

#### Enabling ECN in Multi-Service Multi-Queue Data Centers

#### Wei Bai, Li Chen, Kai Chen, Haitao Wu (Microsoft)

SING Group @ Hong Kong University of Science and Technology

## Background

• Data Centers

- Many services with diverse network requirements



## Background

Data Centers

- Many services with diverse network requirements

• ECN-based Transports

ECN = Explicit Congestion Notification

## Background

- Data Centers
  - Many services with diverse network requirements
- ECN-based Transports
  - Achieve high throughput & low latency
  - Widely deployed: DCTCP, DCQCN, etc.





• ECN-enabled end-hosts

- React to ECN by adjusting sending rates





• ECN-enabled end-hosts

- React to ECN by adjusting sending rates

- ECN-aware switches
  - Perform ECN marking based on Active Queue
     Management (AQM) policies



ECN-enabled end-hosts

- React to ECN by adjusting sending rates

- ECN-aware switches
  - Perform ECN marking based on Active Queue Management (AQM) policies



• Adopt RED to perform ECN marking

#### RED = Random Early Detection

Adopt RED to perform ECN marking
 – Per-queue/port/service-pool ECN/RED

#### Track buffer occupancy of different egress entities

 Adopt RED to perform ECN marking — Per-queue/port/service-pool ECN/RED



Adopt RED to perform ECN marking

 Per-queue/port/service-pool ECN/RED



 Adopt RED to perform ECN marking — Per-queue/port/service-pool ECN/RED



- Adopt RED to perform ECN marking — Per-queue/port/service-pool ECN/RED
- Leverage multiple queues to classify traffic
  - Isolate traffic from different services/applications



- Adopt RED to perform ECN marking — Per-queue/port/service-pool ECN/RED
- Leverage multiple queues to classify traffic
  - Isolate traffic from different services/applications



- Adopt RED to perform ECN marking — Per-queue/port/service-pool ECN/RED
- Leverage multiple queues to classify traffic
  - Isolate traffic from different services/applications



- Adopt RED to perform ECN marking — Per-queue/port/service-pool ECN/RED
- Leverage multiple queues to classify traffic
  - Isolate traffic from different services/applications
  - Weighted max-min fair sharing among queues



- Adopt RED to perform ECN marking — Per-queue/port/service-pool ECN/RED
- Leverage multiple queues to classify traffic
  - Isolate traffic from different services/applications
  - Weighted max-min fair sharing among queues

#### Perform ECN marking in multi-queue context



**RED** Algorithm



• To achieve 100% throughput



• To achieve 100% throughput

$$K \ge C \times RTT \times \lambda$$

Determined by congestion control algorithms

• To achieve 100% throughput



#### Standard ECN marking threshold

• To achieve 100% throughput



The standard threshold is relatively stable in DCN, e.g., 65 packets for 10G network (DCTCP paper)

• Per-queue with the standard threshold

 $-K_{queue(i)} = C \times RTT \times \lambda$ 



• Per-queue with the standard threshold

$$-K_{queue(i)} = C \times RTT \times \lambda$$

Increase packet latency



• Per-queue with the standard threshold

$$-K_{queue(i)} = C \times RTT \times \lambda$$

Increase packet latency



Evenly classify 8 long-lived flows into a varying number of queues

Per-queue with the minimum threshold

$$-K_{queue(i)} = C \times RTT \times \lambda \times w_i / \sum w_j$$

Normalized weight



• Per-queue with the minimum threshold

$$-K_{queue(i)} = C \times RTT \times \lambda \times w_i / \sum w_j$$

– Degrade throughput



• Per-queue with the minimum threshold

 $-K_{queue(i)} = C \times RTT \times \lambda \times w_i / \sum w_j$ 

– Degrade throughput



• Per-port

$$-K_{port} = C \times RTT \times \lambda$$



• Per-port

$$-K_{port} = C \times RTT \times \lambda$$

- Violate weighted fair sharing



• Per-port

$$-K_{port} = C \times RTT \times \lambda$$

#### - Violate weighted fair sharing



Both services have a equal-weight dedicated queue on the switch

## Question

- Can we design an ECN marking scheme with following properties:
  - Deliver low latency
  - Achieve high throughput
  - Preserve weighted fair sharing
  - Compatible with legacy ECN/RED implementation

## Question

- Can we design an ECN marking scheme with following properties:
  - Deliver low latency
  - Achieve high throughput
  - Preserve weighted fair sharing
  - Compatible with legacy ECN/RED implementation

