# Congestion Management – Congestion Isolation

Paul Congdon

Yolanda Yu

Kevin Shen

paul.congdon@tallac.com

yolanda.yu@huawei.com

kevin.shenli@huawei.com

IEEE 802.1 DCB

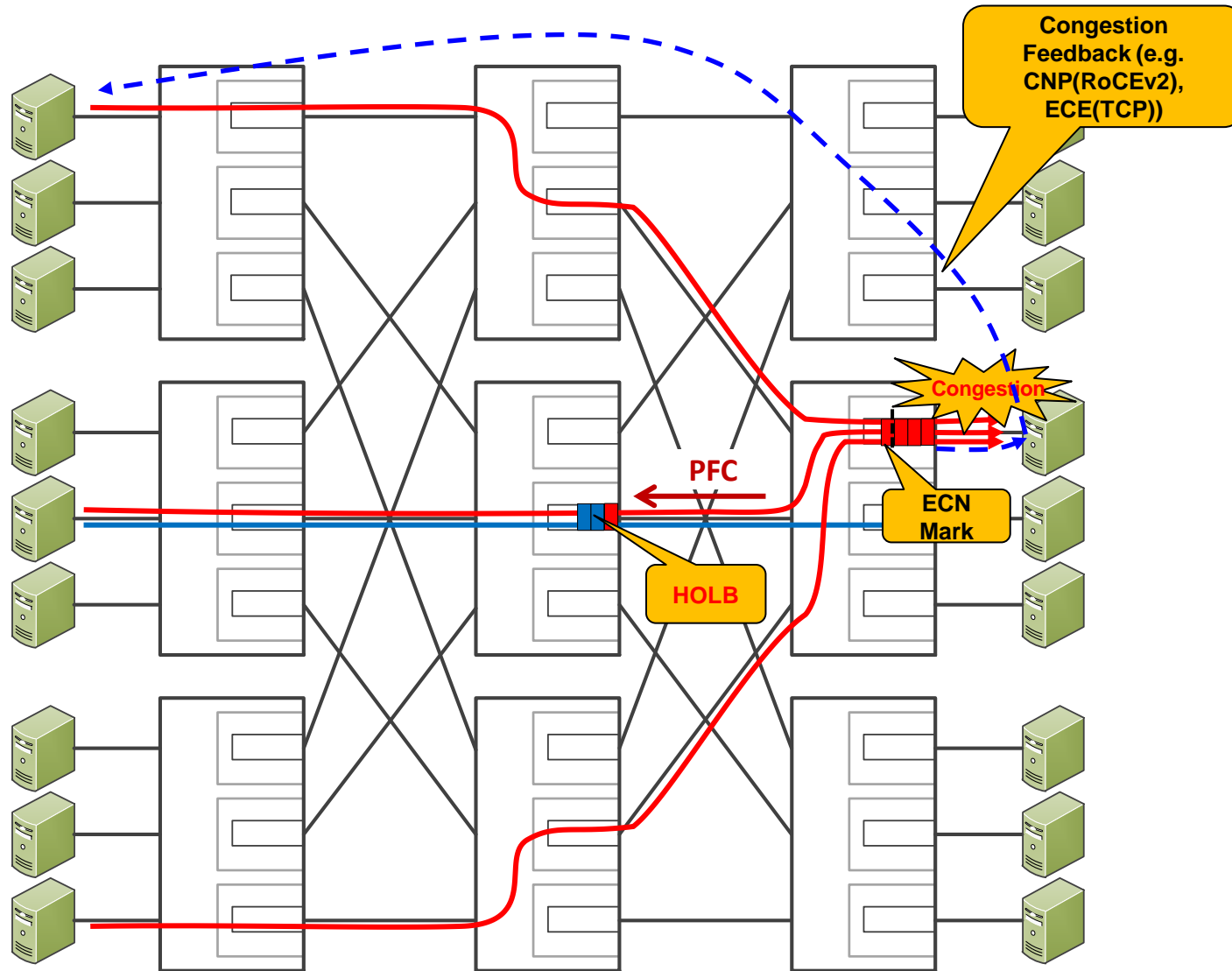St John's Newfoundland

September 2017

# Agenda

- Low-Latency, Lossless, Large-Scale DCNs

- Challenges going forward

- Solution Goals

- Congestion Isolation Details

- Simulation Analysis

- Next Steps

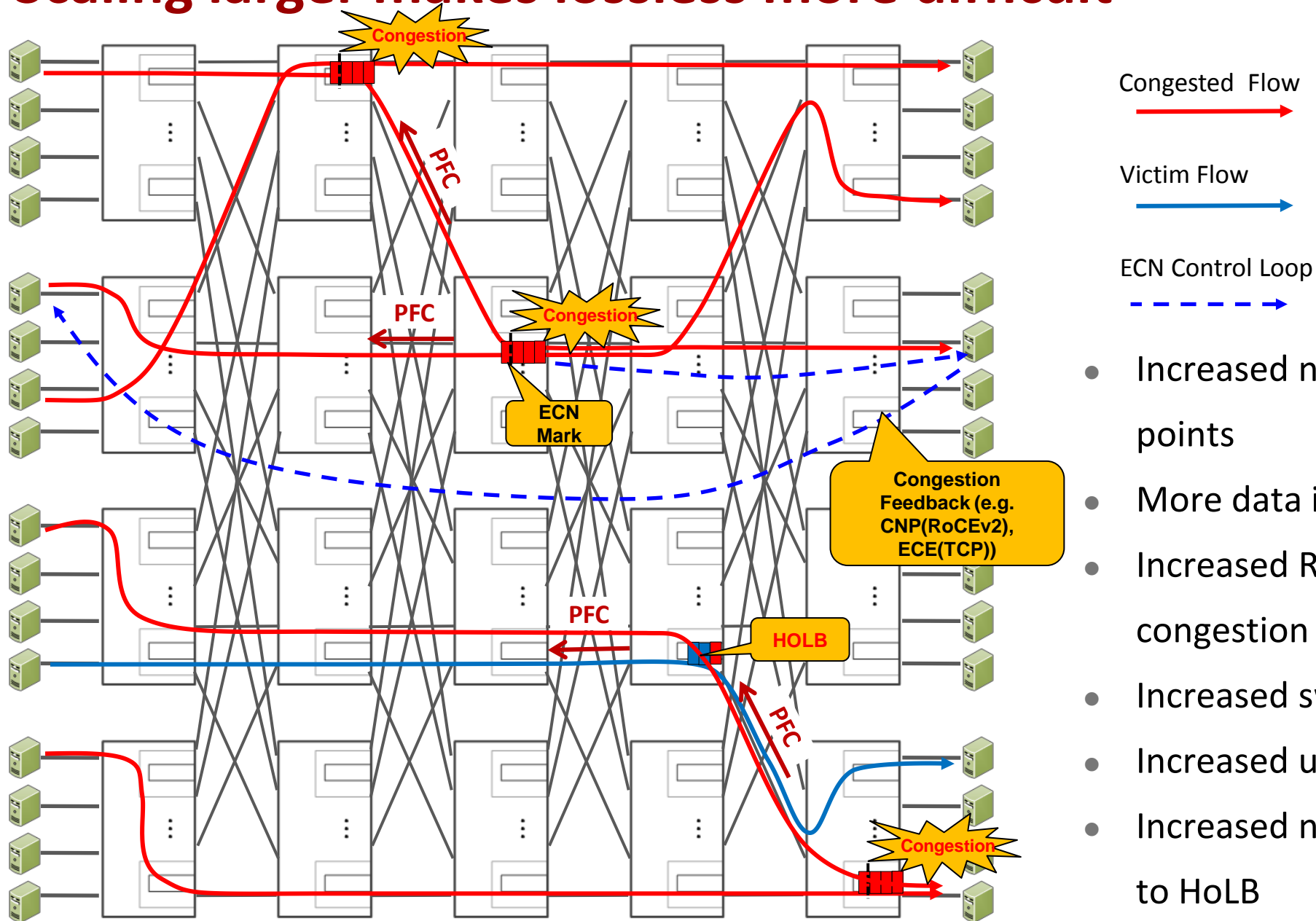# The Case for Low-latency, Lossless, Large-Scale DCNs

- More and more latency-sensitive applications are being deployed in data centers
  - Distributed Storage
  - AI / Deep Learning
  - Cloud HPC
  - High-Frequency Trading
- RDMA is operating at larger scales thanks to RoCEv2
  - Chuanxiong Guo, et. al., Microsoft, "RDMA over Commodity Ethernet at Scale", SIGCOMM 2016
  - Y Zhu, H Eran, et. al., Microsoft, Mellanox, "Congestion control for large-scale RDMA deployments", SIGCOMM 2015
