# Boost HPC application performance thanks to hardware offload

**Bull**
atos technologies

27-28 June, 2017

© Atos

**Bull**
atos technologies

# Disclaimer

Atos, the Atos logo, Atos Consulting, Atos Worldgrid, Worldline, BlueKiwi, Canopy the Open Cloud Company, Yunano, Zero Email, Zero Email Certified and The Zero Email Company are registered trademarks of Atos. June 2017. © 2017 Atos. Confidential information owned by Atos, to be used by the recipient only. This document, or any part of it, may not be reproduced, copied, circulated and/or distributed nor quoted without prior written approval from Atos.
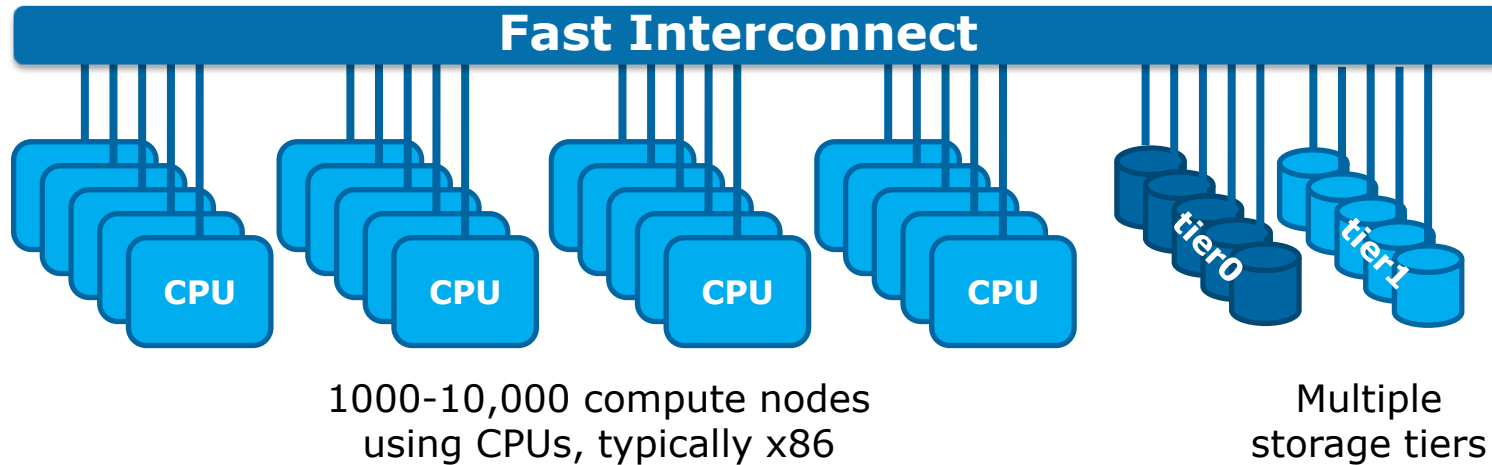
Atos may make changes to specifications and product descriptions at any time, without notice.