#### Our answer: MQ-ECN

## **MQ-ECN'S DESIGN**

• N queues share the link with capacity C



• N queues share the link with capacity C

Input Rate

 $r_1 \\ r_2 \\ \vdots$ 

 $r_N$ 



• N queues share the link with capacity C

#### Input Rate Weight



- N queues share the link with capacity C
- $C = \sum_{i=1}^{N} \min(r_i, w_i \alpha)$  Weighted Fair Share Rate



- N queues share the link with capacity C
- $C = \sum_{i=1}^{N} \min(r_i, w_i \alpha)$
- Use ECN to throttle queue i if  $r_i > w_i \alpha$



- N queues share the link with capacity C
- $C = \sum_{i=1}^{N} \min(r_i, w_i \alpha)$
- $K_{queue(i)} = w_i \alpha \times RTT \times \lambda$



- N queues share the link with capacity C
- $C = \sum_{i=1}^{N} \min(r_i, w_i \alpha)$
- $K_{queue(i)} = w_i \alpha \times RTT \times \lambda$

bit-by-bit round robin



- N queues share the link with capacity C
- $C = \sum_{i=1}^{N} \min(r_i, w_i \alpha)$
- $K_{queue(i)} = w_i \alpha \times RTT \times \lambda$

bit-by-bit round robin



- N queues share the link with capacity C
- $C = \sum_{i=1}^{N} \min(r_i, w_i \alpha)$
- $K_{queue(i)} = w_i \alpha \times RTT \times \lambda$

Time of a round: *T<sub>round</sub>* 



- N queues share the link with capacity C
- $C = \sum_{i=1}^{N} \min(r_i, w_i \alpha)$
- $K_{queue(i)} = quantum_i / T_{round} \times RTT \times \lambda$



•  $K_{queue(i)} = quantum_i / T_{round} \times RTT \times \lambda$ 

•  $K_{queue(i)} = quantum_i / T_{round} \times RTT \times \lambda$ 

## Why does it work?

- $K_{queue(i)} = quantum_i / T_{round} \times RTT \times \lambda$ 
  - Deliver low latency
  - Achieve high throughput

#### $K_{queue(i)}$ adapts to traffic dynamics

- $K_{queue(i)} = quantum_i / T_{round} \times RTT \times \lambda$ 
  - Deliver low latency
  - Achieve high throughput
  - Preserve weighted fair sharing

 $K_{queue(i)}$  is in proportion to the weight

- $K_{queue(i)} = quantum_i / T_{round} \times RTT \times \lambda$ 
  - Deliver low latency
  - Achieve high throughput
  - Preserve weighted fair sharing
  - Compatible with legacy ECN/RED implementation

Per-queue ECN/RED with dynamic thresholds

- $K_{queue(i)} = quantum_i / T_{round} \times RTT \times \lambda$ 
  - Deliver low latency
  - Achieve high throughput
  - Preserve weighted fair sharing
  - Compatible with legacy ECN/RED implementation
- More details
  - Handle inaccurate estimation of  $T_{round}$
  - Apply to round-robin packet schedulers
    - DWRR, WRR, etc.

## **Testbed Evaluation**

- MQ-ECN software prototype
  - Linux qdisc kernel module performing DWRR
- Testbed setup
  - 9 servers are connected to a server-emulated switch with 9 NICs
  - End-hosts use DCTCP as the transport protocol
- Benchmark traffic

Web search (DCTCP paper)

• More results in large-scale simulations

#### **Static Flow Experiment**



#### **Static Flow Experiment**



#### **Static Flow Experiment**



MQ-ECN preserves weighted fair sharing

#### Realistic Traffic: Small Flows (<100KB)



### Realistic Traffic: Small Flows (<100KB)



**MQ-ECN** achieves low latency

#### Realistic Traffic: Large Flows (>10MB)



#### Realistic Traffic: Large Flows (>10MB)



#### MQ-ECN achieves high throughput

## Conclusions

- Identify performance impairments of existing ECN/RED schemes in multi-queue context
- MQ-ECN: simple, yet effective

   High throughput, low latency, weighted fair sharing
   Currently, apply to round-robin packet schedulers
- ECN marking scheme for arbitrary schedulers is an important ongoing effort
- Code: <u>http://sing.cse.ust.hk/projects/MQ-ECN</u>

Thanks!

## Support Arbitrary Packet Schedulers

- Find a general solution to estimate per-queue *effective* draining rate
  - Draining rate should be measured when the queue keeps congested
    - $K_{queue(i)} = w_i \alpha \times RTT \times \lambda$
  - Measurement window is hard to decide