

An abstract graphic on the left side of the image, composed of numerous thin, wavy green lines that swirl and overlap to form a complex, organic shape. The lines are a vibrant green color against the dark blue background.

# Open. Together.



**OCP**  
SUMMIT



# Disaggregation & Data Center TCP: Maximizing OCP Datacenter Efficiency Through the Reduction of Tail Latency

Nicolaas Viljoen, Associate Director Software Engineering, Netronome

Amal Tariq, Network Hardware Engineer, Facebook

Lawrence Brakmo, Kernel Software Engineer, Facebook



Open. Together.



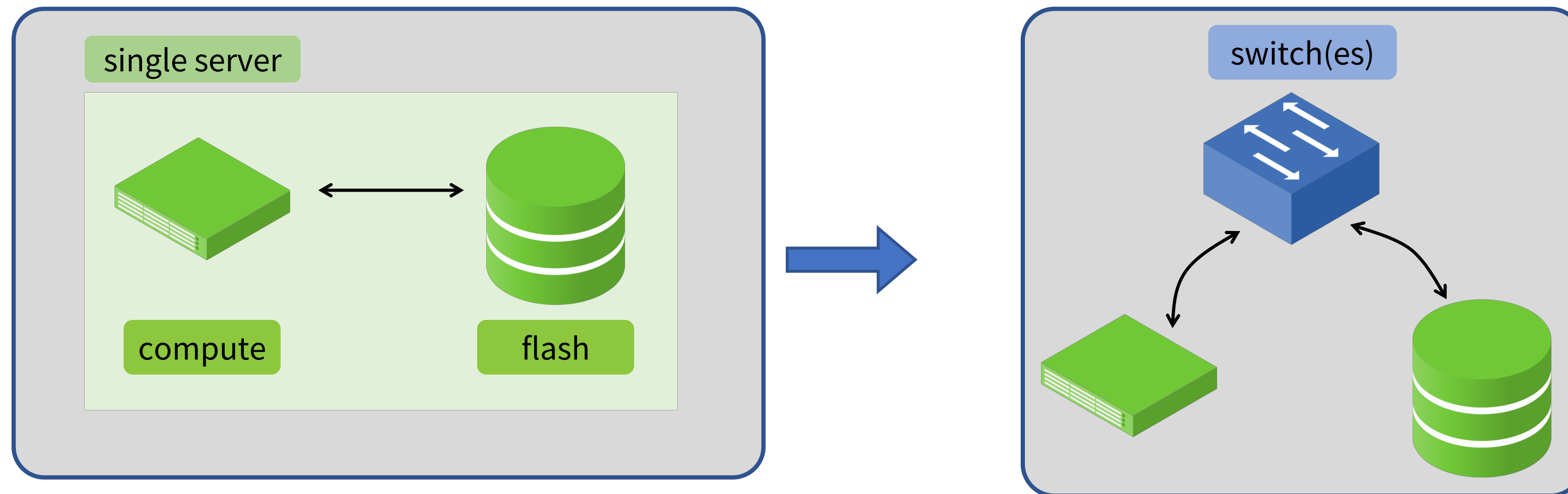
# Motivation for Disaggregation



NETWORKING

**Flash Disaggregation:** moves memory access from within a server to the network

**Stranded Capacity:** Disaggregated flash storage improves utilization by up to 40% [1]



Case Studies

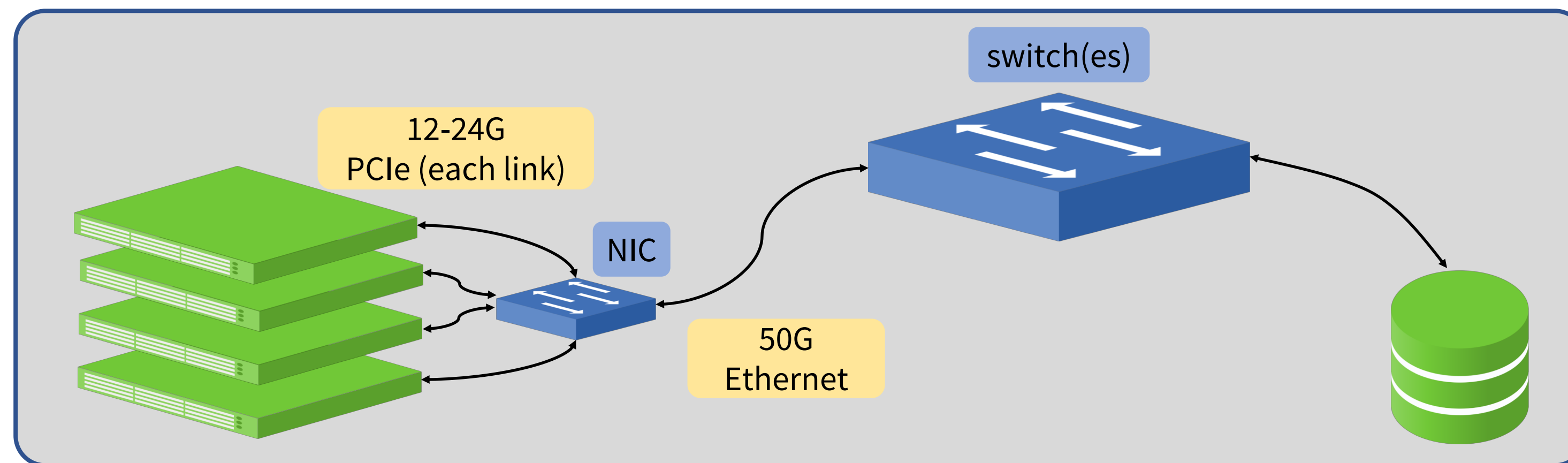
[1] A. Klimovic, C. Kozyrakis, E. Thereksa, B. John, S. Kumar. Flash Storage Disaggregation. Eurosys 2016



# The Yosemite v2 & Disaggregation

**Power efficiency:** The Yosemite v2 increases compute per watt of standard designs

**To combine with disaggregation:** Must handle bursty traffic while minimizing drops to reduce tail latency.  
This is mainly due to in-cast related network stress



## DCTCP + ECN

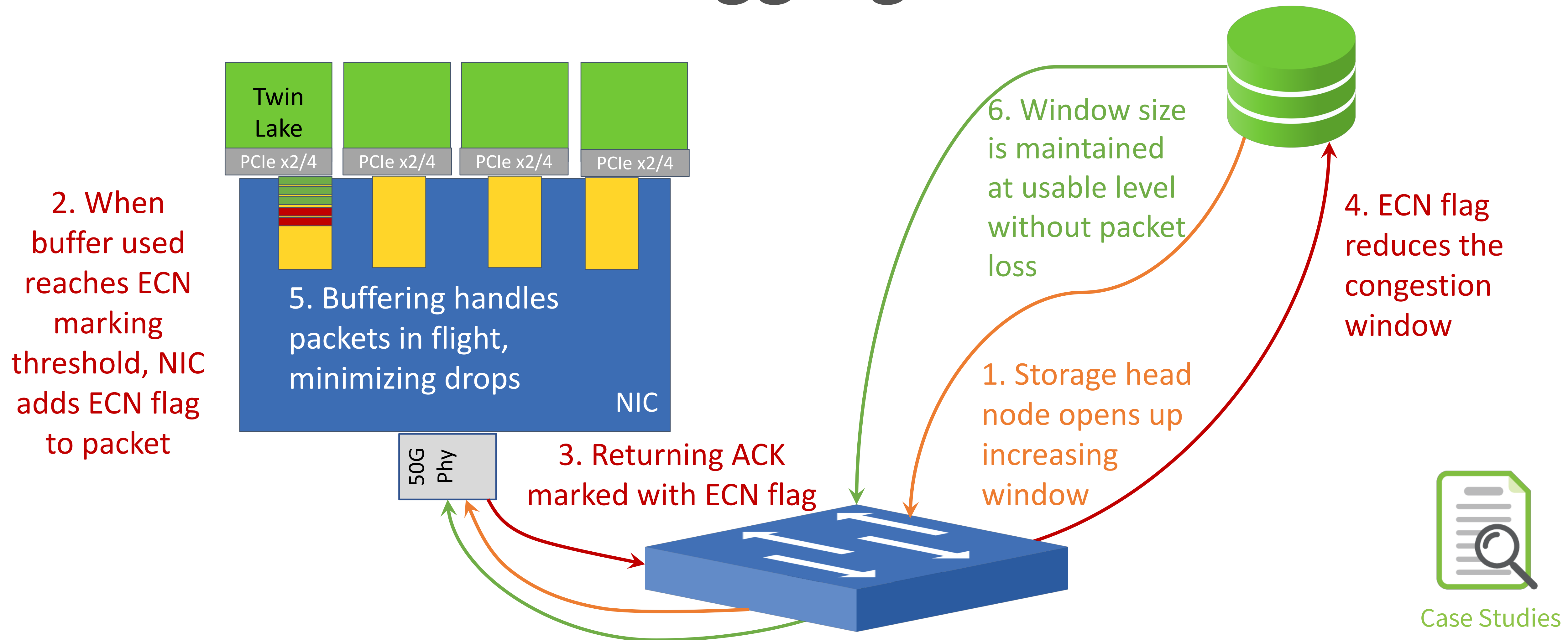
- ECN (explicit congestion notification): aware of *extent* (rather than just presence) of congestion
- Allows sender to respond to congestion before packets are dropped



