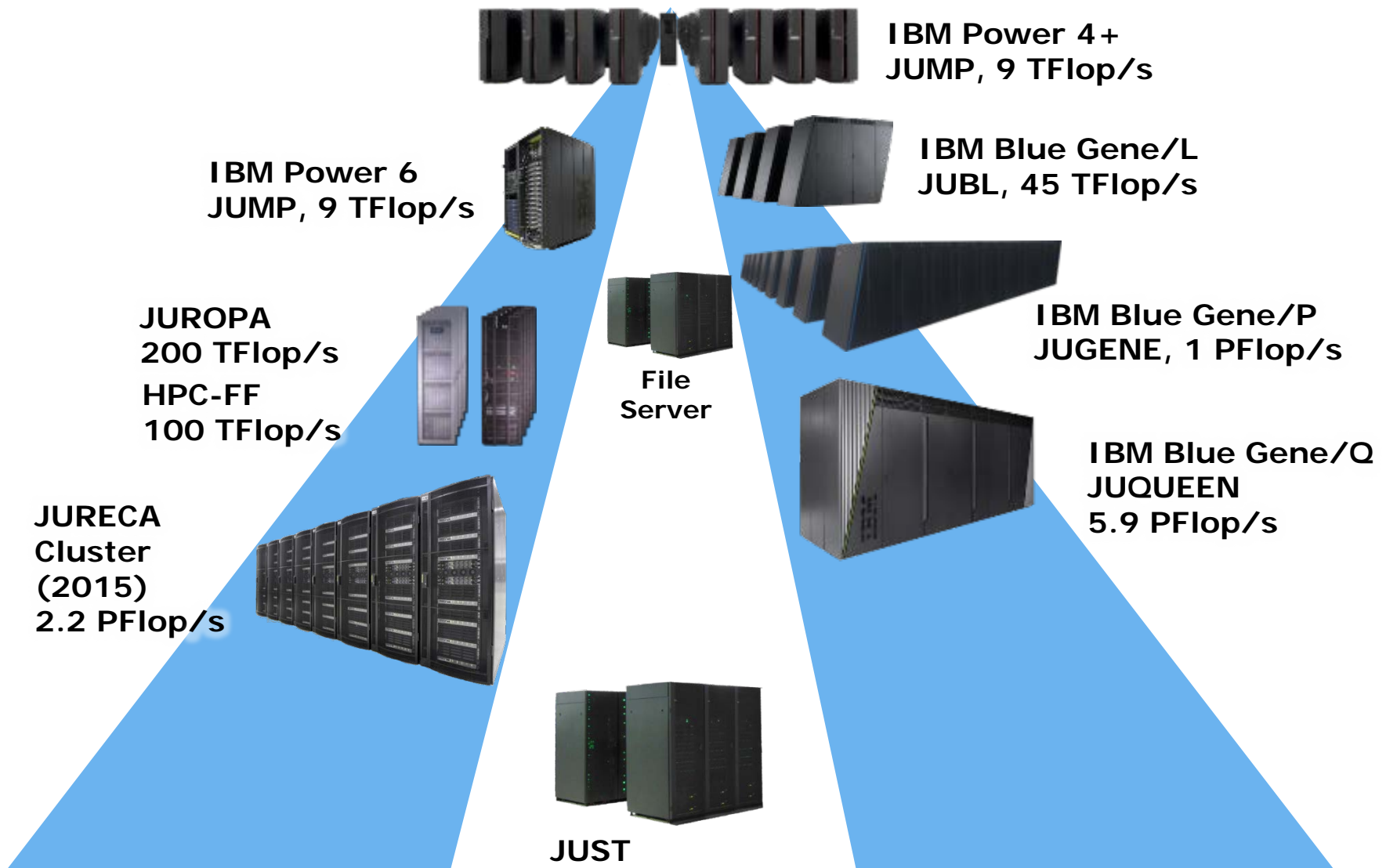


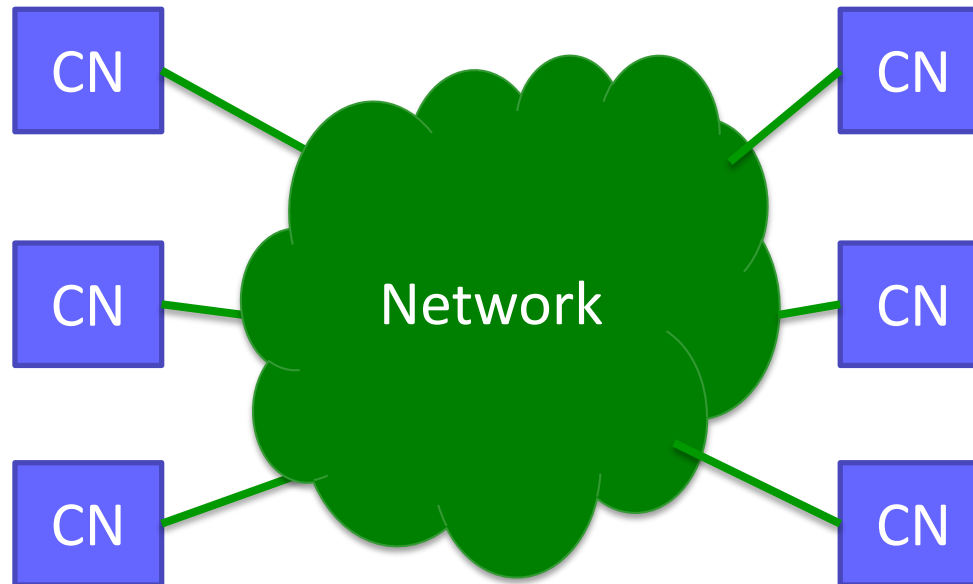
Towards Modular Supercomputing with Slurm

2017-09-25 | Dorian Krause et al.,
Jülich Supercomputing Centre, Forschungszentrum Jülich

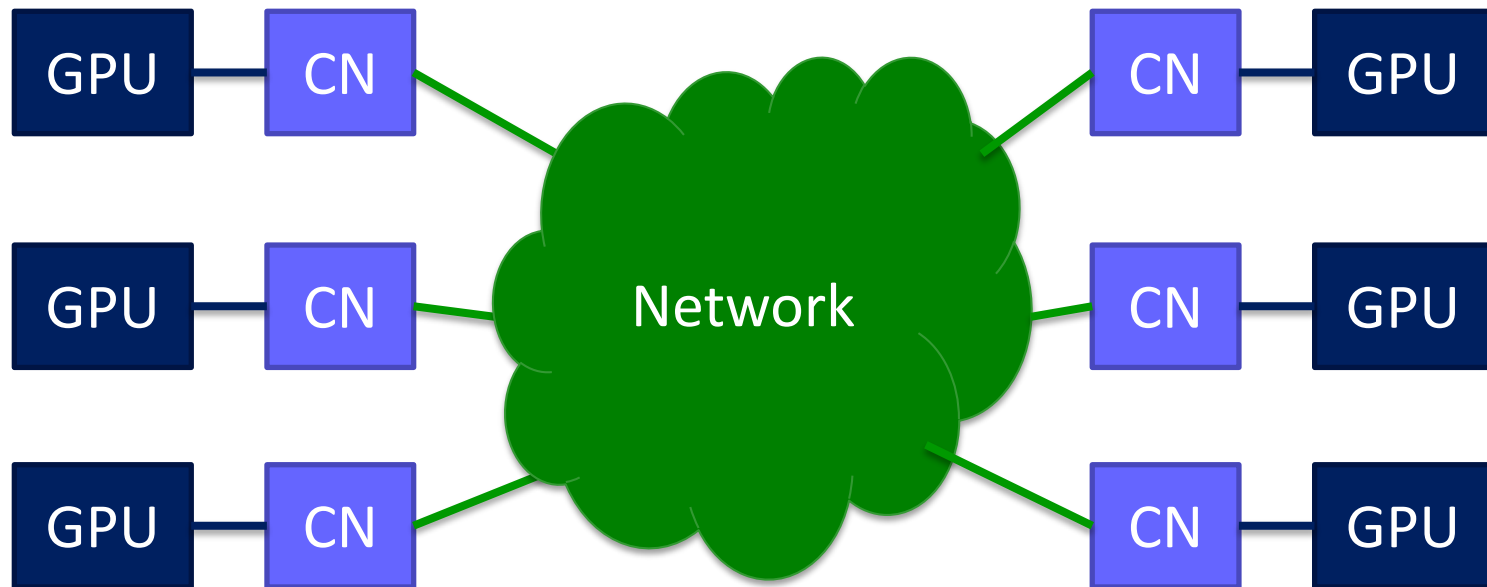
Dual-Architecture Supercomputing Facility



