

Requirement Discussion of Flow-Based Flow Control(FFC)

Yolanda Yu

Nongda Hu

yolanda.yu@huawei.com

hunongda@huawei.com

IEEE 802.1 DCB, Stuttgart, May 2017

Agenda

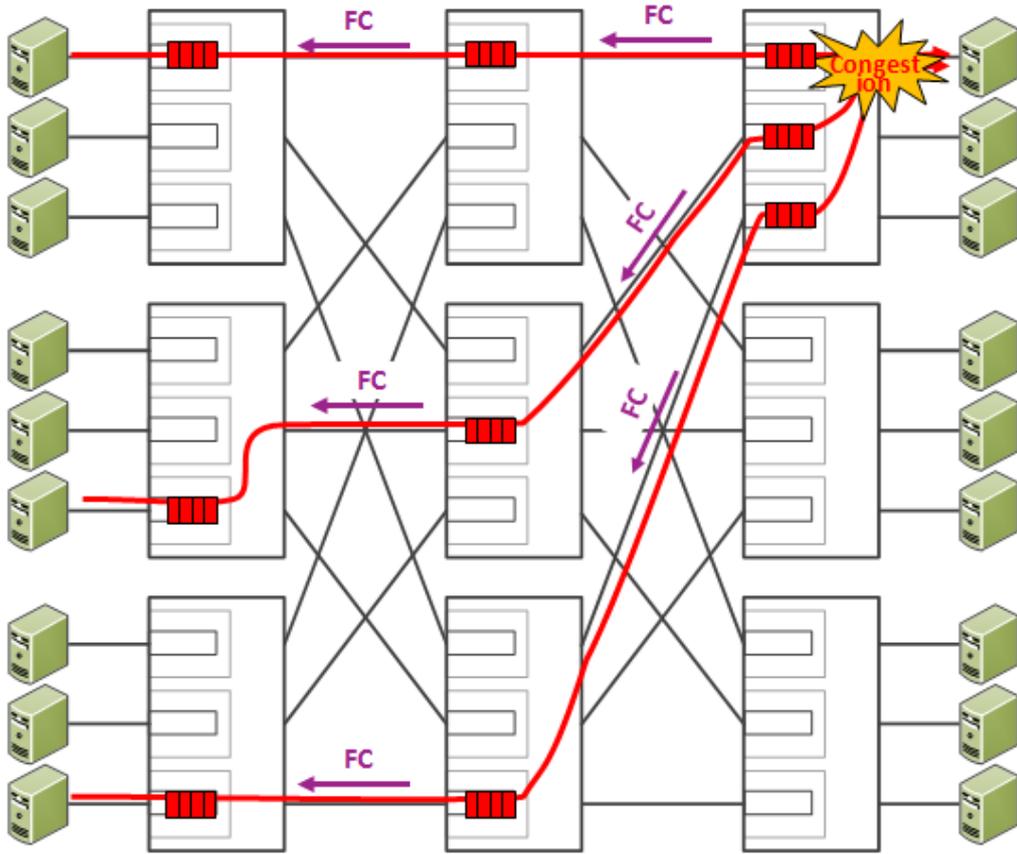
- **Key requirements for the lossless Ethernet network**
- **Some Issues with PFC**
- **Proposal to Flow-based flow control**
- **Summary**

The Requirements of Lossless Ethernet

- NVMe over Fabric (NOF)
 - Represent the requirement of high performance storage and resource pooling.
 - NOF need extremely low latency. Can't accept the latency caused by retransmission.
 - Need lossless Ethernet.
- RDMA technology is more and more used in modern data center
 - The underlying networks for RoCE and RoCEv2 should be configured as lossless.
 - The requirement for an underlying lossless network is aimed at preventing RoCE, RoCEv2 packet drops as a result of contention in the fabric.
 - Currently, IEEE 802.1 Qbb PFC (per-priority link-layer flow-control) is used.
- The convergence of HPC, Storage, LAN
 - Need lossless and low latency Ethernet.