Case Studies



# DCTCP/ECN in Disaggregated Architecture





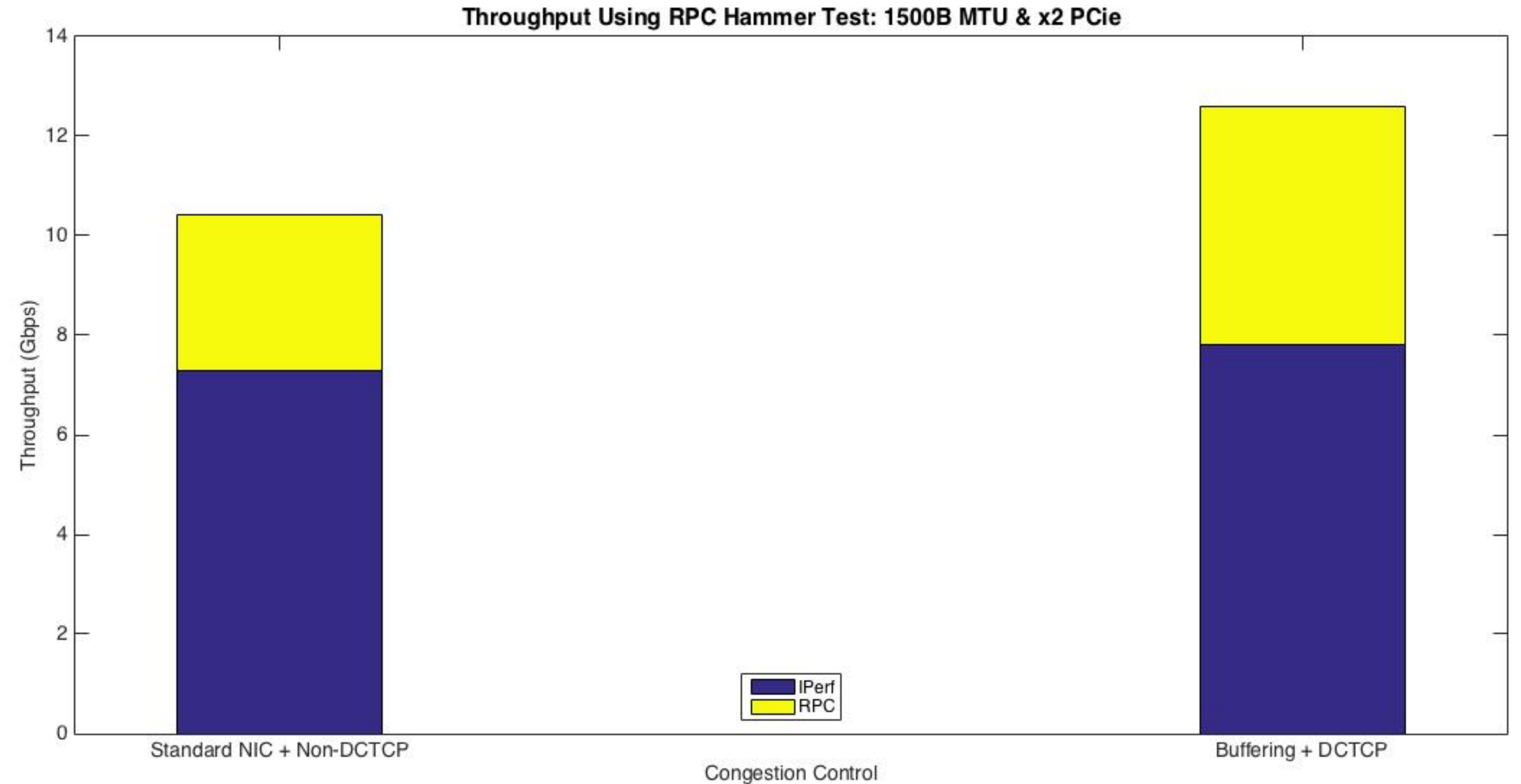
# Throughput with DCTCP

Often a concern raised with DCTCP is that available bandwidth is underutilized.

This is due to the **low ECN threshold** values used.

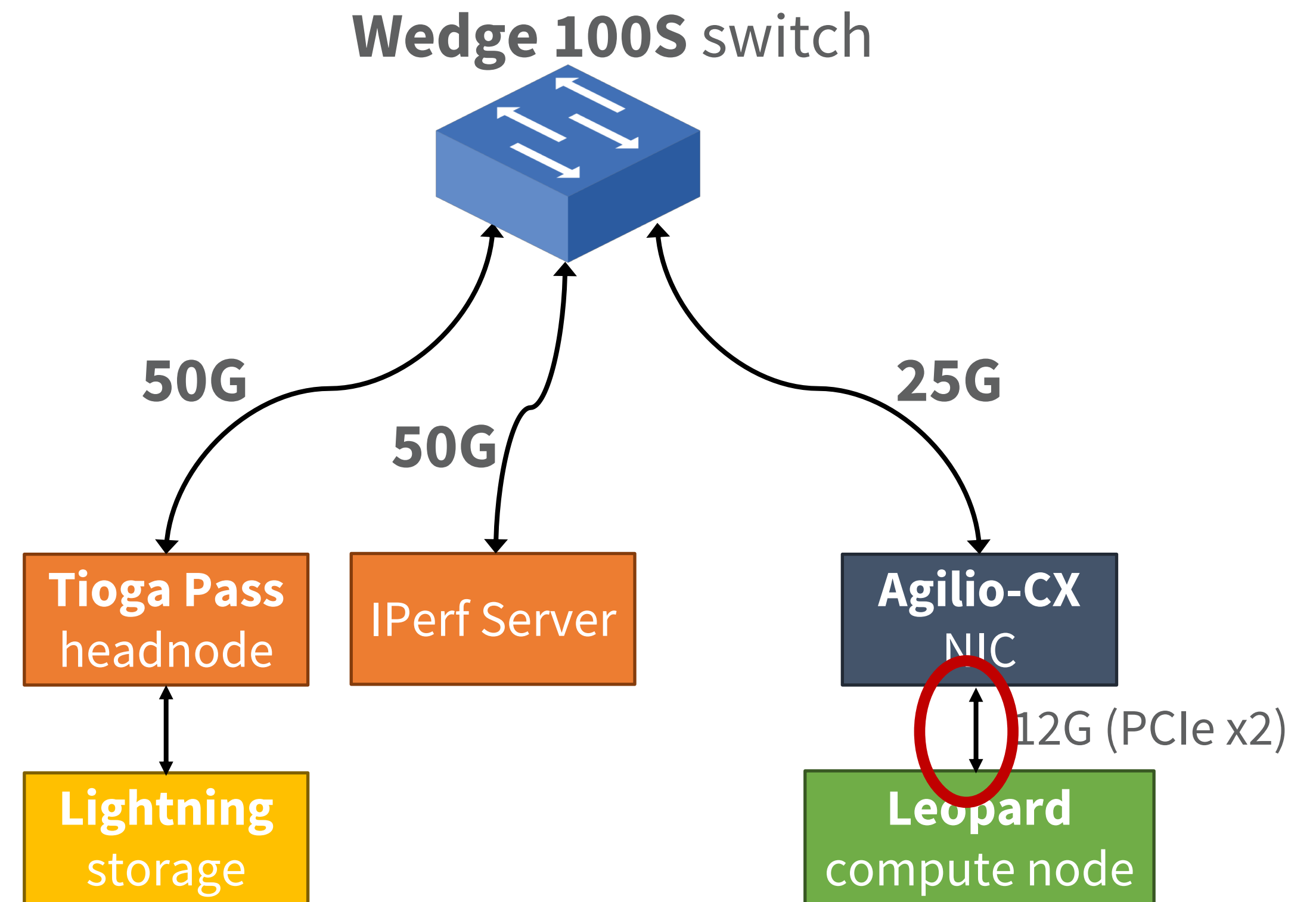
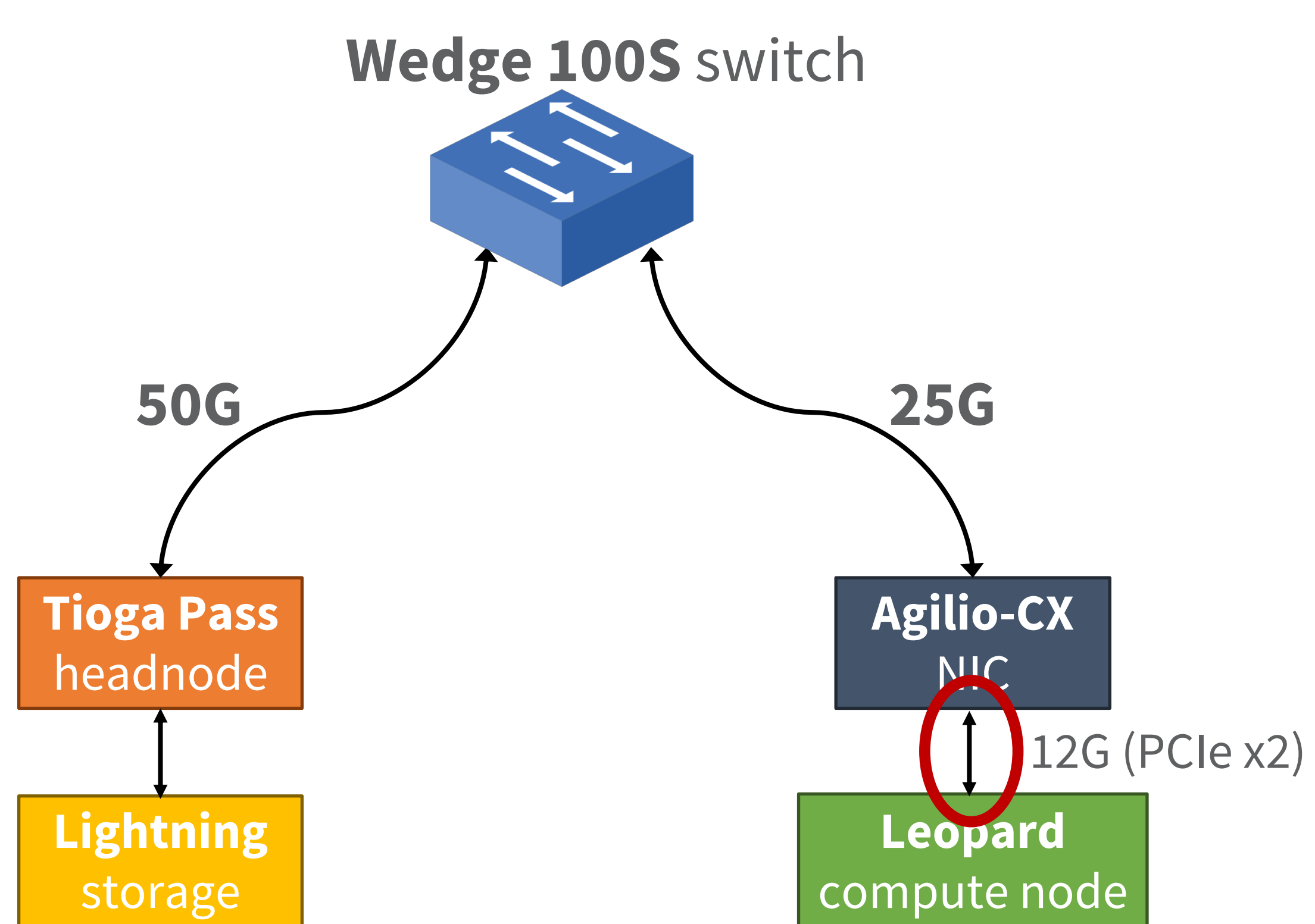
However, with sufficient buffer space to **set higher thresholds**, link underutilization is not an issue.