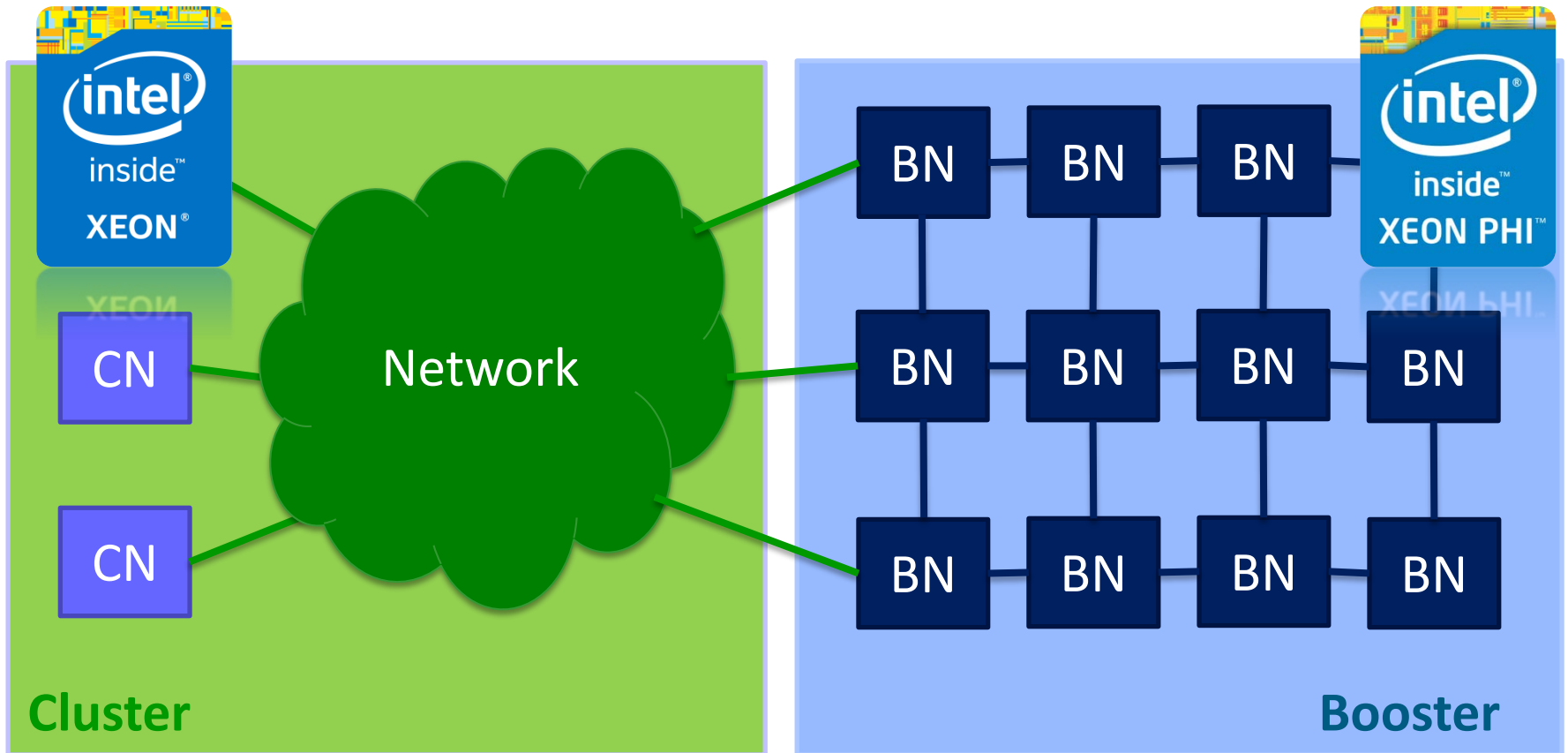
Homogeneous Cluster



Heterogeneous Cluster



Cluster-Booster Architecture



The DEEP projects: DEEP, DEEP-ER, DEEP-EST

www.deep-projects.eu

EU-Exascale projects

27 partners

Total budget: 44 M€

EU-funding: 30 M€

Nov 2011 – Jun 2020

All combine:

- Hardware
- Software
- Applications

in a strong co-design



The DEEP projects: DEEP, DEEP-ER, DEEP-EST



www.deep-projects.eu

EU-Exascale projects

27 partners

Total budget

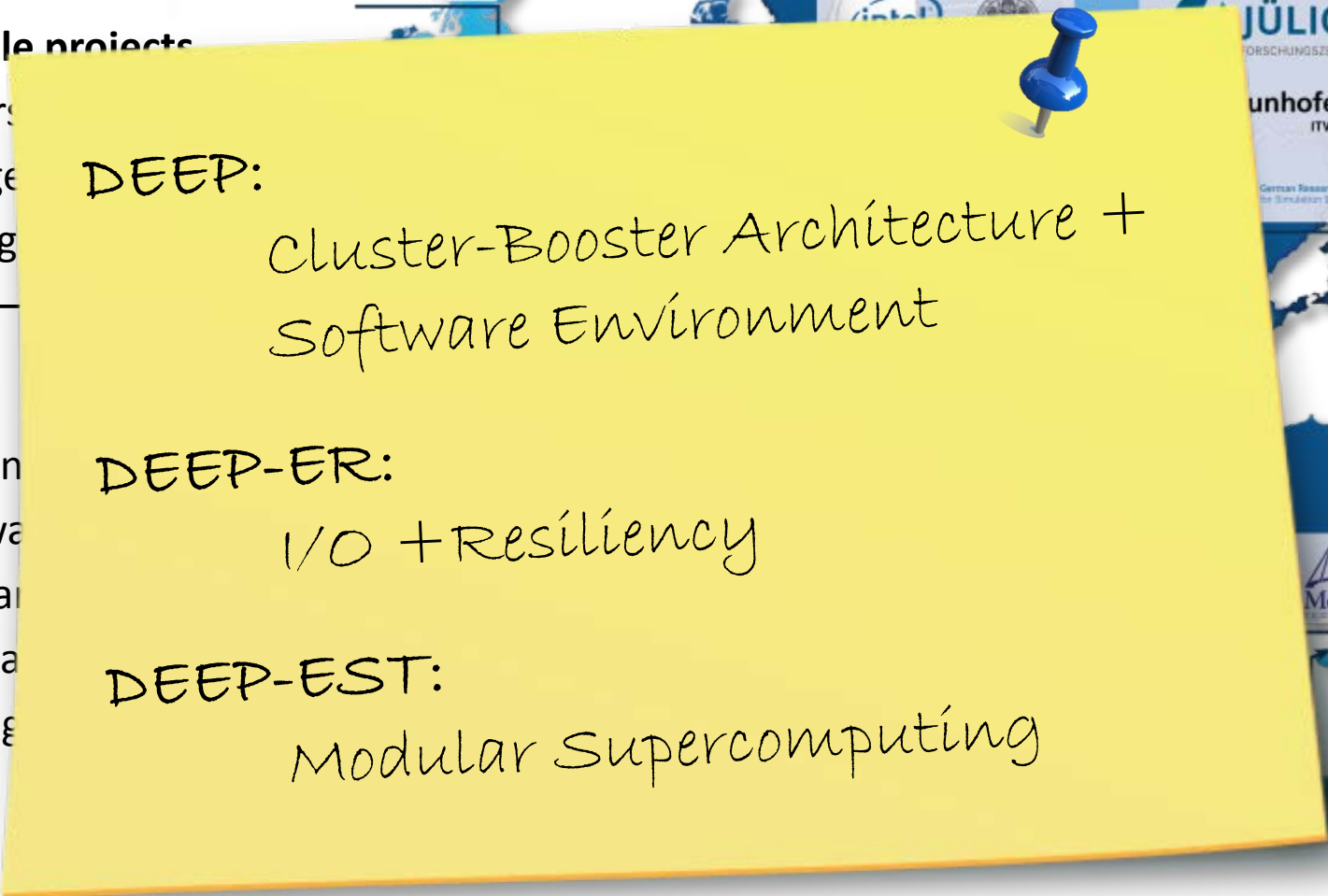
EU-funding

Nov 2011 –

All combin

- Hardwa
- Softwa
- Applica

in a strong



DEEP:
Cluster-Booster Architecture +
Software Environment

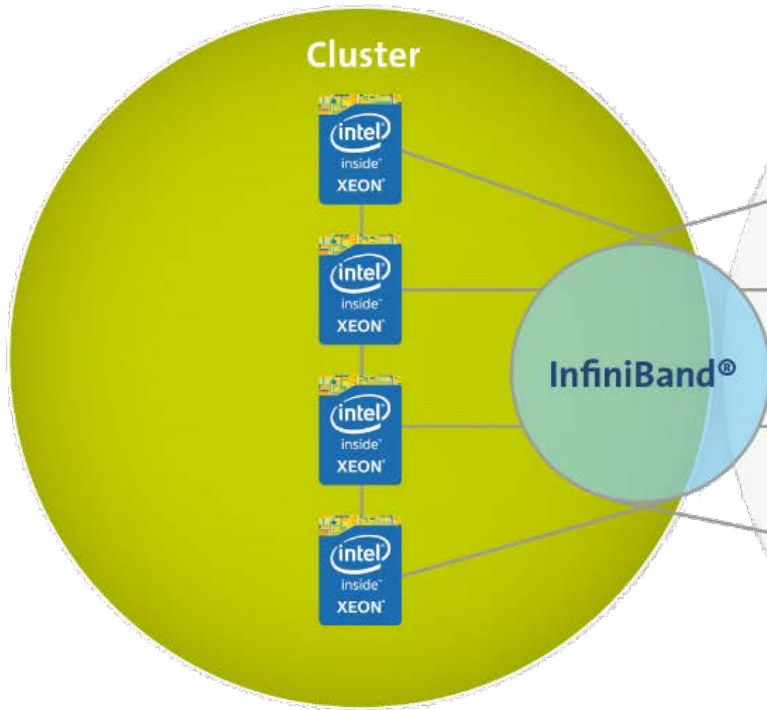
DEEP-ER:
I/O + Resiliency

DEEP-EST:
Modular Supercomputing

DEEP Architecture

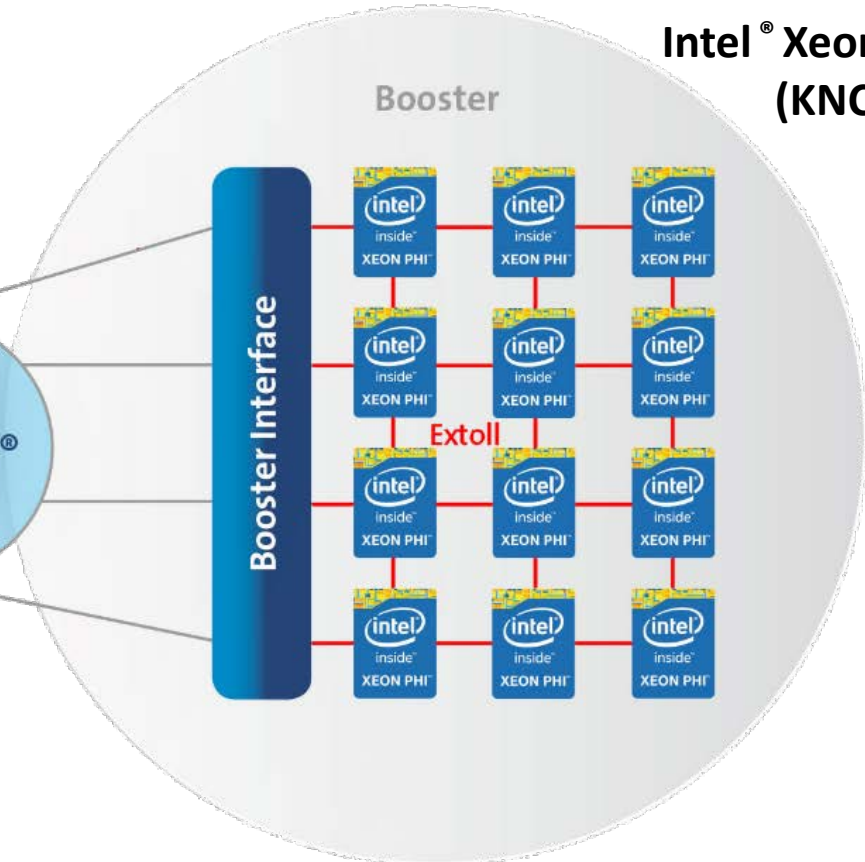


Intel® Xeon®



LOW/MEDIUM
SCALABLE CODE

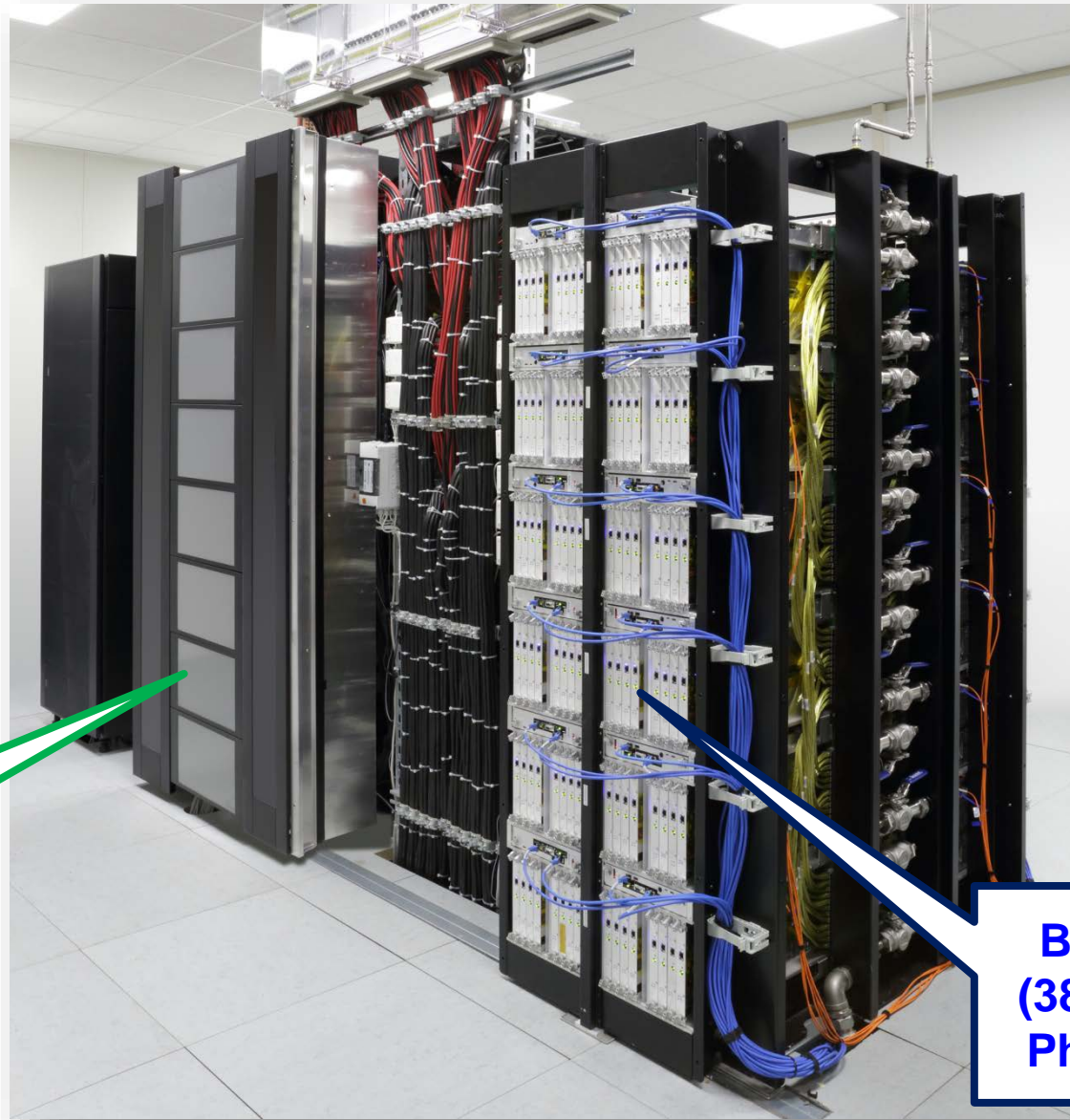
Intel® Xeon Phi™
(KNC)



HIGHLY
SCALABLE CODE

DEEP Prototype

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

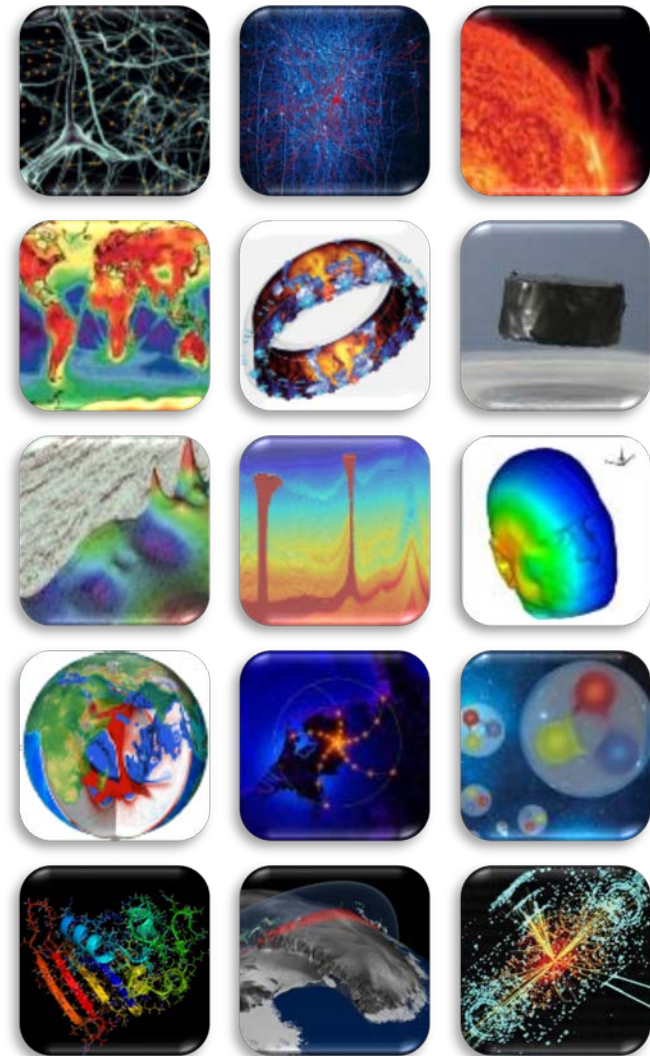
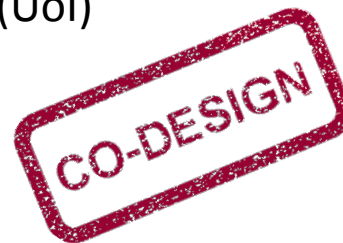