All products, computer systems, dates, and figures specified are preliminary based on current expectations, and are subject to change without notice. This is not a binding offer.
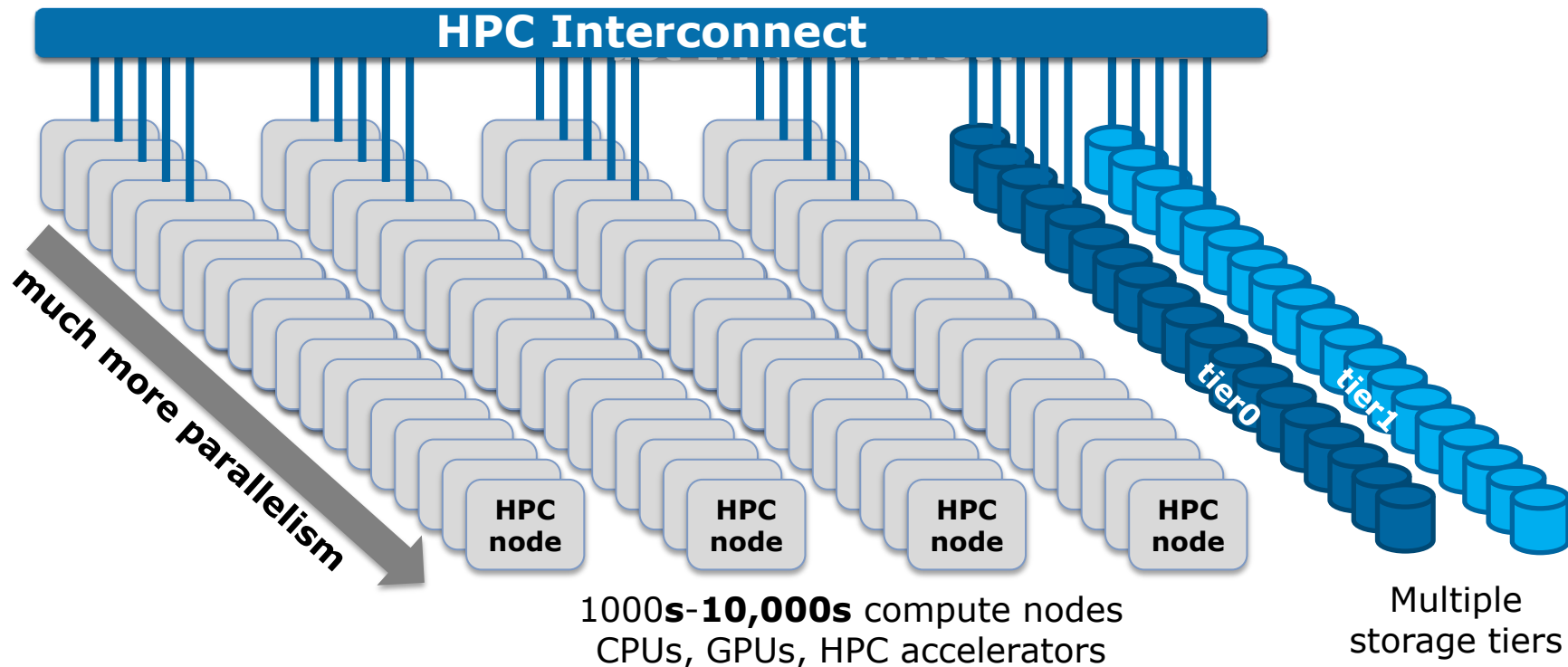
Atos hardware and software products may contain design defects or errors known as errata, which may cause the product to deviate from published specifications. Current characterized errata are available on request.

Copyright © 2017, Atos. All rights reserved.

# HPC systems are highly parallel Petaflops class featuring 1000s CPU nodes

**Fast Interconnect**

**CPU** **CPU** **CPU** **CPU**

tier0 tier1

1000-10,000 compute nodes using CPUs, typically x86

Multiple storage tiers

# 10-100 Pflops systems being deployed … with HPC specific Processing Units

**BXI**

**HPC Interconnect**

much more parallelism

**HPC node** **HPC node** **HPC node** **HPC node**

tier0 tier1

1000**s**-**10,000s** compute nodes
CPUs, GPUs, HPC accelerators

Multiple storage tiers
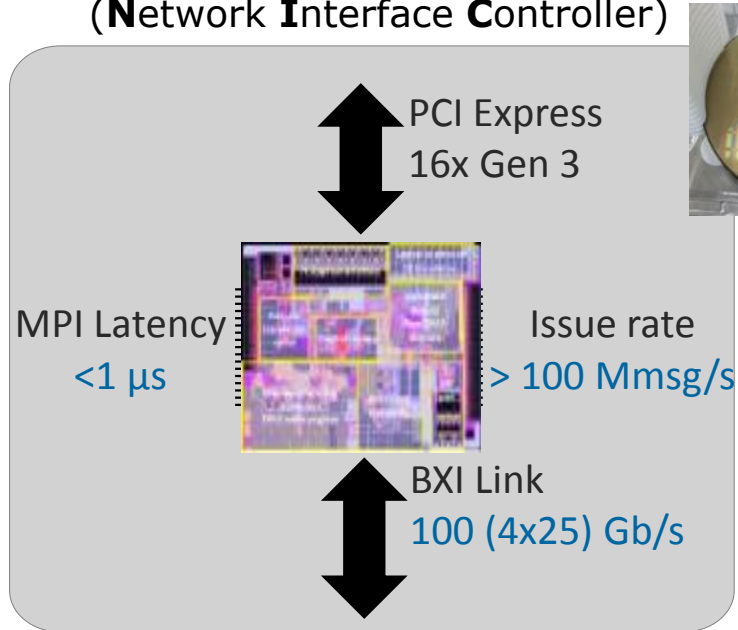
**Bull** atos technologies