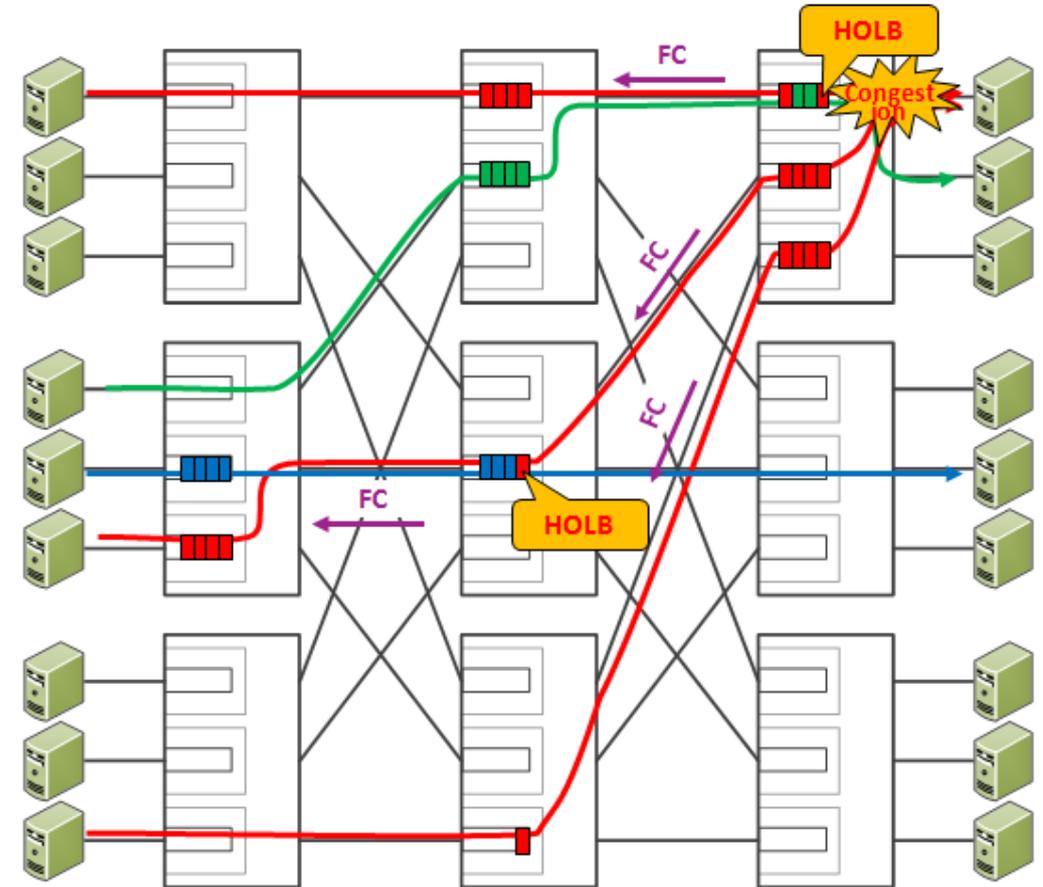
The Issues of Lossless Ethernet

- Congestion Propagation



Congestion can spread progressively through the network, building up the congestion tree.