This removes the throughput regressions otherwise observed.





# Initial tests with DCTCP: setup





# Initial tests with DCTCP

Tested in production environment, shadowing **latency-sensitive** application traffic.

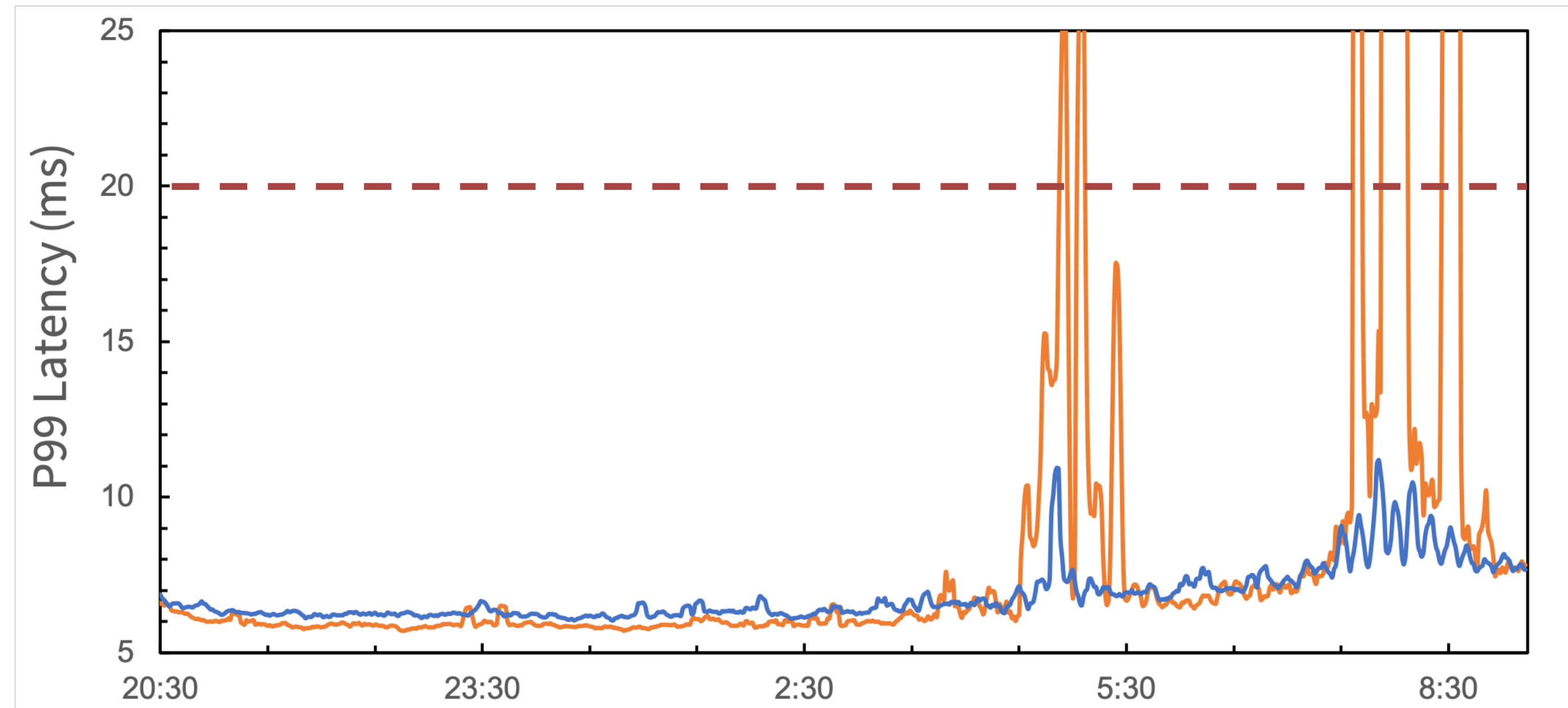
**DCTCP traffic** is the storage reads between compute and head node (*traffic under our control*)

