# Bull Exascale Interconnect (BXI)
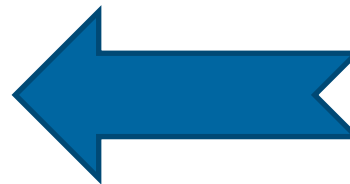
# un nouveau réseau
## pour le calcul de haute performance

Jean-Pierre Panziera

11-10-2016

# HPC applications

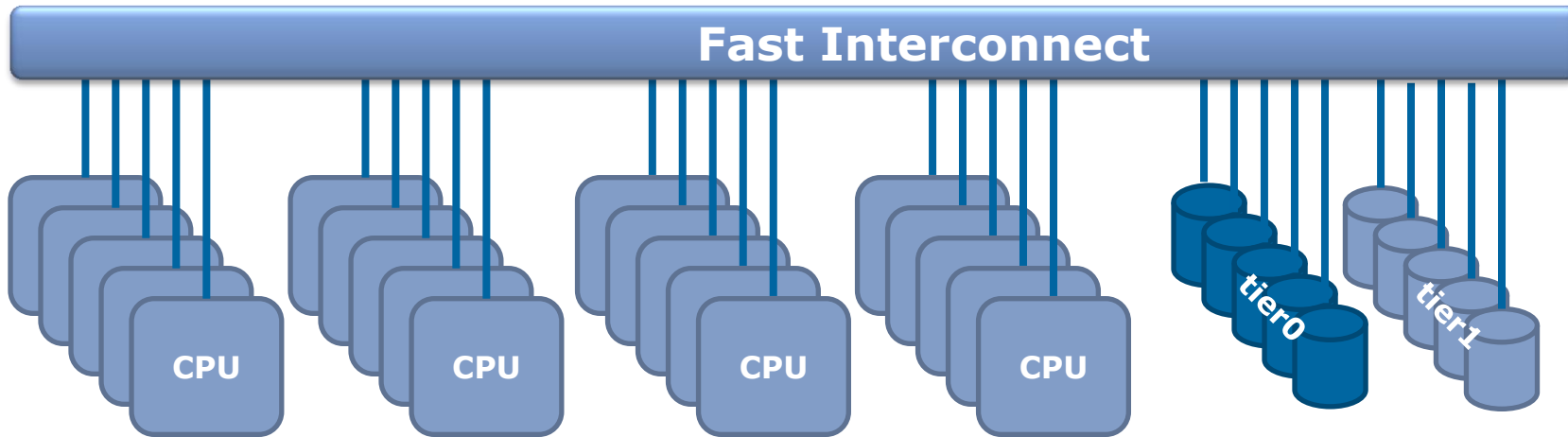HPC applications are characterized by:

▶ X-large computing needs (TeraFlops, PetaFlops, ExaFlops…)

▶ X-large datasets

▶ large number of processors

▶ tight coupling between computing threads

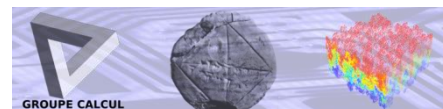▶ Many short MPI messages (latency)

▶ Large IO transfers (bandwidth)

Efficient HPC Interconnect

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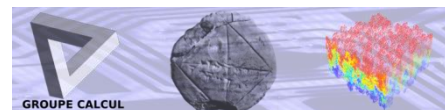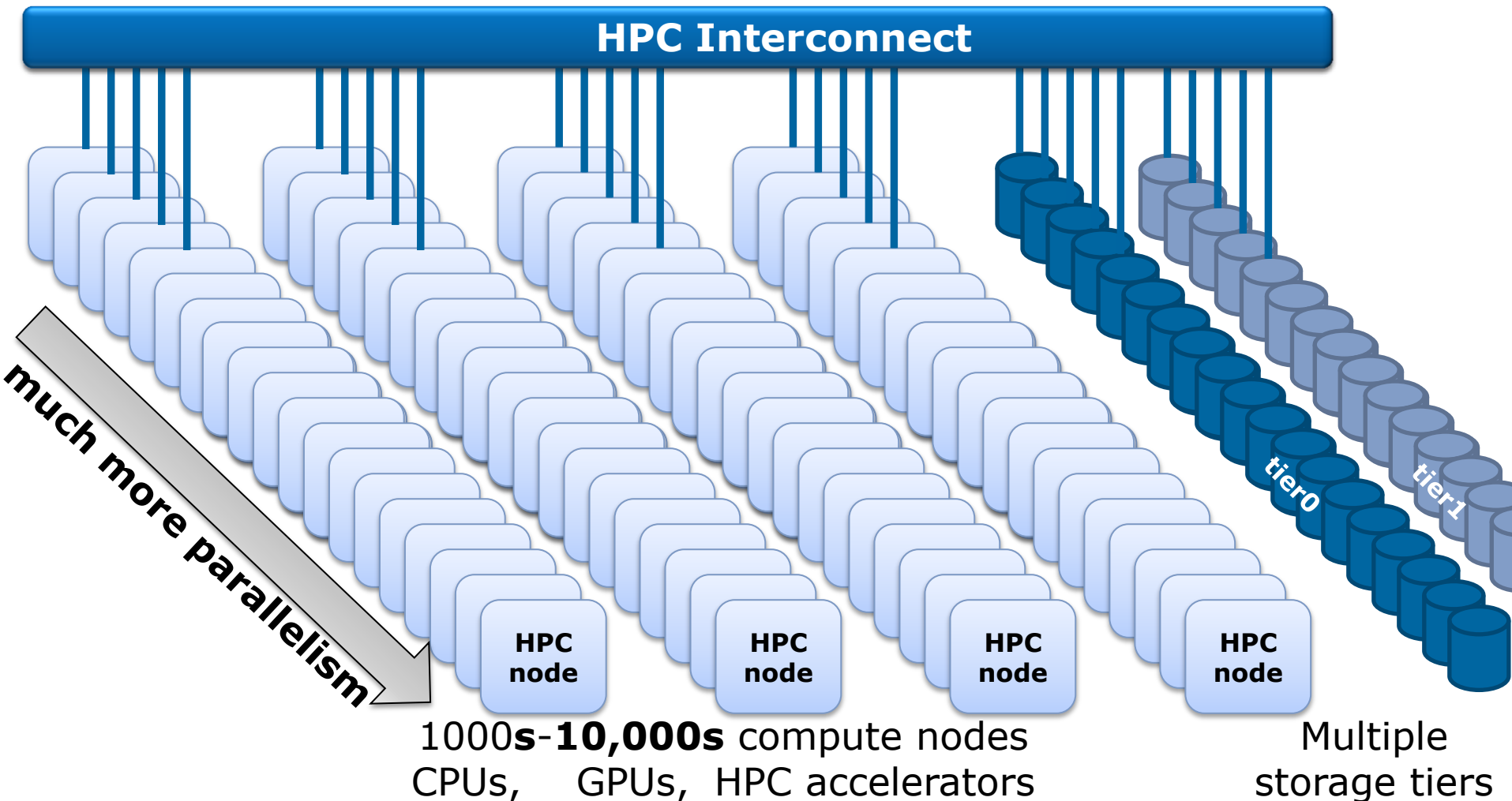# Petaflop class featuring 1000s CPU nodes