**Cluster
(128 Xeon)**

**Booster
(384 Xeon
Phi KNC)**

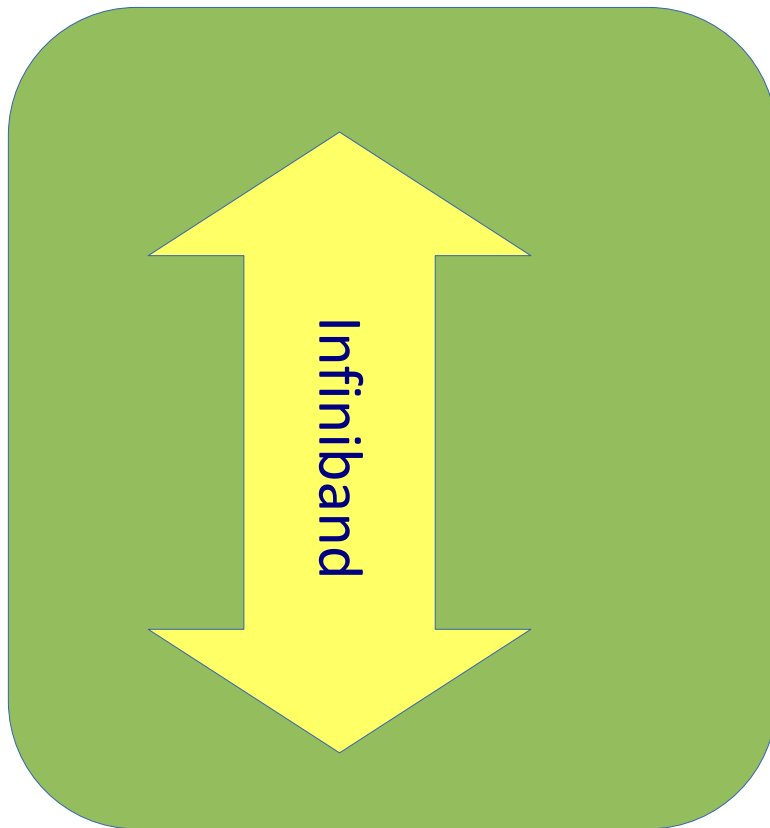
DEEP projects applications (15):

- Brain simulation (EPFL + NMBU)
- Space weather simulation (KU Leuven)
- Climate simulation (Cyprus Institute)
- Computational fluid engineering (CERFACS)
- High temperature superconductivity (CINECA)
- Seismic imaging (CGG + BSC)
- Human exposure to electromagnetic fields (INRIA)
- Geoscience (LRZ)
- Radio astronomy (Astron)
- Lattice QCD (University of Regensburg)
- Molecular dynamics (NCSA)
- Data analytics in Earth Science (Uoi)
- High Energy Physics (CERN)

