

DEEP and DEEP-ER

Estela Suarez
Jülich Supercomputing Centre

10.05.2016

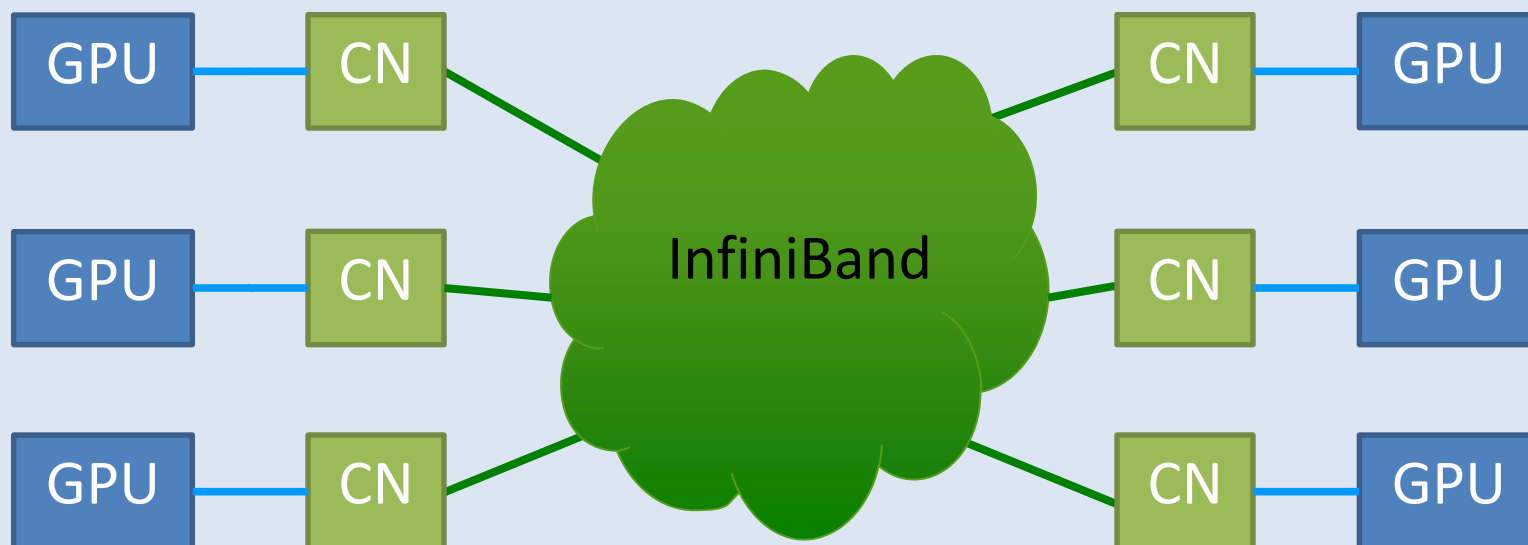
DEEP

- **Cluster-Booster archit.**
- Software stack
- Programming environ.
- Energy efficiency
- Applications:
 - Co-design
 - Evaluation/demonstration
 - Code modernisation

DEEP-ER

- Extend memory hierarchy
- High-performance **I/O**
- Scalable **resiliency**
- Applications:
 - Co-design
 - Evaluation/demonstration
 - Code modernisation

CLUSTER-BOOSTER ARCHITECTURE

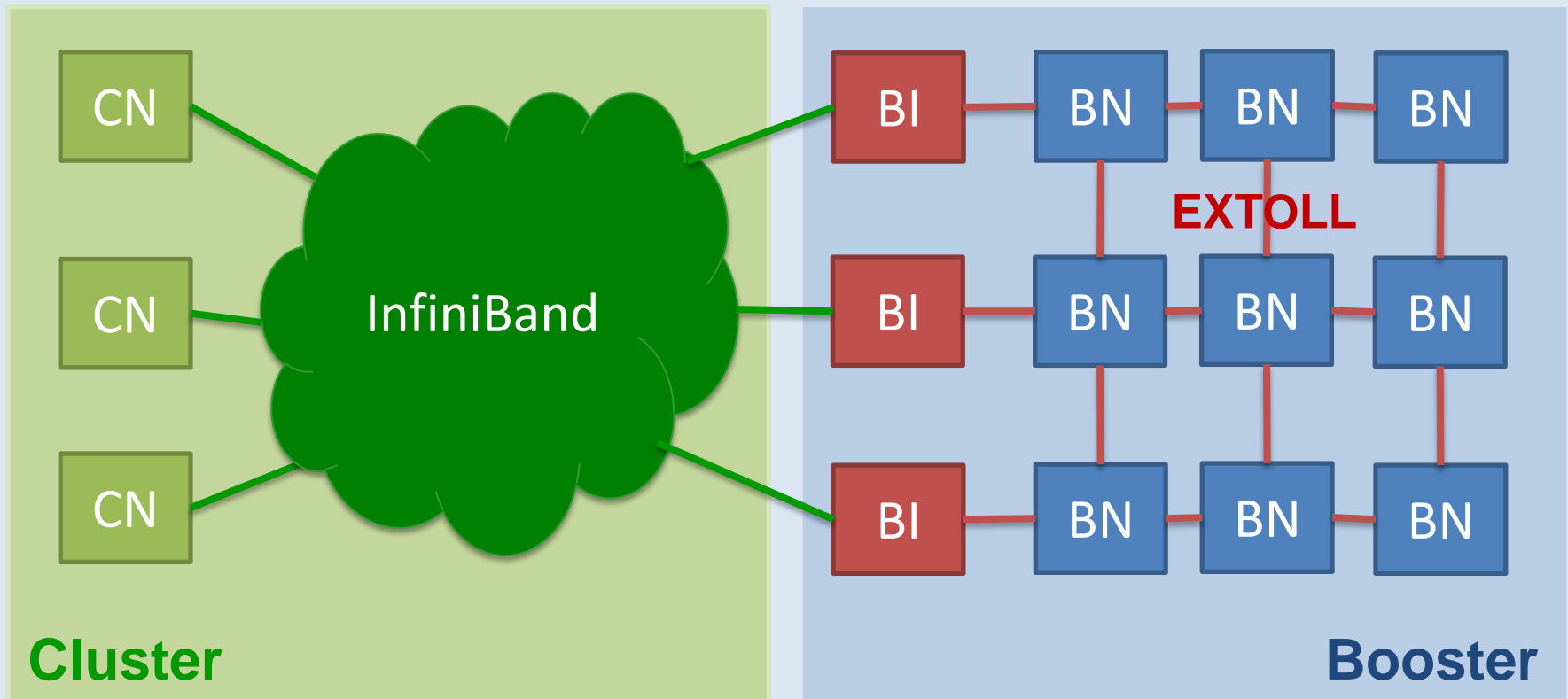


Flat topology

Simple management of
resources

Static assignment of
accelerators to CPUs

Accelerators cannot act
autonomously



Flexible assignment of resources (CPUs, accelerators)
 Direct communication between accelerators
 “Offload” of large and complex parts of applications

Intel® Xeon®

Cluster

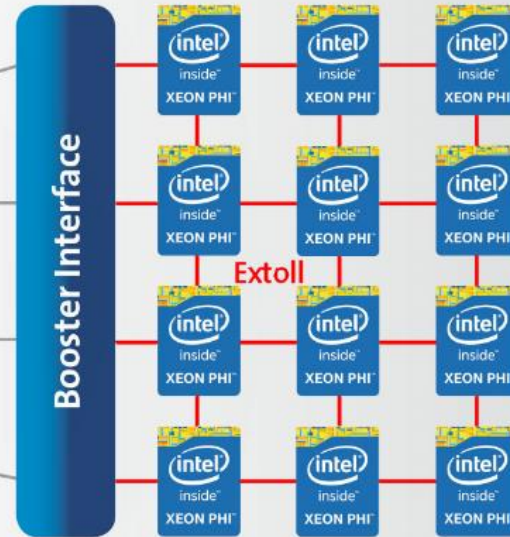


InfiniBand®

LOW/MEDIUM
SCALABLE CODE

Intel® Xeon Phi™

Booster



HIGHLY
SCALABLE CODE

- Installed at JSC
- 1,5 racks
- 500 TFlop/s peak perf.
- 3.5 GFlop/s/W
- Water cooled