**CUBIC traffic** is everything else (e.g. application queries to compute node or external storage writes)

*(traffic not under our control)*

**Latency** as primary performance metric, **other metrics were not to regress** (throughput, QPS, etc.)

PCIe at the NIC was limited to **mimic in-cast** problem



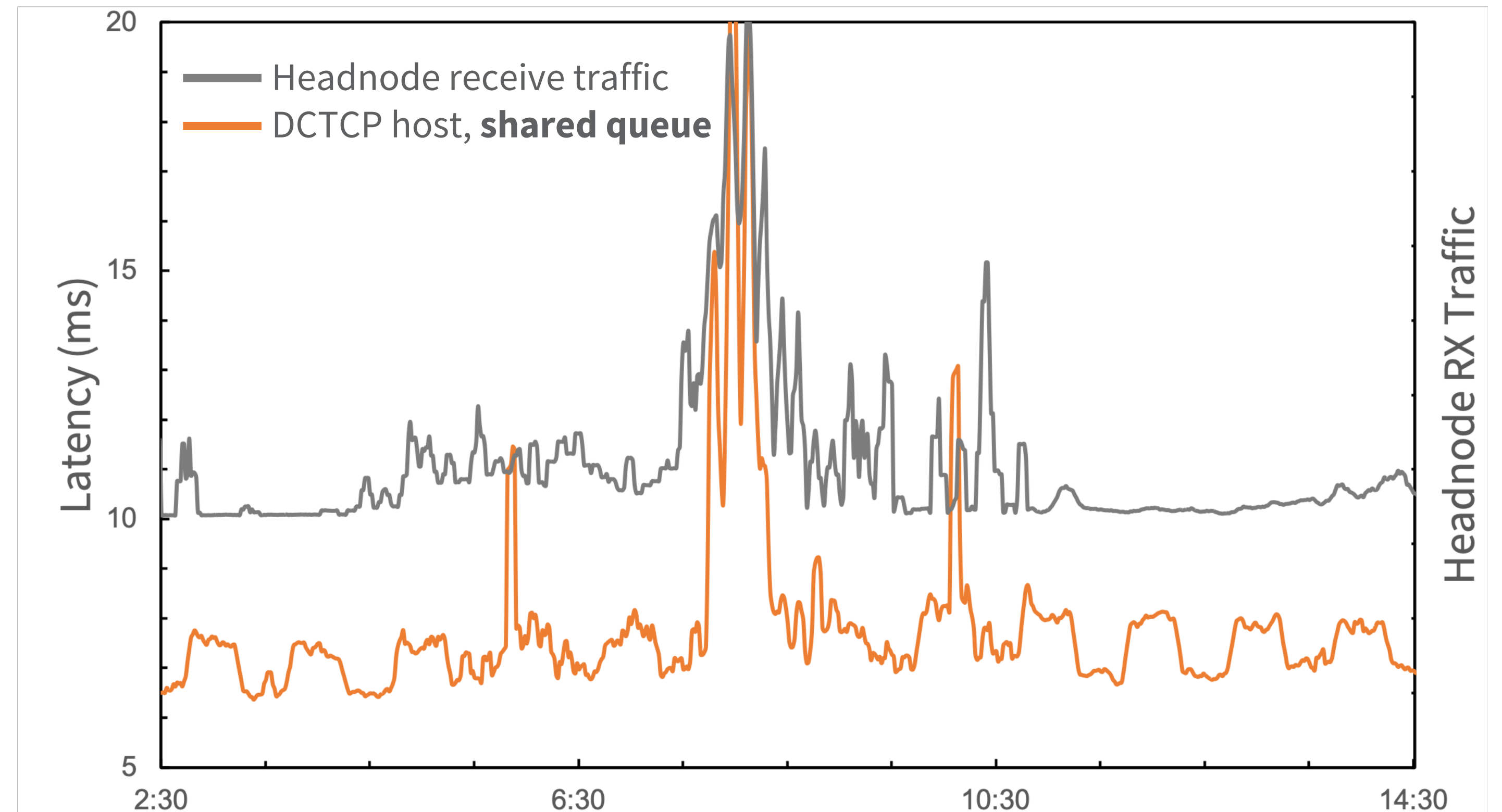
# Split Queues by Congestion Control

Latency peaks correspond to daily partition (CUBIC traffic)

Due to CUBIC traffic dominating, buffer queues are not able to respond to ECN

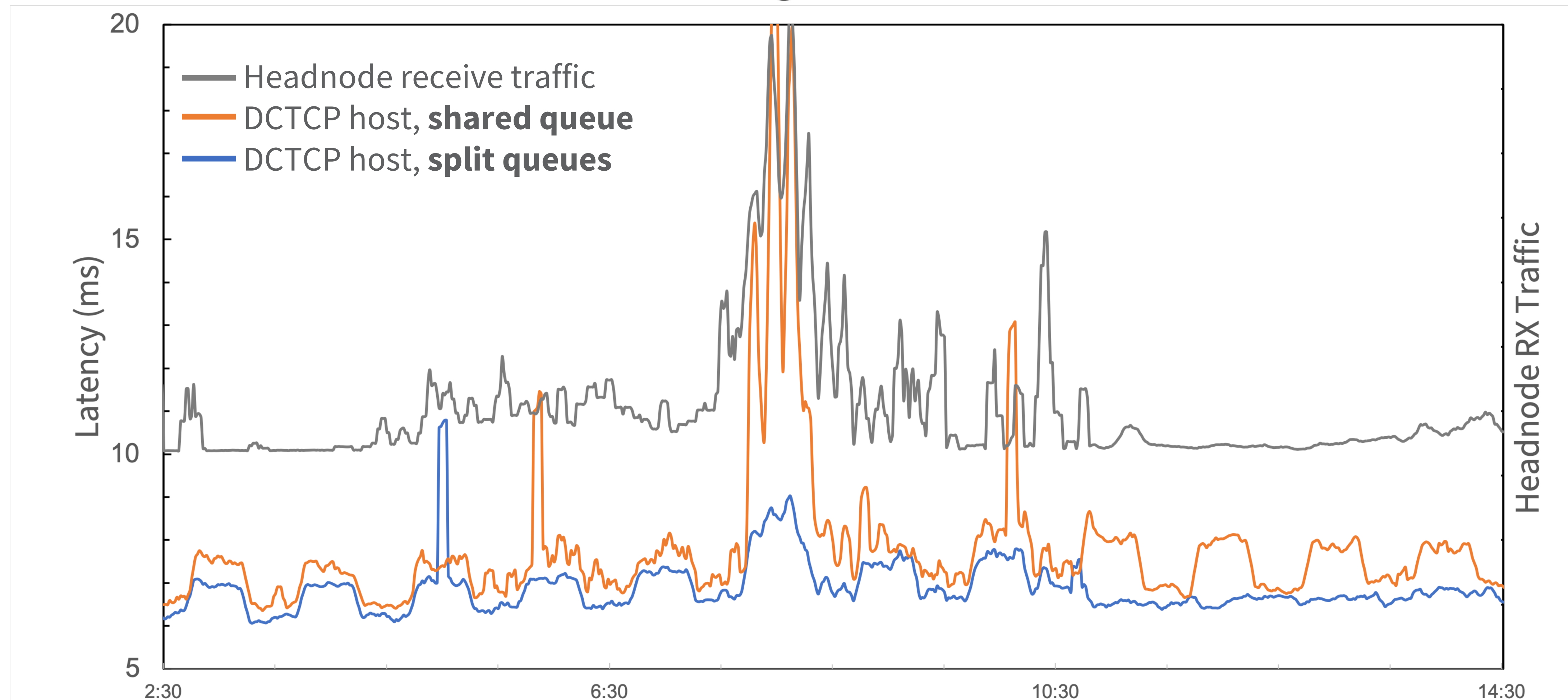
Resolved this by creating separate buffer queues for DCTCP and CUBIC traffic in the NIC