  - Radhika Mittal, et. al., UC Berkeley, Google, "TIMELY: RTT-based Congestion Control for the Datacenter", SIGCOMM 2015
- The scale of Data Center Networks continues to grow
  - Larger, faster clusters are better than more smaller size clusters
  - Server growth continues at 25% - 30% putting pressure on cluster sizes and networking costs

# Lossless DCN state-of-the-art



Congested Flow

Victim Flow

ECN Control Loop

Congestion Feedback (e.g. CNP(RoCEv2), ECE(TCP))

Congestion

PFC

ECN Mark

HOLB

- DCN is primarily an L3 network
- ECN used for end-to-end congestion control
- Congestion feedback can be protocol and application specific
- PFC used as a last resort to ensure lossless environment, or not at all in low-loss environments.
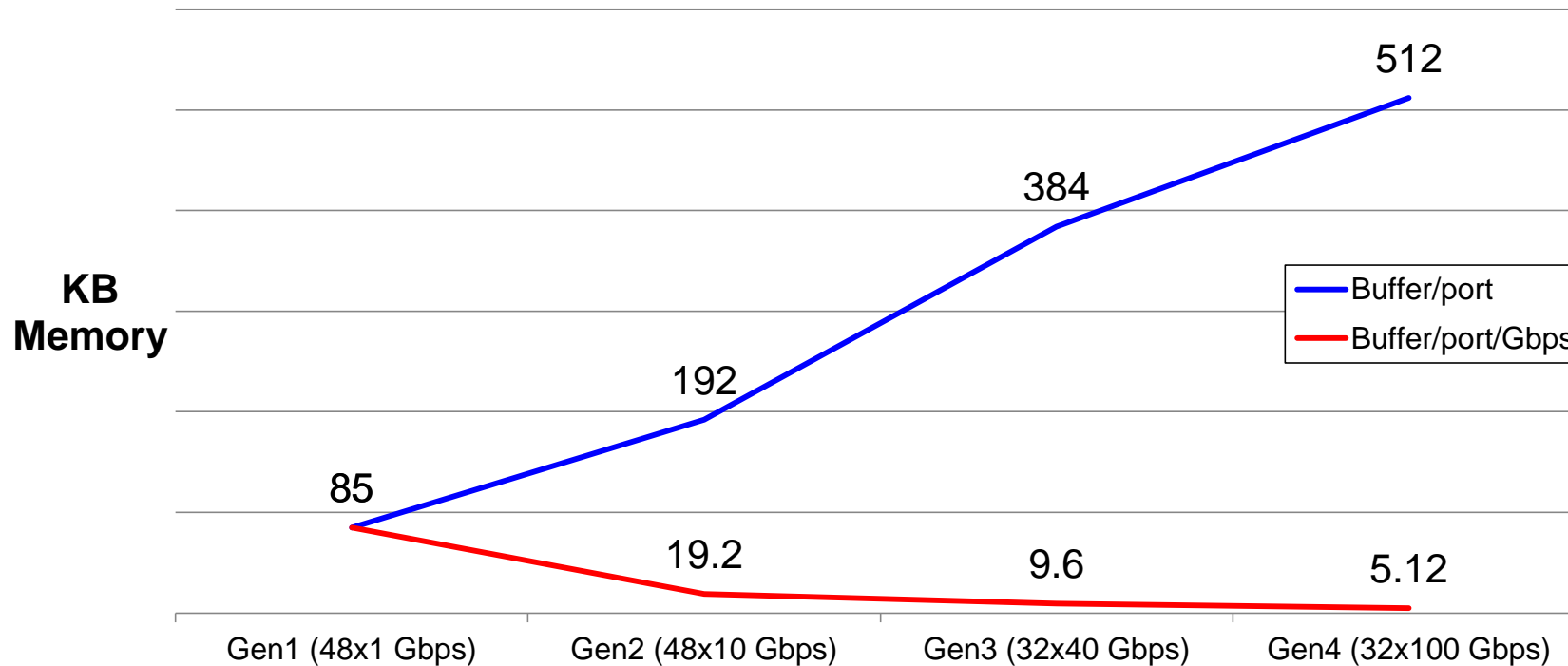- Traffic classes for PFC are mapped using DSCP as opposed to VLAN tags