**BXI**

**Fast Interconnect**

CPU  CPU  CPU  CPU  tier0  tier1

100-10,000 compute nodes
using CPUs, typically x86

Multiple
storage tiers

GROUPE CALCUL

**Bull**
atos technologies

# expecting 10-100 Pflops systems in 2016-17 … 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

GROUPE CALCUL

**Bull**
atos technologies

# HPC Interconnect

▶ **performant**
  – low latency
  – high message rate
  – high bandwidth

▶ **scalable**
  – 10,000**s** nodes

▶ **reliable**
  – fault tolerant
  – redundant

▶ **efficient**
  – handle simultaneously different flow types – small & big messages - MPI & IO
  – Adaptive routing
  – small memory footprint
  – link-level checking & retry, ECC protection

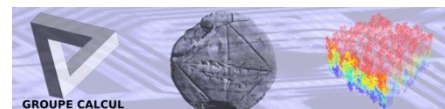▶ Offload communications in **Hardware**
  – HPC cores are many but slow(er)
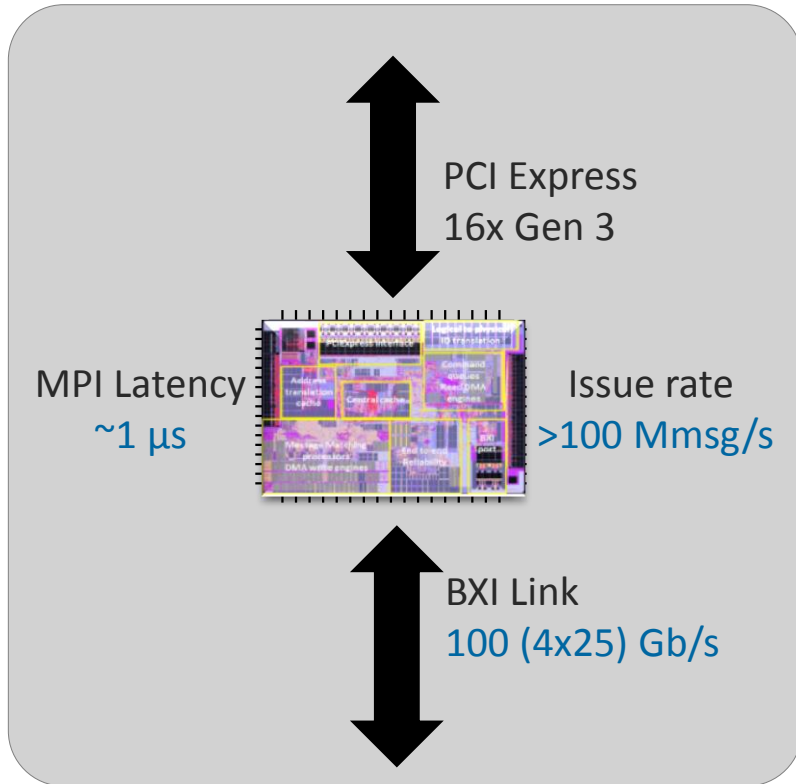
# BXI overview
*High Performance Interconnect for HPC*

▶ **BXI: High Performance Interconnect for HPC**
- Lowest latency, Highest message rate at scale, Highest Bandwidth

▶ **BXI full acceleration in hardware for HPC applications**
- based on Portals 4 (Sandia), BXI provides full HW acceleration for:
  - **MPI** and PGAS communications (send/recv, RDMA)
  - High performance collective operations

▶ **BXI highly scalable, efficient and reliable**
- Exascale scalability → 64k nodes (v1)
- Adaptive Routing, Quality of Service (QoS)
- End-to-end error checking + link level CRC + ECC in ASICs

▶ **BXI co-designed with CEA**
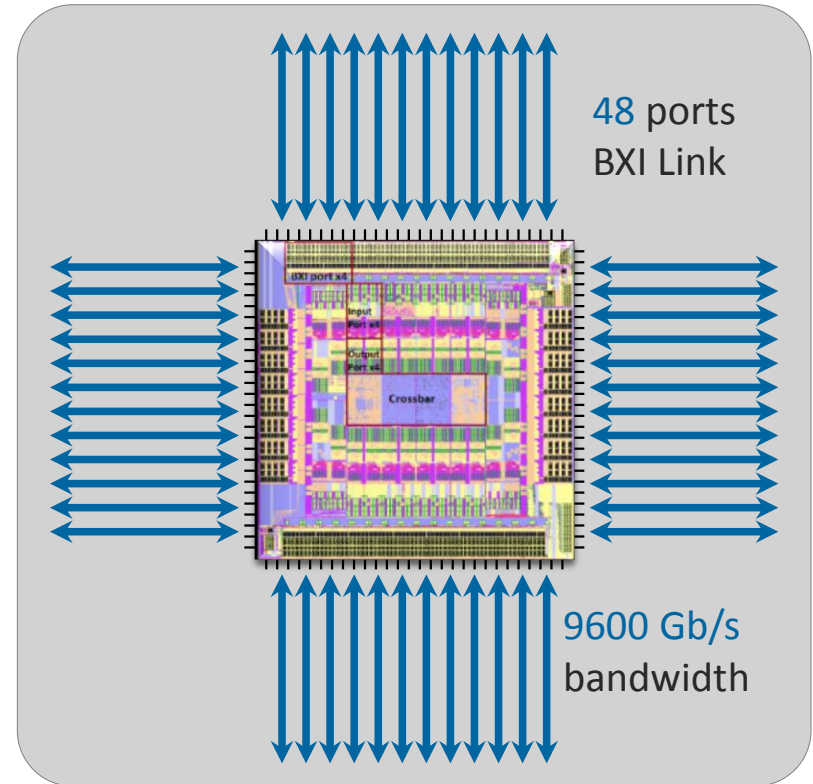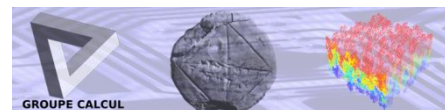
# BXI Network is based on 2 ASICs