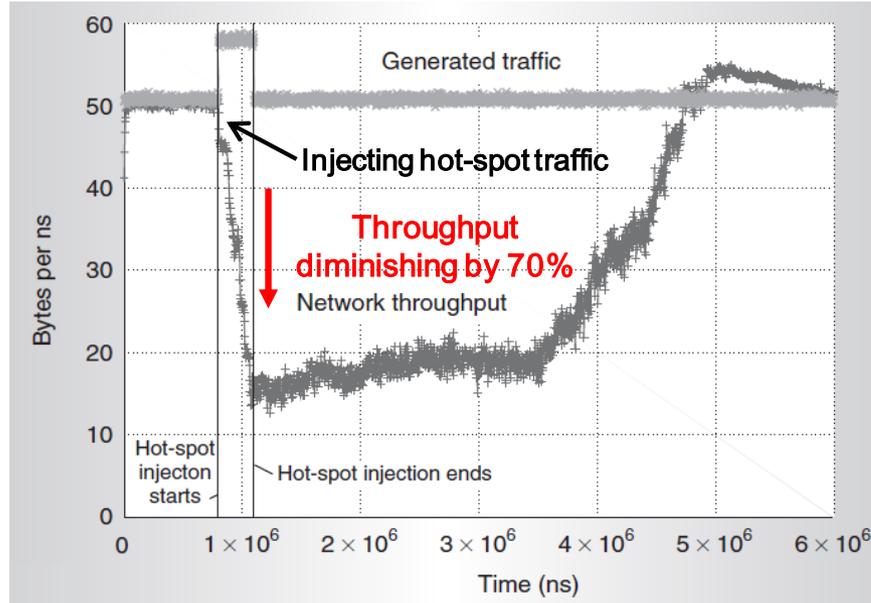
- HOL-Blocking



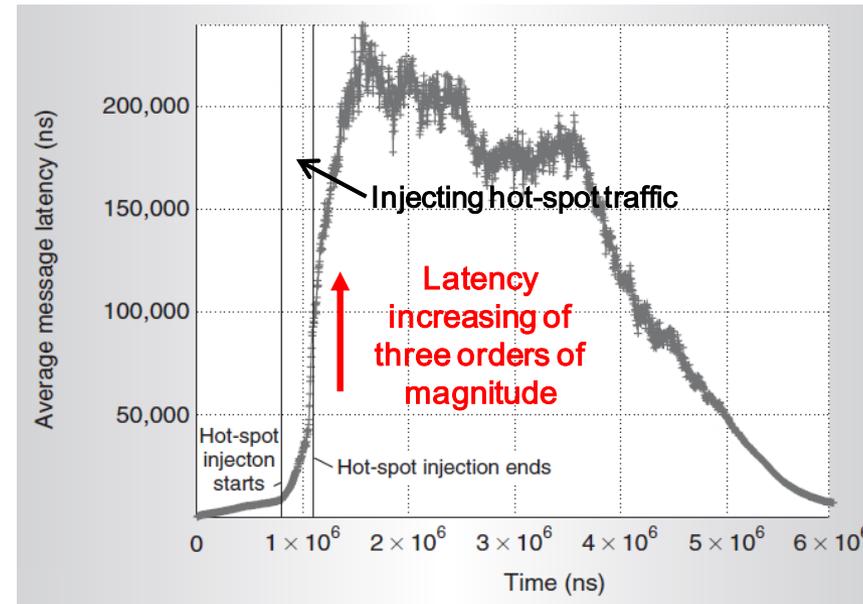
Congested flow can prevent uncongested flow in the same queue, resulting in HOL blocking. The impact of HOL blocking on network performance can be very serious: Network throughput can degrade and packet latency can increase dramatically.

The Issues of Lossless Ethernet

- Issues of Lossless network have been well studied in academic community.
- The impact of HOL blocking on network performance can be very serious.
- As show in paper (Pedro J. Garcia et al, IEEE Micro 2006)^[1]:



Network Throughput and Generated Traffic

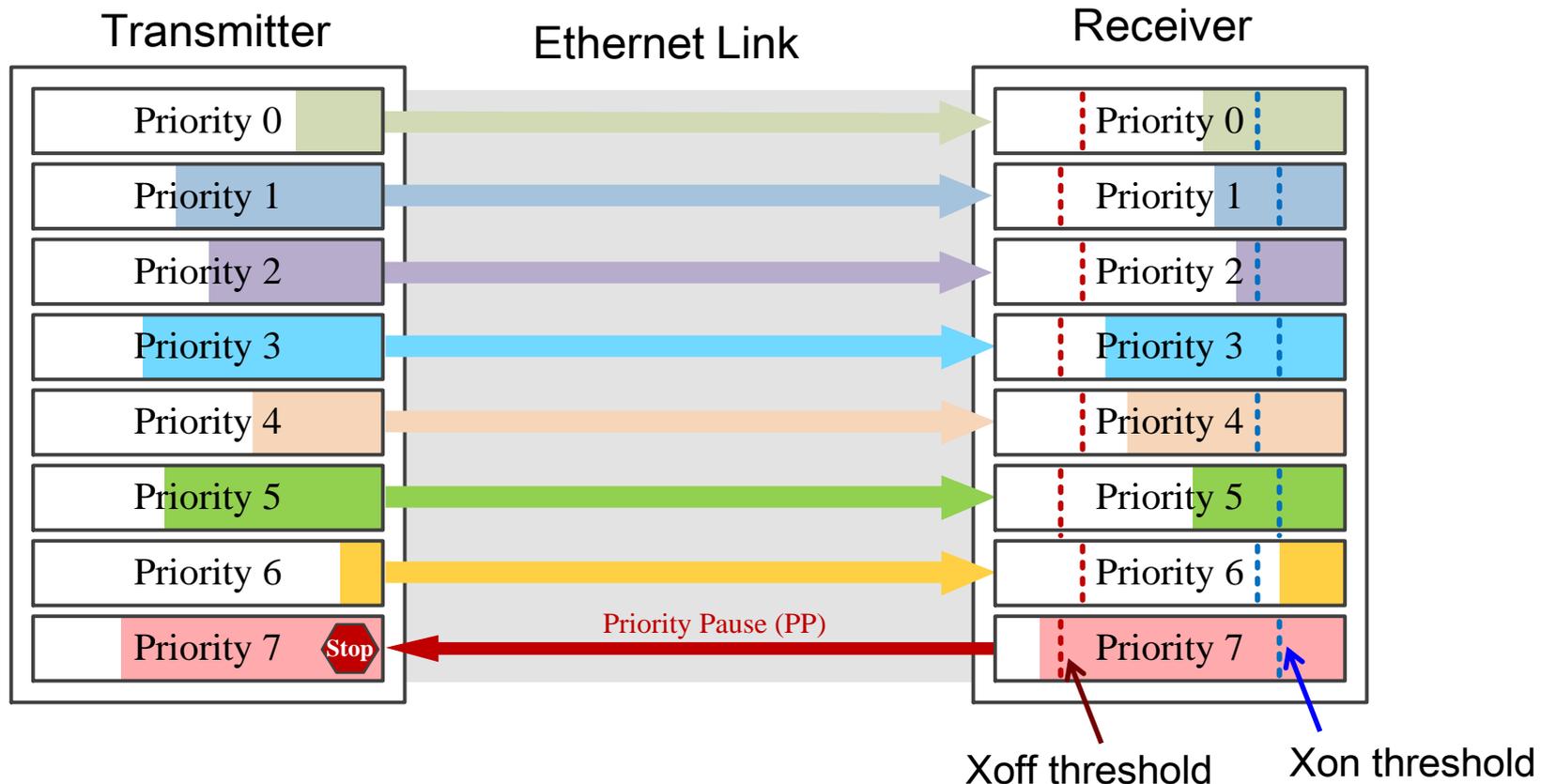


Average Packet Latency

Network Performance Degrades Dramatically after Congestion Appears

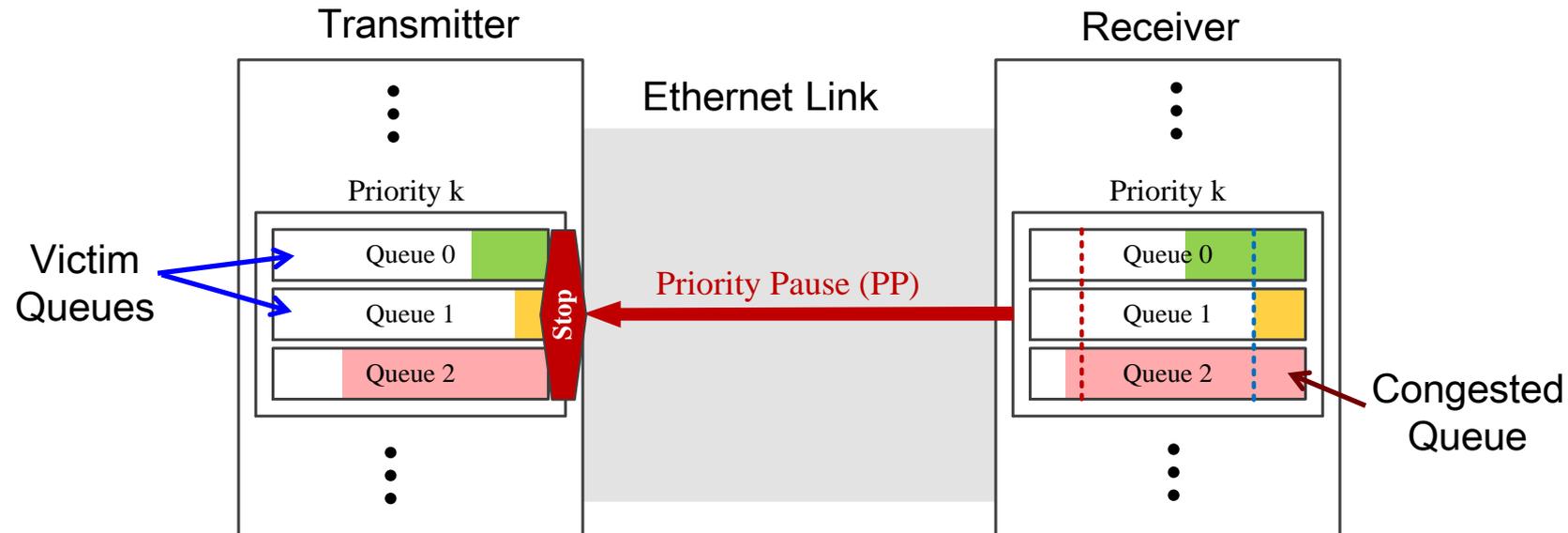
Current PFC Implementation

- IEEE 802.1Qbb PFC implement per-priority flow control, supporting eight priorities.
- Flow control of each priority can be enabled independently with individual Xoff/Xon thresholds.



Issues of PFC