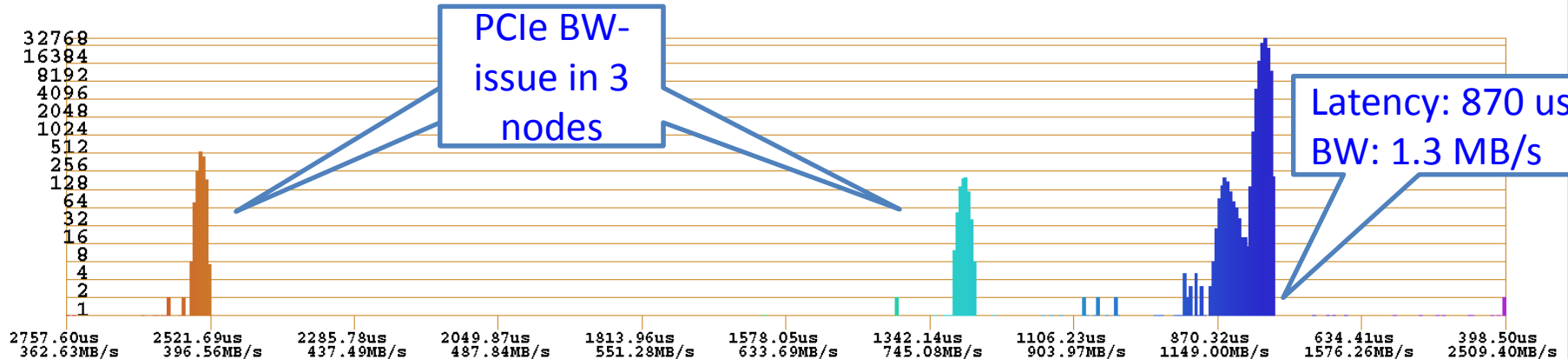
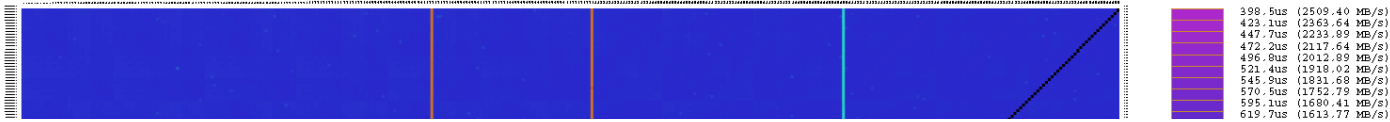


**Cluster
(128 Xeon)**

**Booster
(384 Xeon Phi
KNC)**

Booster measurements

MPI Linktest: ping-pong



length_of_message:	1048576 bytes (1024.00 KBytes)	number_of_tasks:	384
number_of_messages:	25	Execution order:	Serial
Alltoall:	1	Mixing PE rank:	No
Min Value:	398.5us (2509.40 MB/s)	Alltoall Min Value:	616.1us (1 Byte)
Max Value:	2757.6us (362.63 MB/s)	Alltoall Max Value:	5038.0us (1 Byte)
Avg Value:	814.9us (1227.12 MB/s)	Alltoall Avg Value:	795.1us (1 Byte)

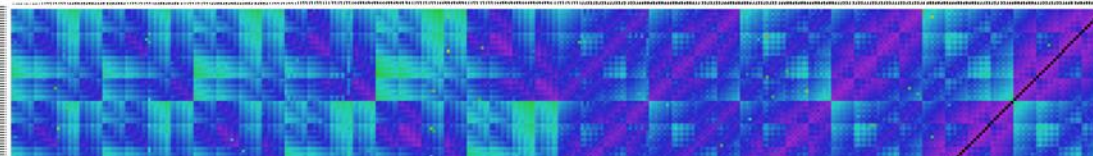
Report generated by FZJ Linktest Result Analyzer, Forschungszentrum Juelich GmbH



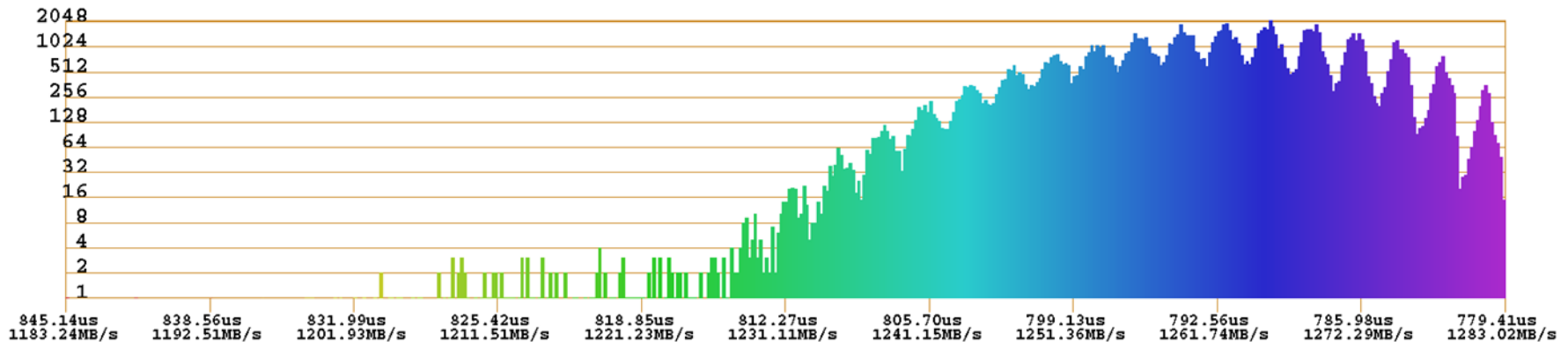
Booster Nodes (KNC) from 1 to 384

Booster measurements

MPI Linktest: ping-pong



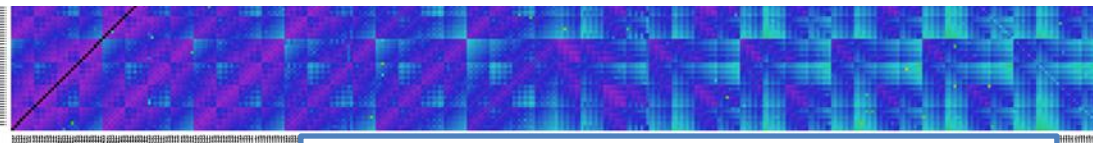
779.4us	(1283.02 MB/s)
780.1us	(1281.89 MB/s)
780.8us	(1280.77 MB/s)
781.5us	(1279.65 MB/s)
782.2us	(1278.53 MB/s)
782.8us	(1277.41 MB/s)
783.5us	(1276.29 MB/s)
784.2us	(1275.18 MB/s)
784.9us	(1274.07 MB/s)
785.6us	(1272.96 MB/s)
786.3us	(1271.85 MB/s)
786.9us	(1270.74 MB/s)
787.6us	(1269.64 MB/s)



length_of_message: 1048576 bytes (1024.00 KBytes) number_of_tasks: 384
 number_of_messages: 40 Execution order: Serial
 Alltoall: 0 Mixing PE rank: No
 Min Value: 779.4us (1283.02 MB/s)
 Max Value: 845.1us (1183.24 MB/s)
 Avg Value: 792.7us (1261.47 MB/s)