Both sets of queues were tuned to ensure fairness



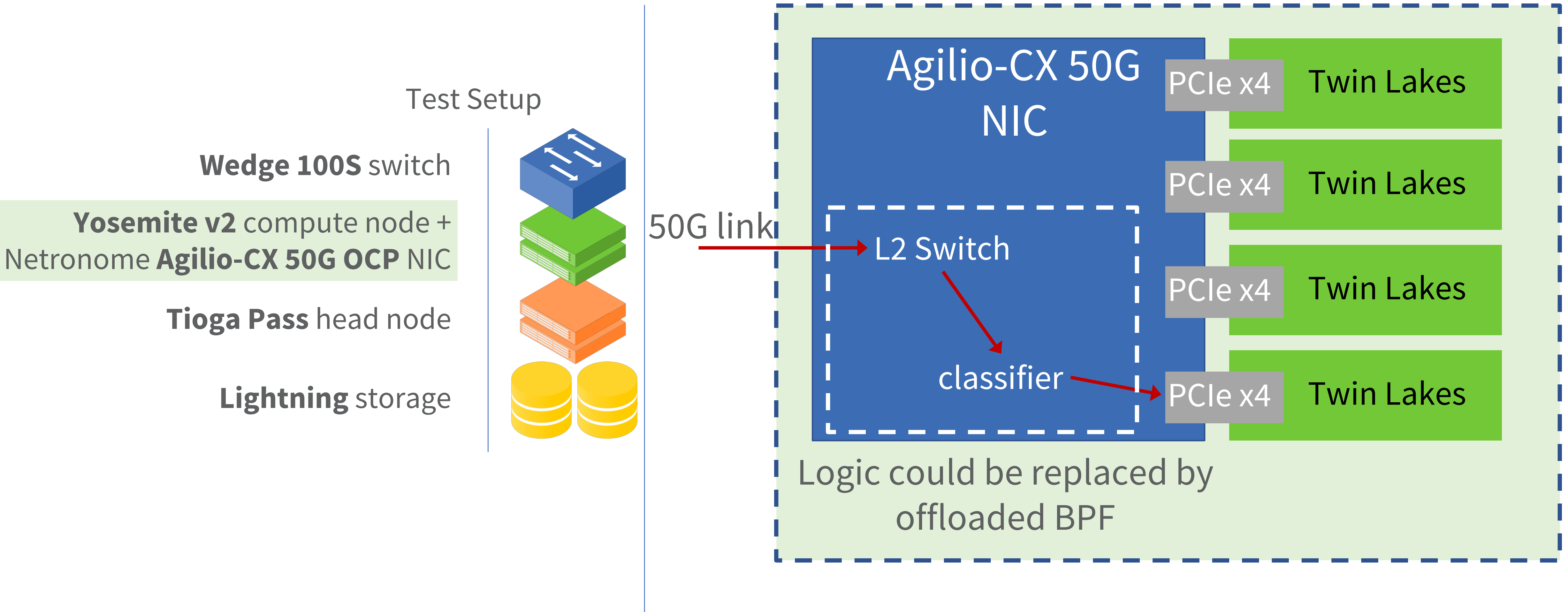


# Split Queues by Congestion Control



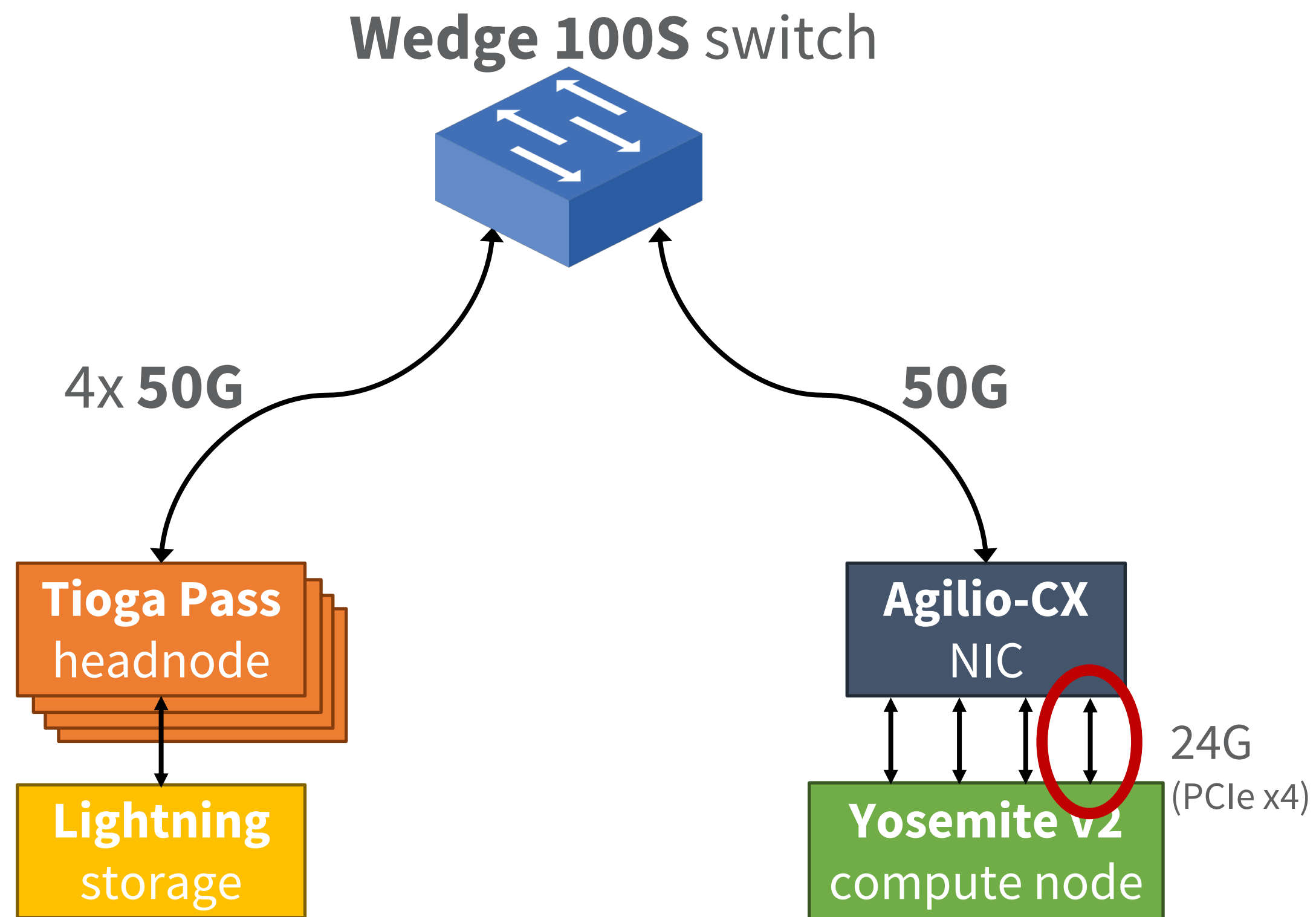
# Multi-Host NIC tests: architecture

## Yosemite v2

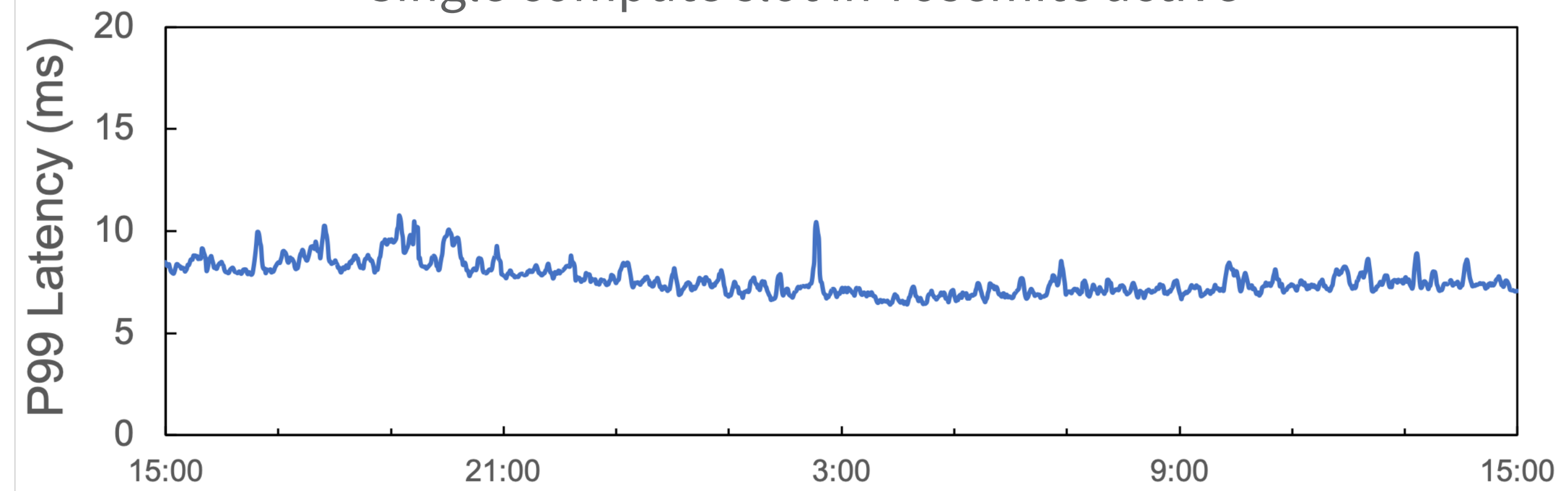




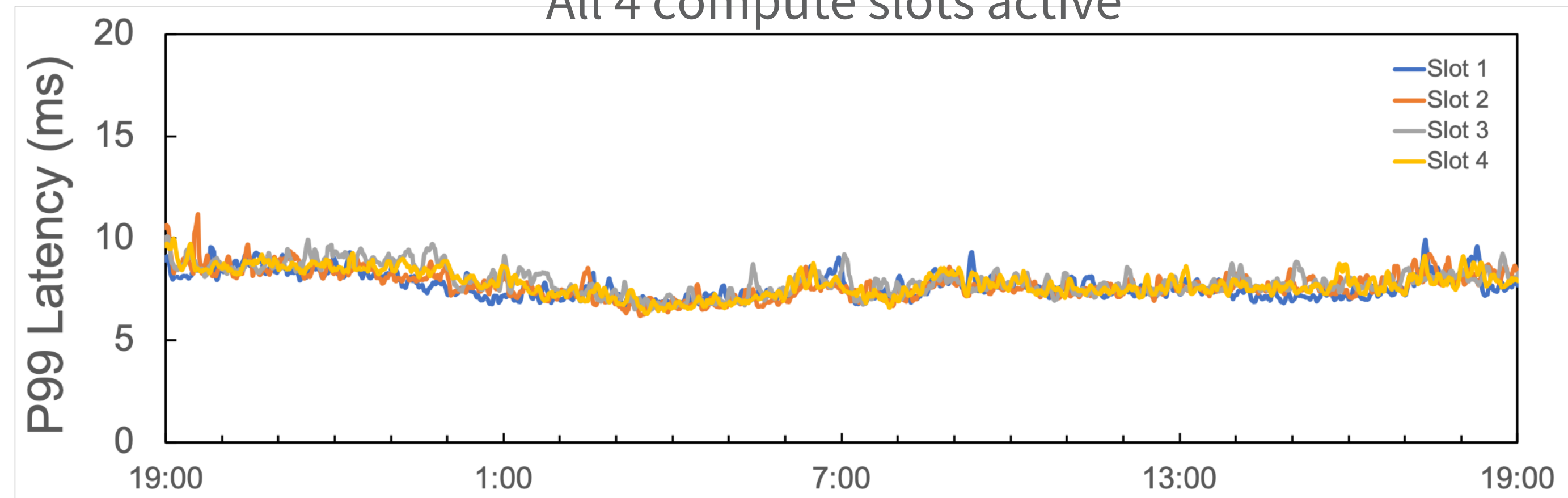