- Since there are hundreds or even thousands of applications in datacenter, a lot of applications have the same priority.
- Traffics of applications in datacenter are very dynamic and unpredictable, and they may affect each other.
- So PFC still suffers congestion propagation and HOL-blocking within each priority.
- More, mismatch between coarse grained flow control (eight priorities in PFC) and rich queues in a port may result in victim queues.

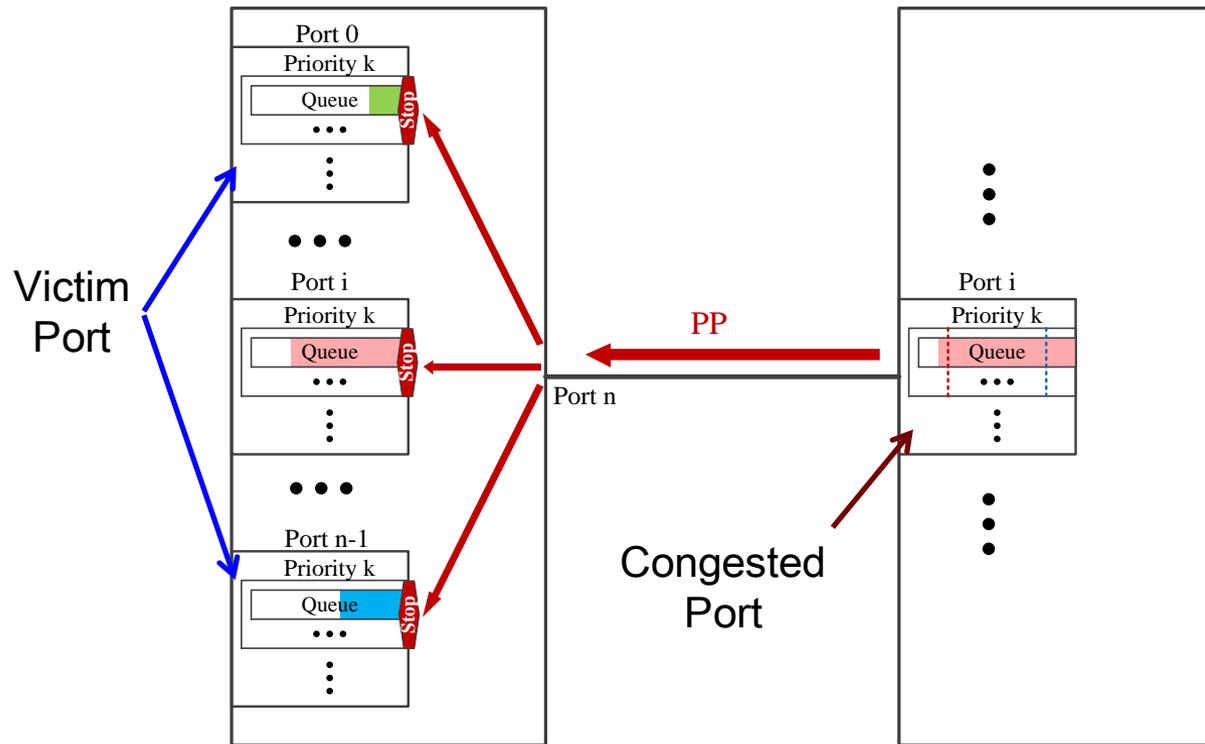


Ideally, only the queue contributing to congestion should be blocked.