Stddev < 10%

Report generated by FZJ Linktest Result Analyzer, Forschungszentrum Juelich GmbH

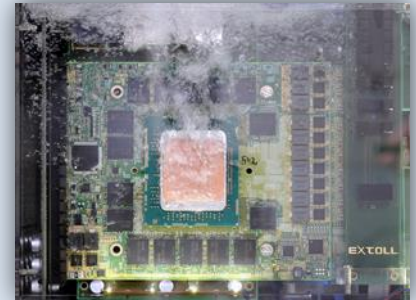


837.0us	(1127.00 MB/s)
838.3us	(1122.90 MB/s)
839.0us	(1121.93 MB/s)
839.7us	(1120.96 MB/s)
840.3us	(1119.99 MB/s)
841.0us	(1118.02 MB/s)
841.7us	(1116.05 MB/s)
842.4us	(1114.08 MB/s)
843.1us	(1112.12 MB/s)
843.8us	(1110.16 MB/s)
844.5us	(1108.20 MB/s)

Booster Nodes (KNC) from 1 to 384

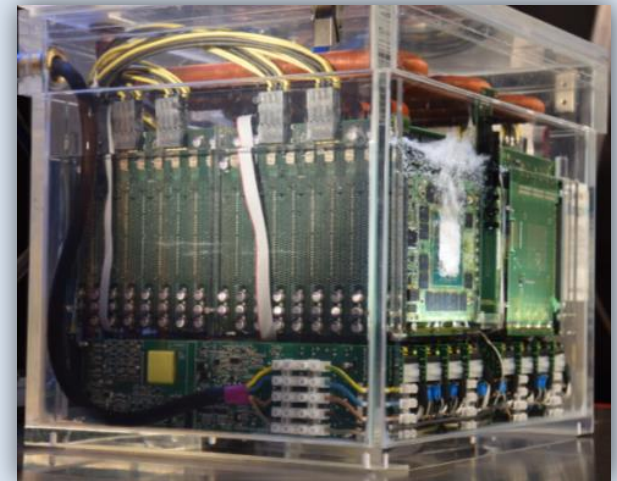
Alternative Booster implementation

- Interconnect EXTOLL ASIC “Tourmalet”
- 32 KNC-node system
- Implement $4 \times 4 \times 2$ topology, with Z dimension open



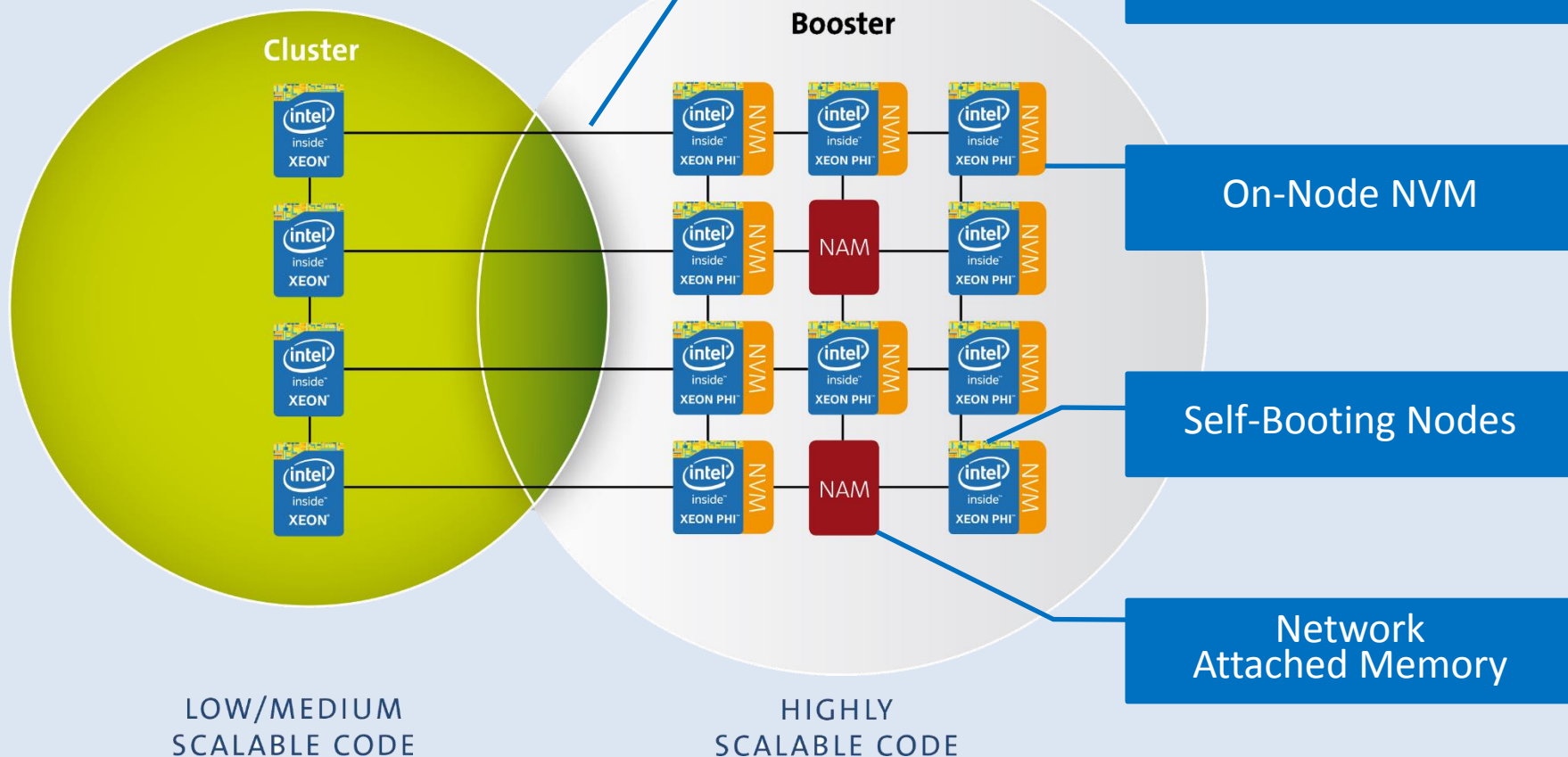
Experiment 2-phase immersion cooling

- NOVEC liquid from 3M
- Evaporates at about 50 degrees
- Condensates again in a water cooling pipe
- Allows very high-density integration



GreenICE Booster

Xeon



DEEP-ER Aurora Blade prototype



Eurotech's Aurora technology

Direct water cooled, high density

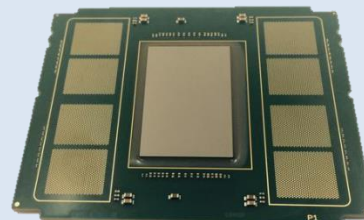
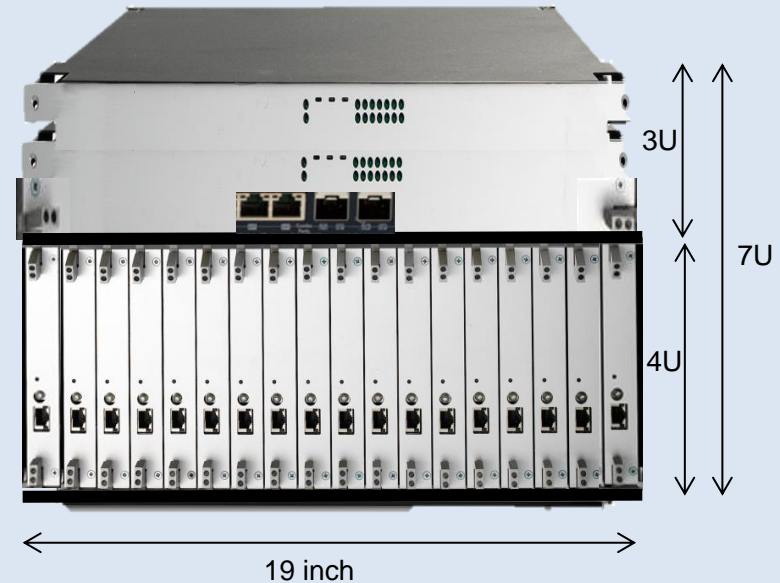