NIC ASIC

switch ASIC

PCI Express
16x Gen 3

MPI Latency
~1 μs

Issue rate
>100 Mmsg/s

BXI Link
100 (4x25) Gb/s

*Lutetia*

48 ports
BXI Link

9600 Gb/s
bandwidth

*Divio*

# NIC main features 1/2

▶ **Implements in hardware the Portals 4 communication primitive**
- – Overlapping communications and computations by offloading to NIC
- – MPI two-sided messaging:
  - • HW acceleration of list management and matching on the NIC
- – PGAS / MPI one-sided messaging:
  - • use fast path inside the NIC

▶ **OS and application bypass**
- – Applications issue commands directly to the NIC, avoiding kernel calls
- – Reception controlled by NIC without OS involvement
- – Reply to a put or a get does not require activity on application side.
  - • Logical to physical ID translation
  - • Virtual to physical memory address translation.
  - • Rendez-vous protocol in HW

# NIC main features 2/2

▶ **Collective Operations offload in HW**

  – using Atomic and Triggered operations units

▶ **End-to-End reliability** recovery mechanism for transient and permanent failures

  – message integrity, 32bits CRC are added to each message (or each message chunk for large transfers).

  – message ordering required for MPI messages is checked with a 16 bit sequence number.

  – message delivery a go-back-N protocol is used to retransmit lost or corrupted messages.

▶ **Allocates Virtual Channels**: Separating different type of messages to avoid deadlocks and to optimize network resources usage (load balancing and QoS)

▶ Offers performance and error **counters** for Applications performance analysis

# BXI
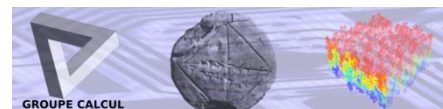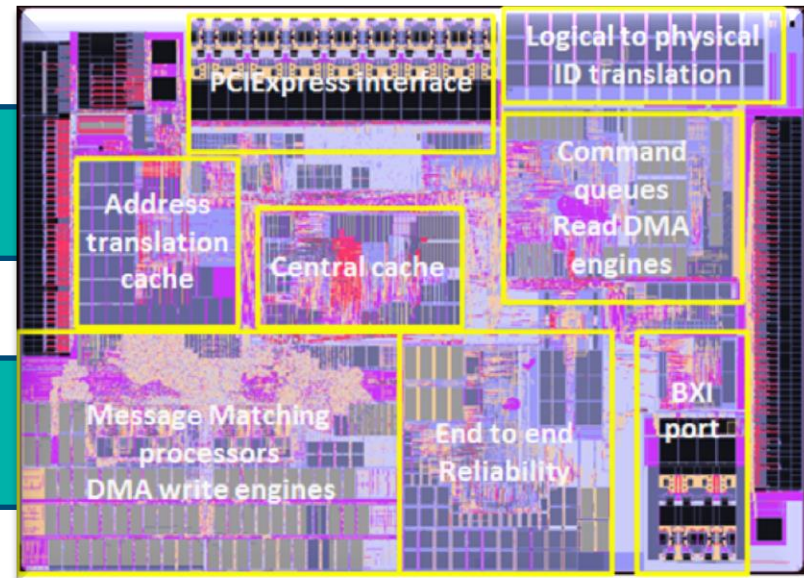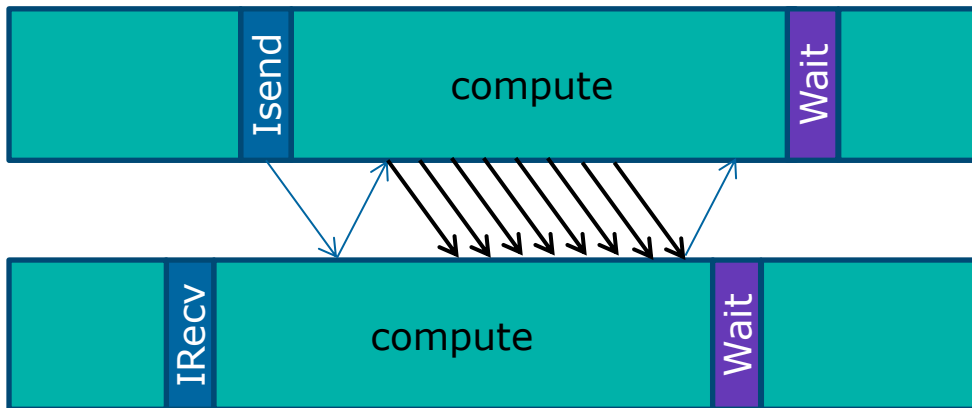# offloading MPI communication in HW



