

# Taming Heterogeneity by Segregation

## The DEEP view on Exascale



Norbert Eicker

Jülich Supercomputing Centre & University of Wuppertal

596. WE-Heraeus-Seminar

Science Applications for Exascale Computing – Exploring New  
Avenues towards Scalability and Fault-Tolerance

7-9 September 2015  
Physikzentrum Bad Honnef, Germany



The research leading to these results has received funding from the European Community's Seventh Framework Programme (FP7/2007-2013) under *Grant Agreements* n° 287530 and n° 610476

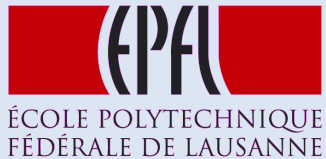
- Motivation
  - Why heterogeneous systems
  - How to organize heterogeneity
- DEEP
  - General concept
  - Hardware architecture
  - Programming paradigm
  - Software stack
  - (Few) Results
- Outlook on DEEP-ER
- Summary



DEEP face-to-face Leuven



DEEP-ER kickoff Jülich



- 16 Partners
  - 4 PRACE hosts
  - 3 Research Centers
  - 5 Industry Partners
  - 4 Universities
  - Coordinator JSC
- 8 Countries
- Duration: 3.5 years
- Budget: 18.3 M€
- EU funding: 8.03 M€

