Daniele De Sensi[1], Salvatore Di Girolamo[1], Kim McMahon[2], Duncan Roweth[2], Torsten Hoefler[1]

[1]ETH Zurich, Switzerland
[2]HPE Cray

# *An In-Depth Analysis of the Slingshot Interconnect*

SPCL
ETH zürich
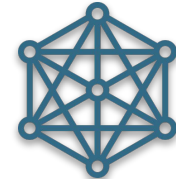
Hewlett Packard
Enterprise

1

All HPC traffic layered over **RoCEv2**

Efficient **software stack**

High-Radix **Switches**

Low-Diameter **Topology**

**Congestion Control**
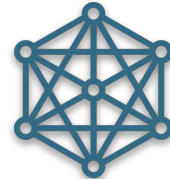
**Adaptive Routing**

**Quality of Service**

All HPC traffic layered over **RoCEv2**

Efficient **software stack**

High-Radix **Switches**

Low-Diameter **Topology**

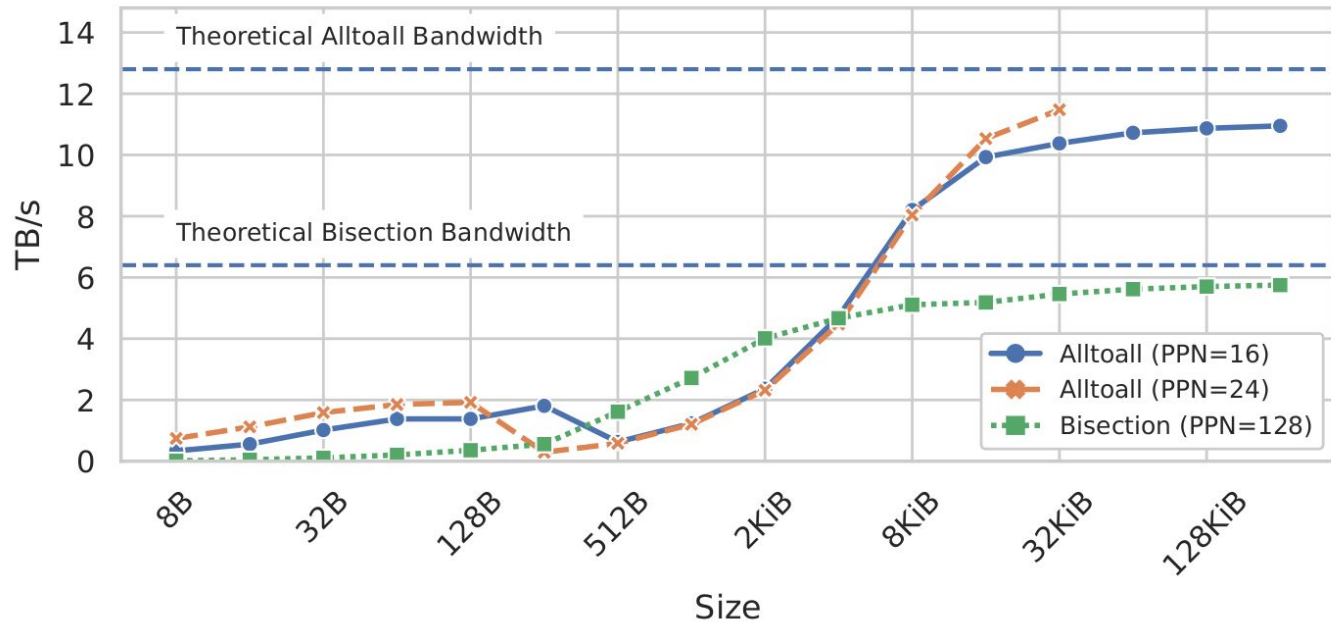**Congestion Control**

**Adaptive Routing**

**Quality of Service**

All HPC traffic layered over **RoCEv2**

Efficient **software stack**

High-Radix **Switches**

Low-Diameter **Topology**

**Congestion Control**

**Adaptive Routing**

**Quality of Service**
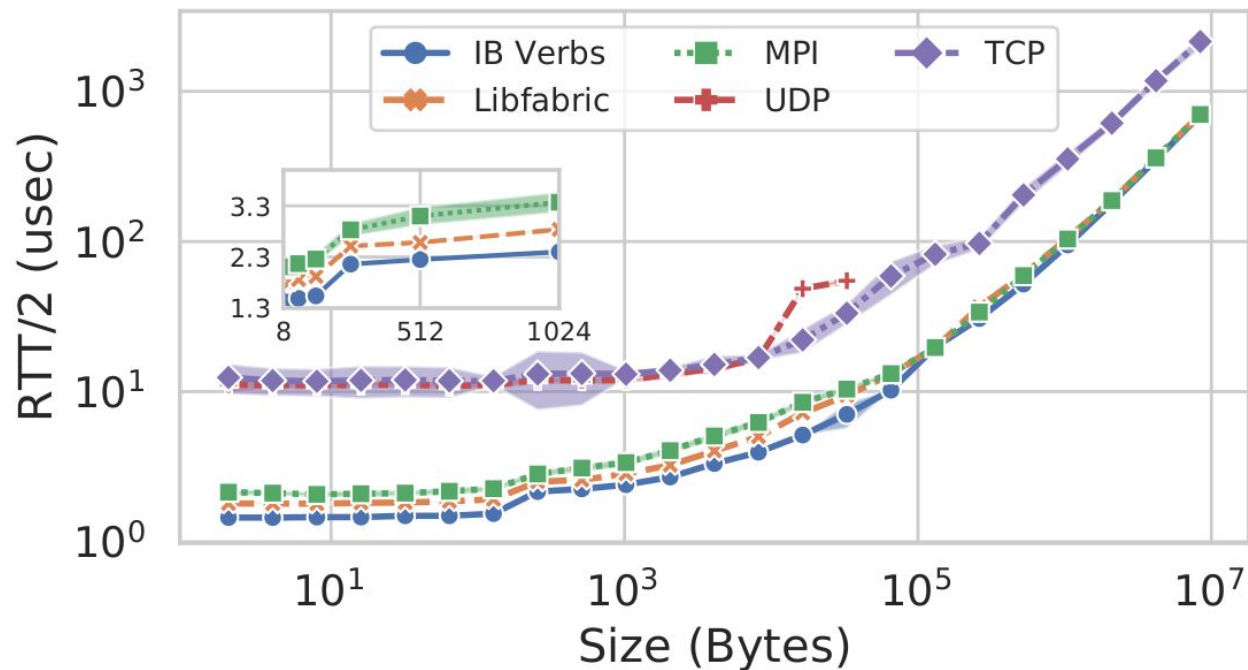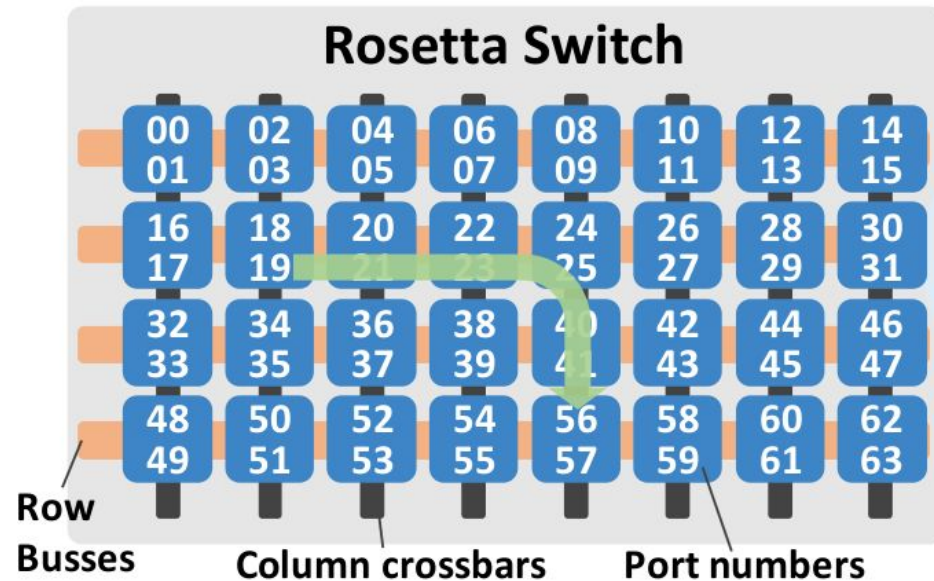
# ETHERNET ENHANCEMENTS