**address V2P**    **size**    **rank L2P**    **message order**
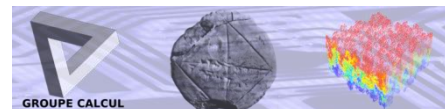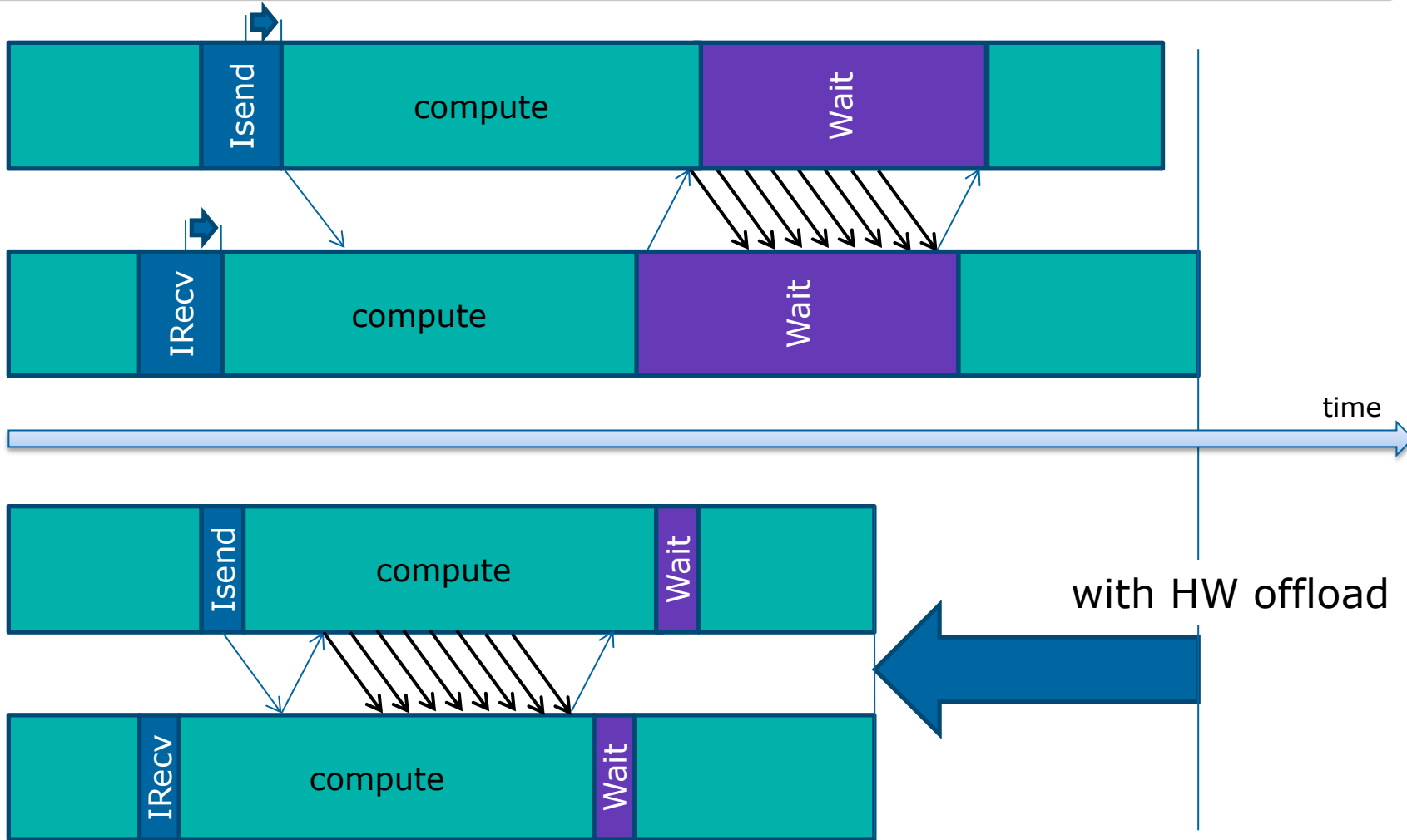
#include <mpi.h>

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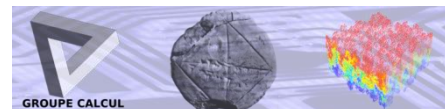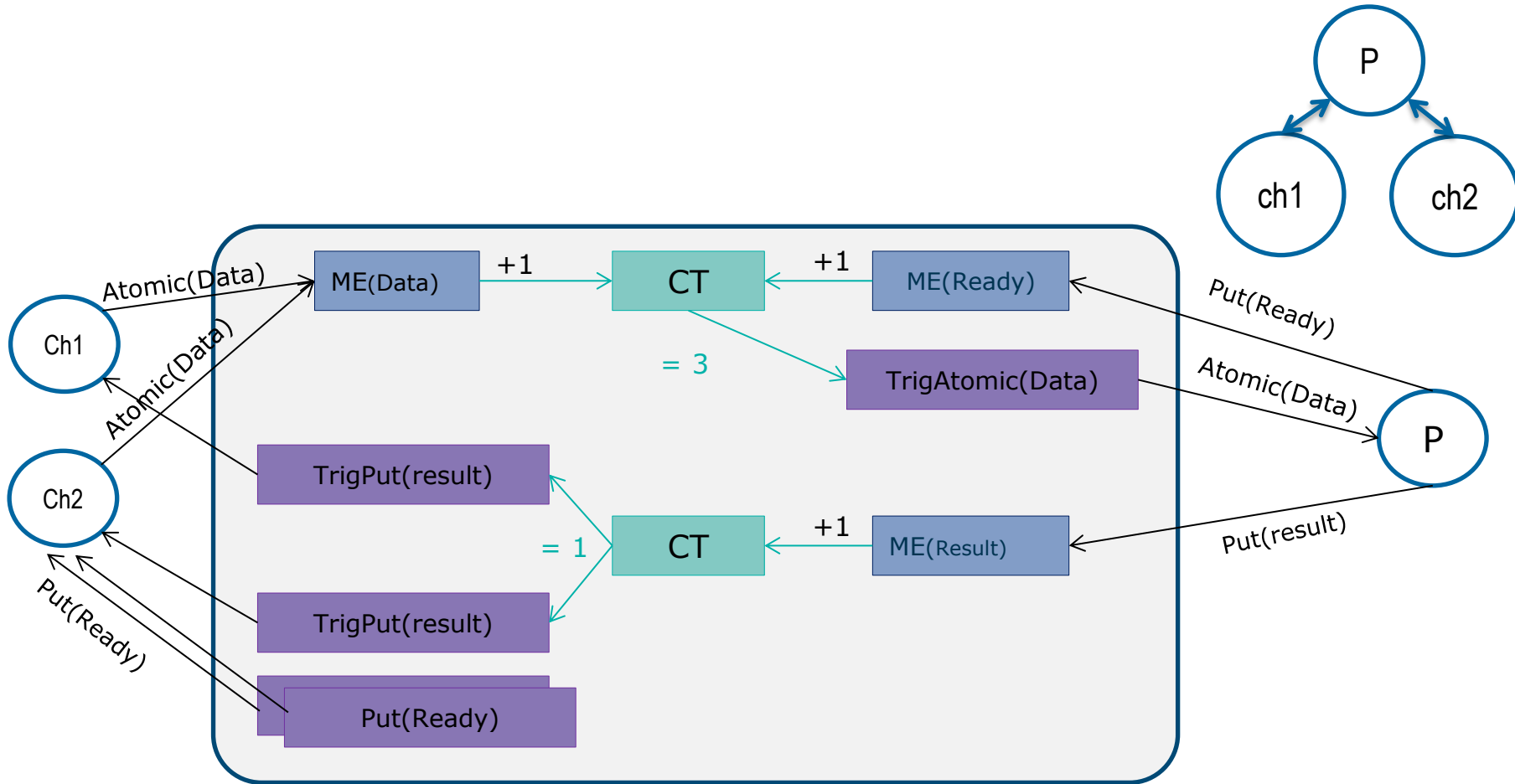
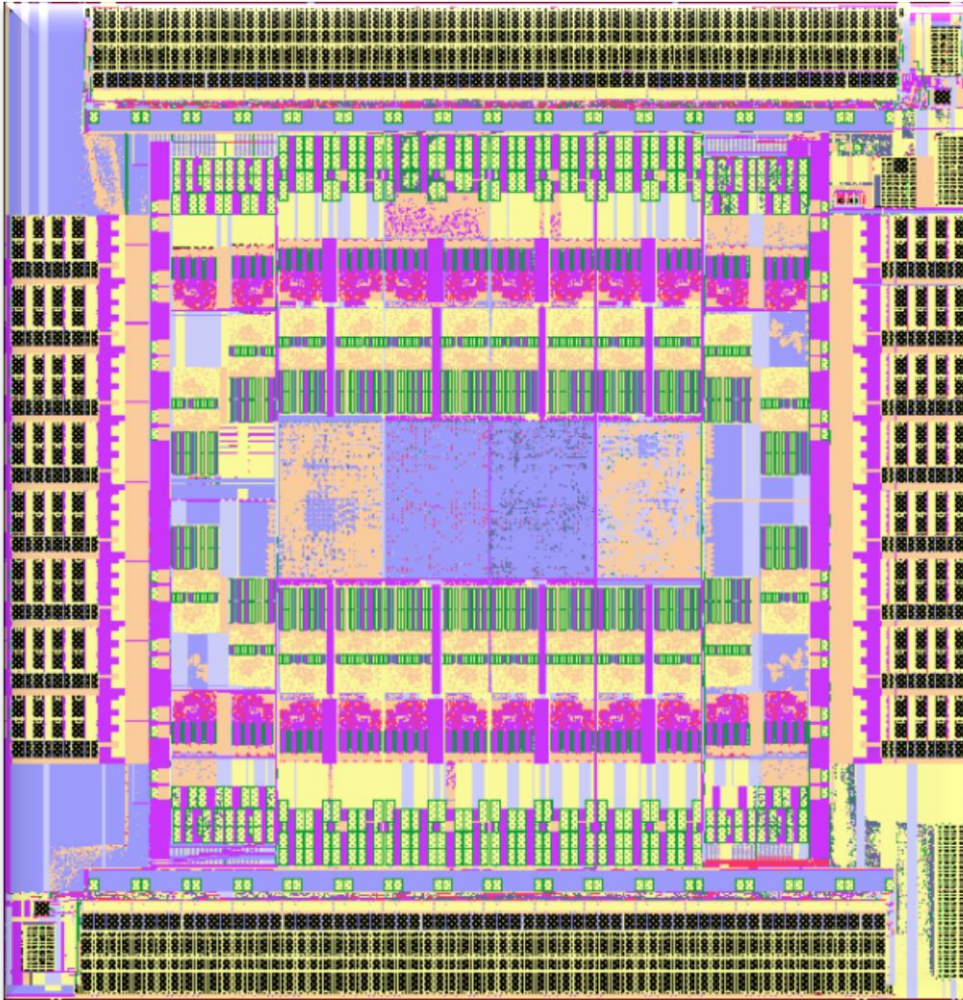int **MPI_Wait**(MPI_Request *_request_, MPI_Status *_status_)