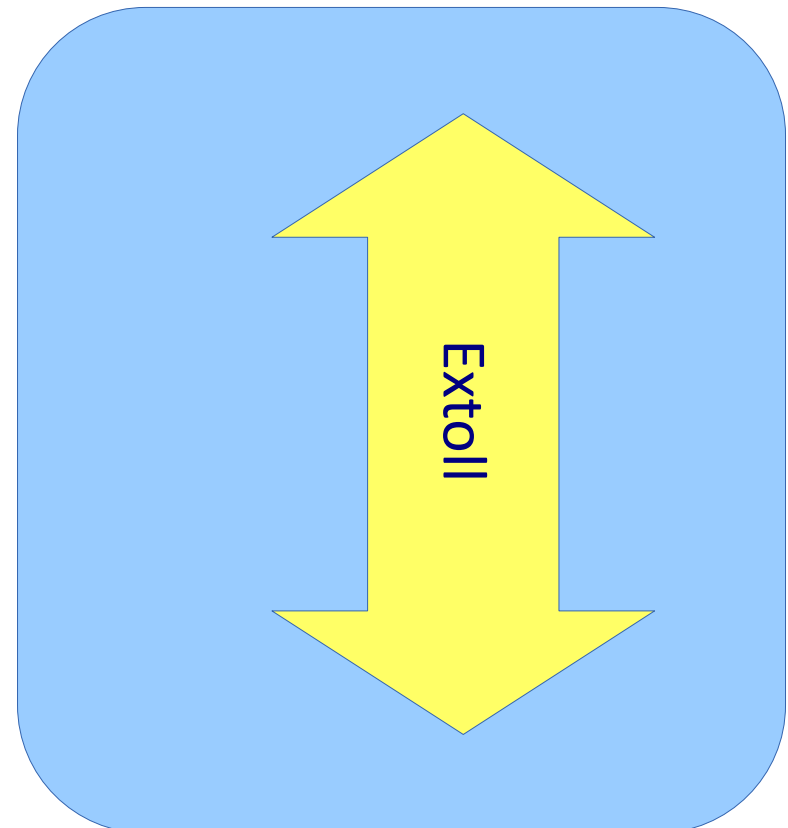


Programming Environment

Cluster



Booster

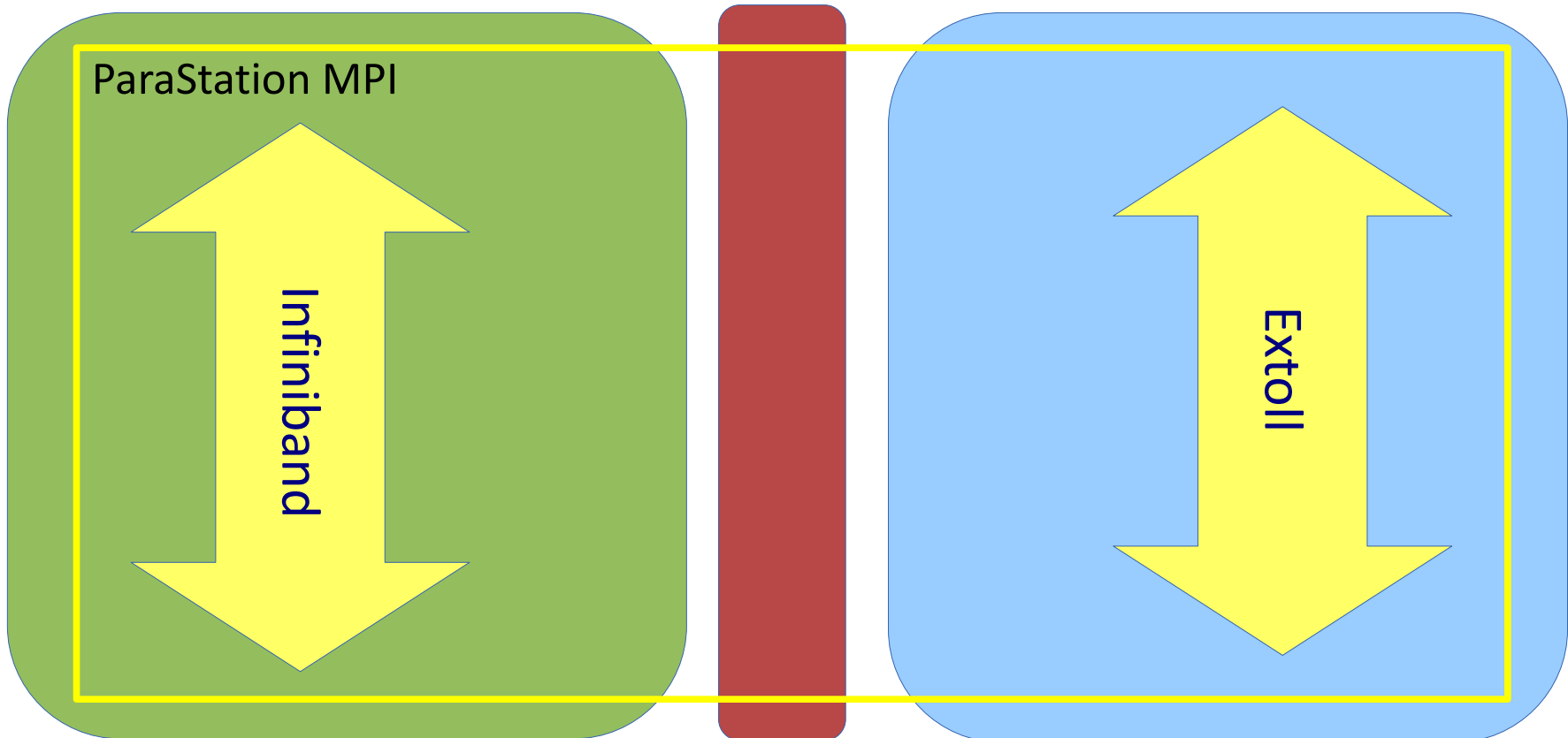


OmpSs on top of MPI provides pragmas to ease the offload process

Programming Environment

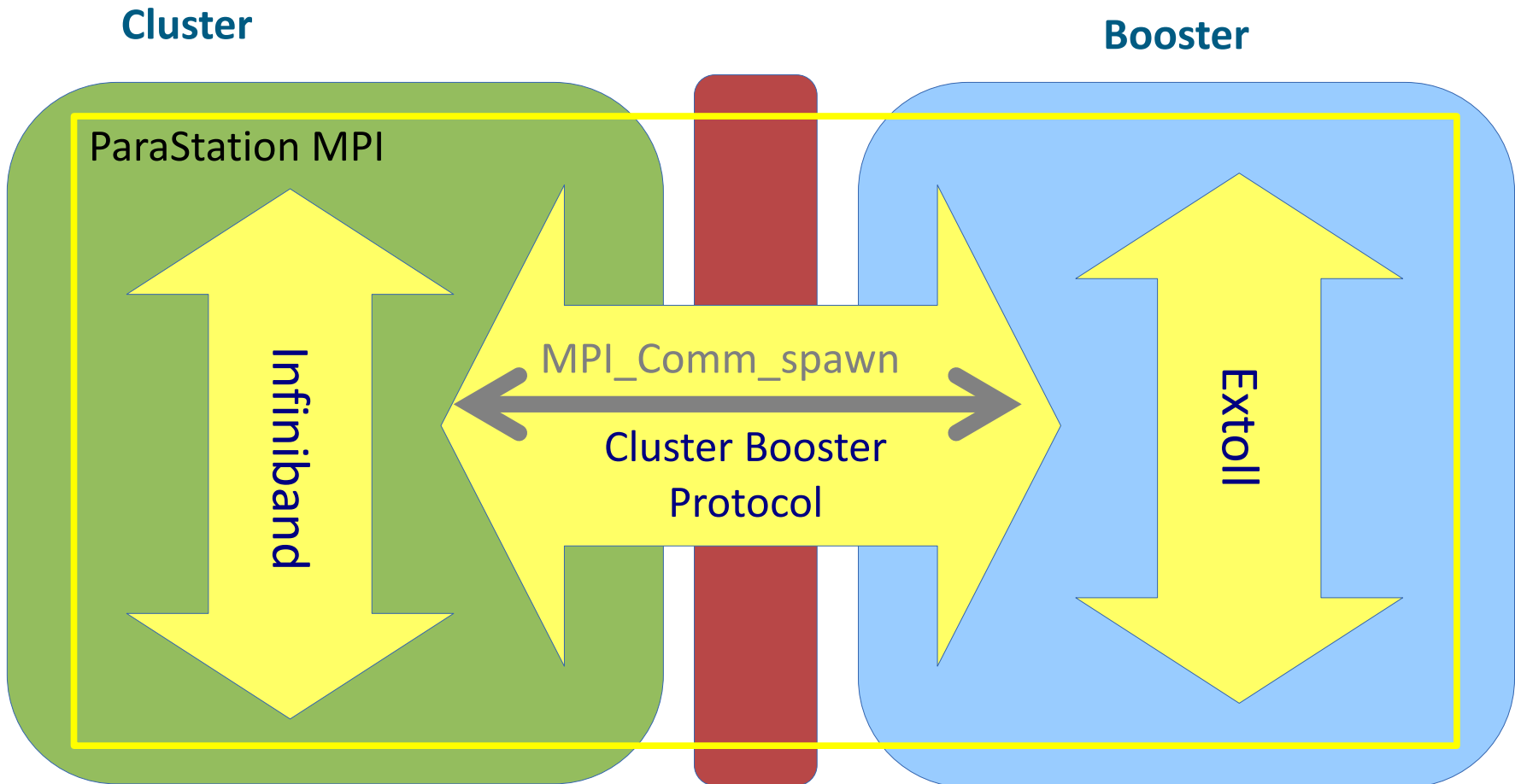
Cluster

Booster



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

Programming Environment

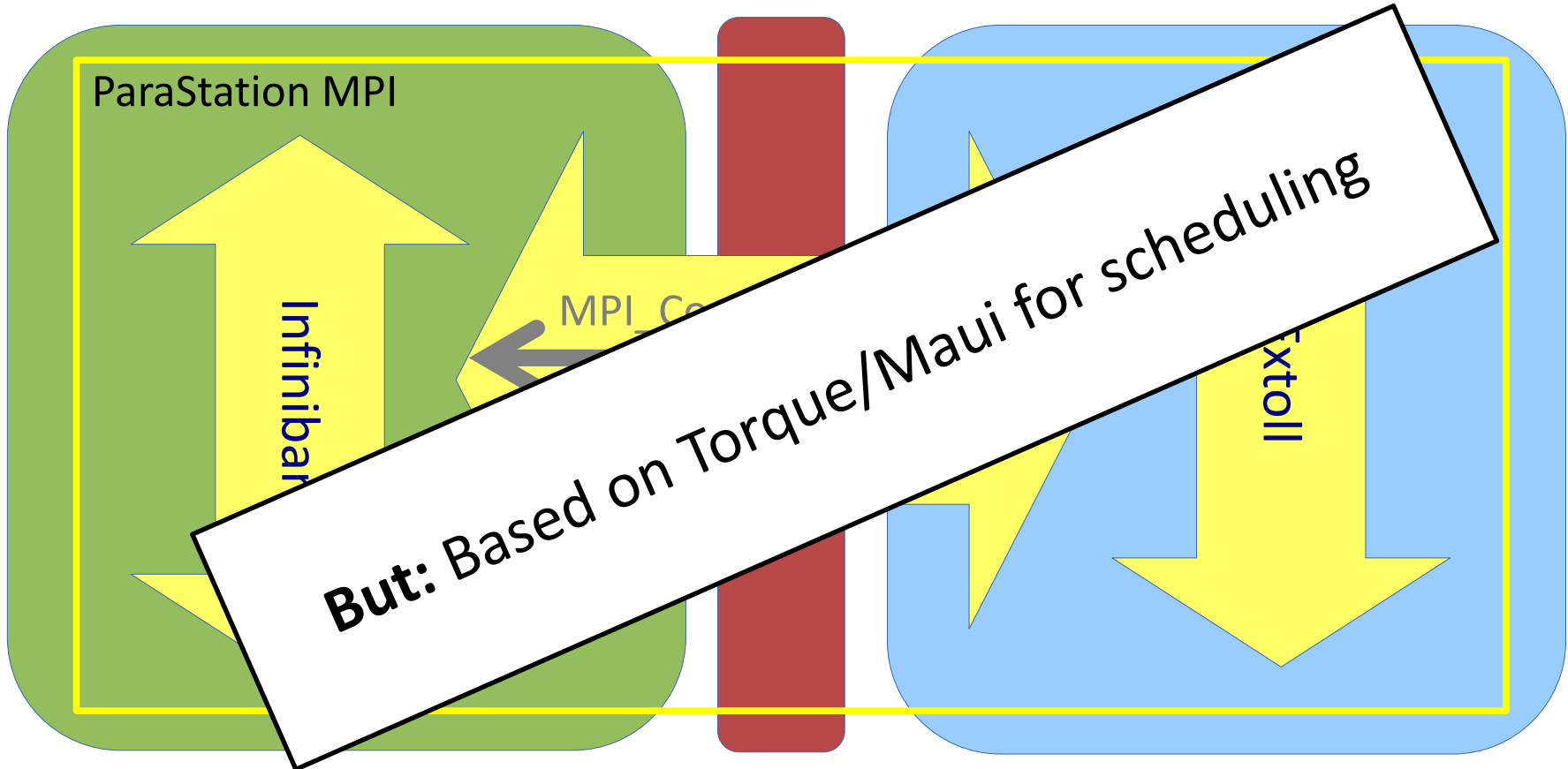