# Network has to become intelligent

**BXI**

▶ **Smart interconnects**, based on the **ability to offload in hardware MPI semantics** from the host CPU, can be translated directly to **greater application performance**.

▶ Communications and computations progress completely independently.

▶ Performance is not impacted by heavy load on the host CPU.

▶ Point to point communication include MPI hardware matching.

▶ Triggered and atomic operations are used for protocol offloading (rendezvous, collectives, etc.).

Bull
atos technologies

# BXI – Interconnect overview

▶ **BXI 1st generation of Bull Exascale Interconnect**

▶ **BXI full acceleration in hardware for HPC applications**

▶ **BXI highly scalable, efficient and reliable**
  – Exascale scalability → 64k nodes,
  – Adaptive Routing,
  – Quality of Service (QoS),
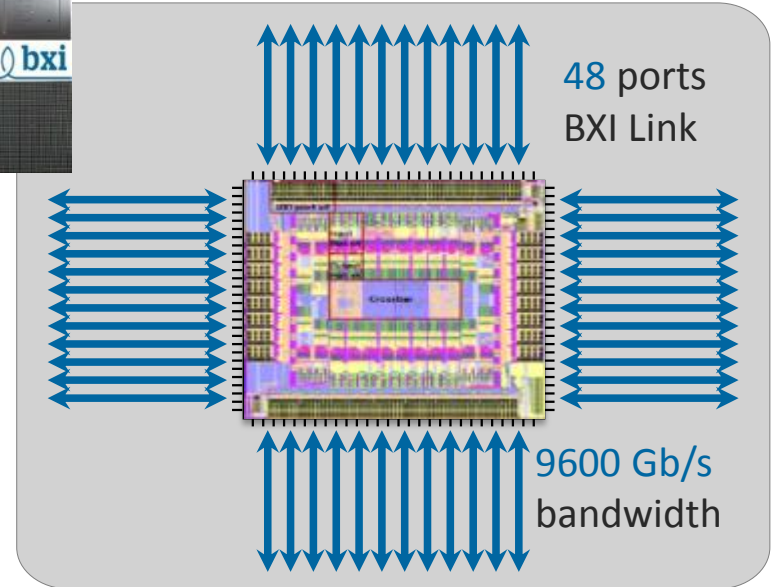  – End-to-end error checking + link level CRC + ASIC ECC.

# BXI Network is based on 2 ASICs

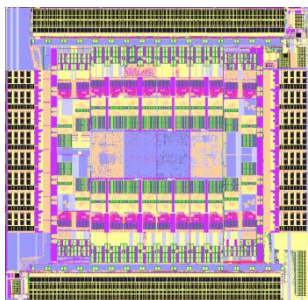**BXI**

► One ASIC, in nodes, for NIC
(**N**etwork **I**nterface **C**ontroller)

► One ASIC for Switch equipment



PCI Express
16x Gen 3

MPI Latency
<1 μs

Issue rate
> 100 Mmsg/s

BXI Link
100 (4x25) Gb/s

48 ports
BXI Link

9600 Gb/s
bandwidth

**Bull**
atos technologies

# BXI – Switch
## Overview

BXI standalone switch: 1U, 48 optical ports

- 48 ports, 192 SerDes @ 25Gb/s
  - Total throughput: 9600 Gb/s
- **Latency**: 130ns
- **Die**: 22 x 23mm
- **Package**: 57.5 x 57.5mm
- **Transistors**: 5.5 billions
- **TDP**: 160W (min 60W)
- **Techno**: TSMC 28nm HPM