# BXI
# offloading MPI communication in HW



time

with HW offload

# BXI Switch overview



- ▶ 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
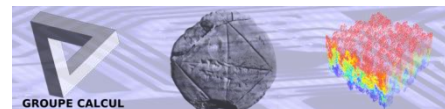  - – Min power : 60W
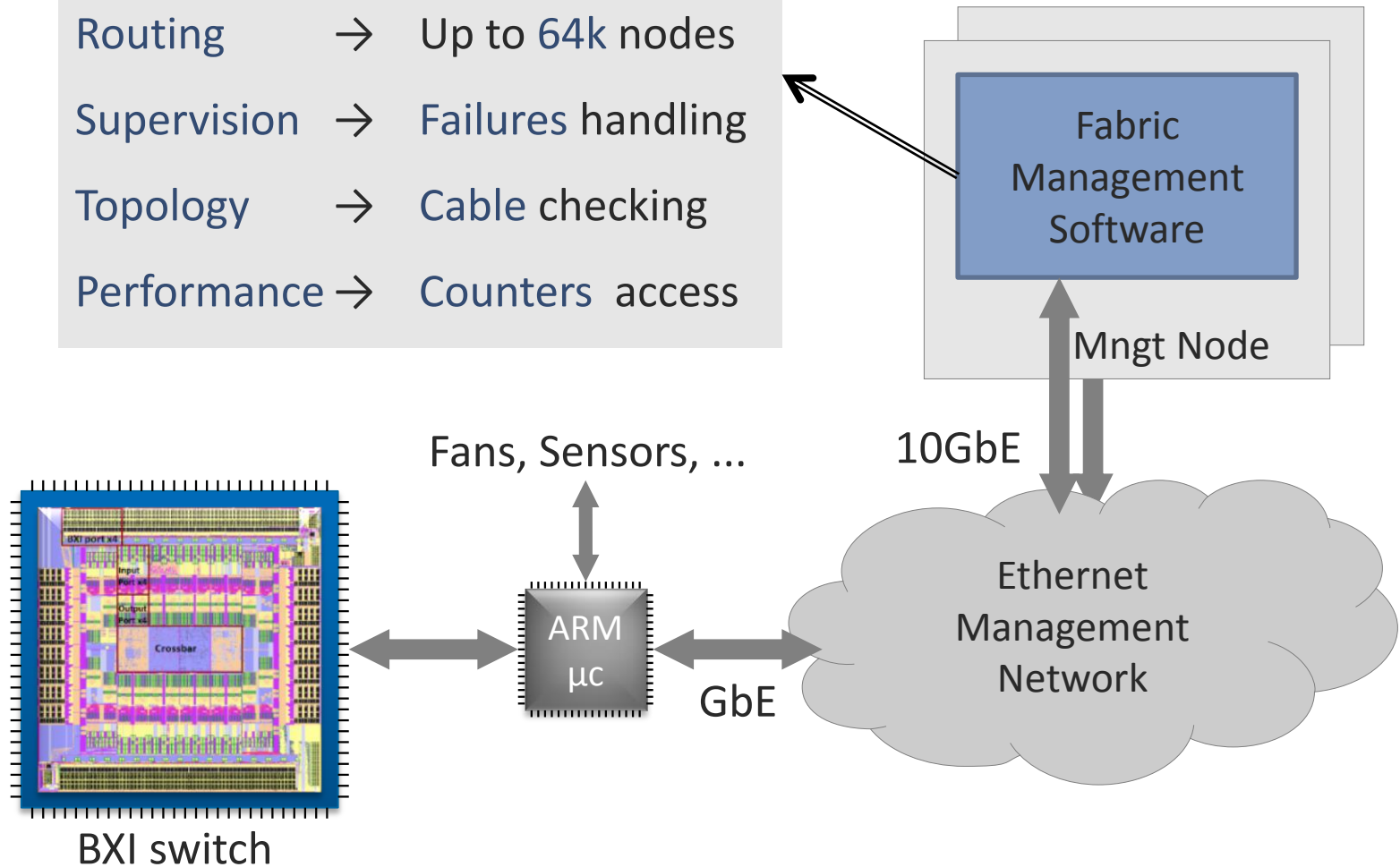- ▶ Techno : TSMC 28nm HPM