Can process both **standard** and **enhanced** Ethernet packets

1024 nodes

# SOFTWARE STACK

Standard **TCP/IP** stack or **libfabric**

All HPC traffic layered over **RoCEv2**

Efficient **software stack**

High-Radix **Switches**

Low-Diameter **Topology**

**Congestion Control**

**Adaptive Routing**

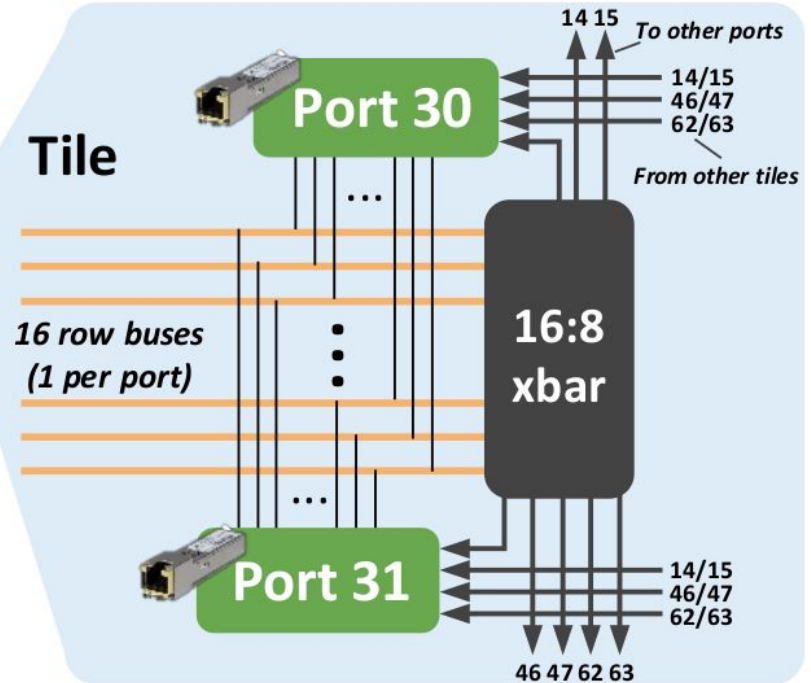