BXI switch ASIC

Redundant fans

Optical modules

Redundant PSUs

# BXI – NIC
## Main features

- **Implements in hardware Portals 4 communication primitives**

- **OS and application bypass**

- **Collective Operations offload in HW**

- **End-to-End reliability**

- **Load balancing & QoS with Virtual Channels**

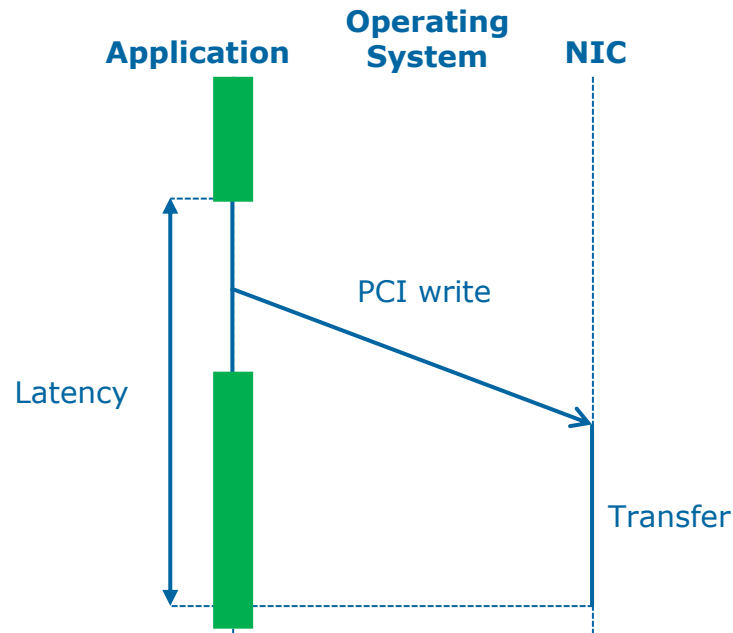- **Performance counters**

# OS Bypass: the first step to offload (1/2)
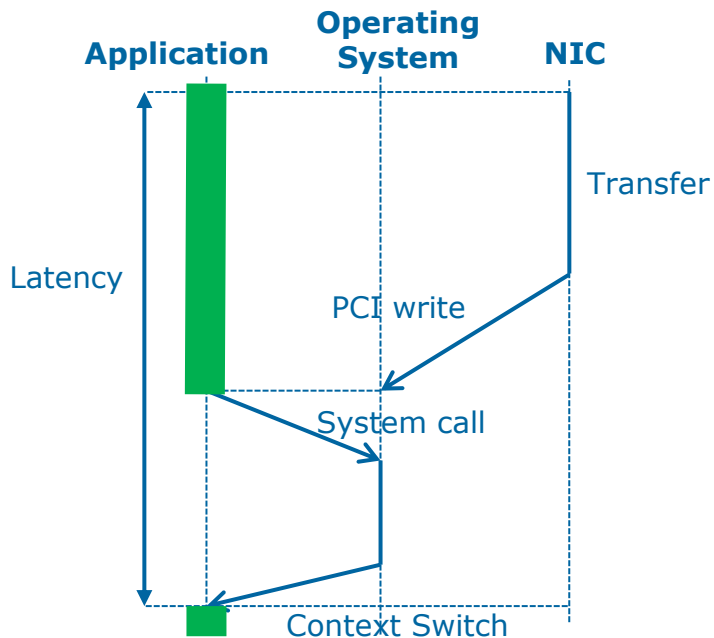## Source node

**Without OS Bypass**

- Application
- Operating System
- NIC

System call

PCI write

Latency

Context Switch

Transfer

**With OS Bypass**

- Application
- Operating System
- NIC

PCI write

Latency

Transfer

# OS Bypass: the first step to offload (2/2)
## Destination node



Without OS Bypass

Application  Operating System  NIC

Transfer

Latency

PCI write

System call

Context Switch

With OS Bypass

Application  Operating System  NIC

Transfer

"zero copy"

Latency

PCI write

# BXI - Offloading MPI communication



**address V2P**   **size**   **rank L2P**   **message order**   *No impact on computation during data transmission*

#include <mpi.h>

int **MPI_Isend**( const void *_buf_, int _count_, MPI_Datatype **datatype**, int **dest**,   int **tag**, MPI_Comm **comm**, MPI_Request *_request_)

int **MPI_IRecv**(void *_buf_, int **count**, MPI_Datatype **datatype**, int **source**, int **tag**, MPI_Comm **comm,** MPI_Request *_request_)