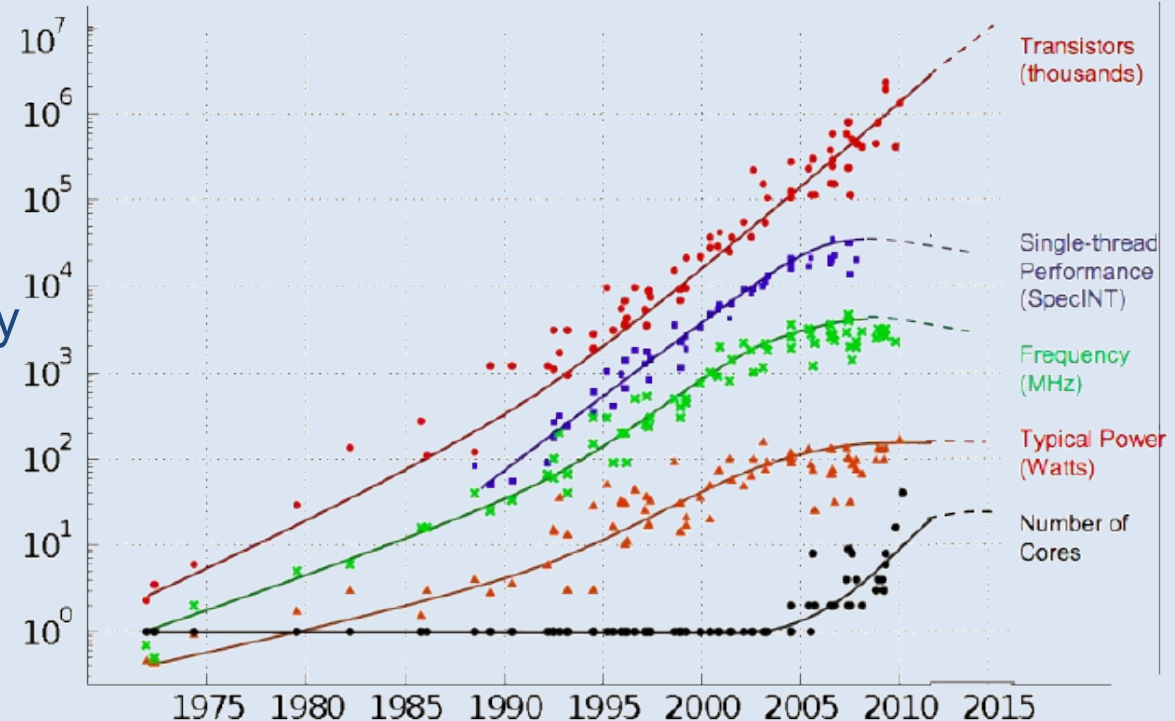


# Heterogeneity



- Observation

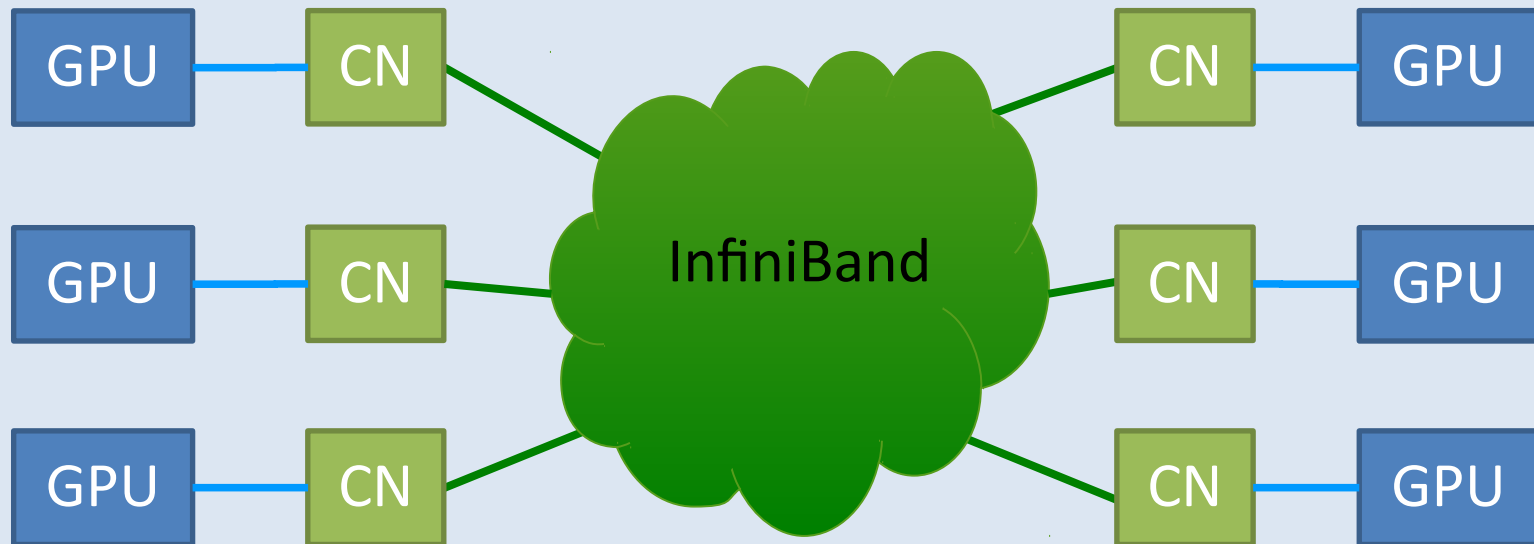
- Clock stagnates since 2002
- Max. clock frequency at about 3 GHz
  - Few exceptions: Power6, Power 7 Gaming
- # transistors still increases



Original data collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond and C. Batten  
Dotted line extrapolations by C. Moore

- Current trends

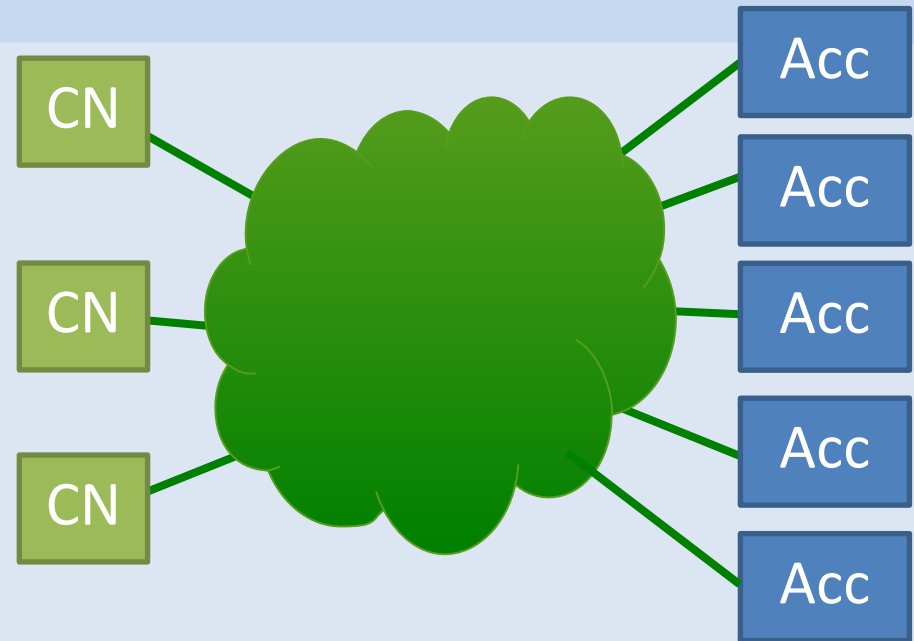
- Multi-Core/Many-Core processors
- Simultaneous Multi Threading (SMT)