# Scaling larger makes lossless more difficult



- Increased number of congestion points
- More data in-flight
- Increased RTT and delay for congestion feedback
- Increased switch buffer requirements
- Increased use of PFC
- Increased number of victim flows due to HoLB

# Switch buffer growth is not keeping up

**KB of Packet Buffer by Commodity Switch Architecture**



Commodity Shallow Buffer Switches in DCNs are desirable:
- Low Latency
- Low Cost

However, packet loss can create performance issues:
- Source: Broadcom, "White Paper: Buffer Requirements for Datacenter Network Switches", DNFAMILY-WP1101, August 25, 2015

Source: "Congestion Control for High-speed Extremely Shallow-buffered Datacenter Networks". In Proceedings of APNet'17, Hong Kong, China, August 03-04, 2017, https://doi.org/10.1145/3106989.3107003

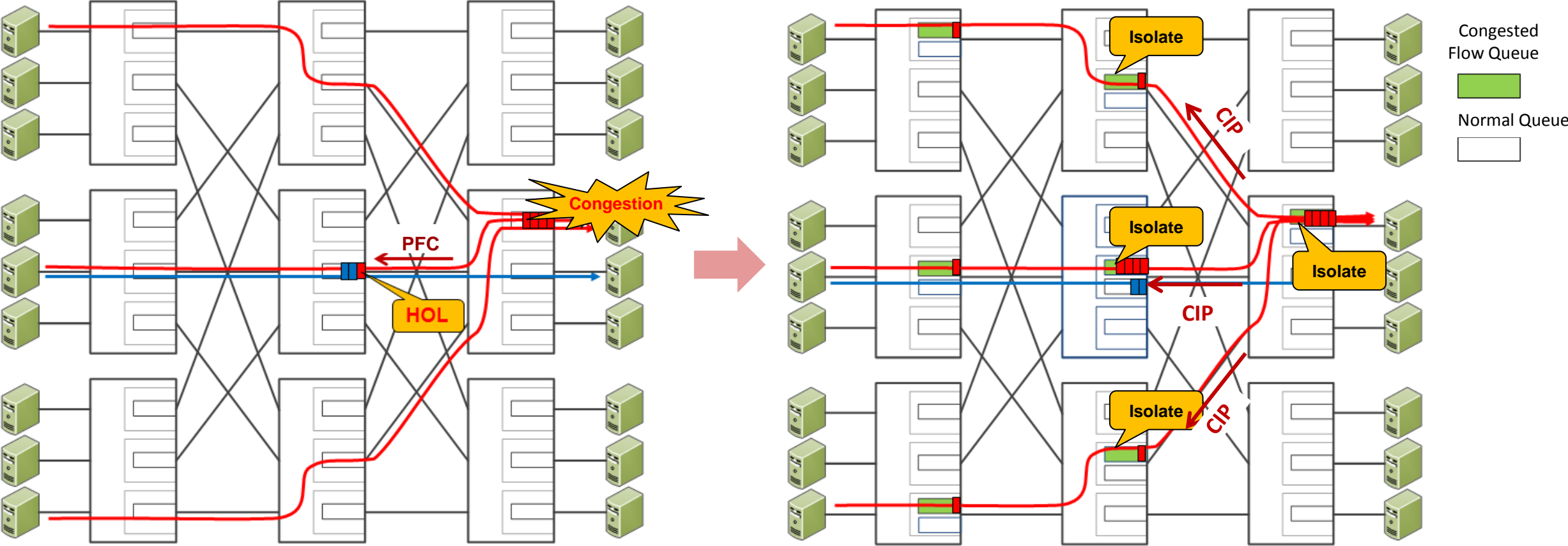# Concerns about over-using PFC

- HoL blocking

- Congestion spreading

- Buffer Bloat, increasing latency

- Increased jitter reducing throughput

- Deadlocks

# Goals

- Support larger, faster data centers (Low-Latency, High-Throughput)

- Support lossless transfers

- Improve performance of TCP and UDP based flows

- Reduce pressure on switch buffer growth

- Reduce the frequency of relying on PFC for a lossless environment


- Eliminate or significantly reduce HOLB caused by over-use of PFC

# Isolate the congestion to mitigate HOLB

# Congestion Isolation

**Definition:** An approach to isolate flows causing congestion and signal upstream to isolate the same flows to avoid head-of-line blocking.

The approach involves:

1.  Identifying the flows creating congestion (e.g. perhaps already done for QCN and/or ECN)
2.  Using implementation specific approaches to dynamically adjust the traffic class of offending flows without packet re-ordering (e.g. DVL – Dynamic Virtual Lanes)
3.  Signaling upstream indications via a Congestion Isolation Packet (CIP)

# Congestion Isolation with Dynamic Virtual Lanes

**Non-Congested Flow Queue：** Normal priority queues. Higher scheduling priority than Congested Flow Queue.
**Congested Flow Queue:** At least one of 8 priority queues. Lower scheduling priority than Non-Congested Flow Queue.
Scheduling assures no out-of-order packets with Non-Congested Flow Queue. There can be multiple congested flow queues
(use 5-tuple hash to map one).

Congested Flow

Non-Congested Flow

Congested Flow Queue

Non-Congested Flow Queue

Transmit Queue

Congested Flow Queue

Non-Congested Flow Queue

Transmit Queue

2

1  Switch B  3

4

2

1  Switch A  3

4

1) When congestion occurs, detect the congested flow, record it in the flow table.

