



# Congestion Control for Large Scale RoCEv2 Deployments

**Hemal V. Shah and Moshe Voloshin**

Data Center Solutions Group (DCSG), Broadcom Inc.

# Agenda

- **Ethernet for HPC**
- **MPI and Communication Topologies**
- **PFC Challenges**
- **RoCEv2 with Congestion Control**
- **Congestion Control Evaluation**
  - OSU Benchmarks
  - HPCG
  - LAMMPS
  - GPCNeT

# Why Ethernet for HPC? – High Level Analysis

Attribute	InfiniBand	Ethernet	Notes
<b>Scalability (# of Nodes)</b>	1000s	1000s	Ethernet as a fabric scales
<b>Throughput</b>	200G	200G	Ethernet network speeds scale
<b>E2E Latency (B2B) – ½ RT</b>	600 ns	1 – 1.2us	Unloaded RoCE latency < 2x IB latency
<b>Switch Latency</b>	~100ns	200-500ns	ETH switch latency in future ~200-300ns
<b>Transport/Network</b>	IBA Transport	RoCEv2/UDP/IP	Ethernet ecosystem not as limited as IB
<b>Verbs API</b>	OFA OFED	OFA OFED	Common host interface
<b>Congestion Management</b>	FECN/BECN	ECN	Standards based congestion control
<b>Availability</b>	Limited sources	Many sources	Large number of Ethernet suppliers

- **Ethernet has advantages in its ubiquitous deployment → Ethernet is Everywhere**
- **Ethernet ecosystem is mature and standards based → Benefits HPC markets**
- **100/200G Ethernet latencies (w/ RoCE) meets HPC requirements**

# RoCEv2 for HPC

## High Performance

- 100 Gbps or higher
- 1-1.2 usec latency
- High packet per second

## CPU Efficient

- RDMA (RoCEv2)
- Kernel bypass
- Hypervisor bypass with SR-IOV

## RoCEv2

- Simple driver model
- Verbs and MPI
- Application offloads

## Software Infrastructure

- Leverages Ethernet
- Converged Infrastructure
- Congestion control

## Network Infrastructure

## MPI and Communication Topologies

- **MPI is widely used in HPC/ML clusters as the communication layer**
- **A process group in MPI represents a collection of processes**
- **The number of processes can be 100s per node**
- **The number of nodes can scale to 1000s in a cluster**
- **The communication pattern is represented by a logical topology**
  - Ring, Binary cube, Tree, etc.
- **Selection of logical topologies depends on apps and communication libs**
- **MPI collectives (Gather, Reduce..) can create congestion in the network**

## Challenges with PFC without Congestion Control

- **Priority Flow Control (PFC) is used for lossless service**
- **PFC is a point-to-point protocol between two Ethernet endpoints**
- **PFC can result in congestion spreading**
- **PFC can create PFC storm due to slow receivers**
- **PFC may result in transport live-lock**

## Congestion Control (CC) with RoCEv2

- **ECN based CC schemes do not require any additional infrastructure support**
- **Congestion control without PFC can be sufficient for most of the workloads**
- **CC with PFC addresses PFC storms/live locks & preserves lossless service**
- **Even w/ large number of competing flows switch egress Q peaks are low**
- **Reaction by sender is quick – few 10s of micro-seconds due to low Q level**
  - Even with low marking threshold, network utilization is high
- **Low marking threshold delivers low end-to-end latency with min interference**
- **Low marking threshold leaves majority of switch buffer for incast absorption**
- **Both probabilistic and deterministic marking are possible**