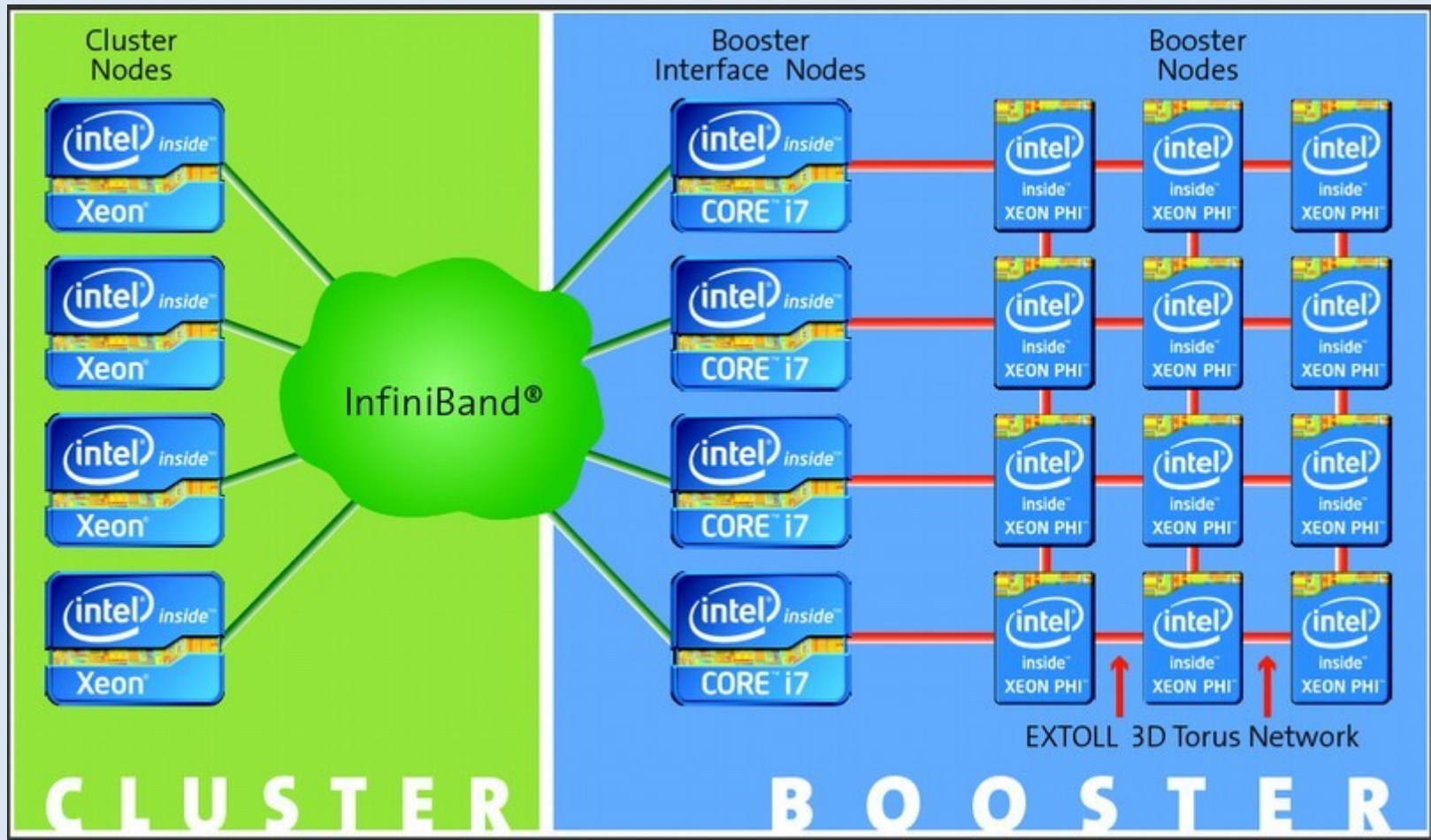
Flat IB-topology  
Simple management of  
resources