# Congestion Isolation with Dynamic Virtual Lanes



**CIP: Congestion Isolation Packet**

3) When Congested Flow Queue exceed the threshold, send CIP (including the flow info, such as 5-tuple info) to upstream to isolate the congested flow.
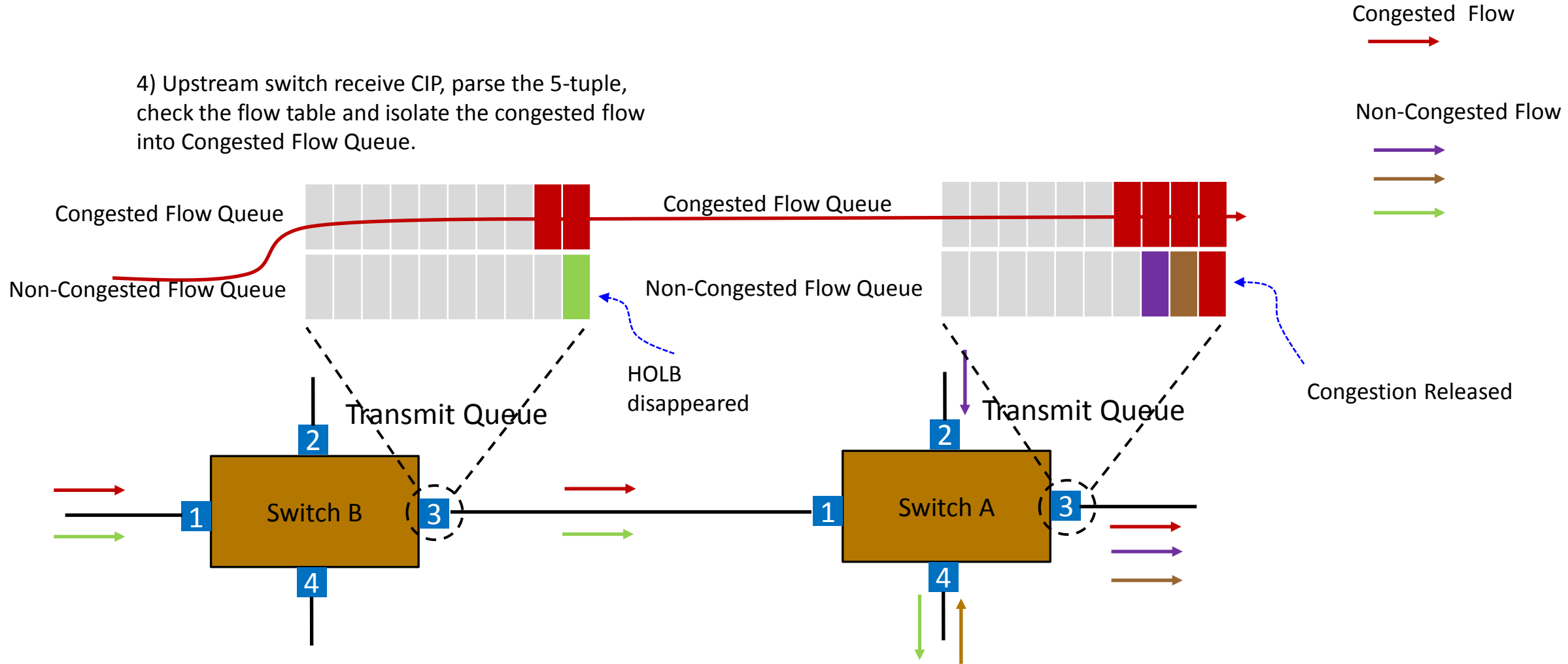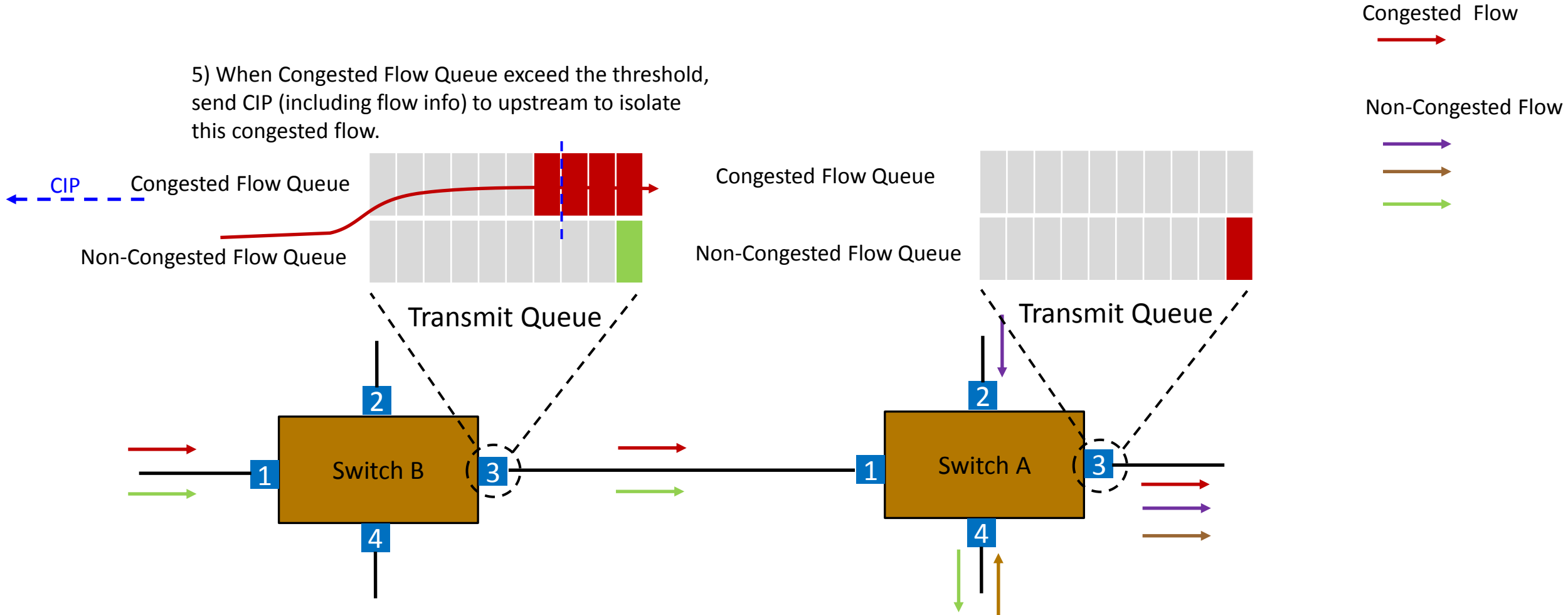
Congested Flow

Non-Congested Flow

CIP

Congested Flow Queue

Non-Congested Flow Queue

Transmit Queue

Congested Flow Queue

Non-Congested Flow Queue

Transmit Queue

2) Isolate the subsequent packets of congested flow to Congested Flow Queue.