# BXI fabric features

► Scalable up to **64K NICs**

► **100 Gb/s** links (4 lanes x25,278125 GT/s)

► **Reliable** and ordered network (end to end  +  Link level)

► **Flexible** with full routing table
  – Many topologies supported (**Fat-Tree**, Torus, Hypercube, **All-to-All**…)
  – Ease routing algorithm optimization

► **Adaptive routing**

► Extensive buffering implementing 16 virtual channels preventing deadlock and efficiently balancing traffic

► Quality of Service (**QoS**) with weighted round robin arbitration
  – highly configurable load balancing
  – Segregation of flows per destination
  – ensuring progress of short messages vs long messages

► High resolution time synchronization
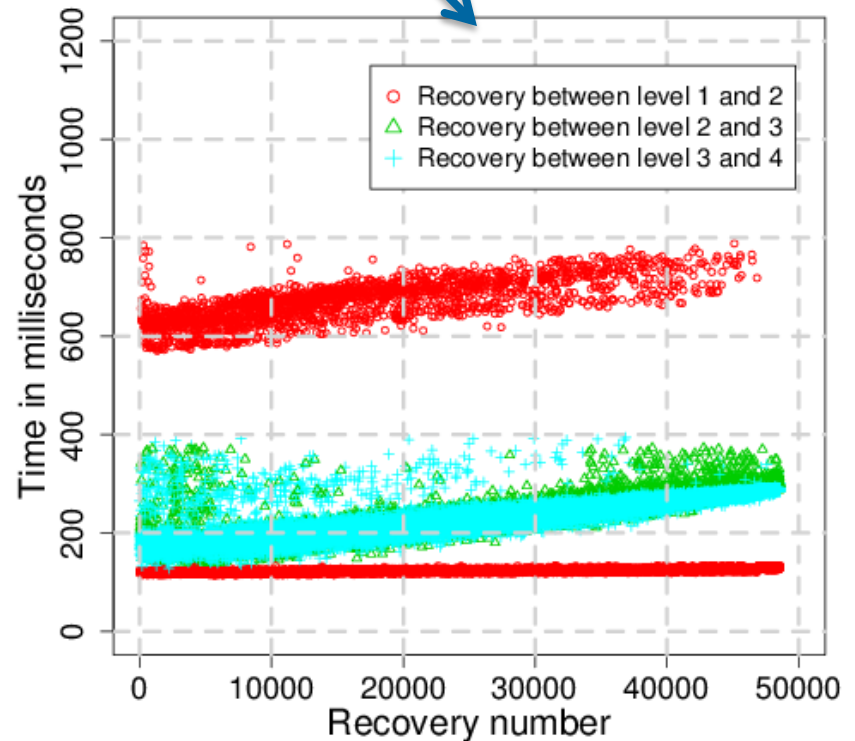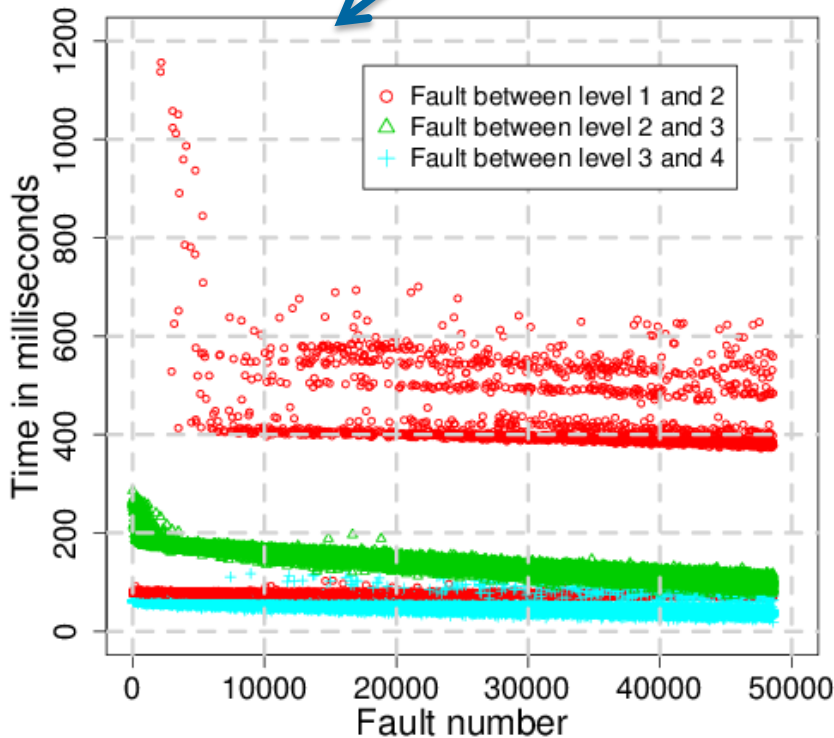
► **Out-of-band management**

