

# The Case for An Alternative Link-Level Flow Control Mechanism at High Speeds

Anoop Ghanwani, Dell

Pat Thaler, Broadcom

Jeff Lynch, IBM

Mitch Gusat, IBM

# Acknowledgement

- For their review
  - Robert Winter
  - Joseph White
  - Shivakumar Sundaram
  - Mark Gravel

# Overview

- Motivation
- Buffer requirements for priority-based flow control (PFC)
- Buffer requirements and the bandwidth delay product
- Other drawbacks of PFC
- Can credits offer a better solution?

# Motivation

- DCB networks offer lossless operation for certain traffic types
- Several protocols either require, or benefit from, a lossless network
  - FCoE, RoCE, iSCSI
- PFC has inherent inefficiencies in buffer utilization
- Can we do better?

# Buffer Requirements for PFC

- PFC requires a certain amount of buffer space to be reserved **for each lossless class**
- The computation of this space is quite involved
  - See Annex N of 802.1Q-2012
  - Depends on cable length, PHY type, use of MACsec, etc.
- The reserved buffers are never used until *after* a PFC event
  - **They remain unused during normal operation, e.g. absorbing bursts**
- This results inefficient buffer utilization

# Buffer Requirements for PFC (2)

- Higher layer/MAC/PHY delays
  - Typically negligible, amounting to  $\ll 1$  MTU
  - Use of MACsec increases this significantly and is proportional to link speed
- Boundary conditions
  - Delay in putting a PFC frame on the wire because the link is busy
  - Delay in effecting PFC once a message is received because a transmission has already begun
  - This is 2 MTU regardless of link speed or cable length
- RTT
  - For a given cable length, the number of bytes required depends on link speed

# Buffer Requirements and the Bandwidth Delay Product

- Consider the following example
  - Cable length = 50 m
  - Speed of light in optical fiber  $\approx 2 \times 10^8$  m/s
  - $RTT = (2 \times 50) / (2 \times 10^8) = 0.5$  usec
  - MTU = 2000 bytes (802.3as) [Ignoring preamble and IFG]

Link Speed	# Bytes in 1 RTT	# MTU in 1 RTT
10 Gbps	$\sim 0.61$ KB	$< 1$
40 Gbps	$\sim 2.44$ KB	$\sim 1.25$
100 Gbps	$\sim 6.1$ KB	$\sim 3.12$
400 Gbps	$\sim 24.4$ KB	$\sim 12.5$
1 Tbps	$\sim 61.0$ KB	$\sim 31.2$

# Buffer Requirements and the Bandwidth Delay Product (2)

- For a given cable length, as link speeds are increased, the amount of buffer that must be set aside for PFC goes up
  - Again, this is not used during normal operation
- The boundary conditions are relatively fixed
  - ~3 MTU regardless of link speed (ignoring MACsec)
- With increasing link speed, RTT is the dominant contributor

# Other Drawbacks of PFC

- PFC is reactive
- Incorrect calculations for buffer space are expensive
  - On the conservative side, they lead to added waste
  - On the aggressive side, they lead to loss
- For a given buffer size and link speed, at some cable length, it becomes impossible to provide lossless operation

# The Case for Credit-based Flow Control

- With credits, there are no buffers set aside for the congestion event
- Instead, the credits must be sized to cover the RTT
  - If sufficient buffers are not available, the link will be underutilized
- Lossless operation is always guaranteed

# Credits Are Not Perfect

- Credit size
  - If too big, could lead to underutilization by fragmentation
- Lost credits
  - If credits are not reliably returned, leads to lower available credits and consequently lower utilization
- Solutions exist for these problems
  - We can debate and refine them during the process of standardization

# Summary

- This presentation discusses some of the drawbacks of PFC as a link level flow control mechanism for high speed links
- It may be worthwhile to look at credit-based flow control
- If developed, there should be a mechanism for negotiating the behavior so as to be backwards compatible

**THANK YOU**