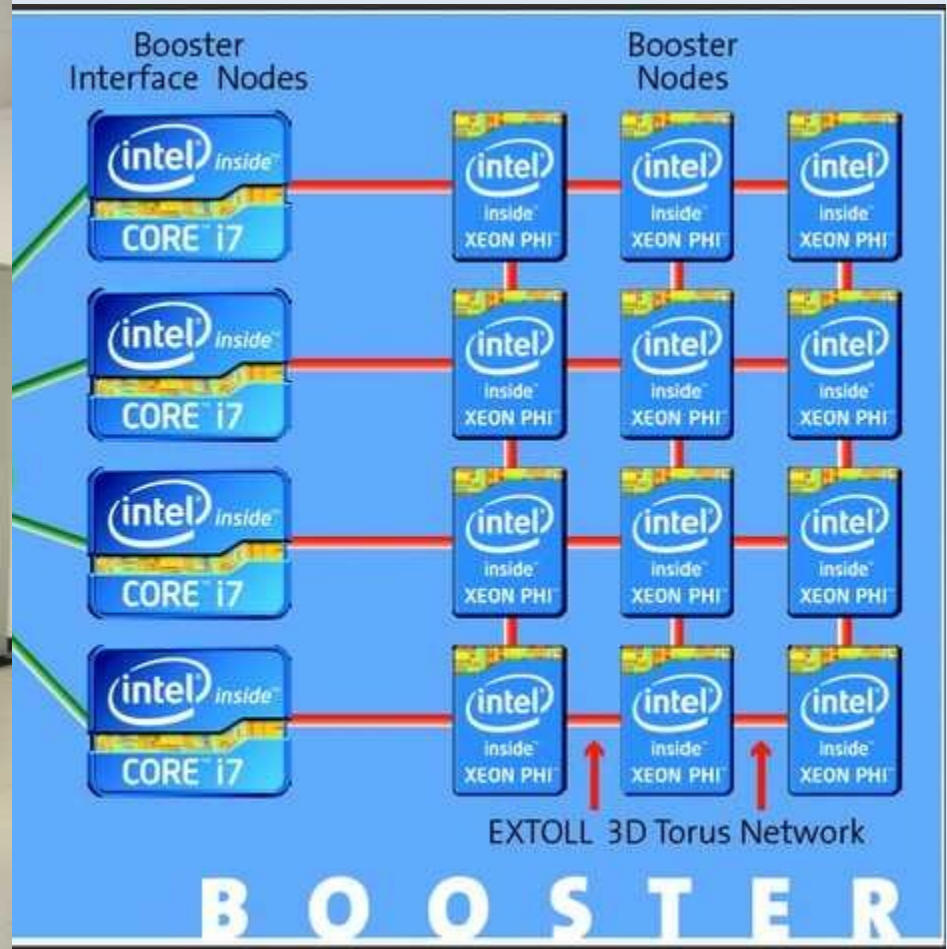
Static assignment of  
CPUs to GPUs  
Accelerators not capable  
to act autonomously

- Go for more capable accelerators (e.g. MIC)
- Attach all nodes to a low-latency fabric
- All nodes might act autonomously
- Dynamical assignment of cluster-nodes and accelerators
  - IB can be assumed as fast as PCIe besides latency
- Ability to off-load more complex (including parallel) kernels
  - communication between CPU and Accelerator less frequently
  - larger messages i.e. less sensitive to latency