Aurora Blade DEEP-ER Booster
(in construction)

Aurora Blade Chassis

Rootcard
-18x EXTOLL
-18x NVMe

Chassis:
-18x KNL
-94GB Mem
-1x backplane



**Intel Xeon Phi
(KNL)**

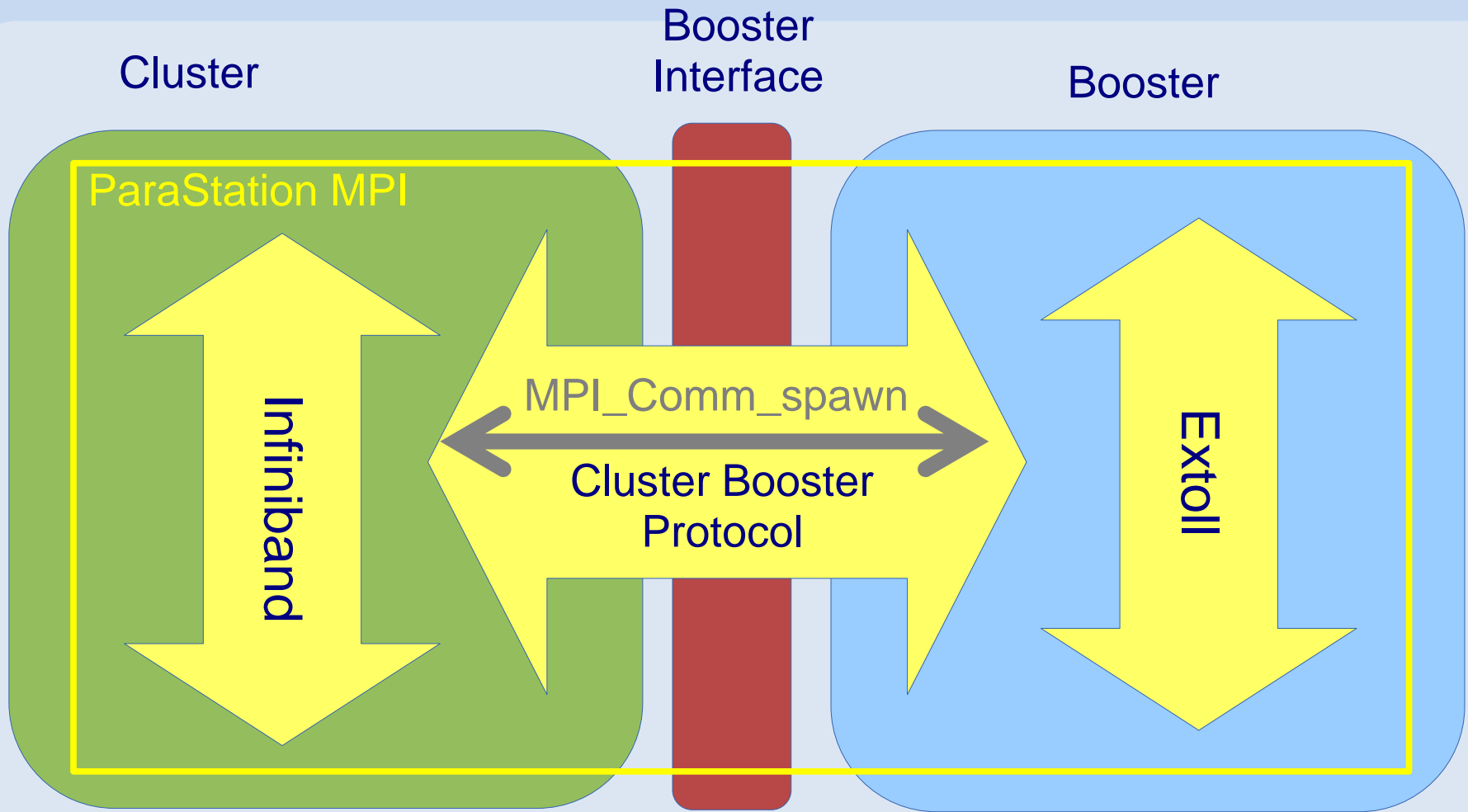


NVMe

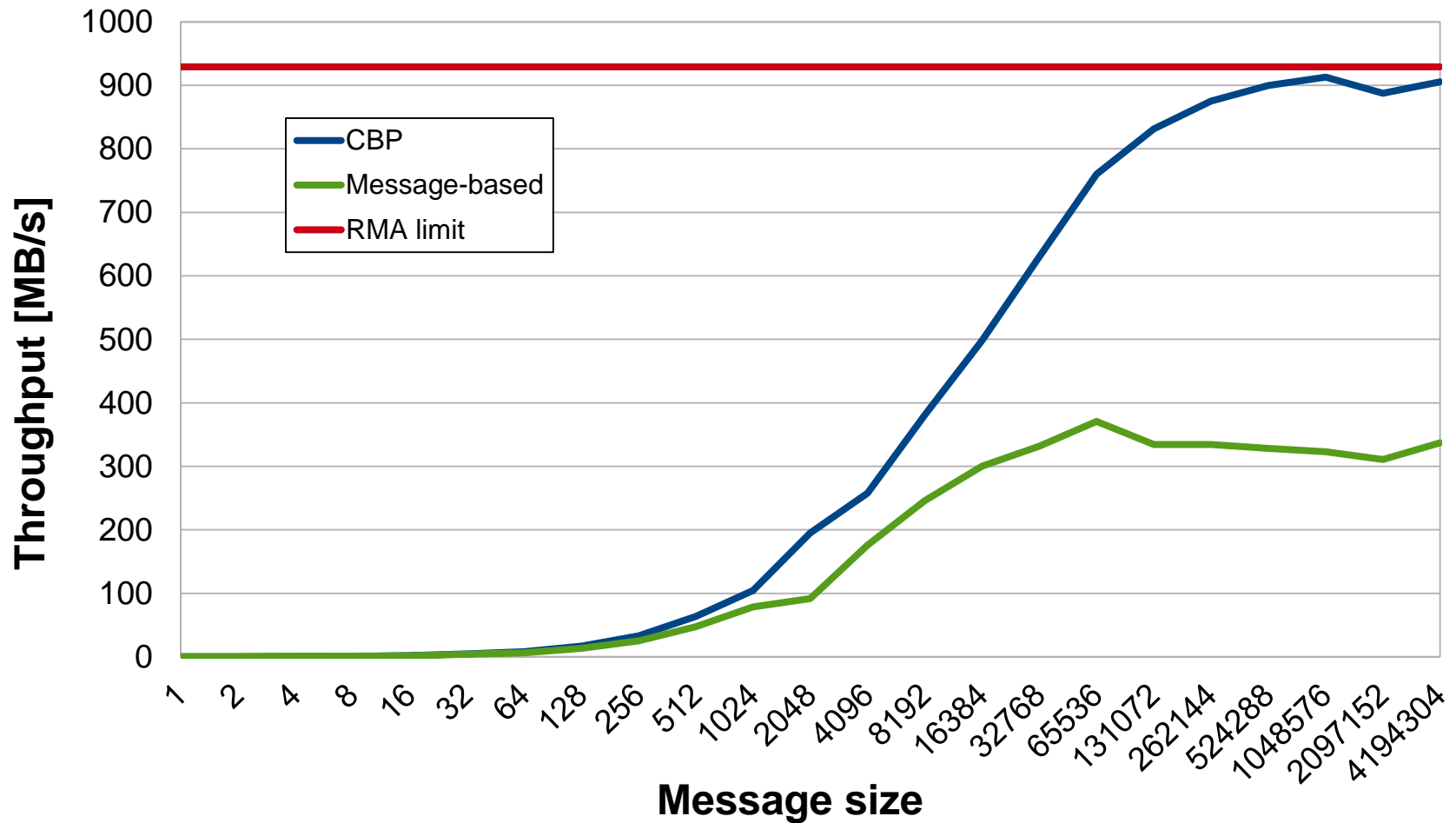


EXTOLL Tourmalet

SOFTWARE



OmpSs on top of MPI provides pragmas to ease the offload process



Source code