Switch B

Switch A

1

2

3

4

1

2

3

4

# Congestion Isolation with Dynamic Virtual Lanes



4) Upstream switch receive CIP, parse the 5-tuple, check the flow table and isolate the congested flow into Congested Flow Queue.

Congested Flow

Non-Congested Flow

Congested Flow Queue

Non-Congested Flow Queue

Congested Flow Queue

Non-Congested Flow Queue

HOLB disappeared

Congestion Released

Transmit Queue

Transmit Queue

Switch B

Switch A

1

2

3

4

1

2

3

4

# Congestion Isolation with Dynamic Virtual Lanes

Congested Flow

Non-Congested Flow

5) When Congested Flow Queue exceed the threshold, send CIP (including flow info) to upstream to isolate this congested flow.

CIP

Congested Flow Queue

Non-Congested Flow Queue

Transmit Queue

Congested Flow Queue

Non-Congested Flow Queue

Transmit Queue

Switch B

Switch A

2

1

3

4

2

1

3

4

# Congestion Isolation with Dynamic Virtual Lanes

Congested Flow

Non-Congested Flow

6) When Congested Flow Queue exceed the high-level threshold, Queue Level Pause is triggered, such as PFC.

Congested Flow Queue

PFC Pause

Non-Congested Flow Queue

Congested Flow Queue

PFC Pause

Non-Congested Flow Queue

Transmit Queue

Transmit Queue

2

Switch B

1

3

4

2

Switch A

1

3

4

# Congestion Isolation Packet

- Objectives/Requirements:
    - Provide upstream neighbor with an indication that a flow has been isolate
    - Provide upstream neighbor with flow identification information
    - No adverse effects of single packet loss
    - Low overhead

**Format of Congestion Isolation Packet**

| |
|---|
| Dest MAC Address |
| Src MAC Address |
| Ethertype = 0x8809 |
| Flow Identification Data (TBD) |
| CRC |

Upstream Port Mac Address

Current Output Port Mac Address

New Ethernet Type

Flow identifying Information
(e.g IP Header, Transport Header,
Virtualization/Tunnel encapsulation).

# Handling the potential out-of-order problem

An instance：Red flow and purple flow are judged as congested flow and are moved to congested flow queue successively.

# Simulation Set-up



- OMNET++ Platform

- 2 Tier CLOS：100G interface with 200ns of link latency 200ns(about 40m)

- Scale：128～1152 servers, 24～72 switches

- Traffic Patterns:
  - Several regional all to all with some persistent incast
  - Flow size distribution is from 5 different real data center applications:
    - Enterprise IT, WebServer, Hadoop, Data Mining, Cache-Follower

- Compared Solutions:
  - PFC+ECN with CI: Congestion Isolation is implemented along with PFC+ECN
  - PFC+ECN without CI: Just PFC+ECN
  - All solutions include small flow prioritization mechanism

# PFC+ECN with CI VS. PFC+ECN without CI



Average flow completion time (all flows) — 25% reduction from Without CI to CI

Average flow completion time (>10MB flows) — 26% reduction

Average flow completion time (1MB~10MB flows) — 26% reduction

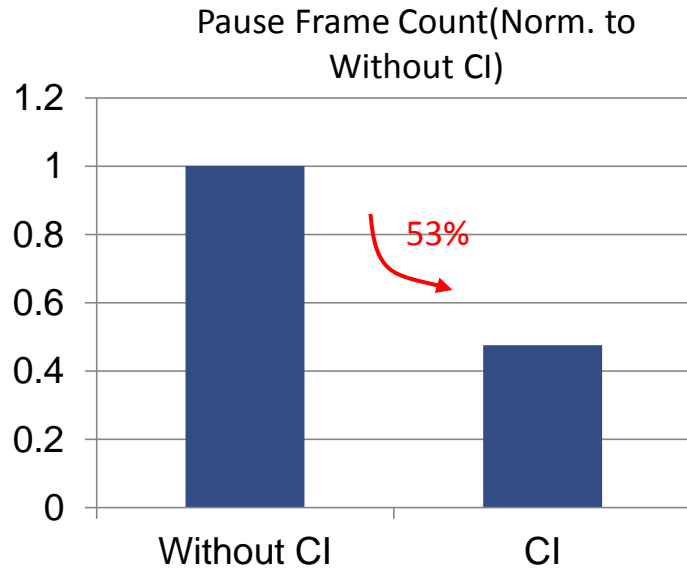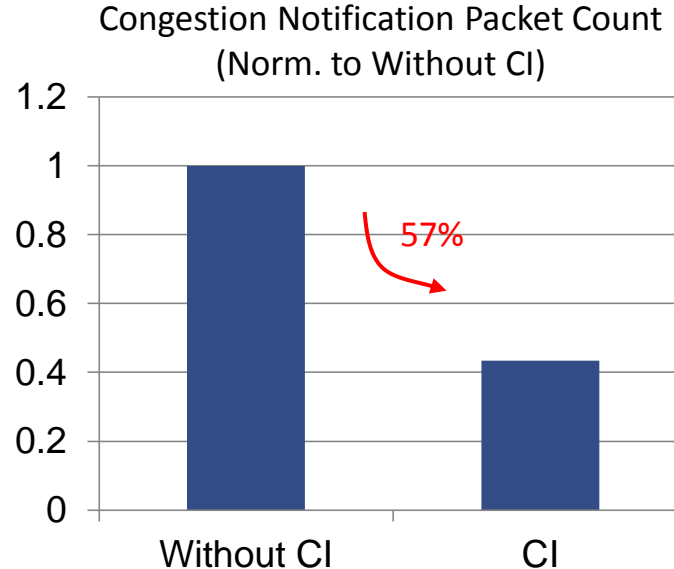Average flow completion time (100KB~1MB flows) — 36% reduction