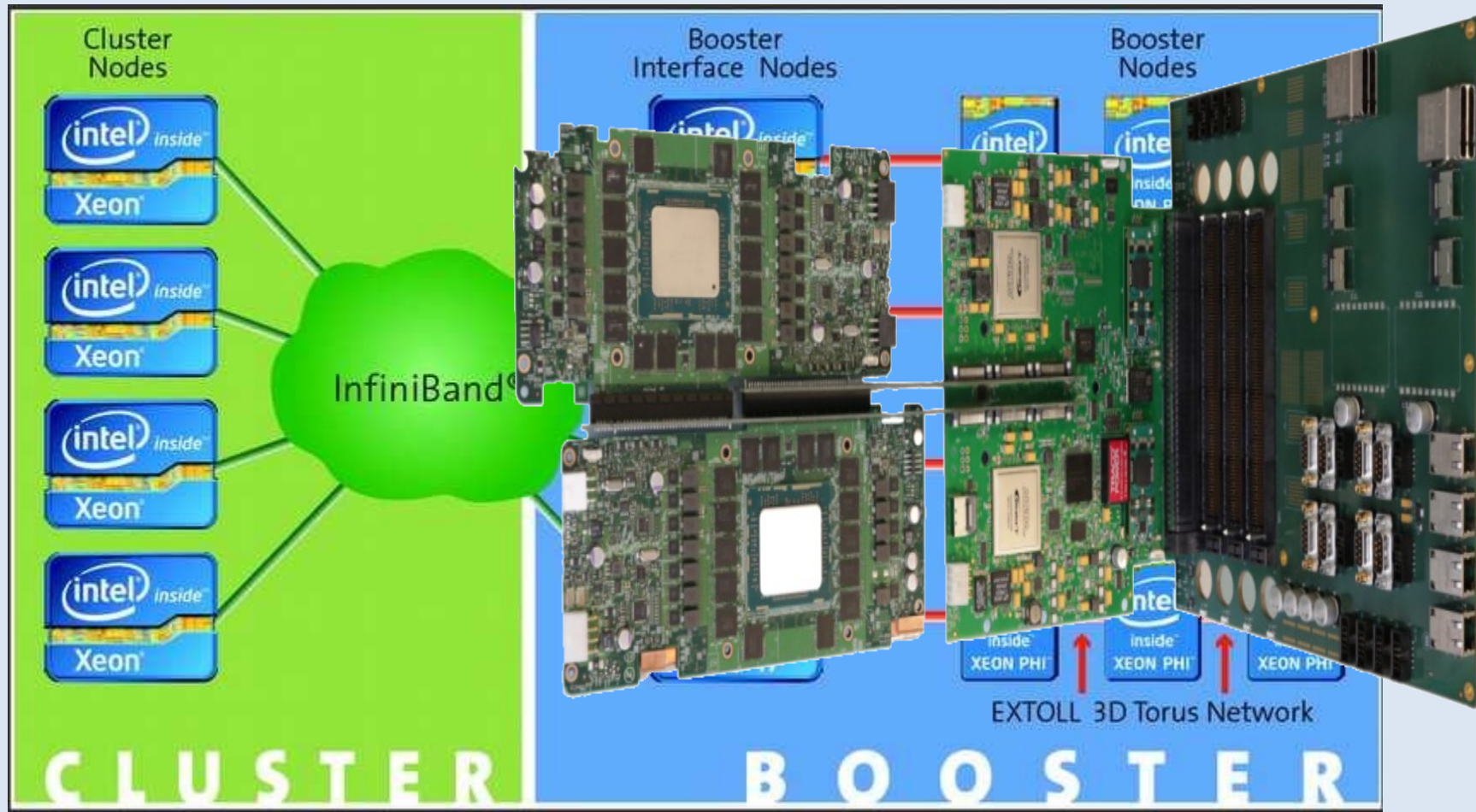


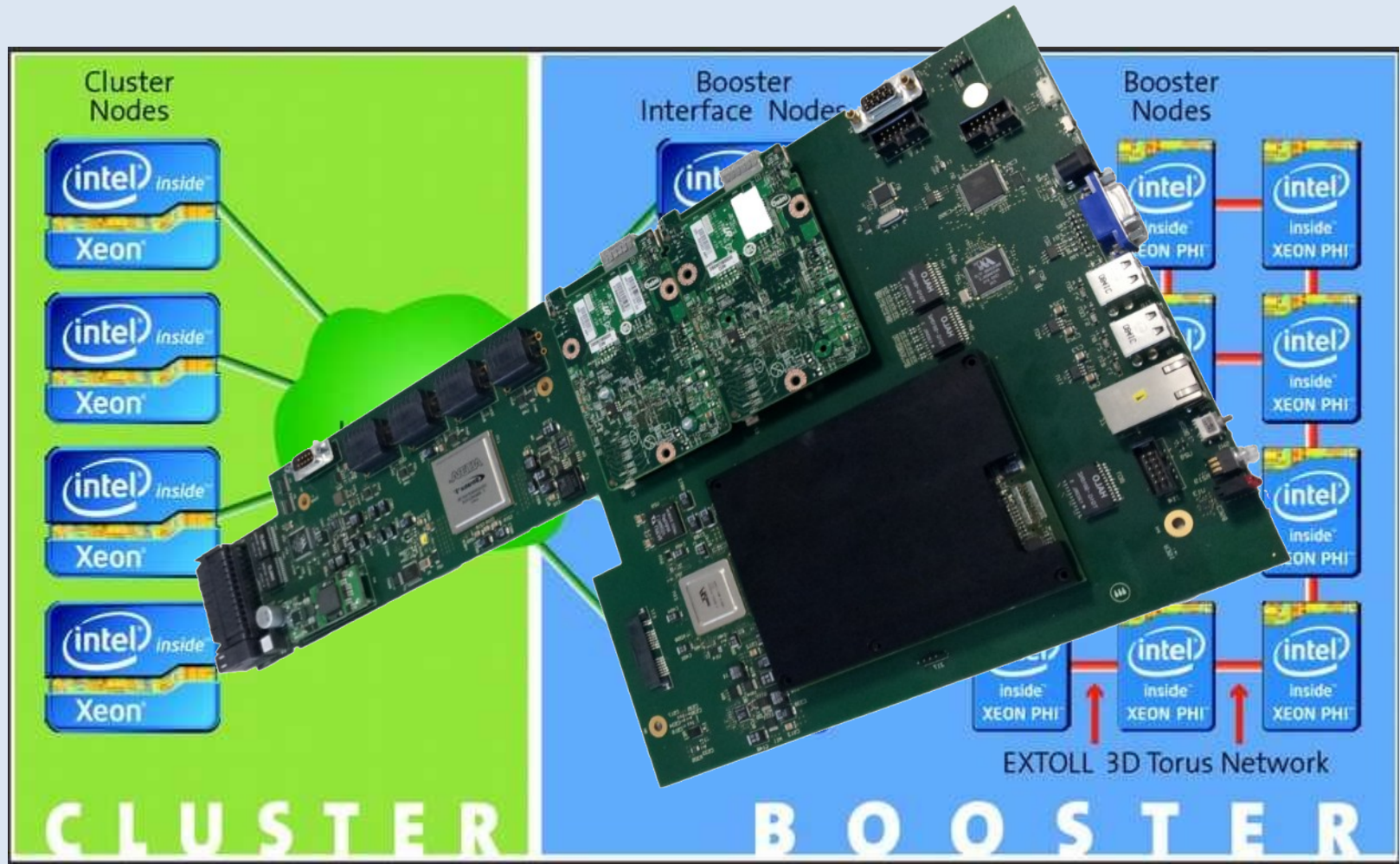