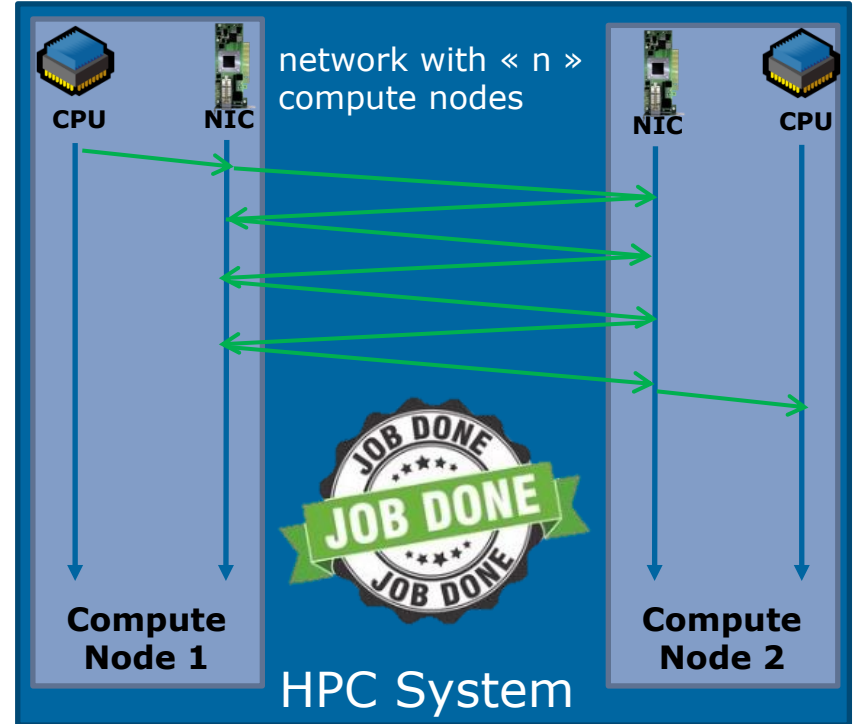int **MPI_Wait**(MPI_Request *_request_, MPI_Status *_status_)

# Offload Mechanism & Benefits



Conventional interconnect: no offload on the NIC

BXI ▶ Faster communications ▶ Less CPU solicitation

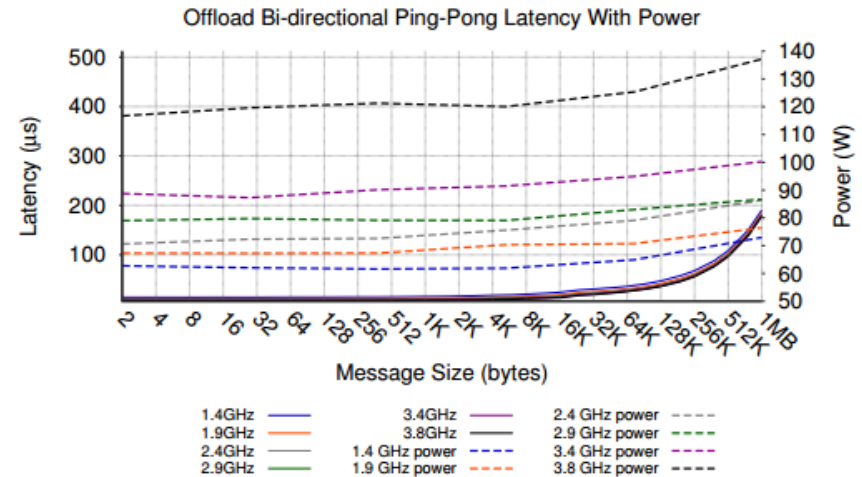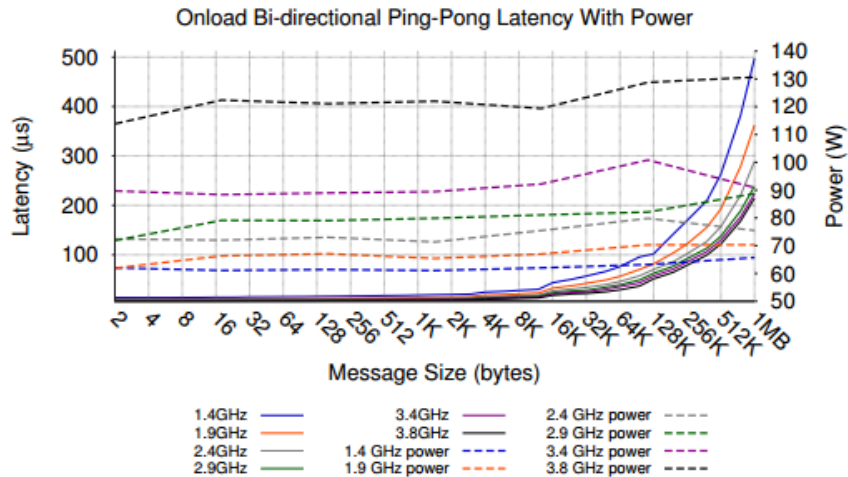network with « n » compute nodes