OmpSs on top of MPI provides pragmas to ease the offload process

Programming Environment

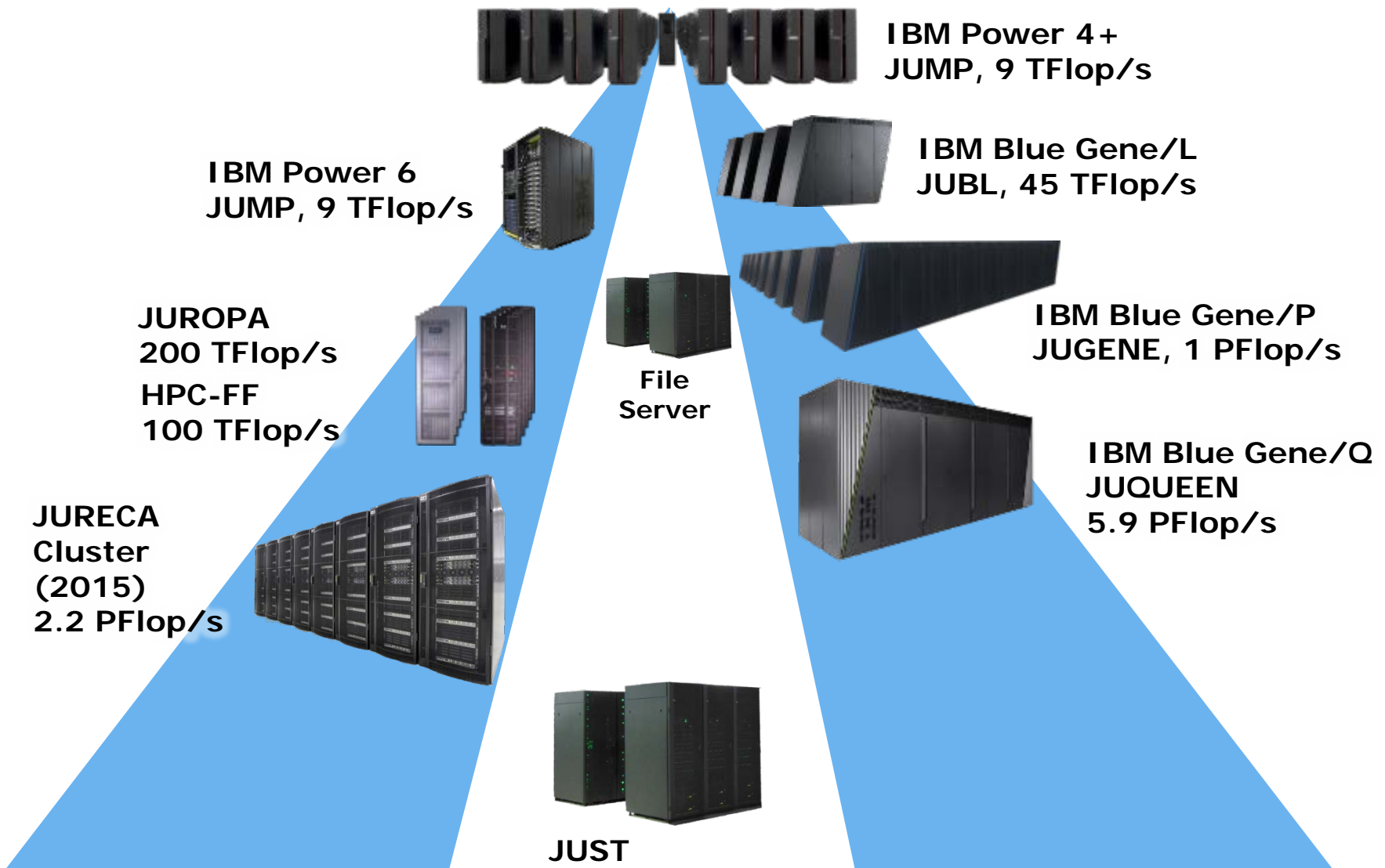
Cluster

Booster

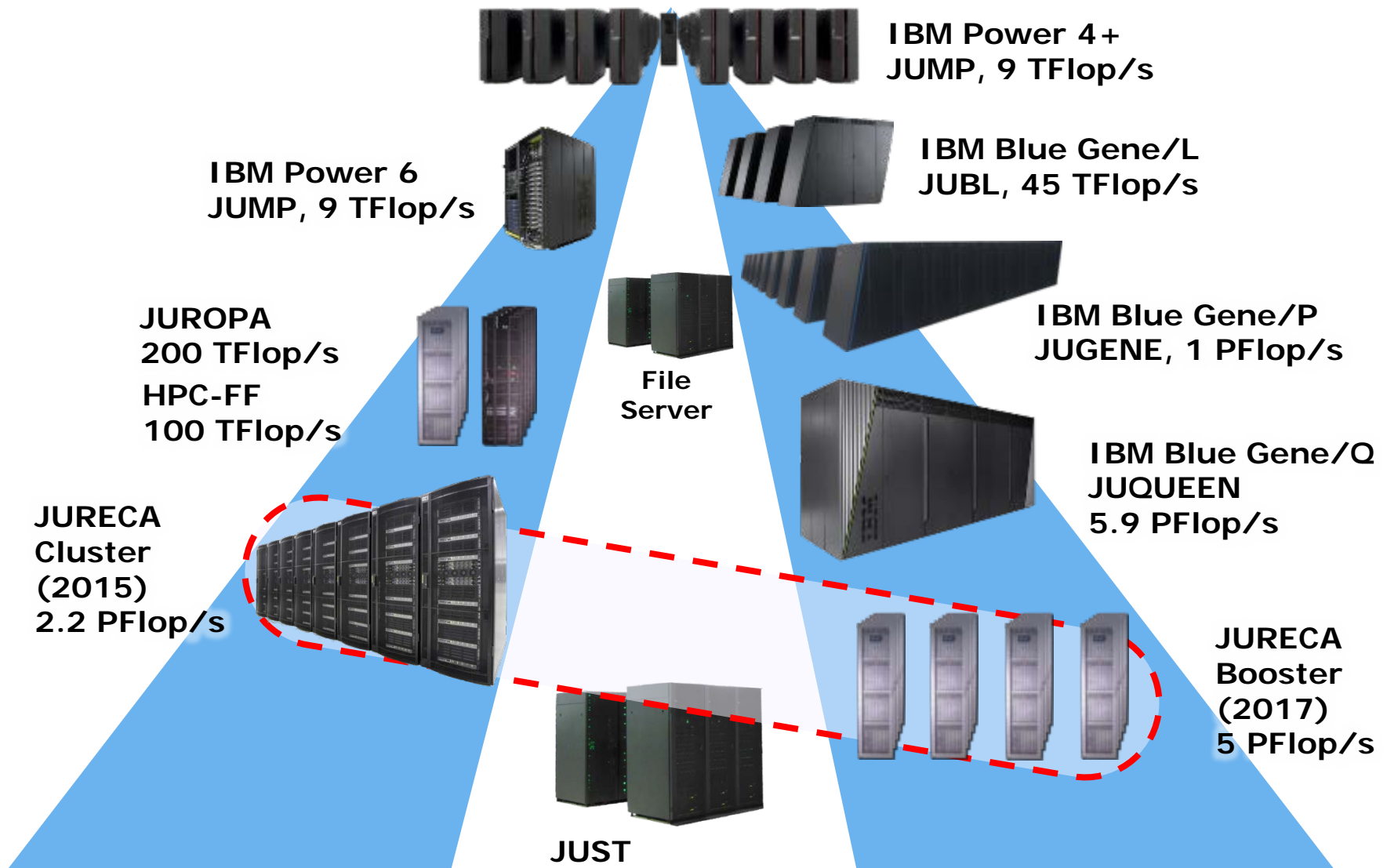


OmpSs on top of MPI provides pragmas to ease the offload process

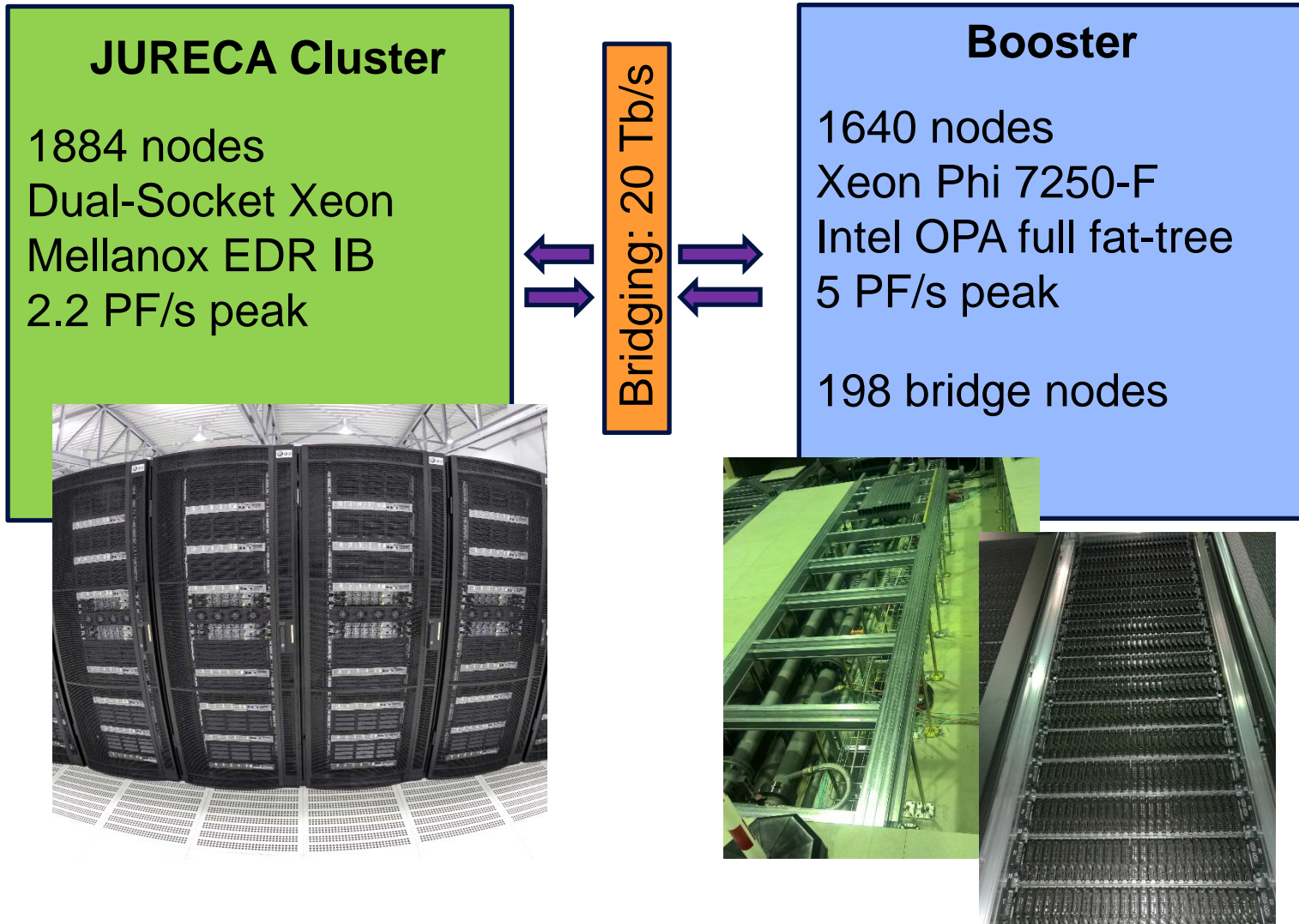
Dual-Architecture Supercomputing Facility



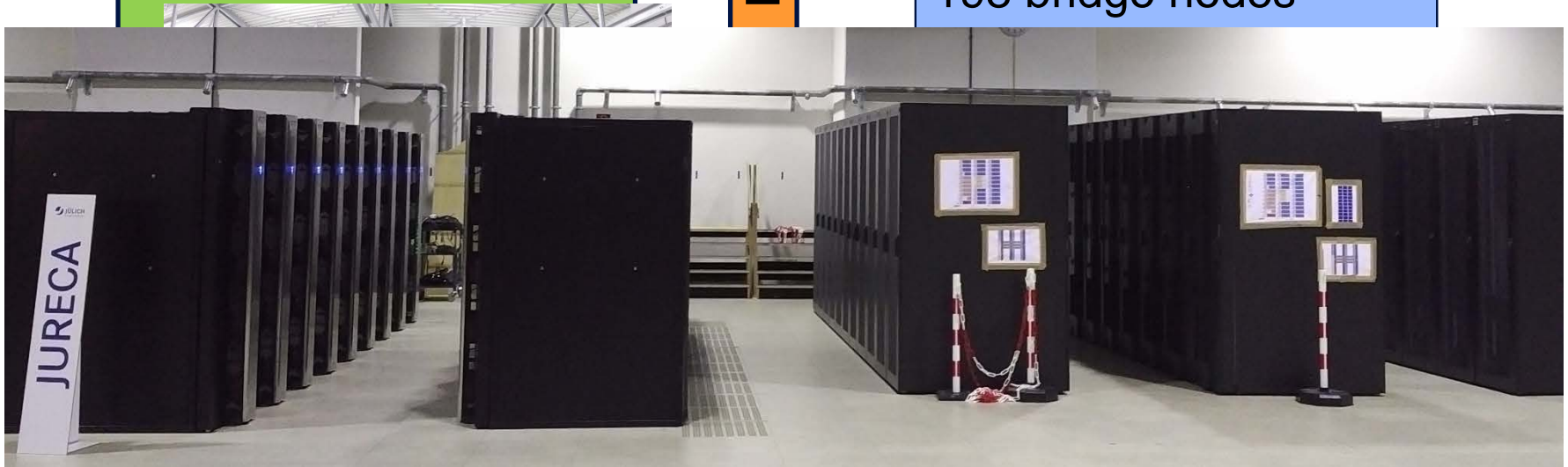