```
int main(int argc, char *argv[]){
    /*...*/
    for(int i=0; i<3; i++){
        #pragma omp task in(...) out (...) onto (com, size*rank+1)
        foo_mpi(i, ...);}
}
```

Compiler

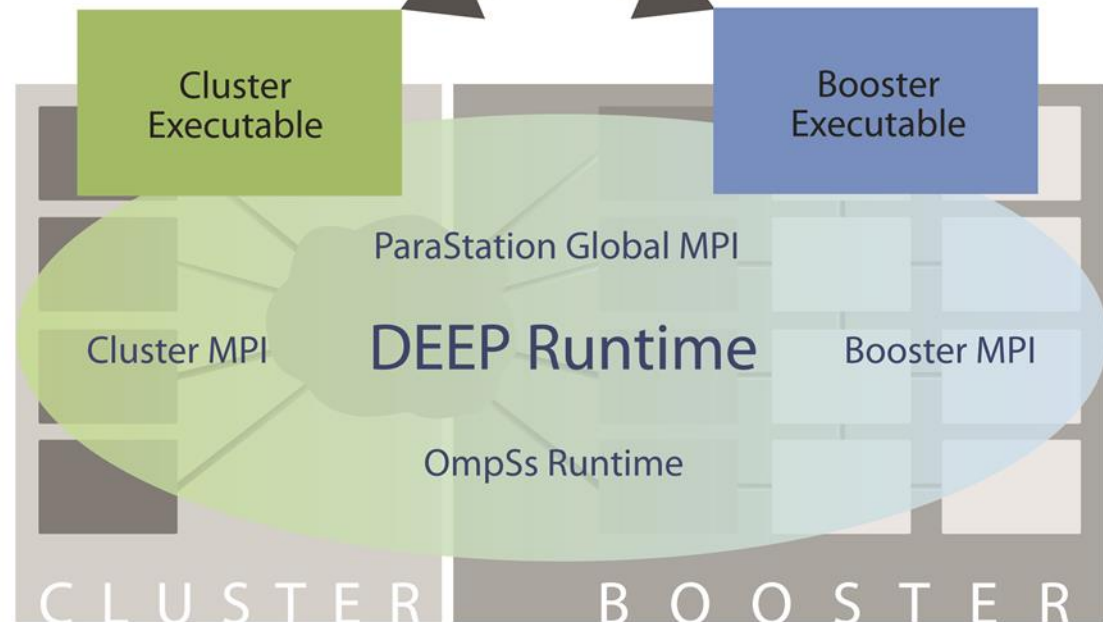
OmpSs Compiler

Application
binaries

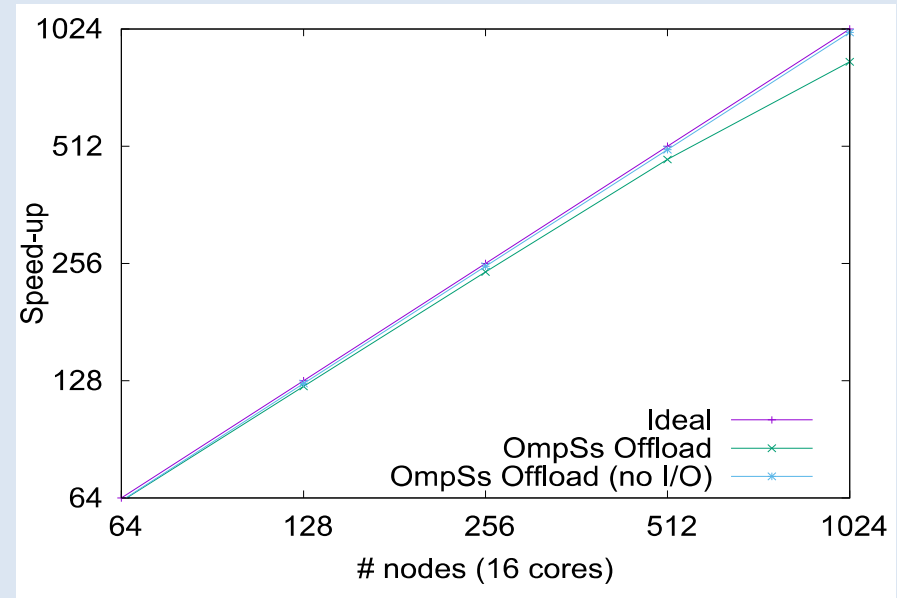
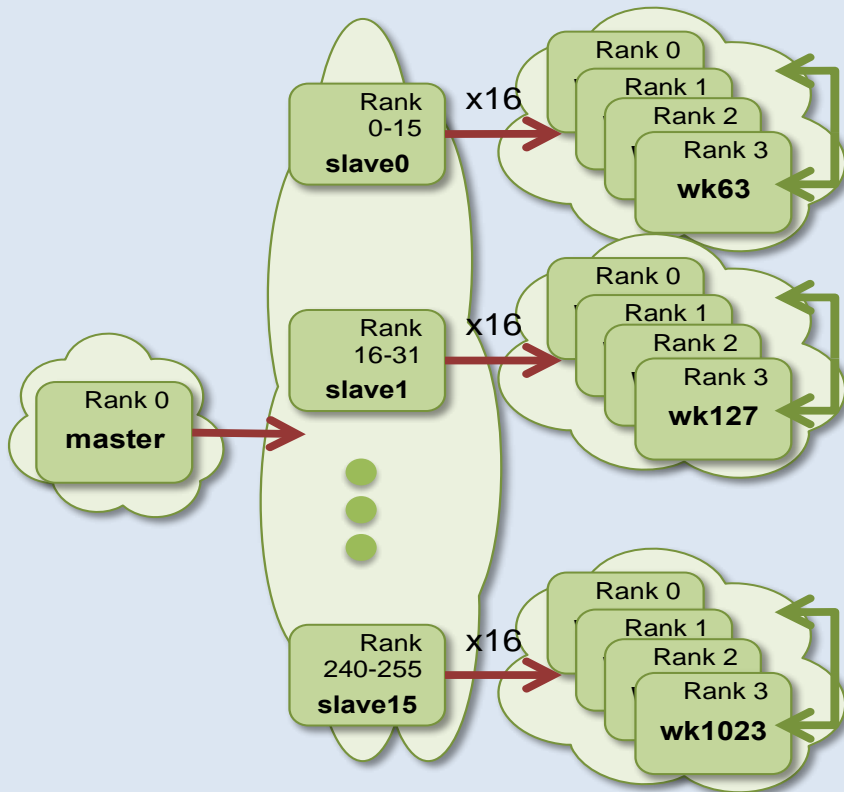
Cluster
Executable

