNX-OS\\_ETH\\_v3\\_6\\_3508\\_UM.pdf](https://www.mellanox.com/related-docs/prod_management_software/MLNX-OS_ETH_v3_6_3508_UM.pdf)
- Tofino Programmable Switch. Accessed: Sep. 11, 2021. [Online]. Available: <https://www.barefootnetworks.com/>
- The Network Simulator ns-2. Accessed: Sep. 11, 2021. [Online]. Available: <http://www.isi.edu/nsnam/ns/>
- M. Alizadeh *et al.*, "pFabric: Minimal near-optimal datacenter transport," in *ACM SIGCOMM Conf. SIGCOMM*, 2013, pp. 435–446.
- G. Kim and W. Lee, "Absorbing microbursts without headroom for data center networks," *IEEE Commun. Lett.*, vol. 23, no. 5, pp. 806–809, May 2019.
- H. Wu, J. Ju, G. Lu, C. Guo, Y. Xiong, and Y. Zhang, "Tuning ECN for data center networks," in *Proc. 8th Int. Conf. Emerg. Netw. Experiments Technol.*, 2012, pp. 25–36.
- M. Alizadeh *et al.*, "Data center TCP (DCTCP)," *ACM SIGCOMM Comput. Commun. Rev.*, vol. 40, no. 4, pp. 63–74, 2010.
- L. Zheng, Z. Qiu, S. Sun, W. Pan, Y. Gao, and Z. Zhang, "Design and analysis of a parallel hybrid memory architecture for per-flow buffering in high-speed switches and routers," *J. Commun. Netw.*, vol. 20, no. 6, pp. 578–592, Dec 2018.
- P. Goyal, P. Shah, N. K. Sharma, M. Alizadeh, and T. E. Anderson, "Backpressure flow control," 2019, *arXiv:1909.09923*.
- H. Zhu *et al.*, "Racksched: A microsecond-scale scheduler for rack-scale computers," in *Proc. 14th USENIX Symp. Operating Syst. Des. Implementation*, 2020, pp. 1225–1240.
- Advanced Congestion & Flow Control With Programmable Switches, Apr. 2020. [Online]. Available: <https://opennetworking.org/wp-content/uploads/2020/04/JK-Lee-Slide-Deck.pdf>
- A. Greenberg *et al.*, "V12: A scalable and flexible data center network," *ACM SIGCOMM Comput. Commun. Rev.*, vol. 39, no. 4, pp. 51–62, 2009.
- A. Roy, H. Zeng, J. Bagga, G. Porter, and A. C. Snoeren, "Inside the social network's (datacenter) network," in *Proc. ACM Conf. Special Interest Group Data Commun.*, 2015, pp. 123–137.
- G. Judd, "Attaining the promise and avoiding the pitfalls of TCP in the datacenter," in *Proc. 12th USENIX Conf. Netw. Syst. Des. Implementation*, 2015, pp. 145–157.
- V. Vasudevan *et al.*, "Safe and effective fine-grained TCP retransmissions for datacenter communication," *ACM SIGCOMM Comput. Commun. Rev.*, vol. 39, no. 4, pp. 303–314, 2009.
- W. Bai, S. Hu, K. Chen, K. Tan, and Y. Xiong, "One more config is enough: Saving (DC)TCP for high-speed extremely shallow-buffered datacenters," in *Proc. IEEE INFOCOM IEEE Conf. Comput. Commun.*, 2020, pp. 2007–2016.
- W. Bai, K. Chen, H. Wang, L. Chen, D. Han, and C. Tian, "Information-agnostic flow scheduling for commodity data centers," in *Proc. 12th USENIX Conf. Netw. Syst. Des. Implementation*, 2015, pp. 455–468.
- M. Park, S. Sohn, K. Kwon, and T. T. Kwon, "MaxPass: Credit-based multipath transmission for load balancing in data centers," *J. Commun. Netw.*, vol. 21, no. 6, pp. 558–568, Dec. 2019.
- D. Shan, W. Jiang, and F. Ren, "Absorbing micro-burst traffic by enhancing dynamic threshold policy of data center switches," in *Proc. IEEE Conf. Comput. Commun.*, 2015, pp. 118–126.
- D. Shan and F. Ren, "Improving ECN marking scheme with micro-burst traffic in data center networks," in *Proc. IEEE Conf. Comput. Commun.*, 2017, pp. 1–9.
- Cisco catalyst 9000 switching platforms: QoS and queuing white paper, San Jose, CA, USA, White Paper, 2020. [Online]. Available: <https://www.cisco.com/c/en/us/products/collateral/switches/catalyst-9000/white-paper-c11-742388.html>
- S. Hu *et al.*, "Aeolus: A building block for proactive transport in datacenters," in *Proc. Annu. Conf. ACM Special Interest Group Data Commun. Appl. Technol. Architectures Protoc. Comput. Commun.*, 2020, pp. 422–434.
- A. K. Choudhury and E. L. Hahne, "Dynamic queue length thresholds for shared-memory packet switches," *IEEE/ACM Trans. Netw.*, vol. 6, no. 2, pp. 130–140, Apr. 1998.
- Arista 7050x3 series switch architecture, Santa Clara, CA, USA, White Paper, 2018. [Online]. Available: [https://www.arista.com/assets/data/pdf/Whitepapers/7050X3\\_Architecture\\_WP.pdf](https://www.arista.com/assets/data/pdf/Whitepapers/7050X3_Architecture_WP.pdf)



**Gyuyeong Kim** (Member, IEEE) received the BS and PhD degrees in computer science from Korea University, South Korea, in 2012 and 2020, respectively. He is currently a research professor with Future Network Center, Korea University. His research interests include networked systems and networking support for AI, big data, cloud, and IoT systems.



**Wonjun Lee** (Fellow, IEEE) received the BS and MS degrees in computer engineering from Seoul National University, Seoul, South Korea, in 1989 and 1991, respectively, the MS degree in computer science from the University of Maryland, College Park, MD, USA, in 1996, and the PhD degree in computer science and engineering from the University of Minnesota, Minneapolis, MN, USA, in 1999. In 2002, he joined the faculty of Korea University, Seoul, South Korea, where he is currently a professor with the School of Cybersecurity. He has authored or co-authored more than 220 papers in refereed international journals and conferences. His research interests include communication and network protocols, optimization techniques in wireless communication and networking, security and privacy in mobile computing, and RF-powered computing and networking. He was on the TPC and an organizing committee member of IEEE INFOCOM from 2008 to 2021, the PC Vice Chair of IEEE ICDCS 2019, and ACM MobiHoc from 2008 to 2009, and more than 130 international conferences.

▷ For more information on this or any other computing topic, please visit our Digital Library at [www.computer.org/csdl](http://www.computer.org/csdl).