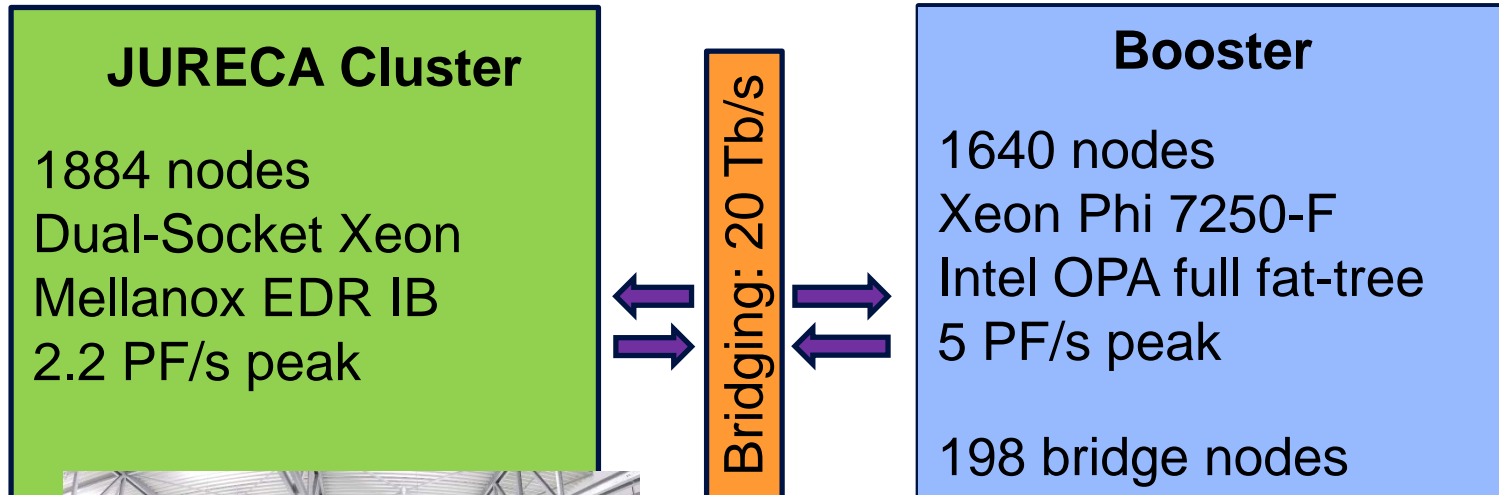
Dual-Architecture Supercomputing Facility



JURECA Cluster-Booster System



JURECA Cluster-Booster System



JURECA Cluster-Booster System

JURECA Cluster

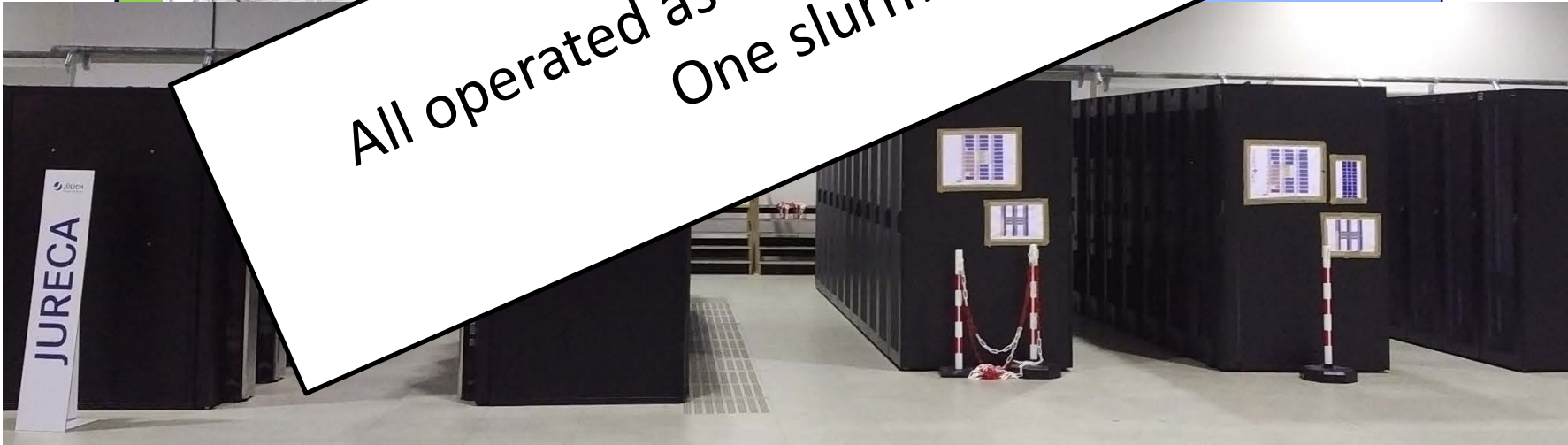
- 1884 nodes
- Dual-Socket Xeon
- Mellanox EDR IB
- 2.2 PF/s peak

20 Tb/s

Booster

1640 nodes

All operated as one system with Slurm:
One slurmctld!