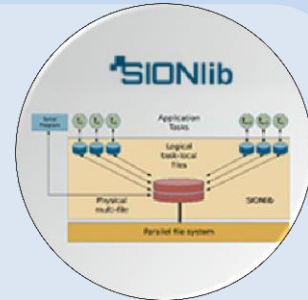
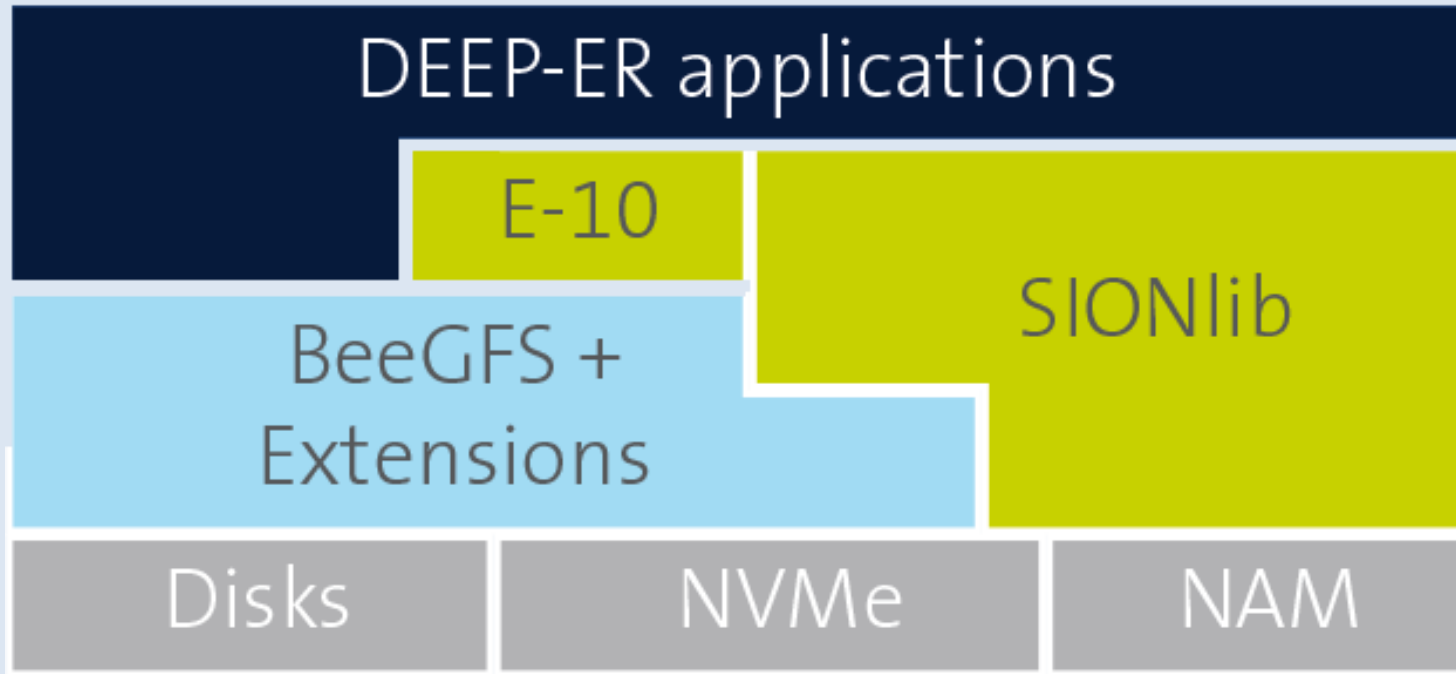
Booster
Executable

DEEP
Runtime

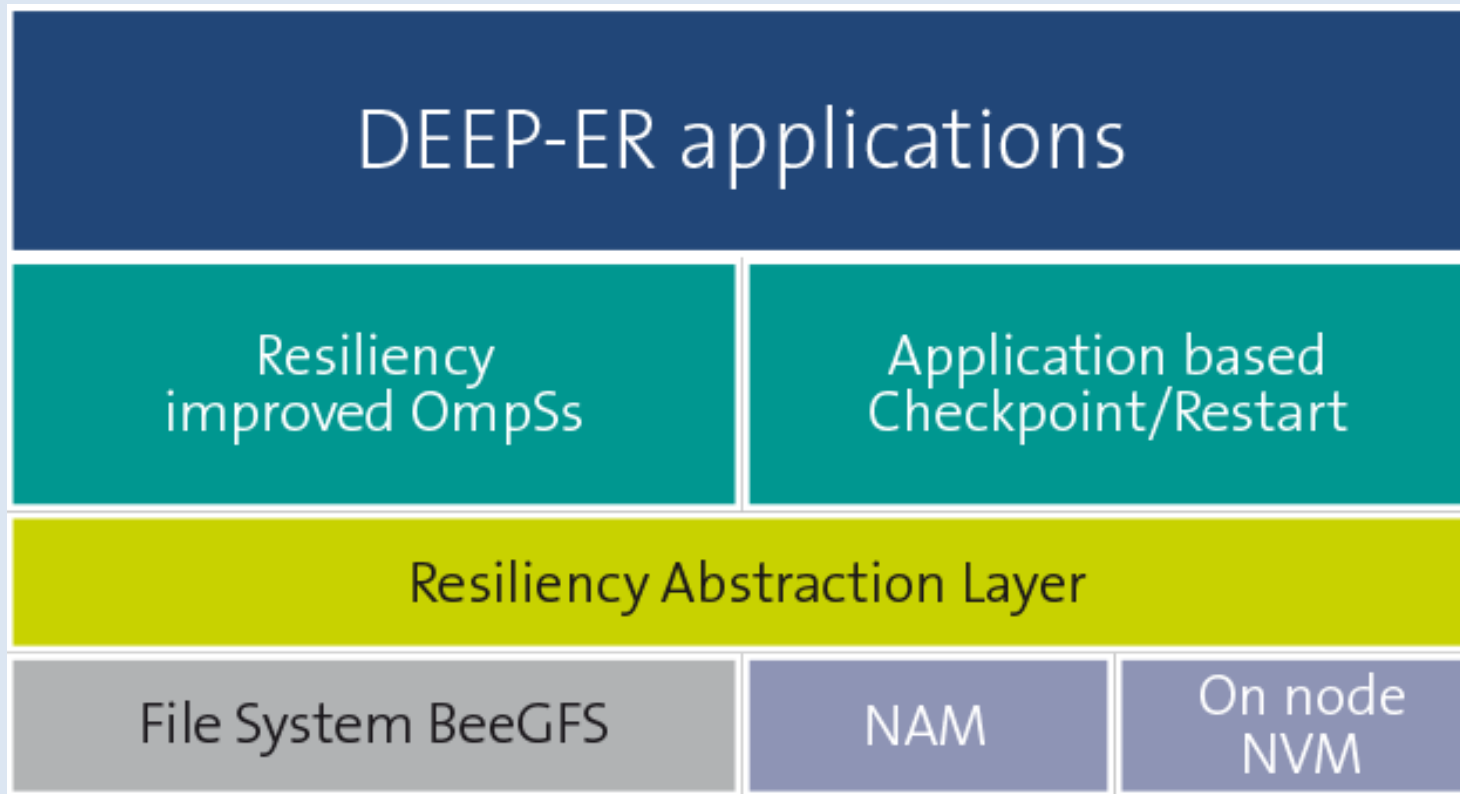


FWI (full wave inversion) code





- Improve I/O scalability on all usage-levels
- Used also for checkpointing



- Develop a hierarchical, distributed checkpoint/restart mechanism leveraging DEEP-ER architecture