## RoCEv2 Application Performance Under Congestion

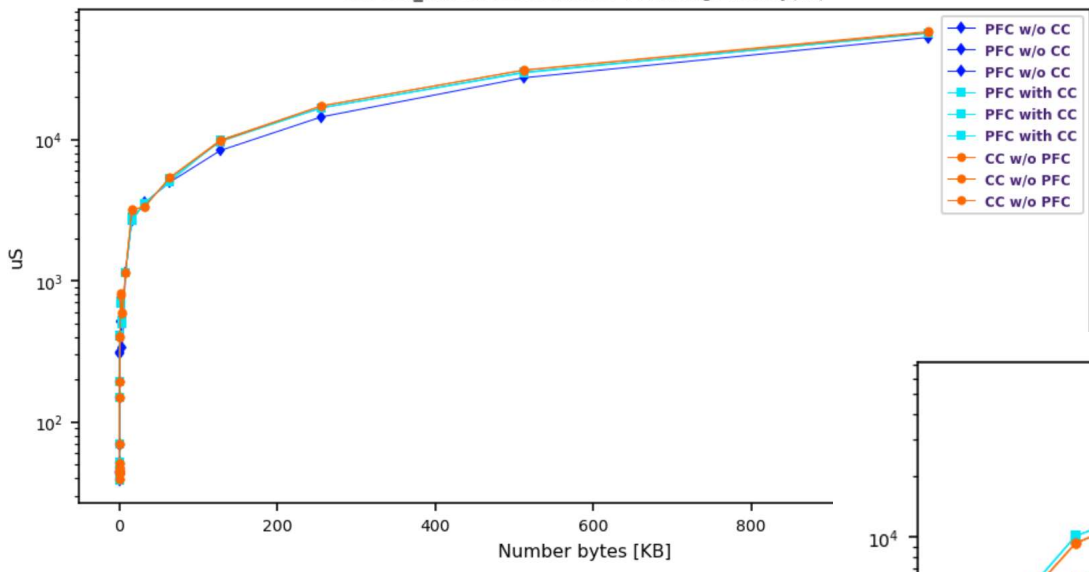
Test Scenario	Overview
<b>OSU Benchmarks</b>	Blocking and non blocking Collective benchmarks for various sizes and with various PPN (Processes Per Node) over 2 to 32 nodes
<b>HPCG</b>	High Performance Conjugate Gradient Benchmark for HPC, with 8, 16, 32 PPN on 8, 16, 32 nodes
<b>LAMMPS</b>	5 benchmarks of Molecular Dynamics with 32,000 atoms per core Scaling efficiency charts relative to CPU time on single core for 32,000 atoms Chart title shows 1 node 1 PPN loop time in seconds for 32,000 atoms