Average flow completion time (<100KB flows) — 15% reduction

- CI reduces the count of PAUSE Frames sent to NICs of servers, so it can alleviate the HOL Blocking of the NIC, which can improve the performance of mice flows.
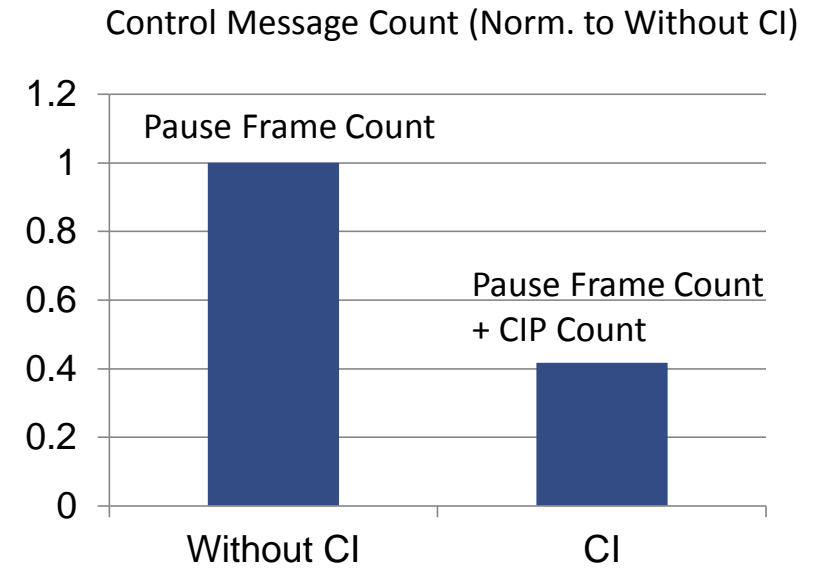- In the PFC+ECN without CI, we also prioritize the mice.

# Why PFC+ECN with CI outperforms PFC+ECN without CI



Pause Frame Count(Norm. to Without CI)
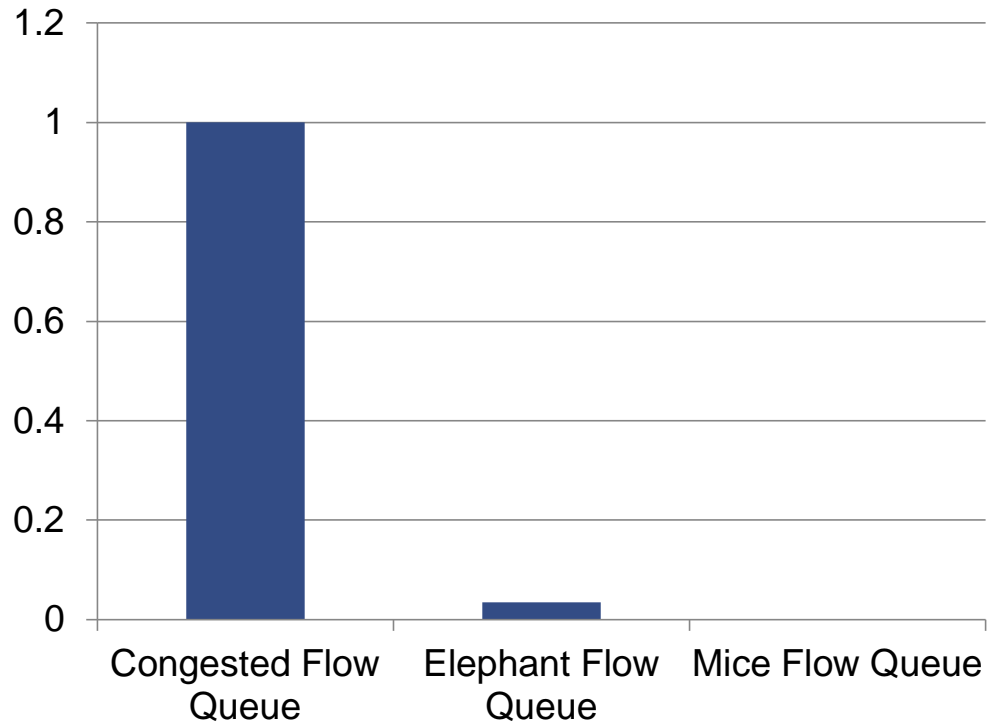
53%

Without CI | CI

Congestion Notification Packet Count (Norm. to Without CI)

57%

Without CI | CI

Control Message Count (Norm. to Without CI)

Pause Frame Count

Pause Frame Count + CIP Count

Without CI | CI

- CI reduces the pause frame count by 53%.

- CI reduces the CNP count by 57%.

- The count of new control message generated by CI is much less than the count it reduces the count of Pause frames.
- It has the same order-of-magnitude with large flow count.

# Why PFC+ECN with CI outperforms PFC+ECN without CI



Pause Frame Count Generated by Different Queues(Norm. to Congested Flow Queue)



Different flow count(Norm. to All Flow)