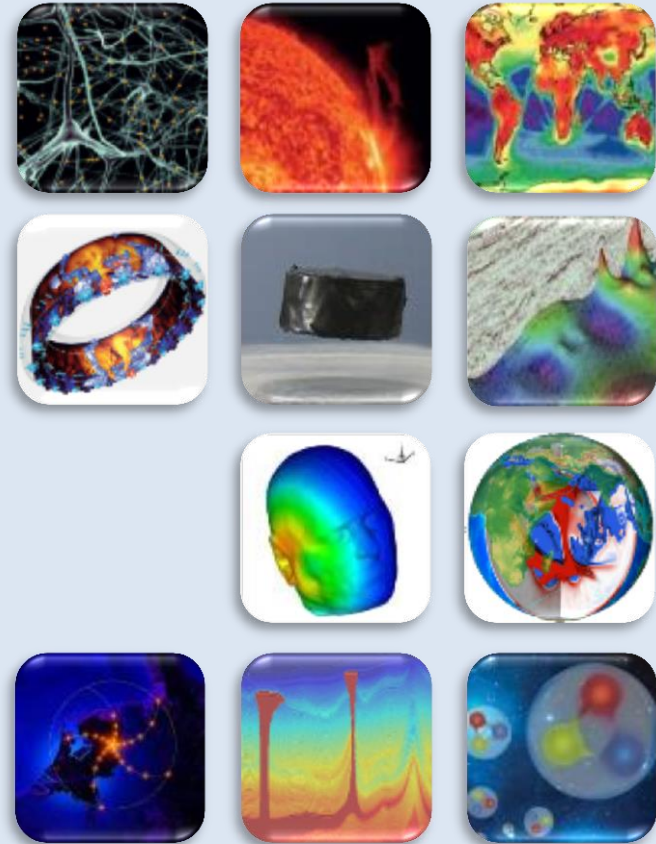
APPLICATIONS

- **DEEP+DEEP-ER applications:**

- Brain simulation (EPFL)
- Space weather simulation (KULeuven)
- Climate simulation (Cyprus Institute)
- Computational fluid engineering (CERFACS)
- High temperature superconductivity (CINECA)
- Seismic imaging (CGG)
- Human exposure to electromagnetic fields (INRIA)
- Geoscience (LRZ Munich)
- Radio astronomy (Astron)
- Oil exploration (BSC)
- Lattice QCD (University of Regensburg)

- **Goals:**

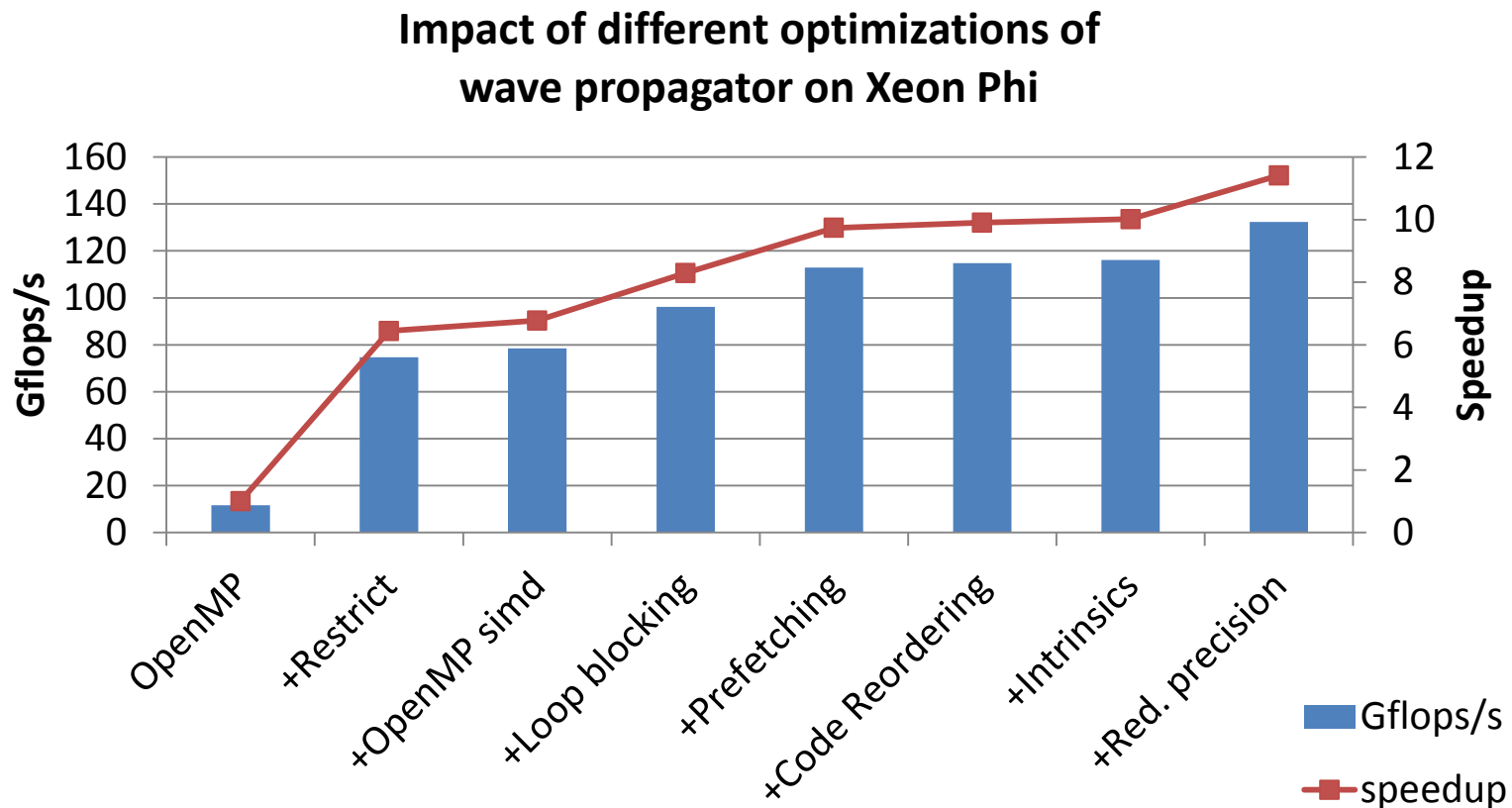
- Co-design and evaluation of architecture and its programmability
- Analysis of the I/O and resiliency requirements of HPC codes



- More **flexible** than a standard architecture
 - This enables different use models:
 1. Dynamic ratio of processors/coprocessors
 2. Use Booster as pool of accelerators (globally shared)
 3. Discrete use of the Booster
 4. Discrete use + I/O offload
 5. Specialized symmetric mode
- Enables a **more efficient use of system resources**
 - Only resources actually needed are blocked by applications
 - Dynamic allocation further increases system utilization