# Multi-Host NIC tests



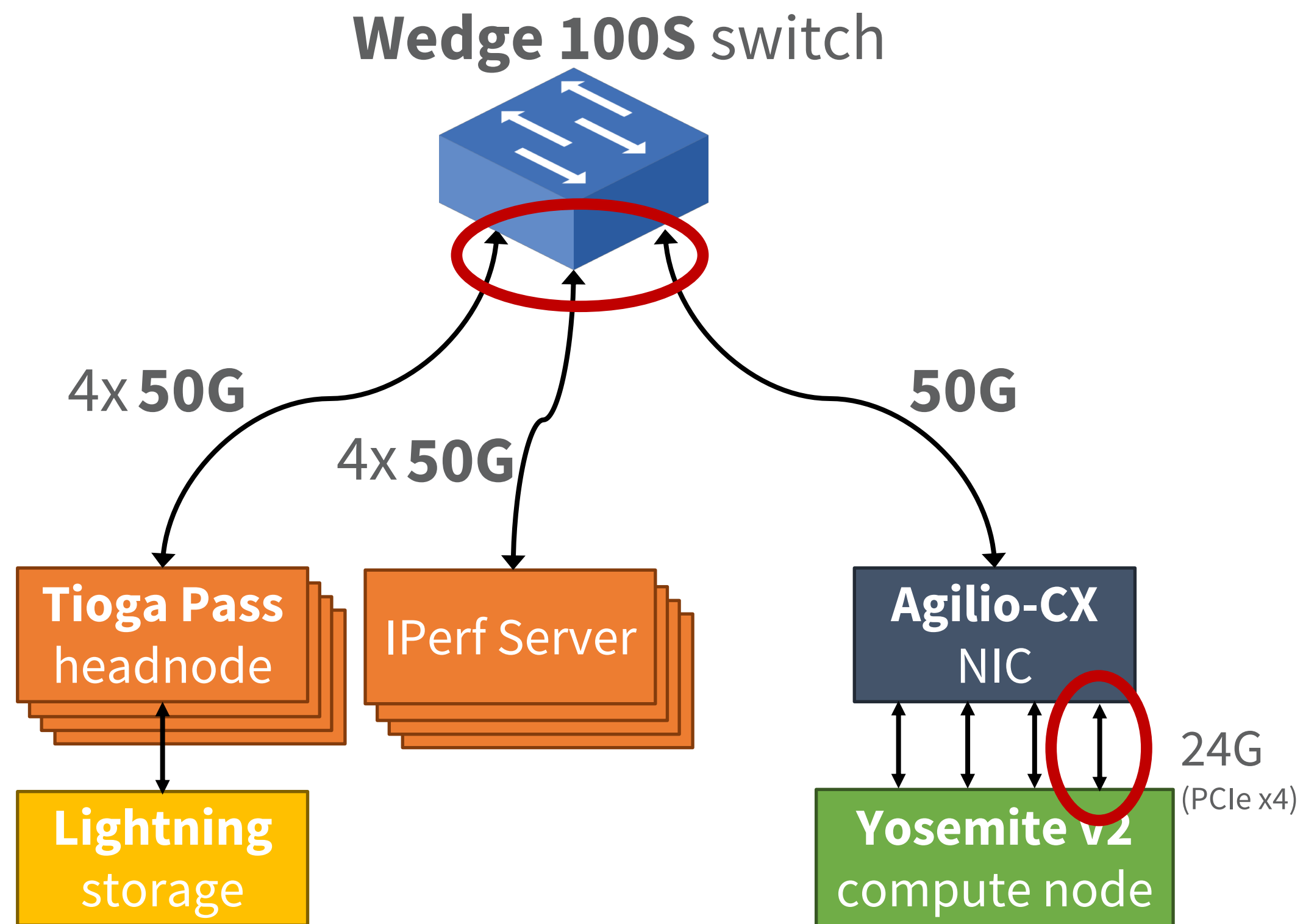
Single compute slot in Yosemite active



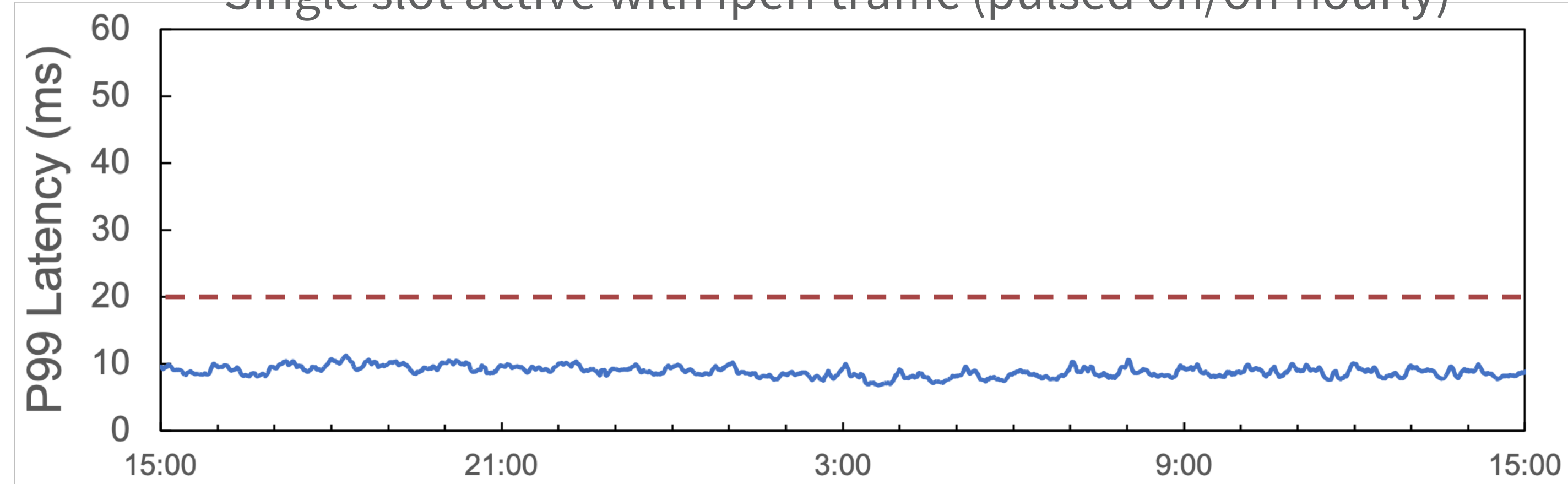
All 4 compute slots active



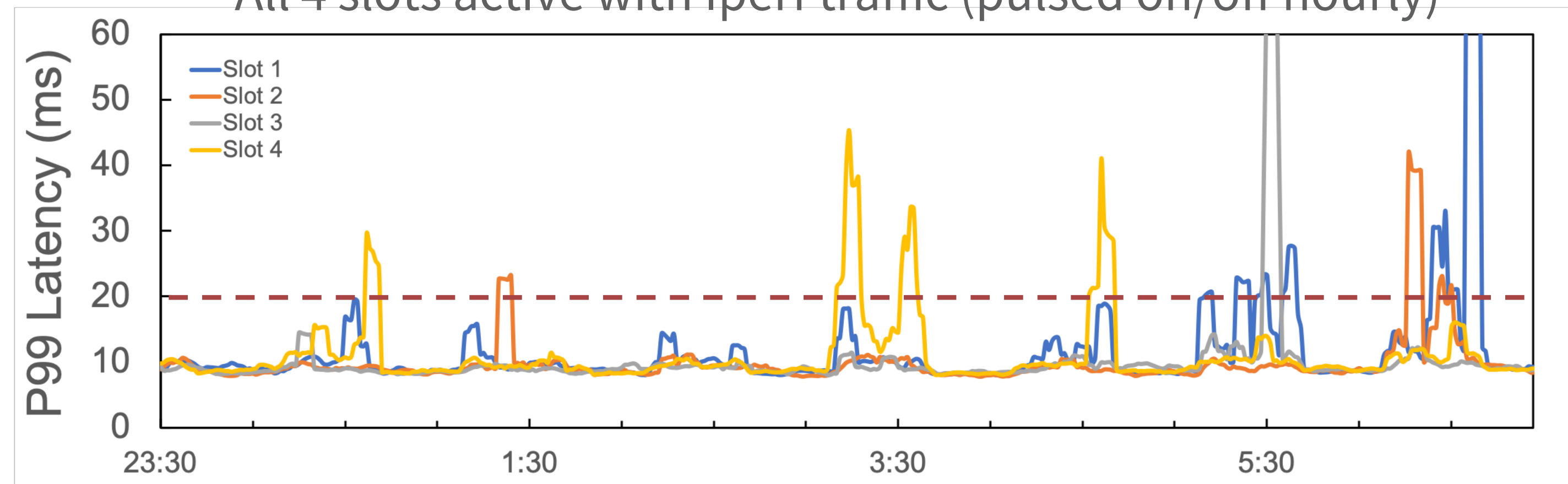
# Multi-Host NIC tests: added stress



Single slot active with Iperf traffic (pulsed on/off hourly)