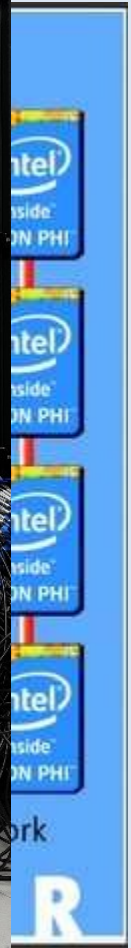
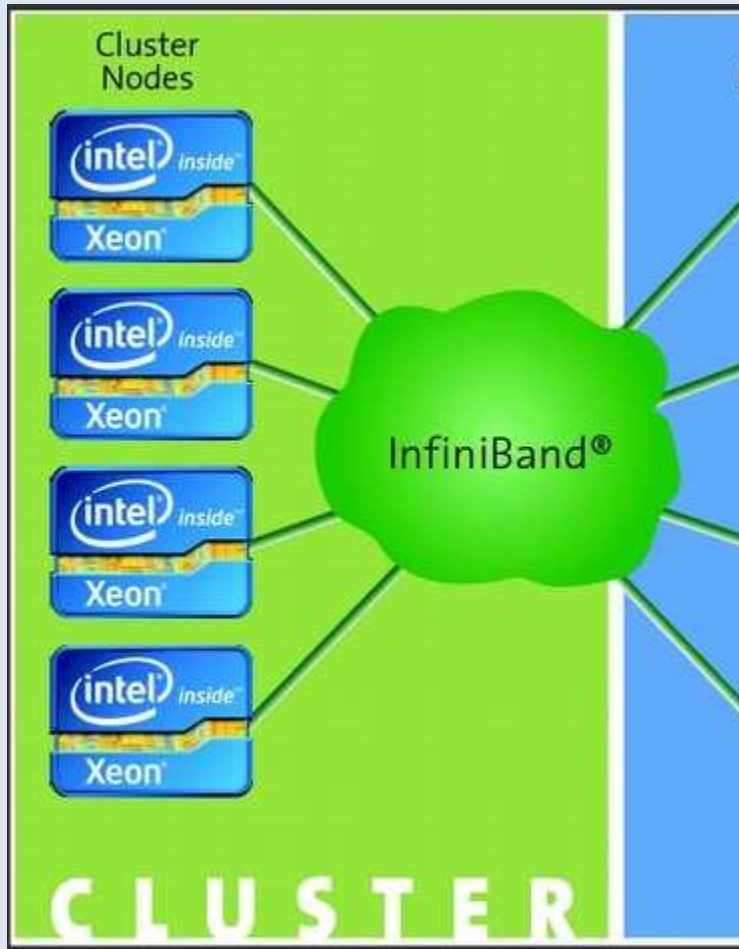