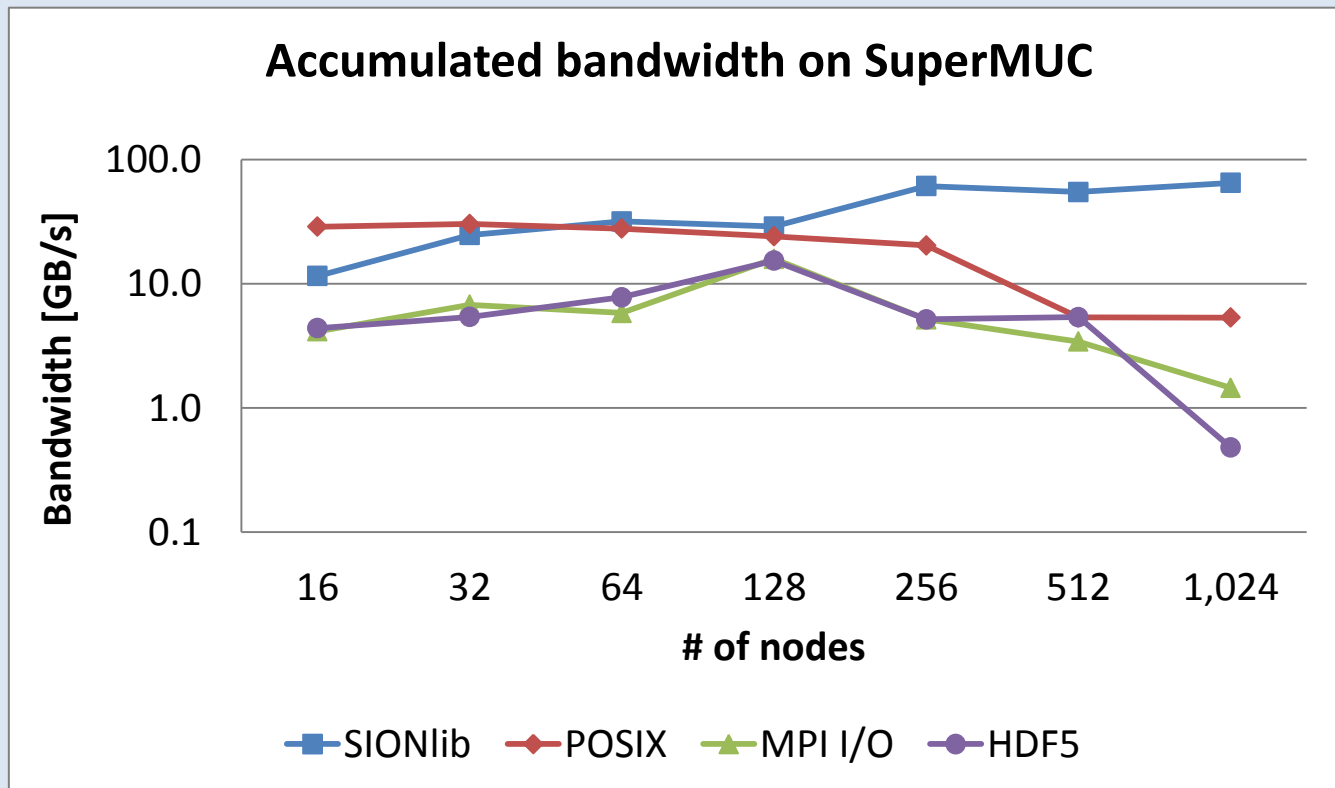
BSC: Enhancing Oil Exploration (FWI, wave propagator)

1 XeonPhi (60 cores), 180 OpenMP threads



LRZ: Rapid crustal deformation & earthquake source equation (Seisol)

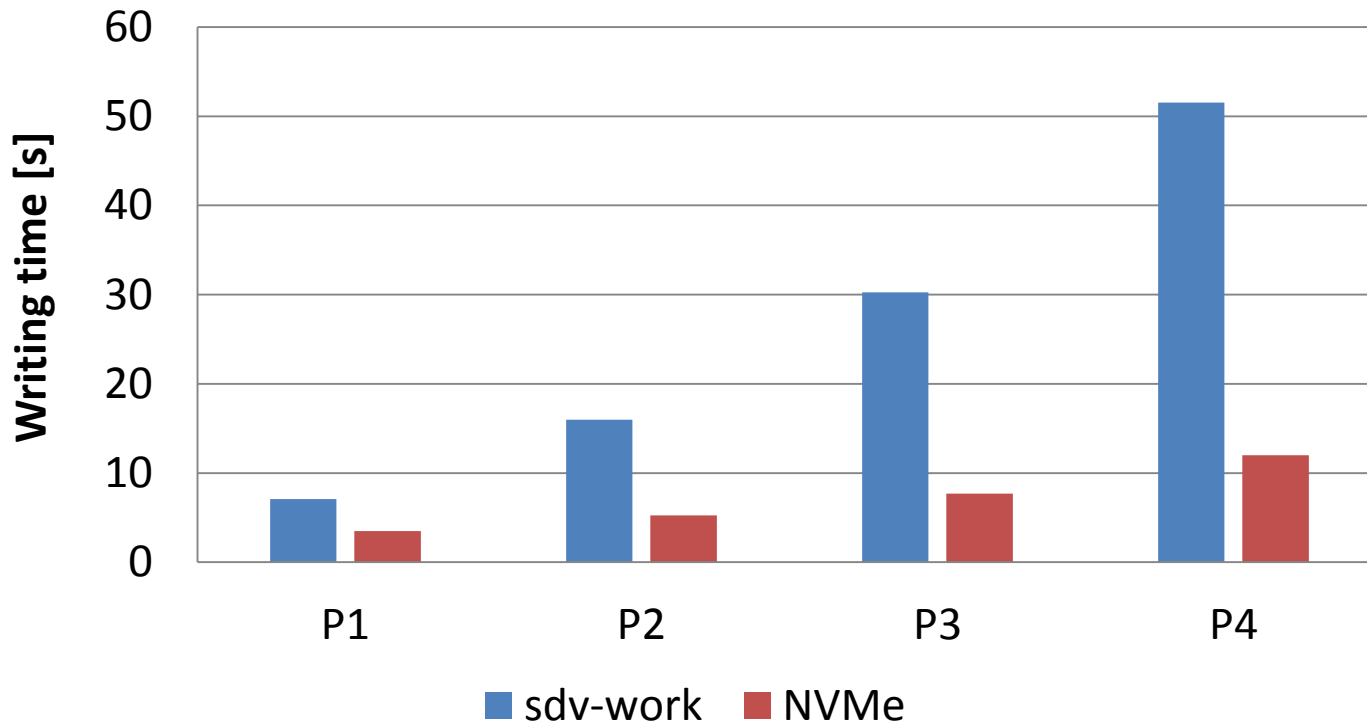
1 process per node, 16 threads per process
writing 20 checkpoint files (4GB/checkpoint)



Inria: Assessment of Human exposure to EM fields

24 MPI processes, 1 thread per process

I/O performance of MAXW-DGTD



Increasing
model precision
 $P1 < P2 < P3 < P4$

- DEEP:
 - Holistic project (HW+SW +App.) in strong co-design
 - Showed the potential of an alternative approach to heterogeneous computing
 - Prototype available also to external users
 - If interested, please contact pmt@deep-project.eu
- DEEP-ER:
 - HW: Prototype under development (impacted by KNL delay)
 - SW + APPs: Very good progress in software and applications
 - APPs: Benefit from new memory technologies shown.
Strong co-design within the project.

- **To SRA:**
 - HPC System Architecture and component (Cluster-Booster concept, EXTOLL network, memory, energy efficiency, etc.)
 - System Software management: (Interconnect management, ParaStation cluster management, etc.)
 - Programming environment (Heterogeneous programming model (MPI + OmpSs), optimisations for resiliency, etc.)
 - Balance compute, I/O and storage performance (BeeGFS and I/O software to efficiently manage multilevel memory hierarchy, etc.)
 - Energy and resiliency: (Warm water cooling, RAS monitoring, DEEP-ER resiliency, etc.)
 - Application code modernisation
- **To international collaboration:**
 - European Exascale Projects (EEP) initiative

EU-Exascale projects
20 partners
Total budget: 28,3 M€
EU-funding: 14,5 M€
Nov 2011 – Mar 2017

Visit us @
ISC'16, Frankfurt
(Germany)
20.-22.06.2016

-Booth
-BoF
-Workshop