- 96.6% of the pause frames are generated by congested flow queues
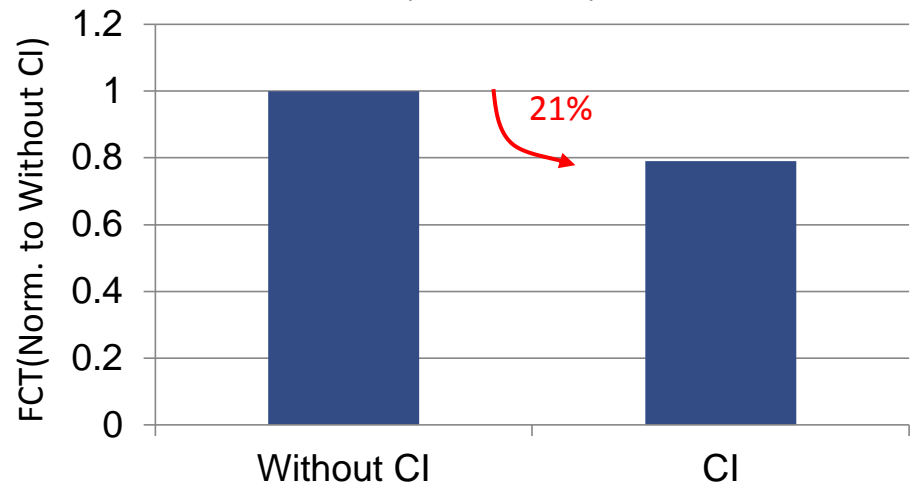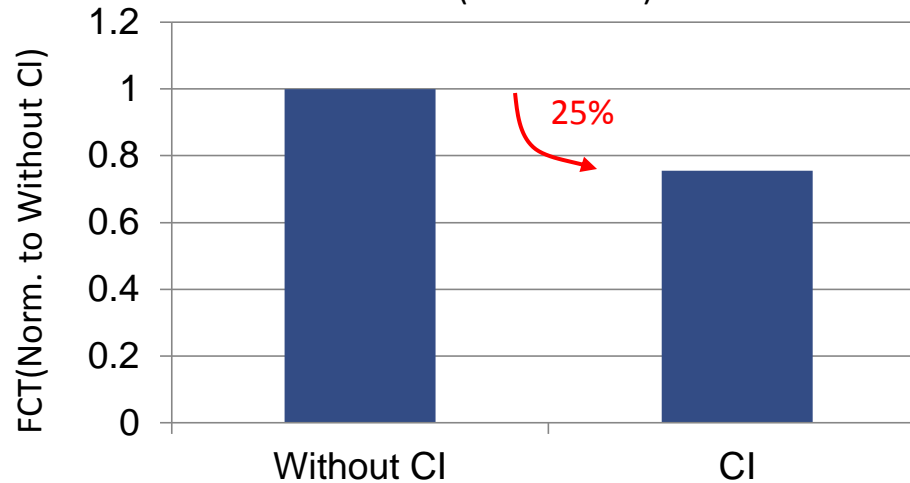
- The count of isolated flows is quite small. In our simulation with 22188 flows and 1152 server nodes. The proportion is 2% for total flows , and 12% for large flows.
- So the HOLB only occurs among the congested flows
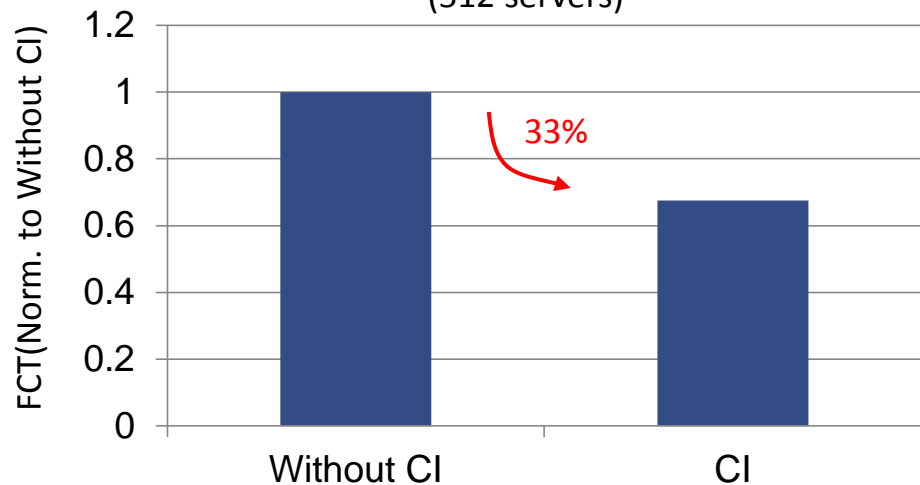
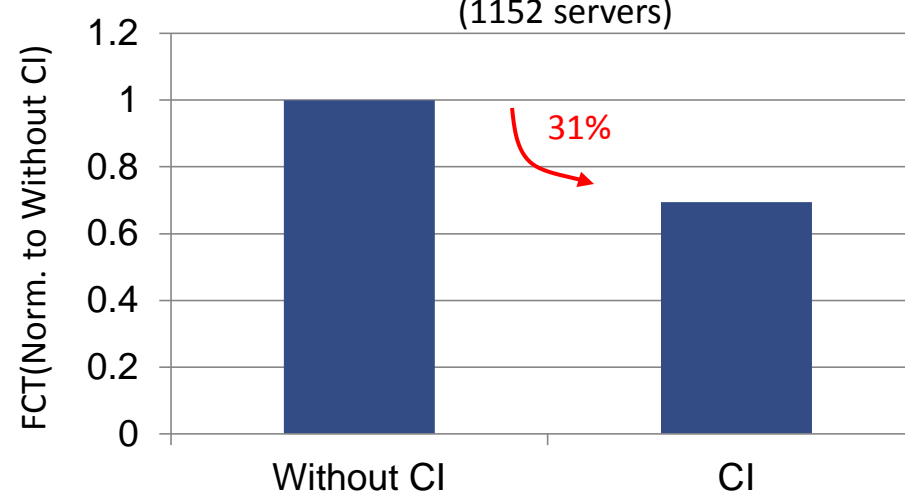# Comparison for different scale