CPU   NIC   NIC   CPU

Compute Node 1   Compute Node 2   HPC System

network with « n » compute nodes

CPU   NIC   NIC   CPU

JOB DONE   JOB DONE

Compute Node 1   Compute Node 2   HPC System

# BXI - Offloading MPI communication



Task A    B    C    D          A    B    C    D          A    B    C    D

compute

AllReduce   AllReduce   AllReduce   AllReduce

MPI-1 Collective          Non-triggered operations          BXI "triggered operations"
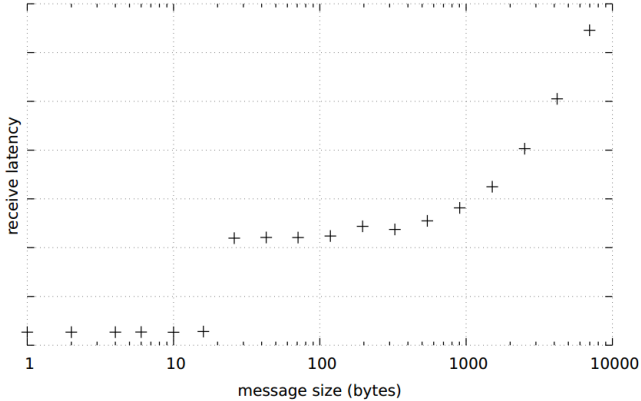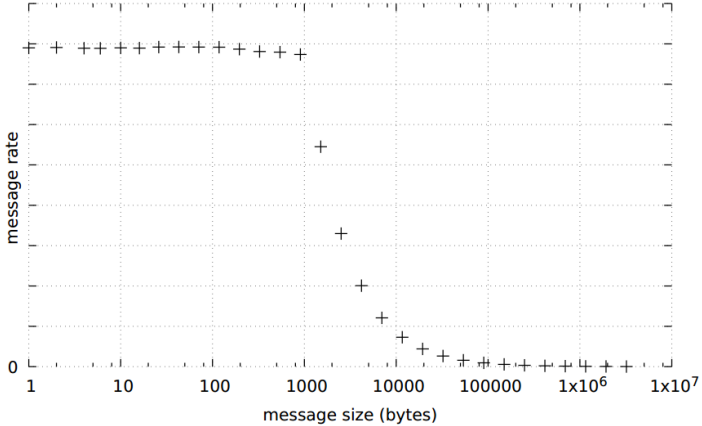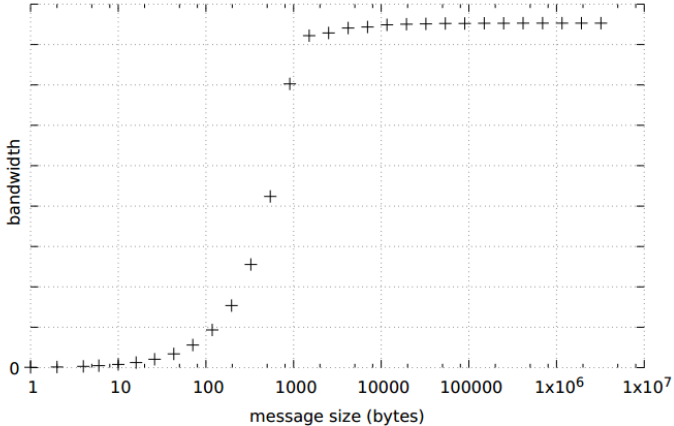
# Offload offers other benefits

▶ Offloaded networking approach:

– is **less frequency sensitive** than onloaded networking,

– provides **major power reductions**, particularly significant with large-scale systems.