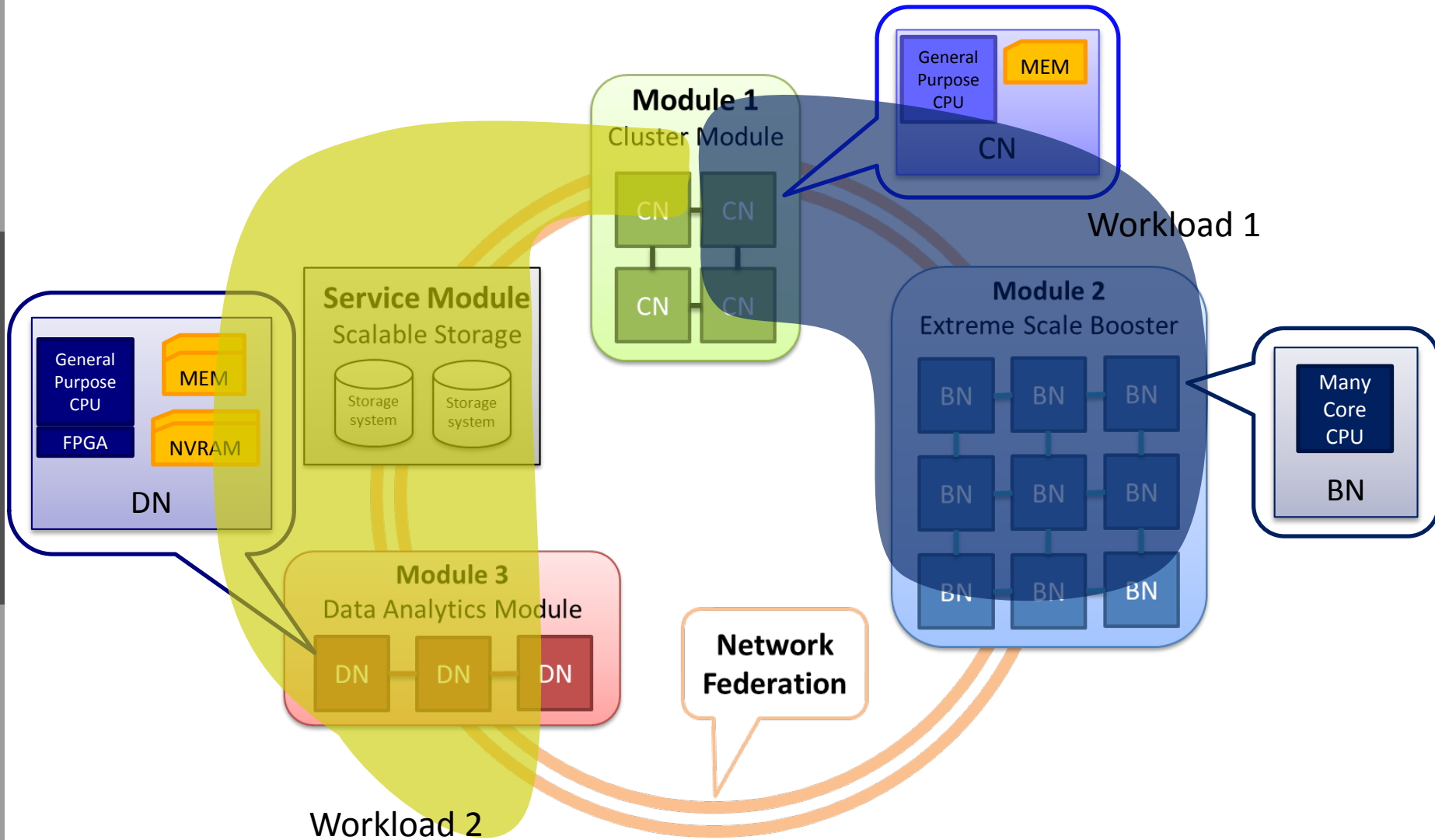
Building blocks

Cluster Hardware	✓
Cluster Software	✓
Booster Hardware	✓
Booster Software	✓
Cluster-Booster Bridging for MPI	✓
Workload Management	✓

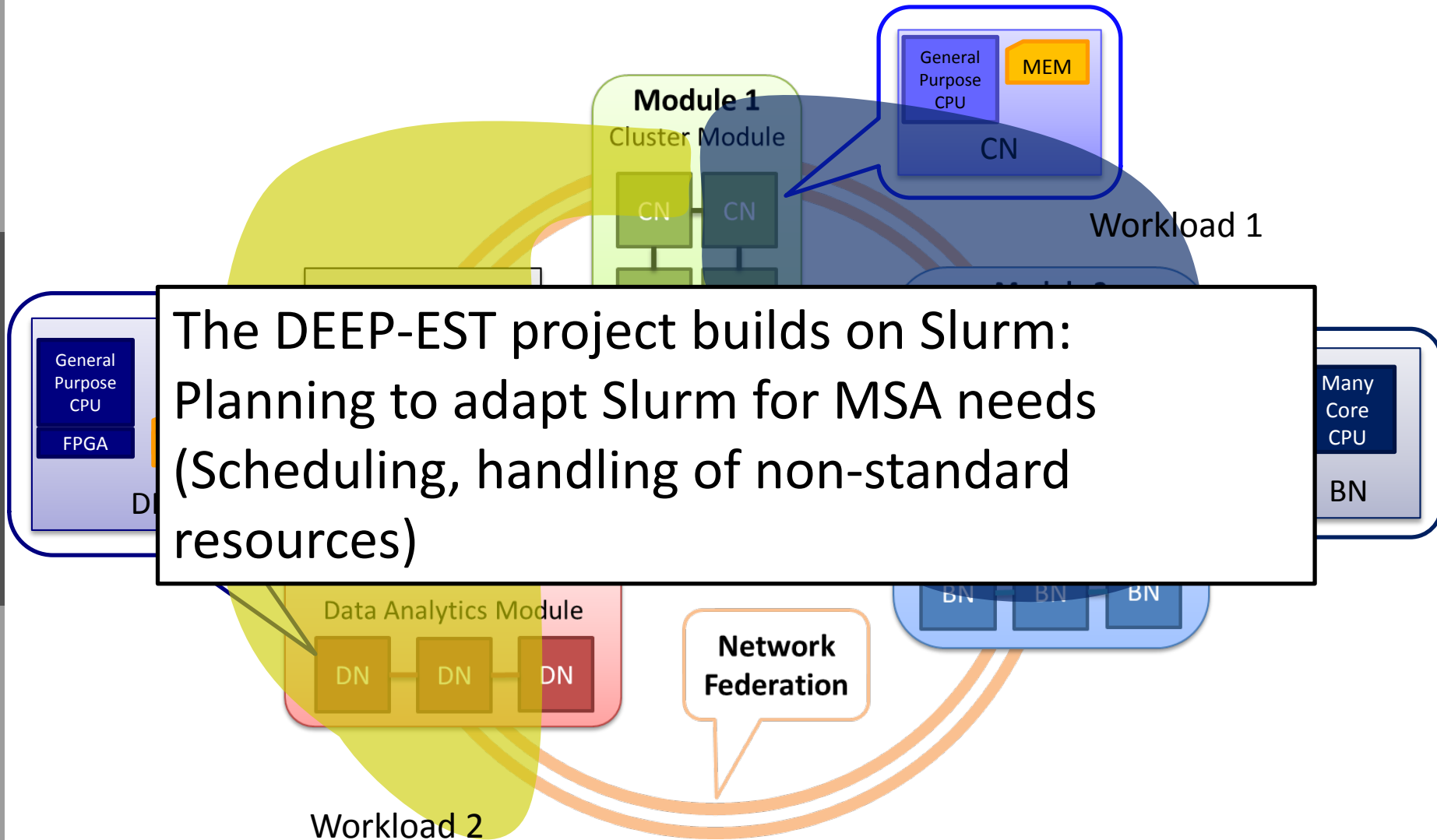
Slurm requirements

- Support for heterogeneity, QoS per partition ✓
- Support for flexible heterogeneous jobs across partitions

Next step: Modular Supercomputing



Next step: Modular Supercomputing



- Cluster-Booster Architecture combines multi- and many-core (or other accelerator) technologies in one system and enables a flexible resource selection.
- JURECA is currently extended to realize the Cluster-Booster Architecture in a 7.2 PF/s production system.
- Slurm is a key building block of our software stack. The community and SchedMD help us achieve our goal.
- Working towards Modular Supercomputing Architecture in DEEP-EST using Slurm. We hope to be able to contribute back.