# Fabric Management Software

Routing       →      Up to 64k nodes

Supervision    →      Failures handling

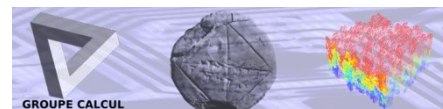Topology      →      Cable checking

Performance →      Counters access

Fabric Management Software

Mngt Node

Fans, Sensors, …

10GbE

ARM µc

GbE

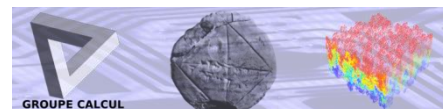Ethernet Management Network

BXI switch

# BXI Routing Online Mode Processing Time e.g. 64k nodes
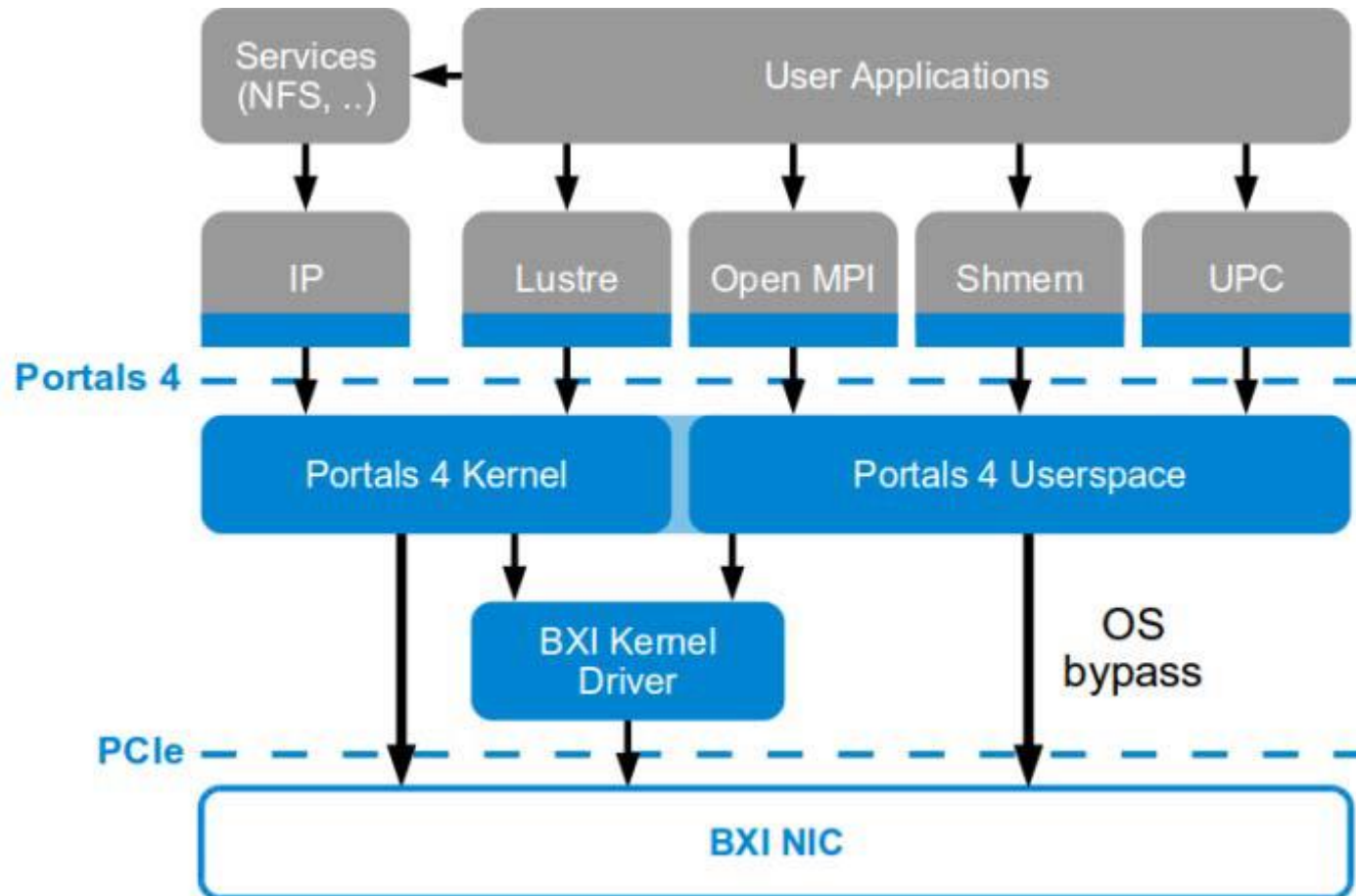
**BXI**

routing table updates computed in < 1s

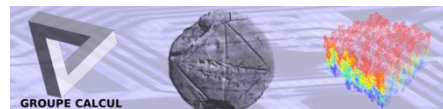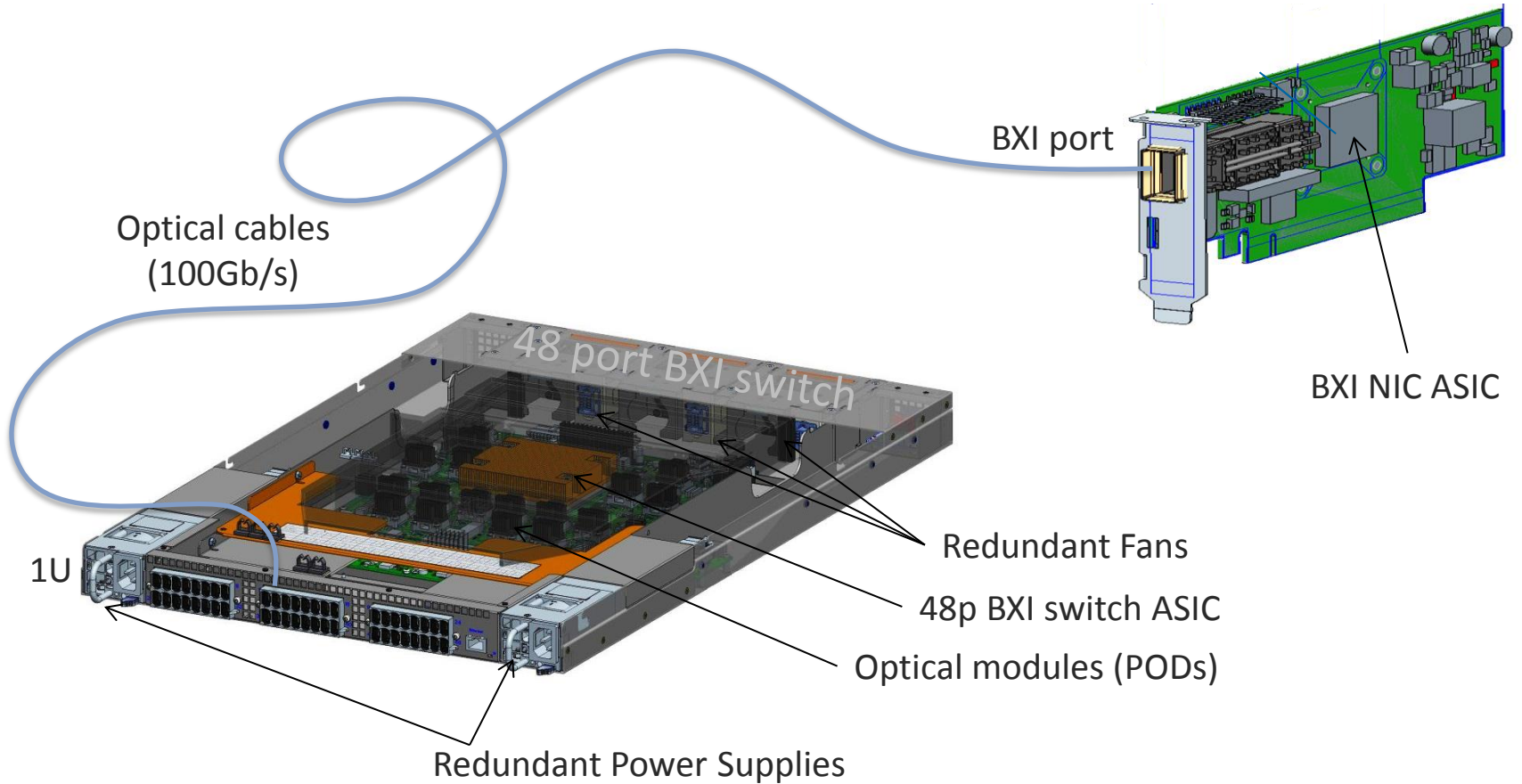on link failure