Onload Bi-directional Ping-Pong Latency With Power



Offload Bi-directional Ping-Pong Latency With Power

▶ Paper co-authored by Matthew Dosanjh, Ryan Grant, and Ron Brightwell entitled "**Re-evaluating Network Onload vs. Offload for the Many-Core Era**".

# Performance

# Conclusion and Next Steps

▶ BXI system 1st prototype on Top/Green500,

▶ BXI full system (> 8k nodes) in 2017,

▶ More installations coming,

▶ Growing BXI ecosystem,

▶ Next generation of BXI in preparation.

# Thanks

For more information please contact:
T+ +33  (1) 30 80 74 94
M+ +33  (6) 86 49 33 21
fabien.locussol@atos.net

**Bull**
atos technologies

# Backup

**Bull**
atos technologies

# BXI Software Suite
## Overview

# BXI Fabric Management

| Supervision | Performance Analysis | Topology Manager | Routing | CLI |
|---|---|---|---|---|

**BXI Fabric Tools**

**BXI Switch**

► BXI switches are managed through a **distributed and out-of-band fabric management suite** allowing to scale up to 64K nodes.

► **Out-of-band management** eliminates any interference of the management traffic with the applications traffic.
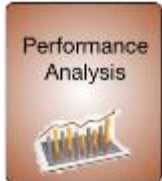
# BXI fabric management
## Fabric Monitoring/Profiling

► Administrators, support and users need information about fabric status

► Each class of users needs different kind of information with different purposes
  - **Administrator**: what happened? What is the status?
  - **Support**: is it working? Why not? Low level debug
  - **User**: how does my application use the interconnect? Reproducibility

► BXI provides counters and sampling:
  - Probes (set of counters + frequency) can be configured
  - 4 probes per switch maximum at the same time
  - Frequency up to **1Hz**
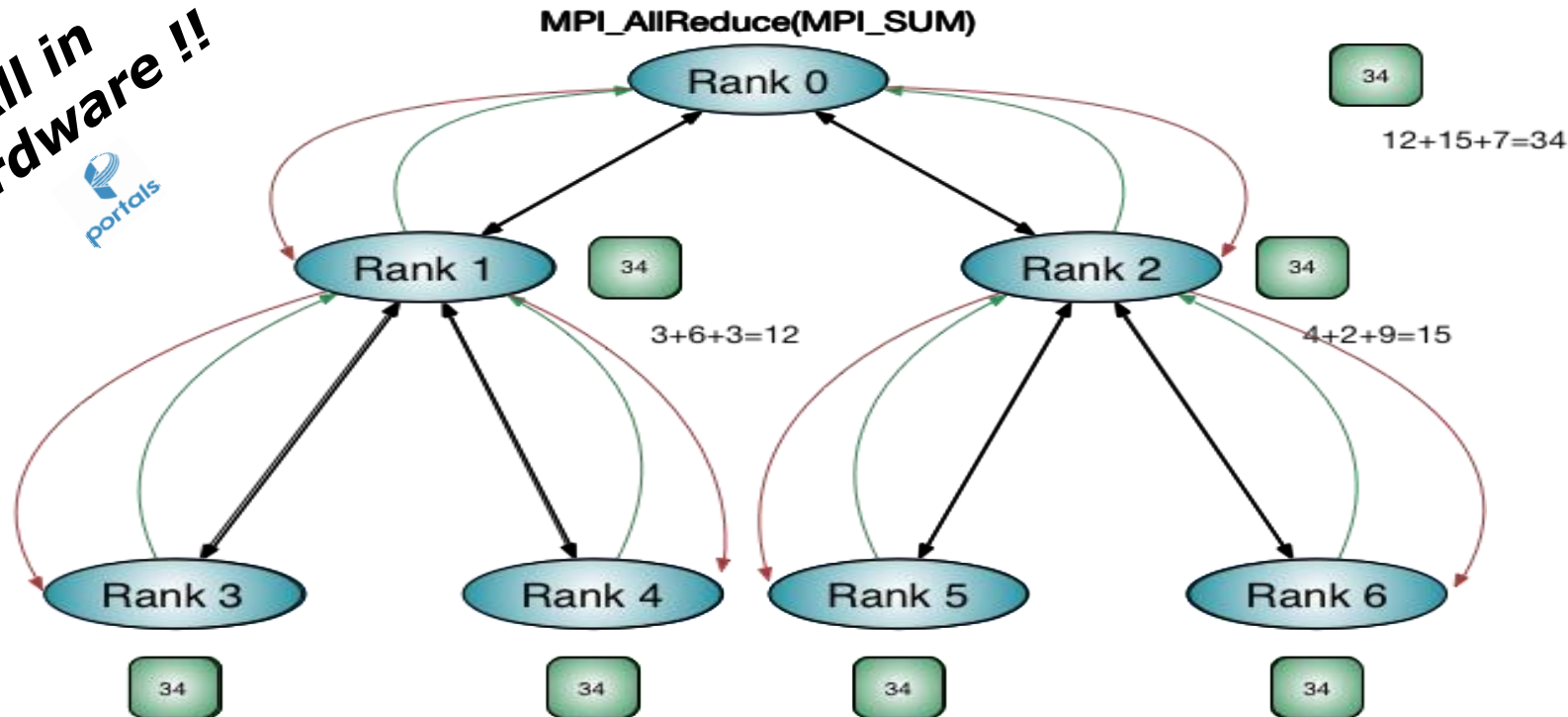  - Threshold setting for monitoring warning