Average flow completion time (128 servers) — 21%



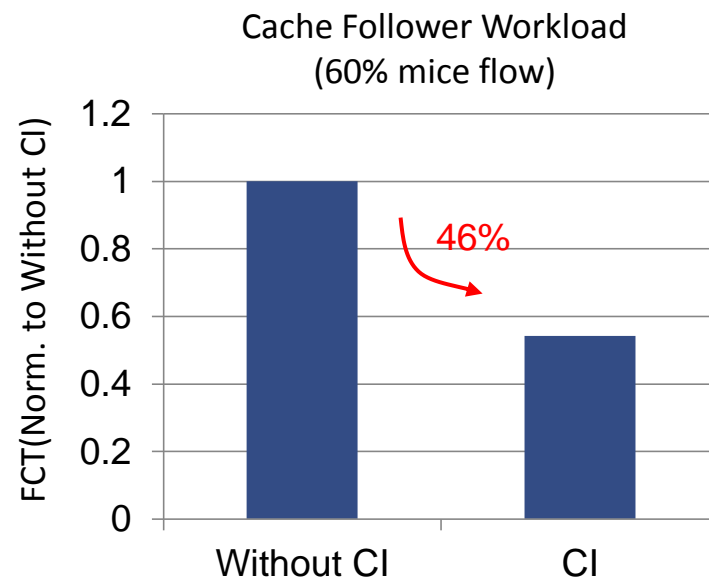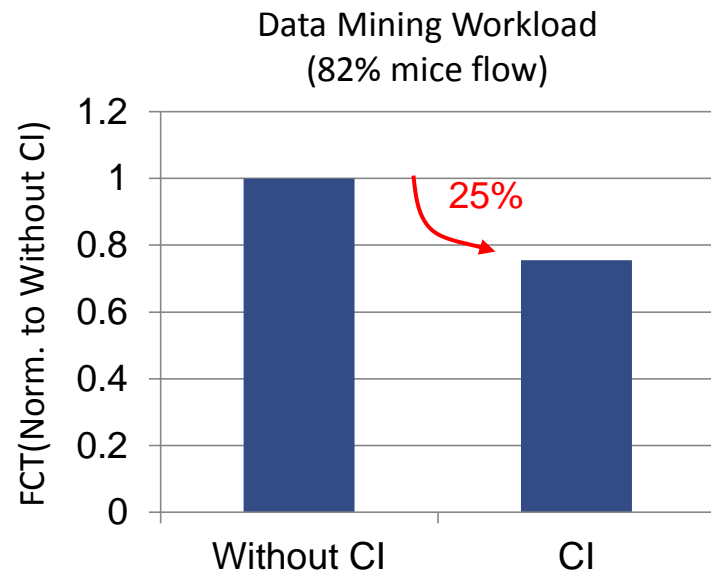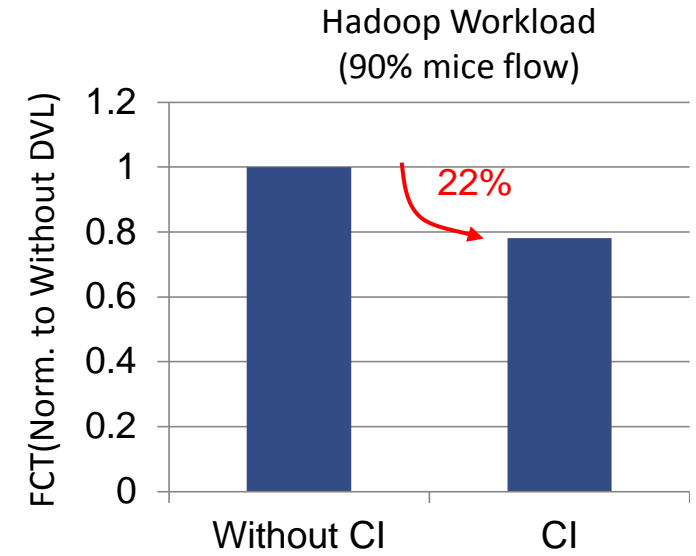Average flow completion time (288 servers) — 25%
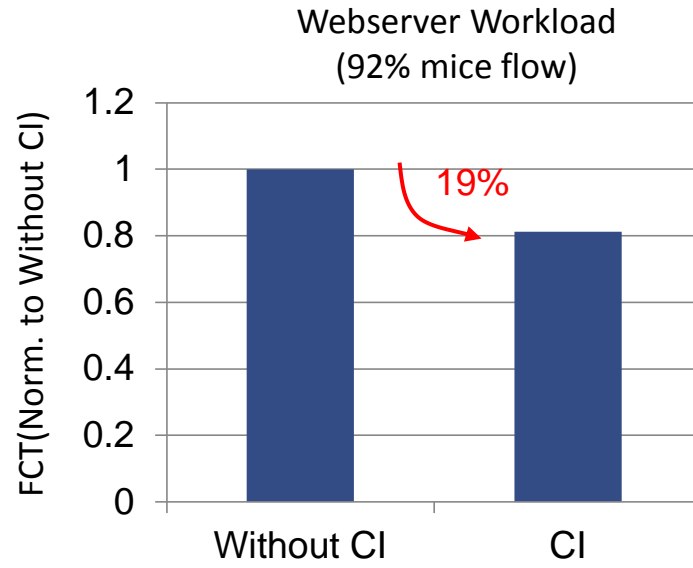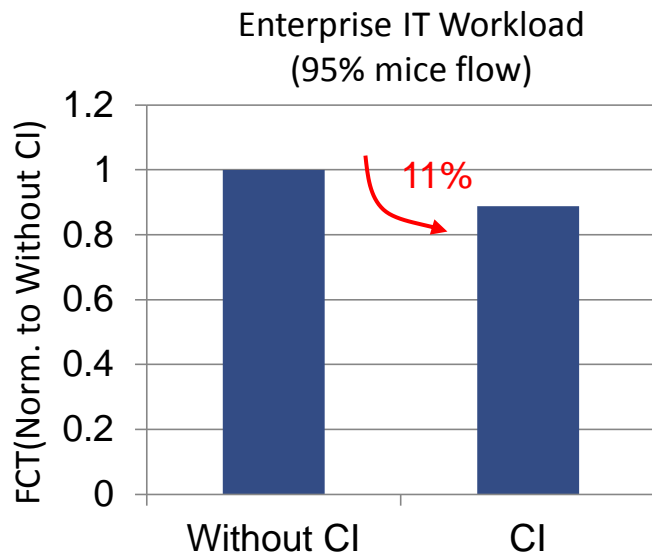


Average flow completion time (512 servers) — 33%



Average flow completion time (1152 servers) — 31%

# Comparison for different workload – Flow Completion Times



Enterprise IT Workload (95% mice flow), Webserver Workload (92% mice flow), Hadoop Workload (90% mice flow), Data Mining Workload (82% mice flow), Cache Follower Workload (60% mice flow)

# Summary

- Current data center design will be challenged to support the needs of large scale, low-latency, lossless networks.

- Congestion Isolation provides the following benefits:

    - Supports lossless as well as low-latency

    - Mitigates Head-of-Line blocking caused by PFC

    - Improves average flow completion times

    - Reduces or eliminates the need for PFC on non-congested flow queues

- Next Steps

    - Call for interest in creating a project

    - Respond to comments and feedback

# Thank you

www.huawei.com