**All tests ran with NIC link BW of 100 gbps in 3 configurations:  
PFC without CC, PFC with CC and CC without PFC**

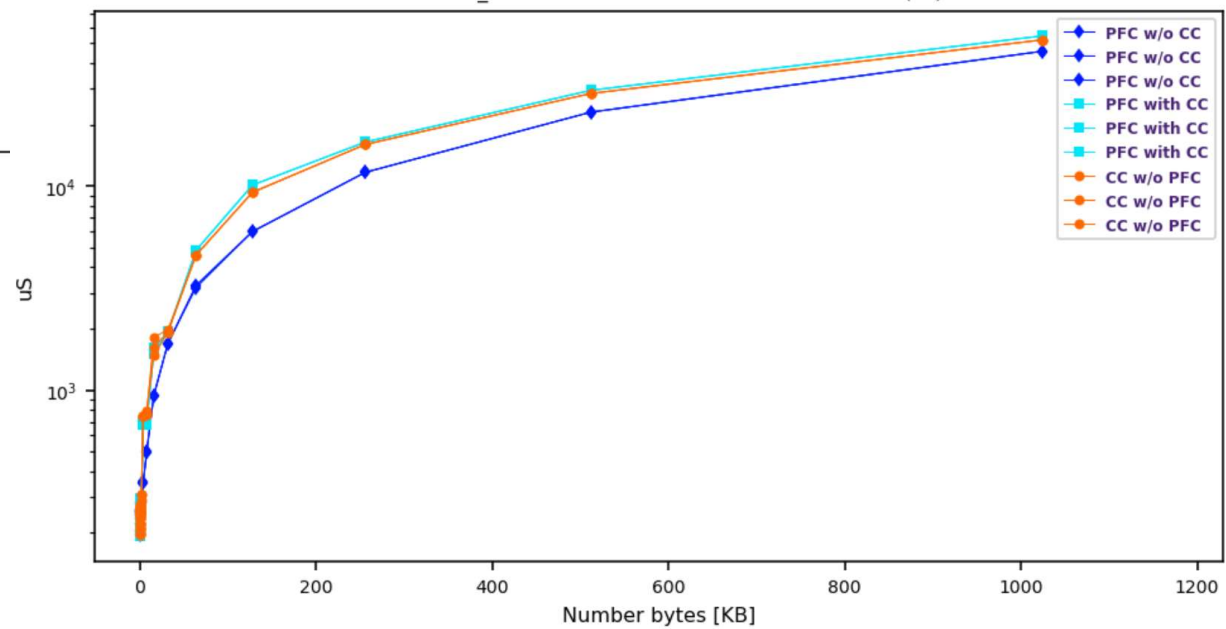


# OSU Benchmarks Results – completion time [uS]

osu osu\_Alltoall test 32 nodes 4 PPN Avg Latency(us)

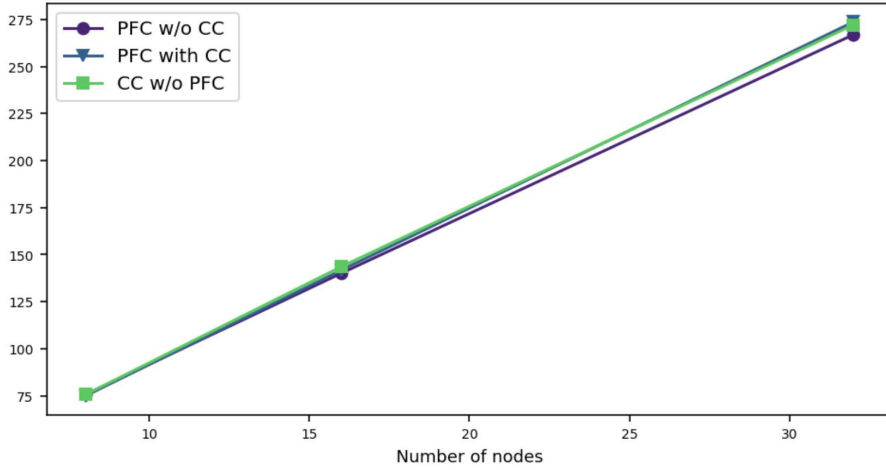


osunb osu\_iAlltoall test 32 nodes 4 PPN Pure Comm.(us)