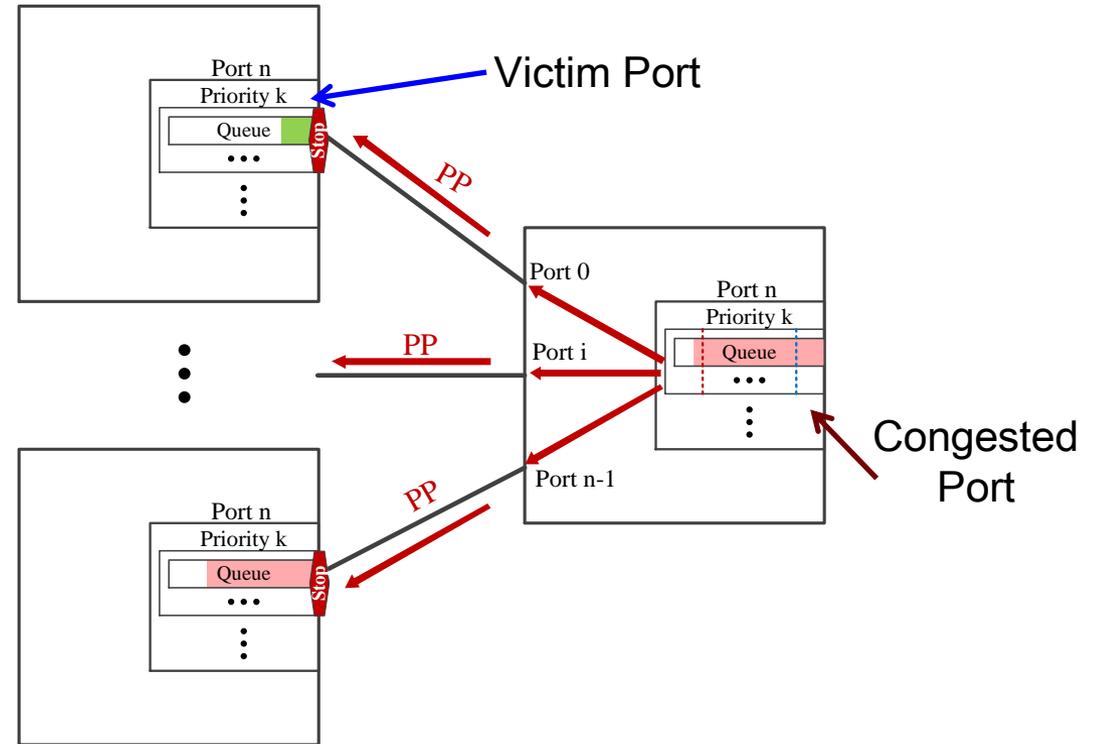
Issues of PFC

- Even more, one congested port in downstream side may block several ports in upstream side, because PFC is static (i.e. one priority k to one or several priority k).