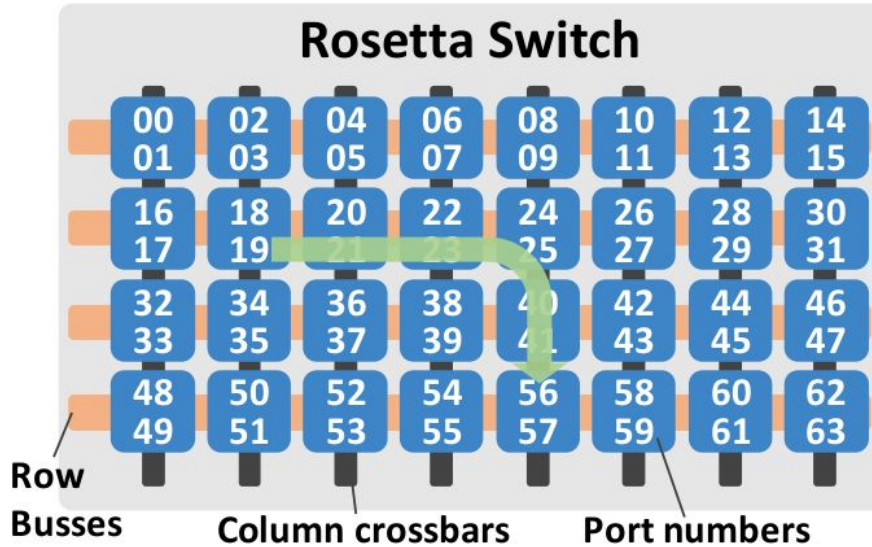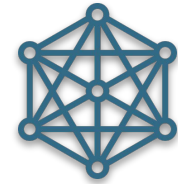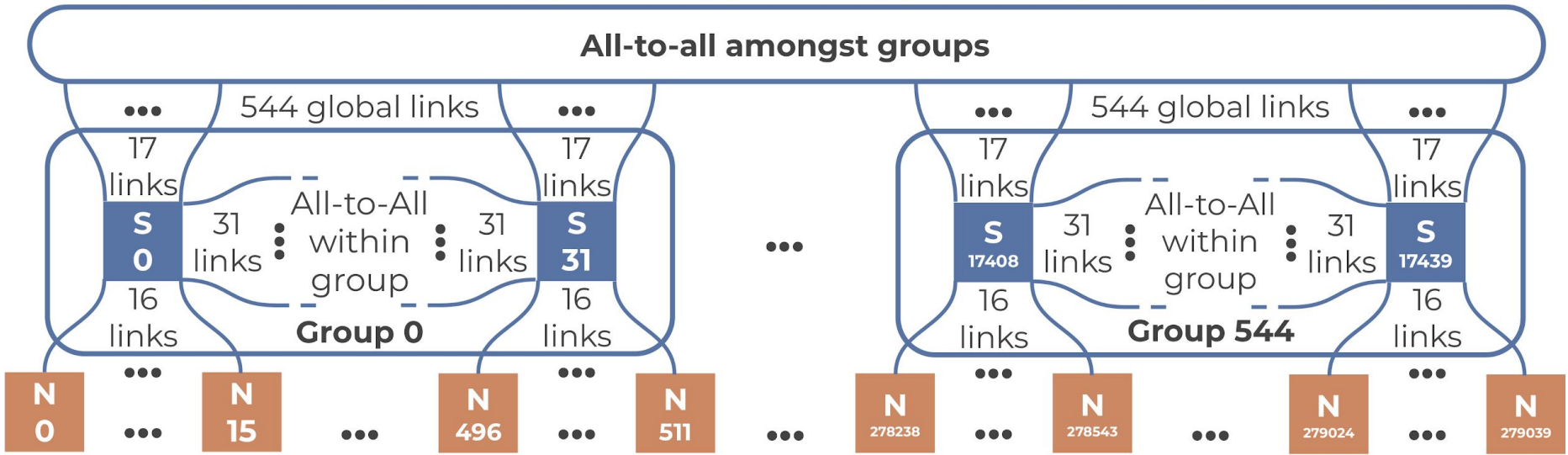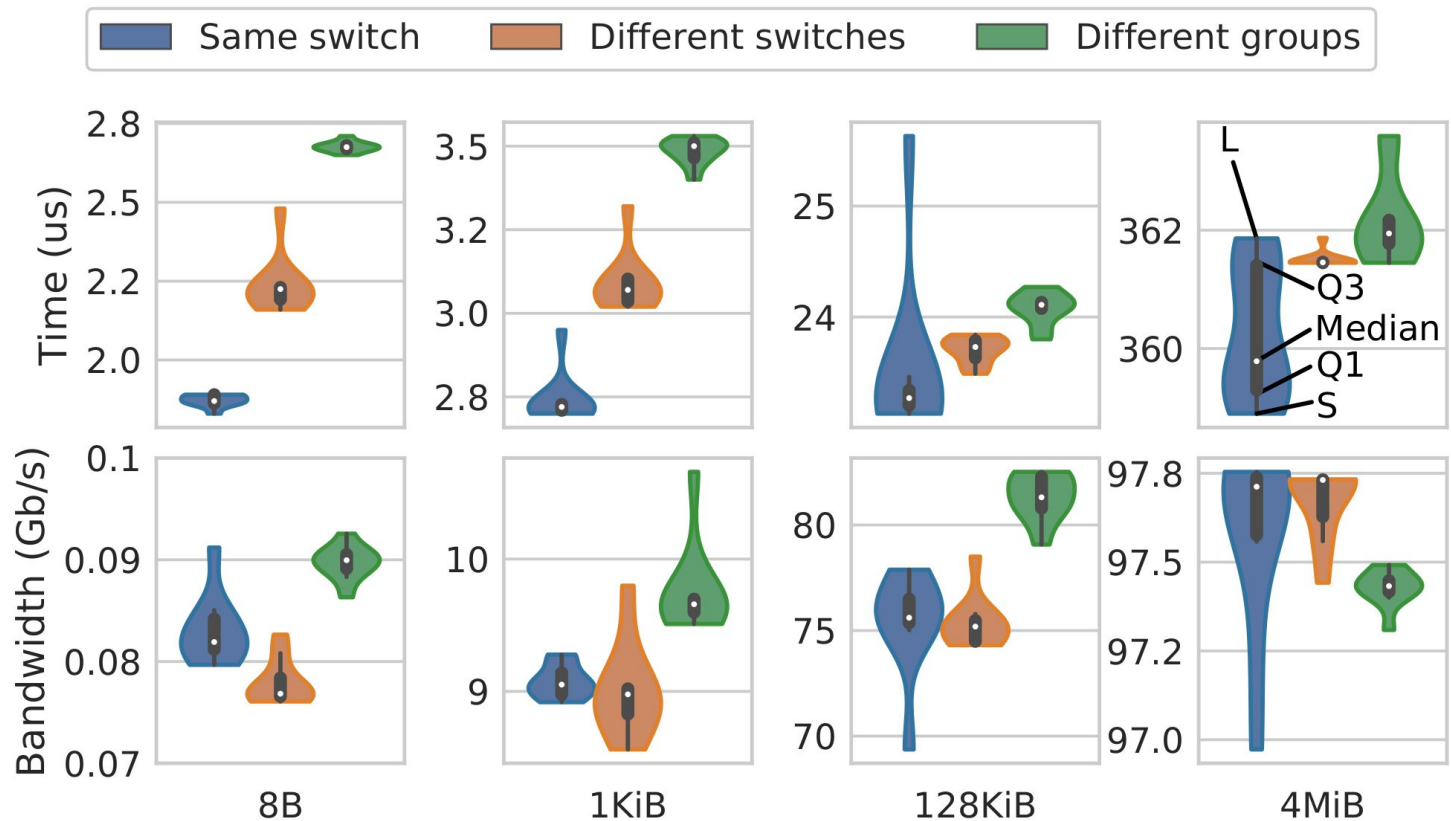**Quality of Service**