► Dynamic routing solutions ensure interconnect's reliability and performances:
  - Topology specific and topology agnostic routing
  - Distributed routing between management nodes and switch embedded solutions

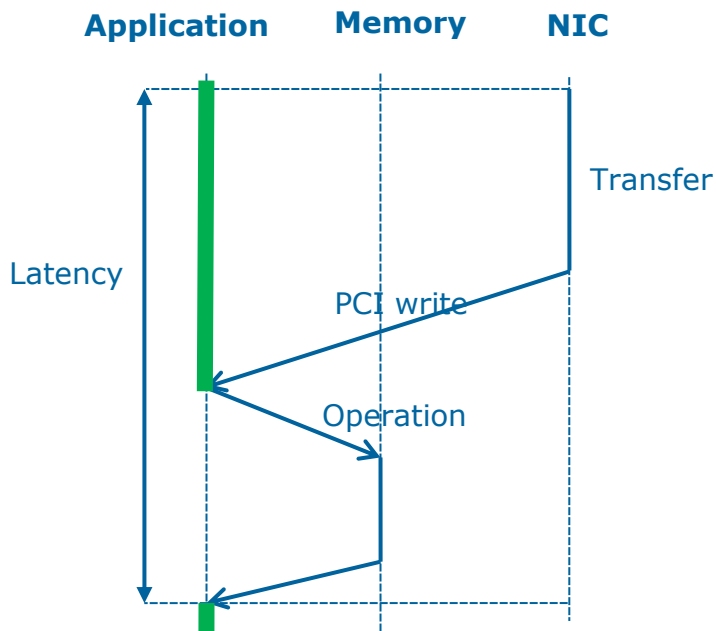# BXI - Offloading MPI communication
## Ex. MPI_AllReduce principle



**MPI_AllReduce(MPI_SUM)**

All in Hardware !!

Rank 0

34

12+15+7=34

Rank 1    34    Rank 2    34

3+6+3=12    4+2+9=15

Rank 3    Rank 4    Rank 5    Rank 6

34    34    34    34

# Application Bypass

## Without Application Bypass

**Application**    **Memory**    **NIC**

Transfer

Latency

PCI write

Operation

## With Application Bypass (BXI)

**Application**    **Memory**    **NIC**

Transfer
"zero copy"

Latency

PCI write

Atomic operations