Input Queue Model



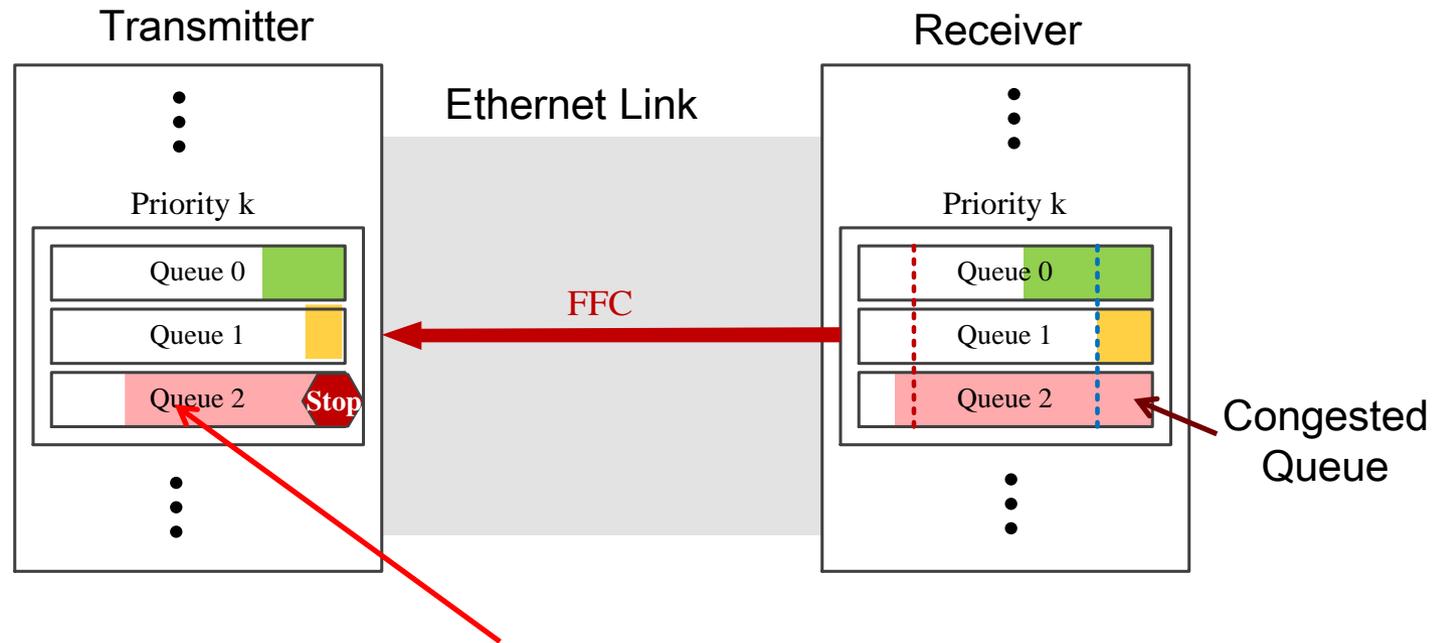
Output Queue Model



Ideally, only the port contributing to congestion should be blocked.

Proposal

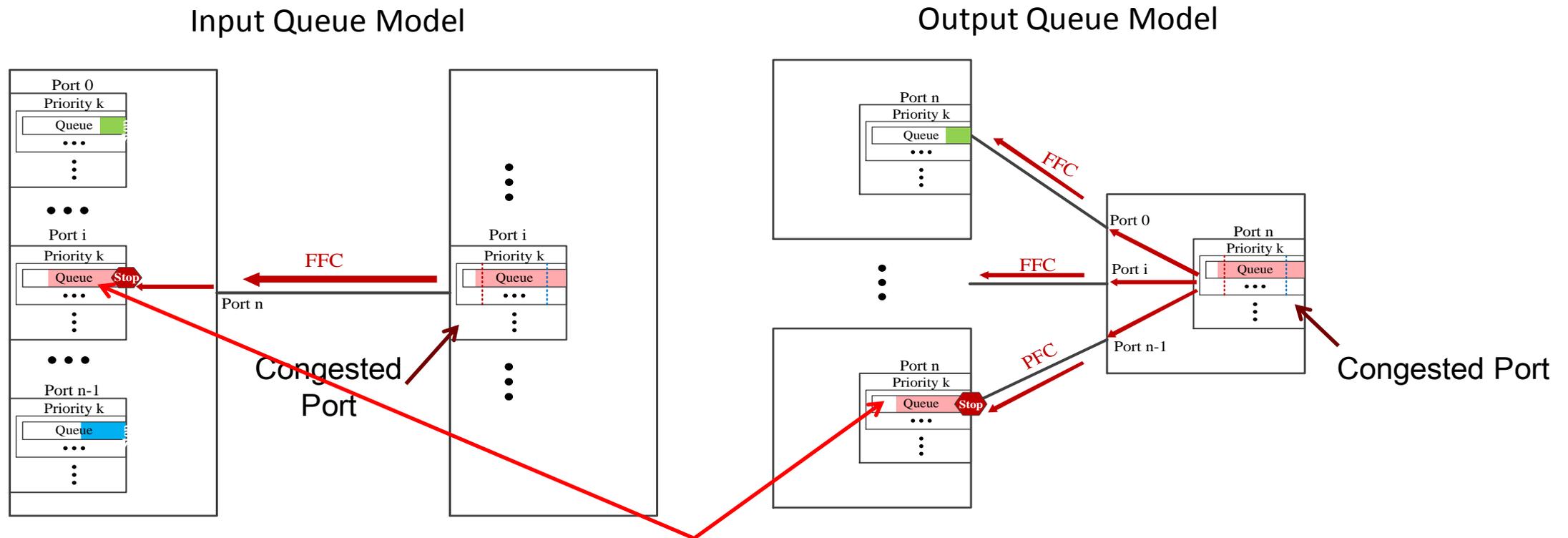
- We propose Flow-based Flow control (FFC):
 - FFC frame carrying flow information to indicate which flows to be paused.
 - On the downstream side, when queue occupancy reaches a threshold, a FFC frame indicating flows entering this queue is generate and sent upward.
 - On the upstream side, when a FFC frame is received, the flow information is parsed from the FFC frame and used to determine which queue to be paused.