on link recovery



Quintin, Vignéras; Fault-Tolerant Routing for Exascale Supercomputer: The BXI Routing Architecture. HiPINEB'15
Quintin, Vignéras; Transitively Deadlock-Free Routing Algorithms . HiPINEB'16
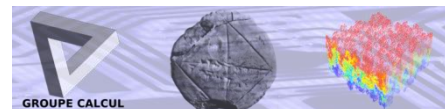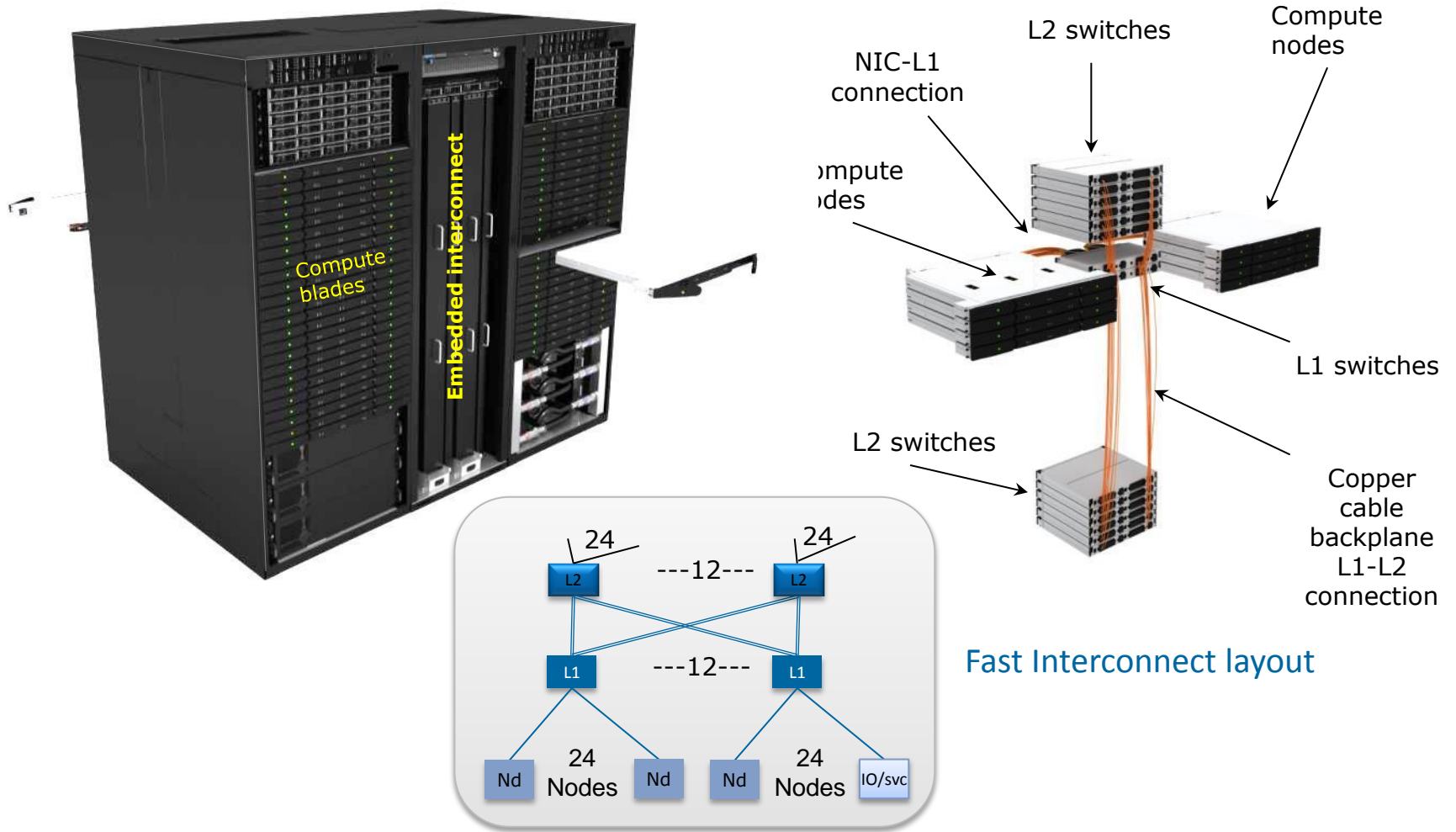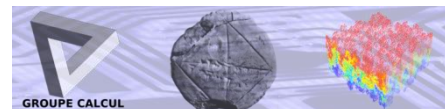
# BXI Software compute stack

# BXI
## PCI adapter card and 48p standalone switch



BXI port

Optical cables (100Gb/s)

48 port BXI switch

BXI NIC ASIC

1U

Redundant Fans

48p BXI switch ASIC

Optical modules (PODs)

Redundant Power Supplies

# "Sequana" – Embedded interconnect

**BXI**

Compute blades

**Embedded interconnect**

L2 switches

Compute nodes

NIC-L1 connection

Compute nodes

L1 switches

L2 switches

Copper cable backplane L1-L2 connection

### Fast Interconnect layout

24

24

---12---

L2            L2

L1   ---12---   L1

24 Nodes

24 Nodes

Nd        Nd        Nd        IO/svc

GROUPE CALCUL

**Bull** atos technologies

# Sequana cells interconnection

L3 switches

... 48x ...

**Direct connections**

**Fat-Tree**

# BXI wrap up

- ▶ BXI is Atos new High Performance Interconnect for HPC

- ▶ BXI offloads communication primitives into the NIC

- ▶ BXI boosts MPI communications in Hardware

- ▶ Highly scalable, up-to 64k nodes

- ▶ Fist BXI system installed in Q4-2016

- ▶ Large BXI deployment (8+K nodes system) in 2017