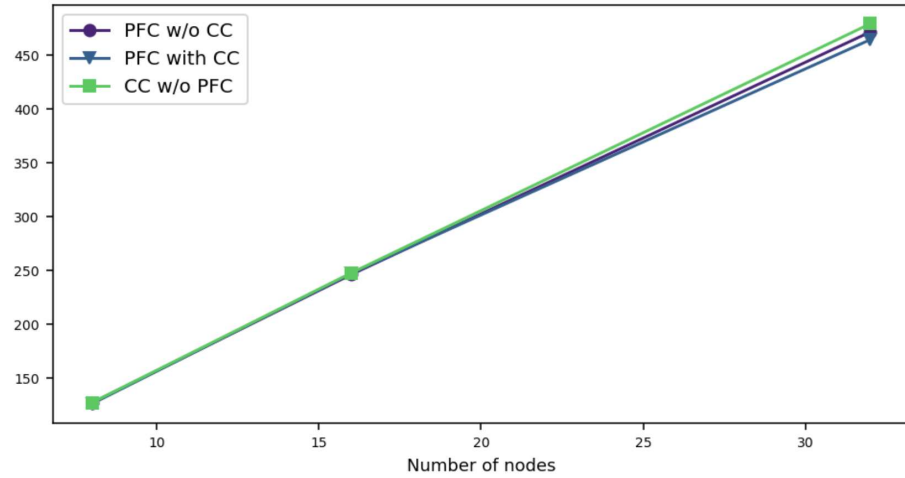


# HPCG Results

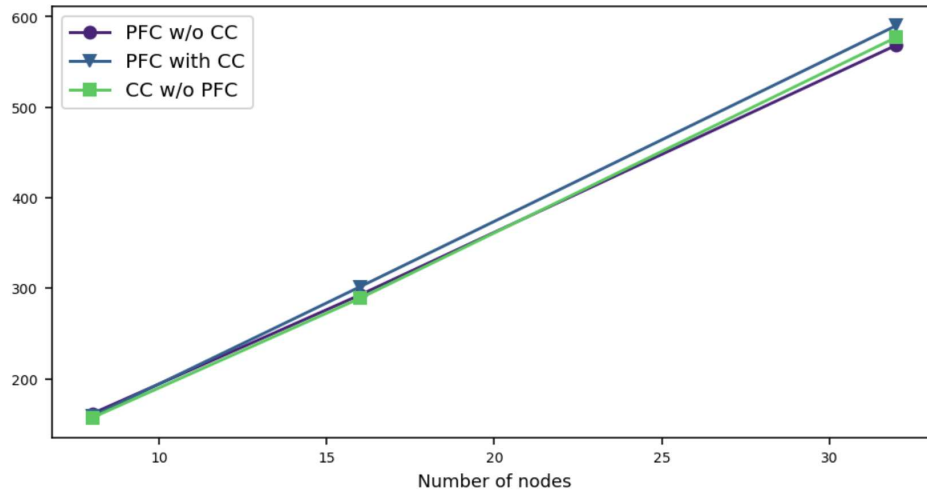
GFLOPs/s for 8 PPN



GFLOPs/s for 16 PPN



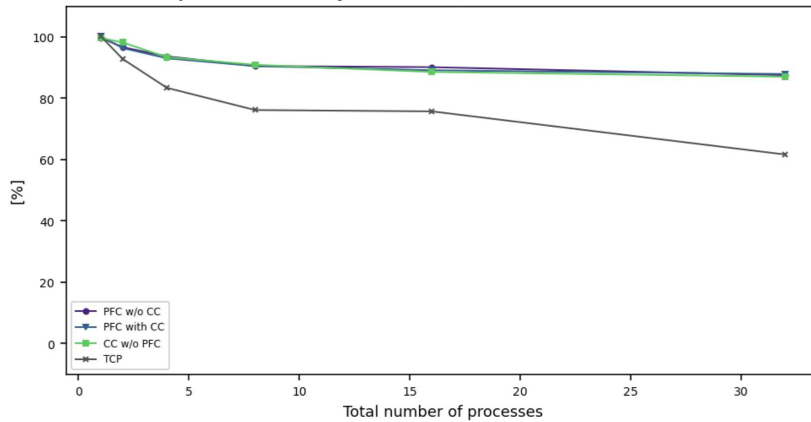
GFLOPs/s for 32 PPN



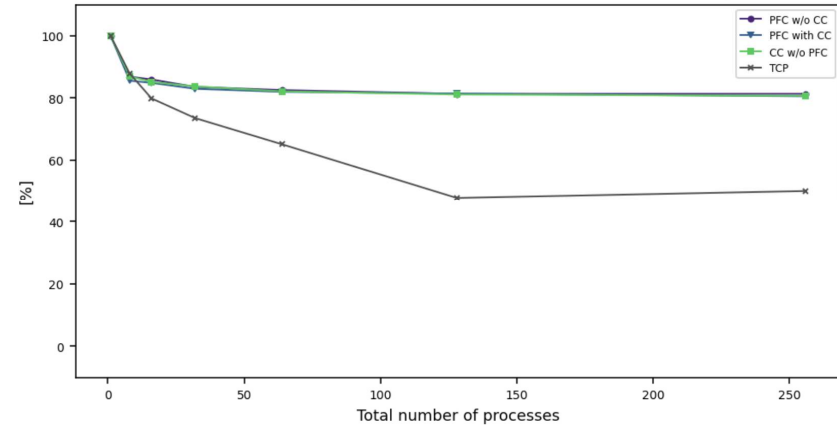
- HPCG application scales well in all configurations with varying PPN