# SWITCH - ROSETTA

**64 x 200Gb/s ports**

**32 tiles**

# SWITCH - ROSETTA

All HPC traffic layered over **RoCEv2**

Efficient **software stack**

High-Radix **Switches**

Low-Diameter **Topology**

**Congestion Control**

**Adaptive Routing**

**Quality of Service**

# SLINGSHOT TOPOLOGY

Switches can be connected into **arbitrary topologies**

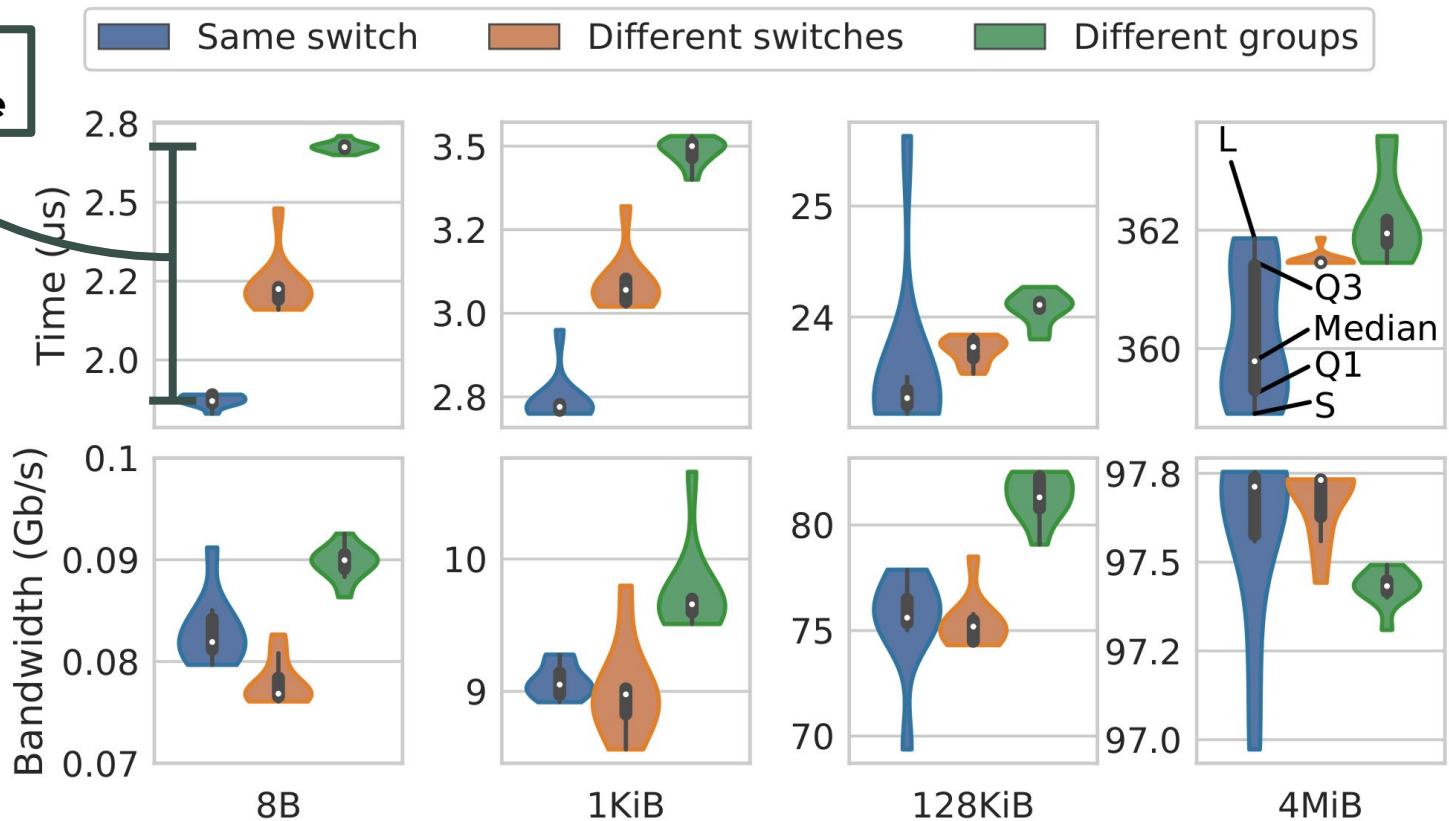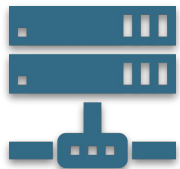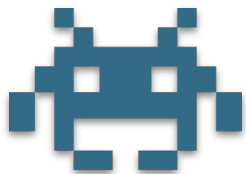**Dragonfly** is the default topology