Only the queue contributing to congestion is blocked.

Proposal

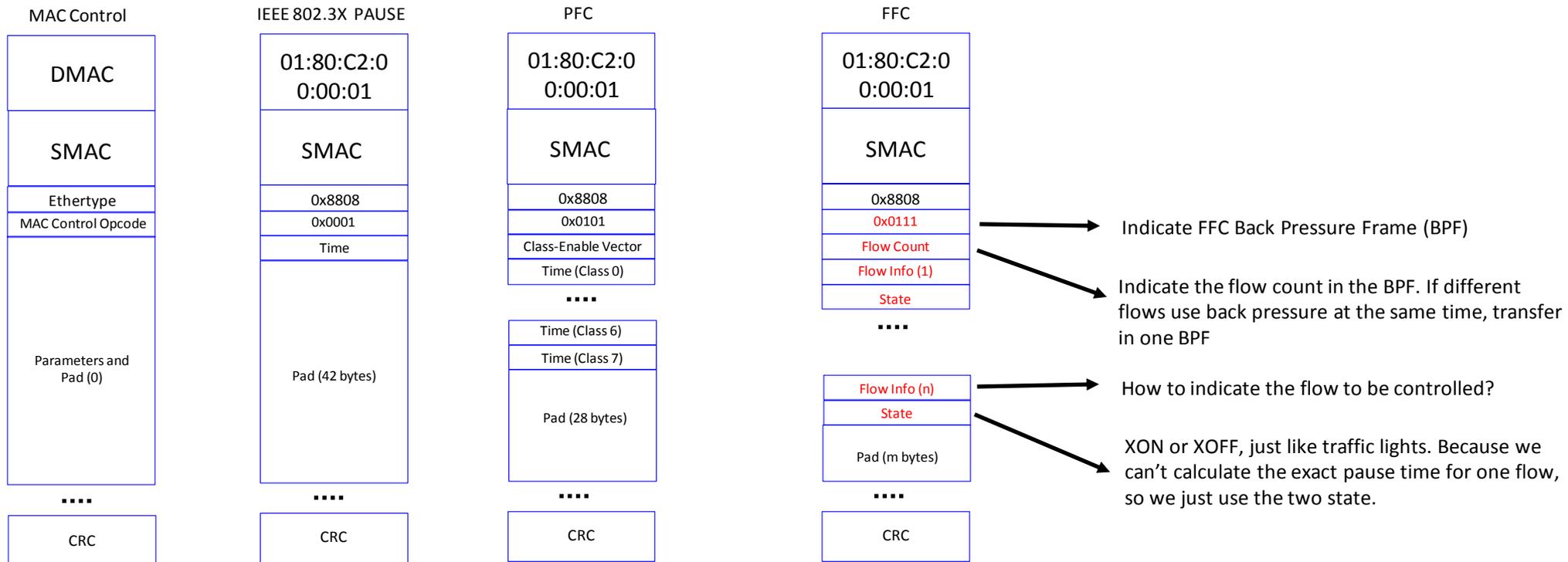
- As a result, two kind of dynamic abilities are acquired:
 - The granularity of flow control is dynamic: as the number of queue increases, finer grained flow control can be supported.
 - The queues to be blocked are dynamically determined: only the upstream queues including the flows aiming to the downstream congested queue (i.e. the flows whose information is carried in the FFC frame) will be blocked. => no victim queues.



Only the queue contributing to congestion is blocked.

Proposal

- Define FFC frame
 - Flow information
 - Action: Xoff/Xon
- Requirement for switch chips
 - One or more aggregated flows will be mapped into the queues per requirement (PFC map flows to the queue per priority).



Summary

- PFC is a coarse-grained flow control method and may suffer serious HOL blocking which will degrade the network performance dramatically.
- A fine-grained flow control mechanism should be considered in order to resolve the mismatch between the queue and the service priority.
- Flow-based flow control(FFC) can provide fine-grained and dynamic congestion management.
- Need to consider how to mapping the aggregated flows to the queues to acquire the explicit flow control.
- Need to consider the structure of FFC message frame and the inherited relationship with PFC PP(Priority Pause) frame.

Thank you

www.huawei.com