# LAMMPS efficiency vs. # processes

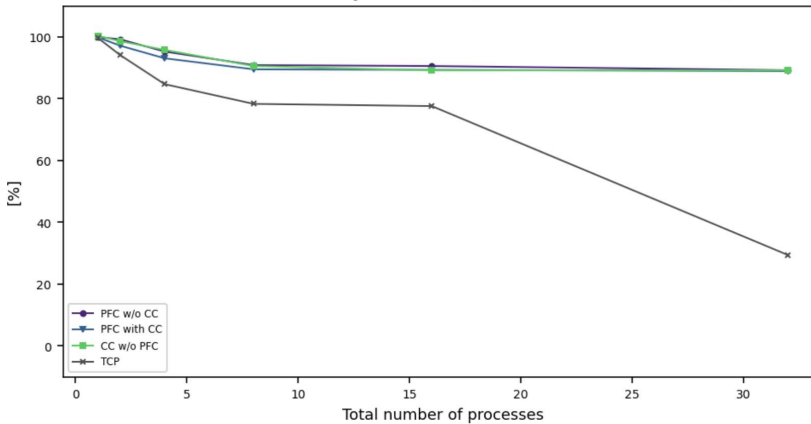
lj test - Efficiency for 1 PPN (1 Node 1 PPN time 1.61)



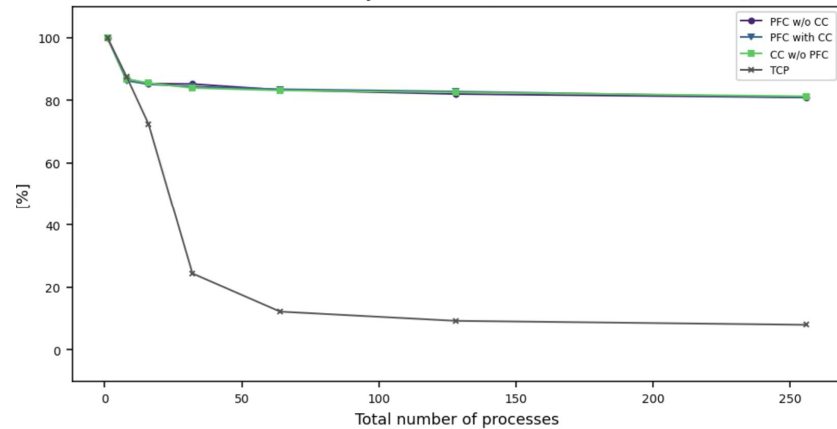
lj test - Efficiency for 8 PPN (1 Node 1 PPN time 1.61)