All 4 slots active with Iperf traffic (pulsed on/off hourly)





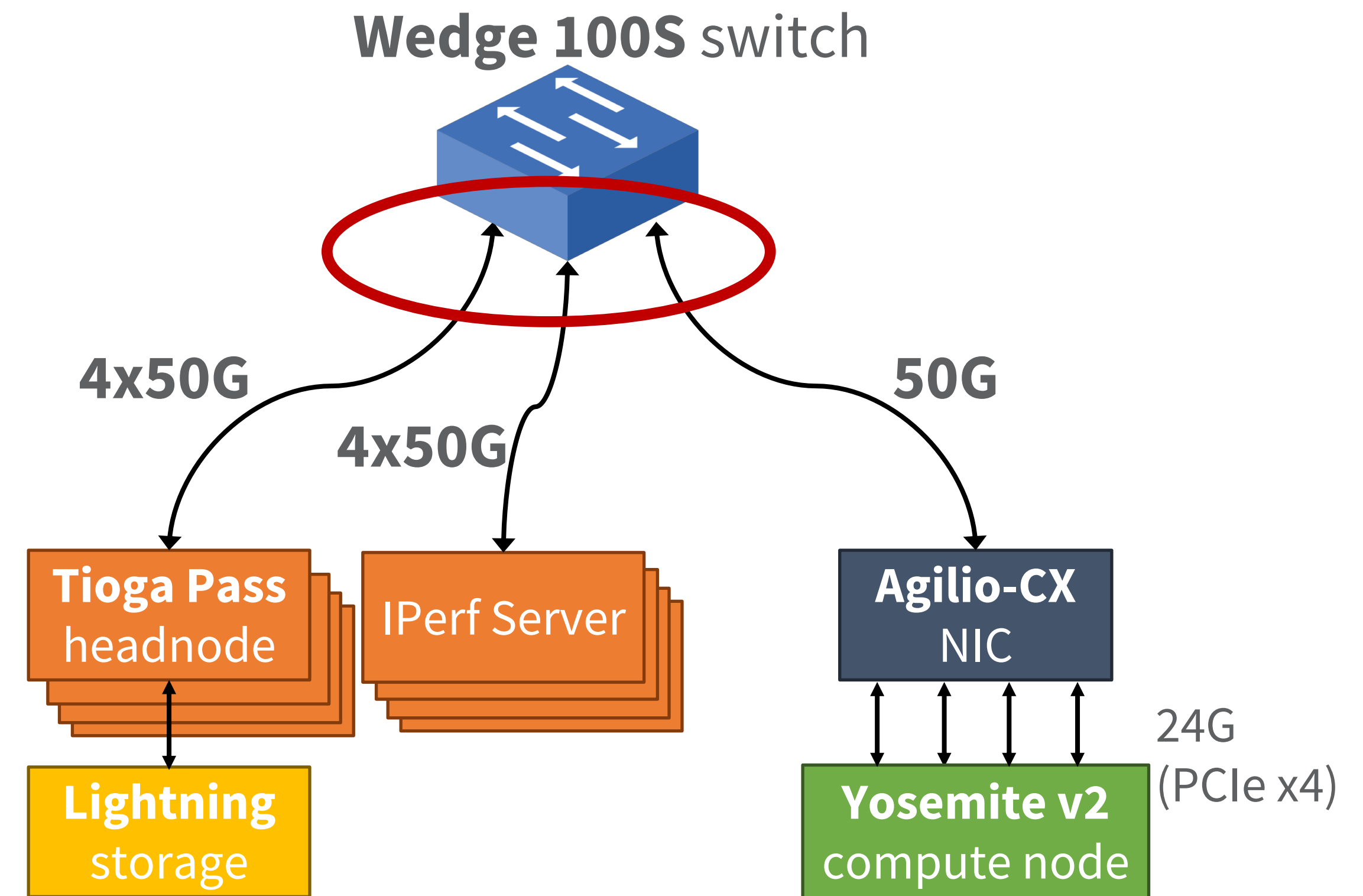
# ECN in Switch

High latencies due to **congestion in switch**.

Need a solution for congested switches in conjunction with NICs.

Had previously tested **enabling ECN in switches** for use with DCTCP, these tests confirmed the advantages of enabling ECN and using DCTCP across rack.

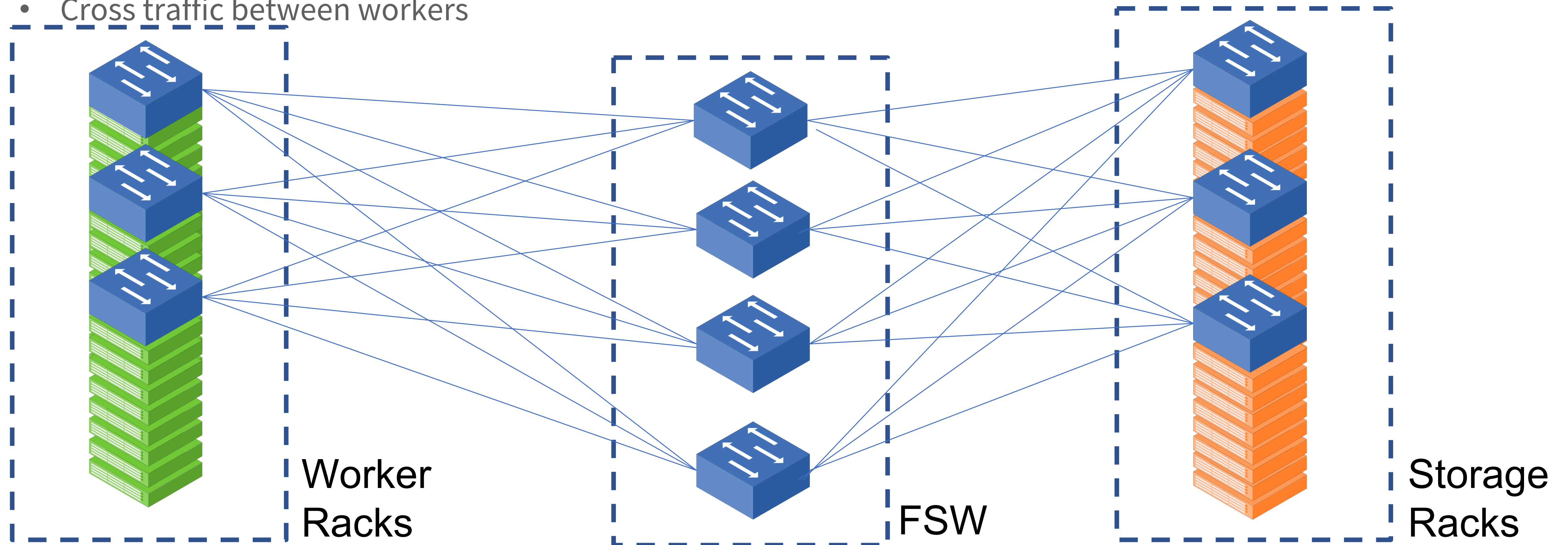
We will go over these tests in the following slides.



# DCTCP in Switch: Topology

## 6 rack tests

- 3 racks are store servers
- 3 racks (workers) read data from store servers
- Cross traffic between workers





# DCTCP in Switch: Benchmarks

	CUBIC	DCTCP
FSW to Worker Avg Link Util %	69.9	69.8
Storage CPU (%)	X	X
Worker CPU (%)	Y	Y+1%
FSW Discards (bits)	89M	235K (0.3%)
Worker rack discards (bits)	417M	0
Storage Retransmits	0.020 %	0.000 %
Worker Retransmits	0.173 %	0.078 %
Storage ECN CE Marked (%)		6.5
Worker ECN CE Marked (%)		12.8