All-to-all amongst groups

544 global links

17 links

**S 0**  31 links  All-to-All within group  31 links  **S 31**  17 links

16 links  **Group 0**  16 links

544 global links

17 links

**S 17408**  31 links  All-to-All within group  31 links  **S 17439**  17 links

16 links  **Group 544**  16 links

**N 0**   **N 15**   **N 496**   **N 511**   **N 278238**   **N 278543**   **N 279024**   **N 279039**

= Endpoints       = Switches

11

# SLINGSHOT TOPOLOGY - LATENCY & BANDWIDTH

# SLINGSHOT TOPOLOGY - LATENCY & BANDWIDTH

All HPC traffic layered over **RoCEv2**

Efficient **software stack**

High-Radix **Switches**

Low-Diameter **Topology**

**Congestion Control**

**Adaptive Routing**

**Quality of Service**

# CONGESTION CONTROL

ECN/QCN **hard to tune** and **slow to converge**

**Tracks** the traffic between **every pair of endpoints**

Slows down **offending traffic** only

**Improves** average and tail **latencies**

15

# CONGESTION CONTROL TESTS

Run two concurrent jobs: **victim** and **aggressor**

| Microbenchs. | silo |
| MILC | sphinx |
| HPCG | xapian |
| LAMMPS | img-dnn |
| FFT | resnet-proxy |

**Tailbench**

**incast**
(heaviest congestion)

**all-to-all**
(intermediate congestion)

# CONGESTION IMPACT - 512 NODES



17

# CONGESTION IMPACT - ADDITIONAL ANALYSIS

All HPC traffic layered over **RoCEv2**

Efficient **software stack**

High-Radix **Switches**
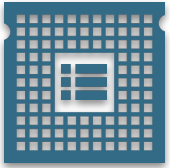
Low-Diameter **Topology**

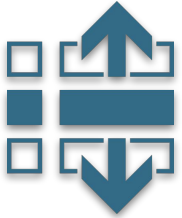**Congestion Control**

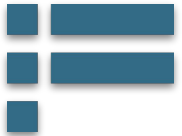**Adaptive Routing**

**Quality of Service**

# QUALITY OF SERVICE

Each traffic class occupies **hardware resources** in the switches

**Tunable** priority, ordering, minimum/maximum bandwidth, …

Jobs can be assigned to a small number **traffic classes**
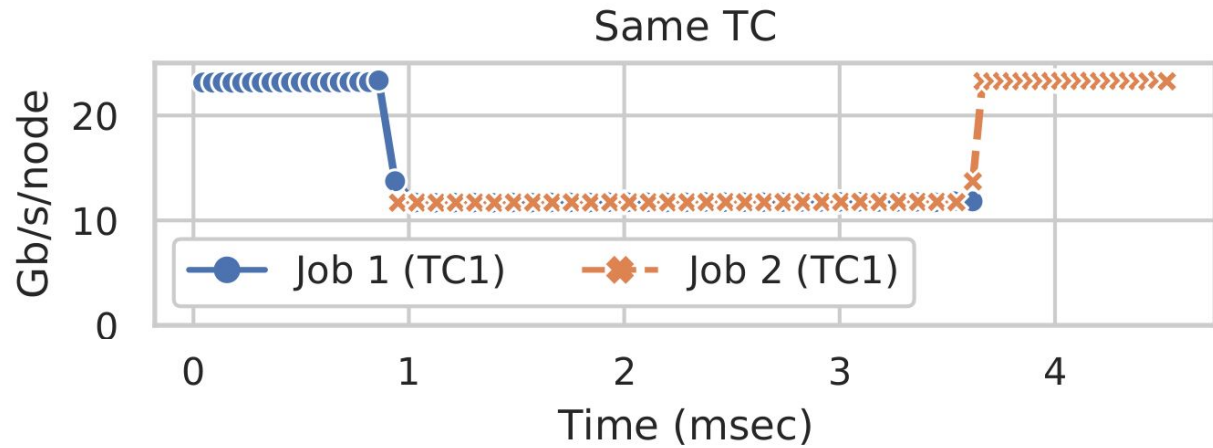
Traffic class can be changed on a **per-packet** basis