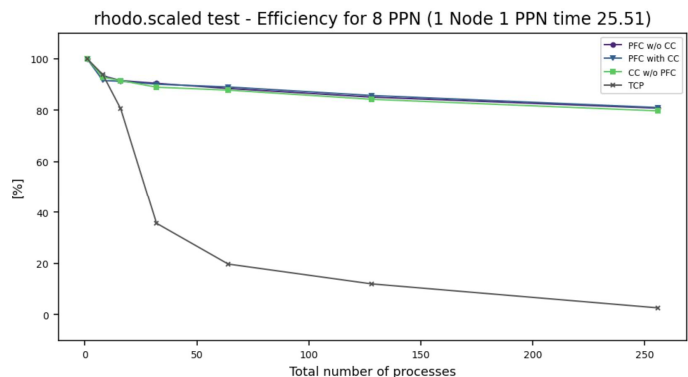
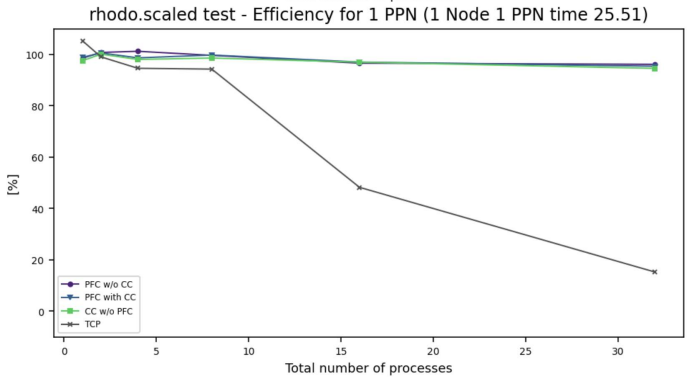
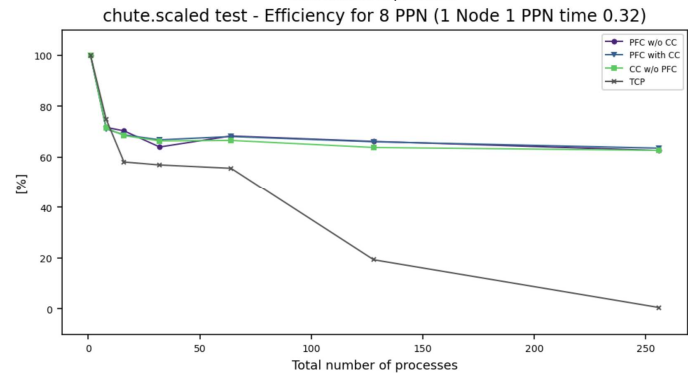
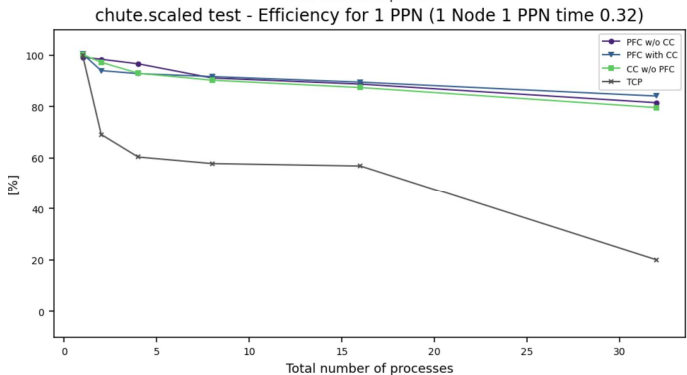
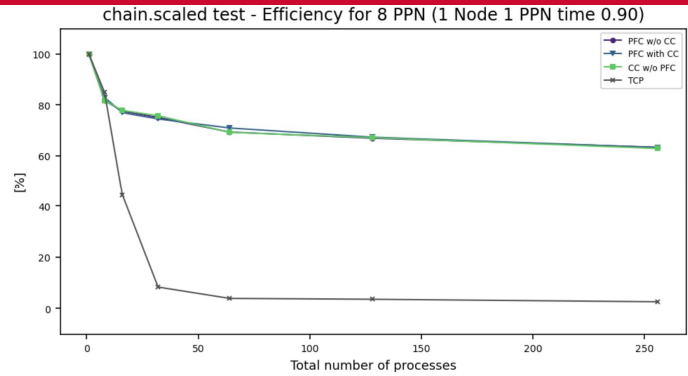
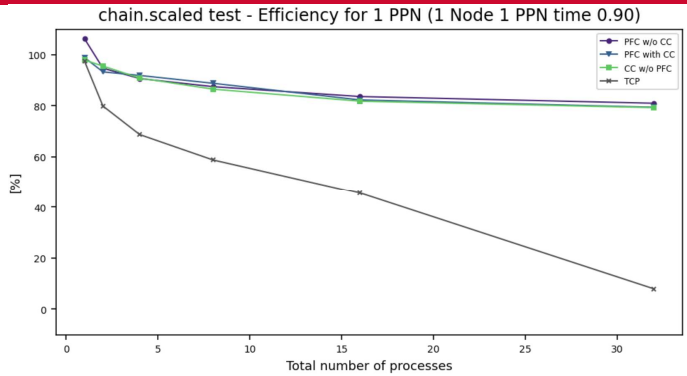
eam test - Efficiency for 1 PPN (1 Node 1 PPN time 4.43)



eam test - Efficiency for 8 PPN (1 Node 1 PPN time 4.43)