**Note:** If we increase load until link utilization is 99%:

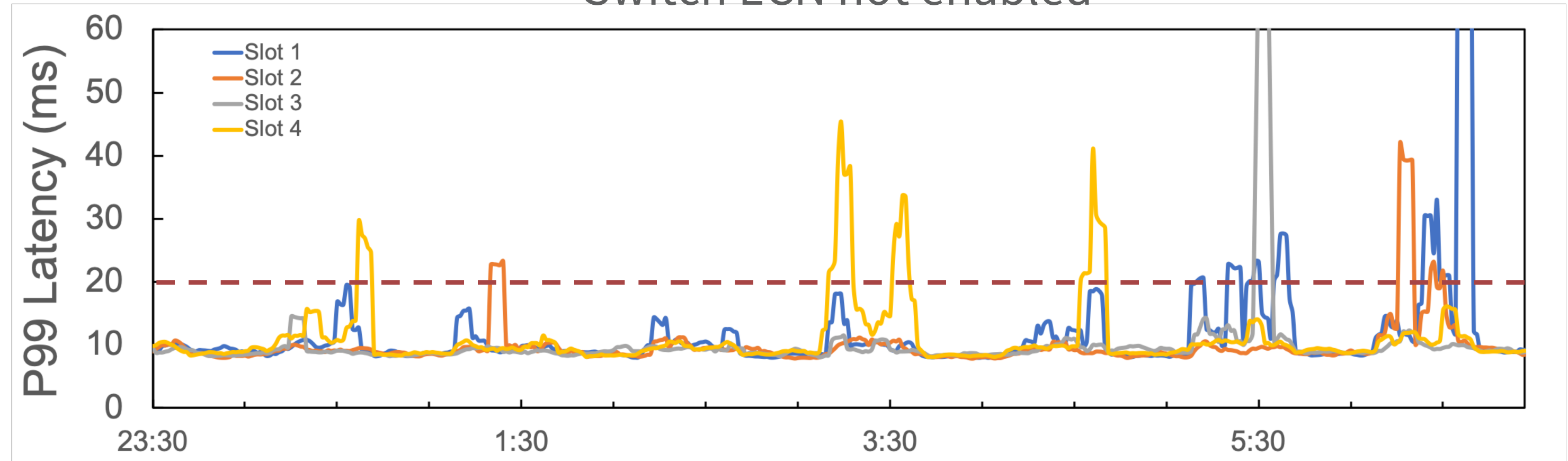
- FSW discards in CUBIC are 160B vs. 157M (0.1%) under DCTCP
- Storage retransmits are .6% under CUBIC vs. 0.001% under DCTCP

# Multi-Host NIC + ECN in Switch

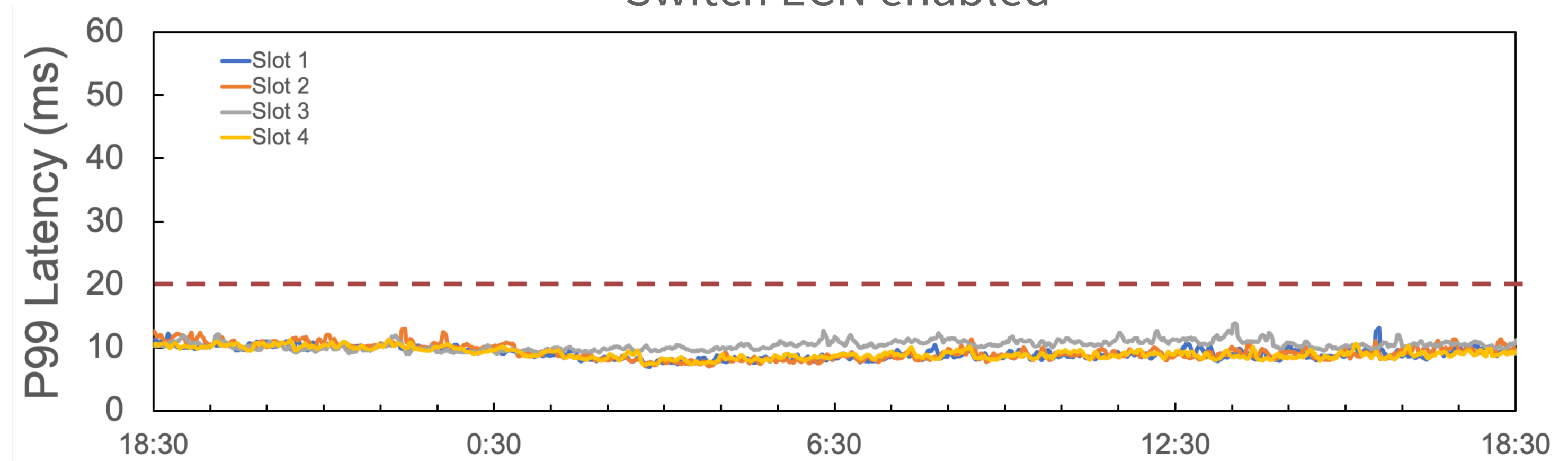
Suspected switch bottlenecks because all hosts used (including iperf servers) were in the same rack.

Enabled ECN marking in the switch and were able to significantly reduce tail latency.

Switch ECN not enabled



Switch ECN enabled





# Summary



NETWORKING

Linking advances in congestion control with OCP based SmartNICs reduces tail latency significantly

This allows OCP Yosemite v2 systems to be used in a wider variety of use cases, significantly improving efficiency

Without also implementing ECN/DCTCP in the switches it is possible to construct cases with high latency

Combining ECN/DCTCP in the multihost NICs and in RSWs, it is possible to 'guarantee' low tail latency



Case Studies



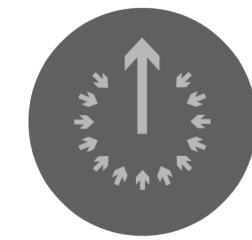
# Product/Facility Info



OPEN  
ACCEPTED™

## Wedge100S

[https://www.opencompute.org/wiki/Networking/SpecsAndDesigns#Facebook\\_Wedge\\_100S\\_32x100G](https://www.opencompute.org/wiki/Networking/SpecsAndDesigns#Facebook_Wedge_100S_32x100G)  
<https://www.opencompute.org/products/190/edgecore-networks-wedge100s-100gbe-data-center-switch>



OPEN  
INSPIRED™

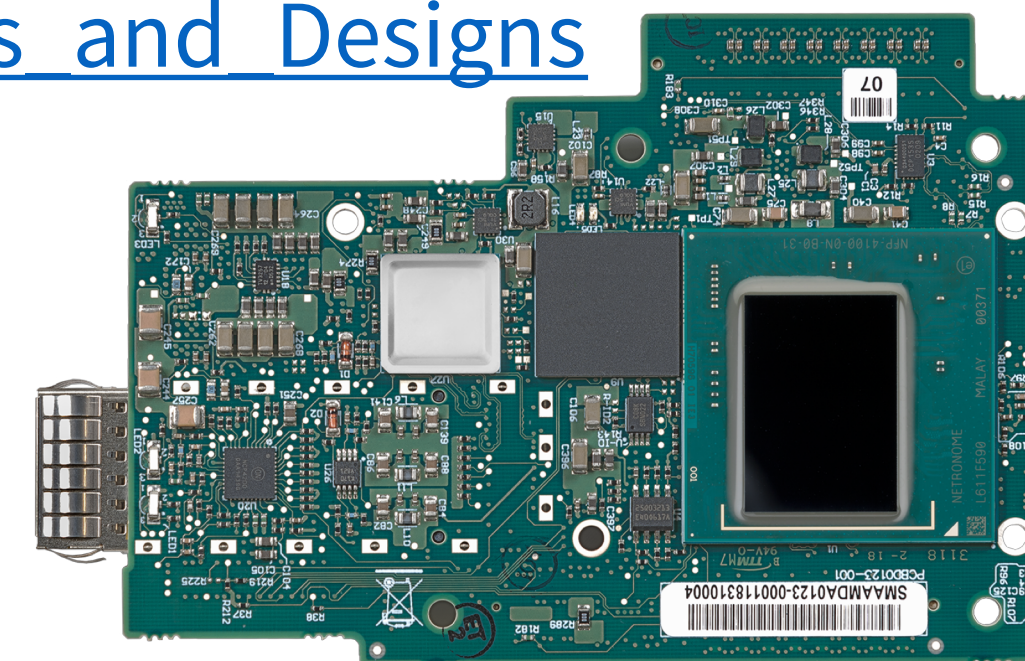
## Yosemite v2

<https://www.opencompute.org/products/275/wiwynn-yosemite-v2>



## Agilio-CX 50G OCP NIC

[https://www.opencompute.org/wiki/Server/Mezz#Specifications\\_and\\_Designs](https://www.opencompute.org/wiki/Server/Mezz#Specifications_and_Designs)





# Call to Action

**Netdev (kernel):** [netdev@vger.kernel.org](mailto:netdev@vger.kernel.org)

**Mezz:** [opencompute-mezz-card@lists.opencompute.org](mailto:opencompute-mezz-card@lists.opencompute.org)

**Server:** [opencompute-server@lists.opencompute.org](mailto:opencompute-server@lists.opencompute.org)

**Switch:** [opencompute-networking@lists.opencompute.org](mailto:opencompute-networking@lists.opencompute.org)

## Additional Information:

1. Flash Disaggregation: <http://csl.stanford.edu/~christos/publications/2016.flash.eurosys.pdf>
2. DCTCP: <https://web.stanford.edu/~balaji/papers/10datacenter.pdf>





# Open. Together.

OCP Global Summit | March 14–15, 2019