20

# QOS TESTS

**25%** bandwidth **tapering**

2 jobs running **bisection bandwidth** tests

**TC1**: 80% minimum bandwidth

**TC2**: 10% minimum bandwidth



Same TC

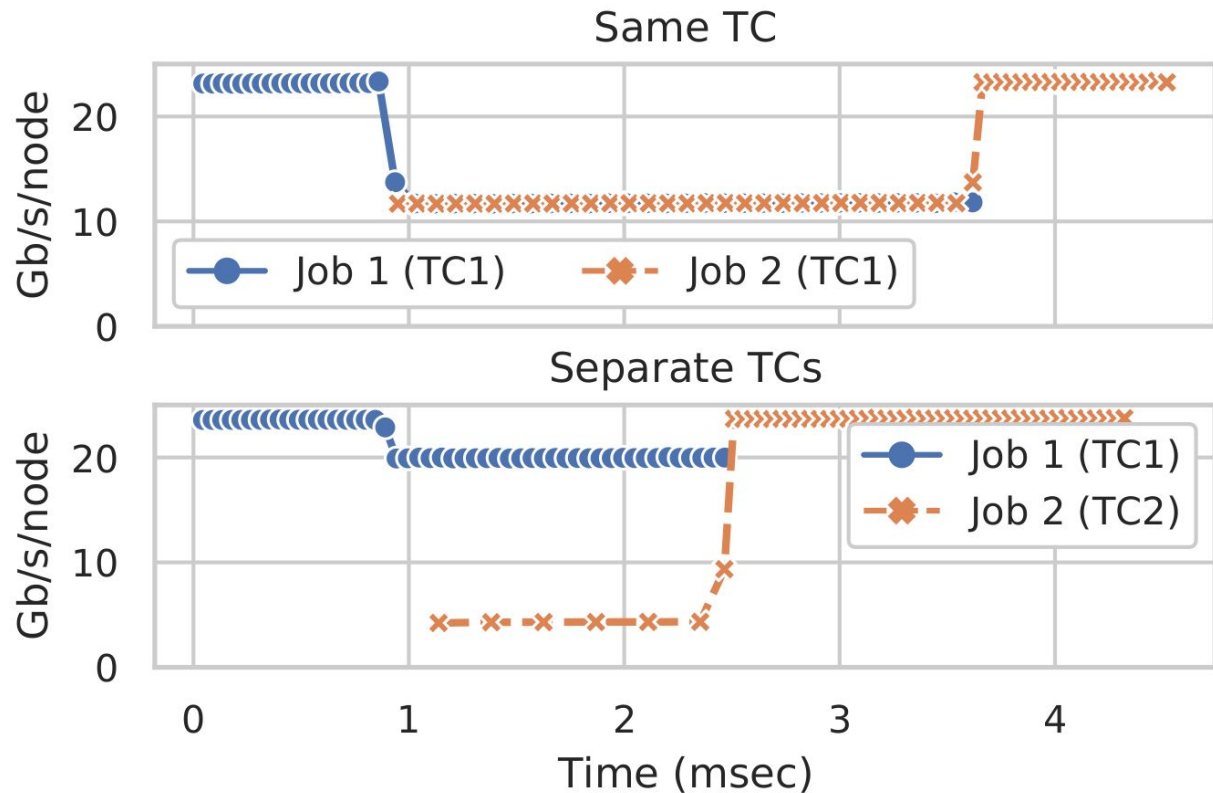Job 1 (TC1)    Job 2 (TC1)

Gb/s/node vs Time (msec)

# QOS TESTS

**25%** bandwidth **tapering**

2 jobs running **bisection bandwidth** tests

**TC1**: 80% minimum bandwidth

**TC2**: 10% minimum bandwidth

CONCLUSIONS

# CONCLUSIONS