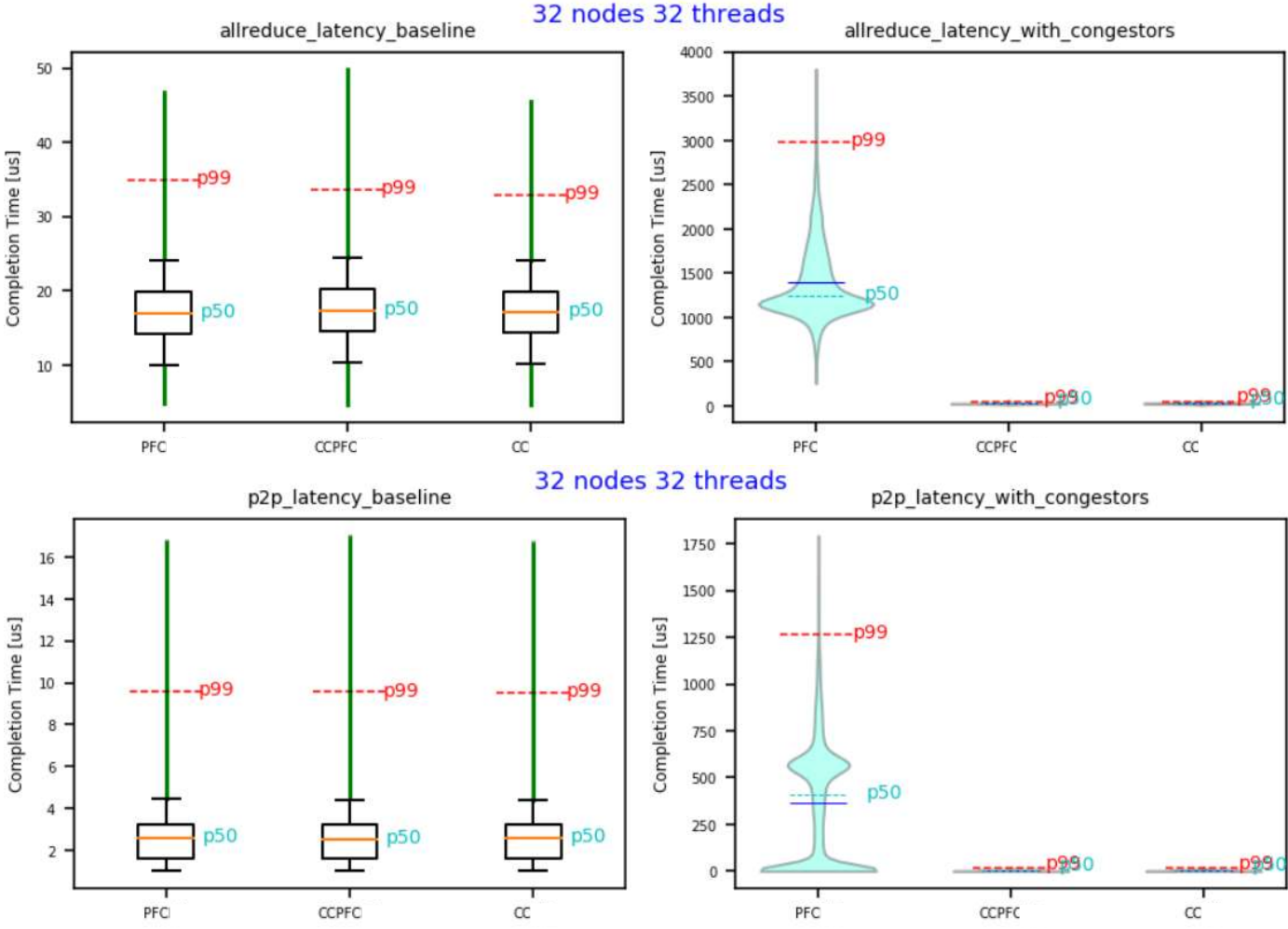
- ROCE significantly outperform TCP in all configurations



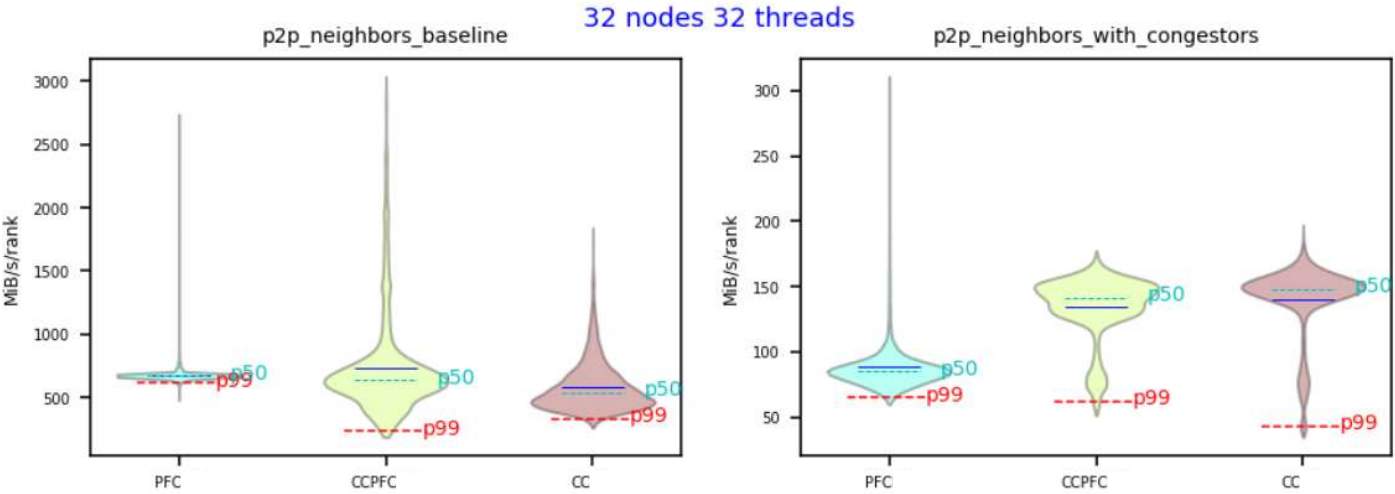
## GPCNeT benchmark

<b>GPCNeT</b>	<p>Global Performance and Congestion Network Test</p> <p>MPI test designed to measure relative performance under load and congestion</p> <p>Designed for large multi-layer switch network</p> <ul style="list-style-type: none"><li>• 20% of nodes w/ test tasks: allreduce, p2p latency, random ring neighbor exchange</li><li>• 80% of nodes assigned with congestor tasks: All2All, incast, RMA put and get</li><li>• Nodes would share switch buffering resources and will cross path between switches</li></ul>
---------------	--

# GPCNeT results – Results on 32 Nodes



# GPCNeT results – Results on 32 Nodes



**CCONLY**

**CC+PFC**

**PFC**

Network Tests running with Congestion Tests - Key Results		with Congestion Tests - Key Results	
Name	Congestion	Congestion Impact	Congestion Impact Factor
	Avg	Avg	Avg
RR Two-sided Lat (8 B)	1.8X	1.8X	120.9X
RR Two-sided BW+Sync (131072 B)	3.9X	4.6X	7.9X
Multiple Allreduce (8 B)	1.7X	1.5X	79.8X

## Congestion Control (CC) for RoCEv2

- **RoCEv2 is designed to scale – no inherent limitation at the protocol level**
- **ROCEv2 demonstrate significant performance advantage over TCP**
- **PFC provides lossless service with a significant interference/blocking**
  - PFC without congestion control should be avoided
- **CC with or without PFC is essential for node and process scaling**
  - ECN marking in switches enable Congestion Notifications to minimize congestion
  - CC algorithms are evolving to provide better congestion avoidance & faster congestion reaction
  - CC algorithm that maintains low switch queue level reduces interference/blocking
- **CC enhancements in NICs further improve performance & scalability**
  - ECN marking/CNP generation
  - Hardware-based congestion control
  - Deterministic marking policy (DCTCP style)



Thank You





**BROADCOM**®

connecting everything®