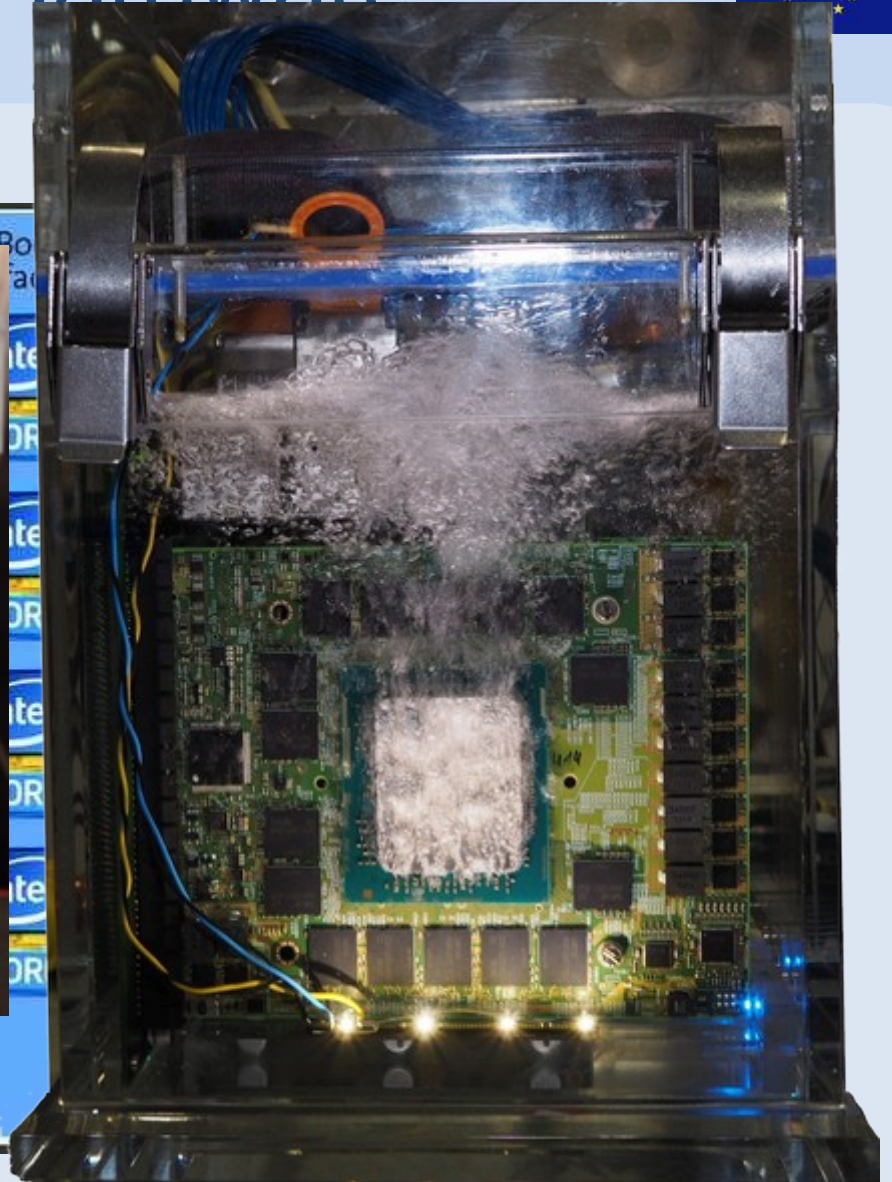




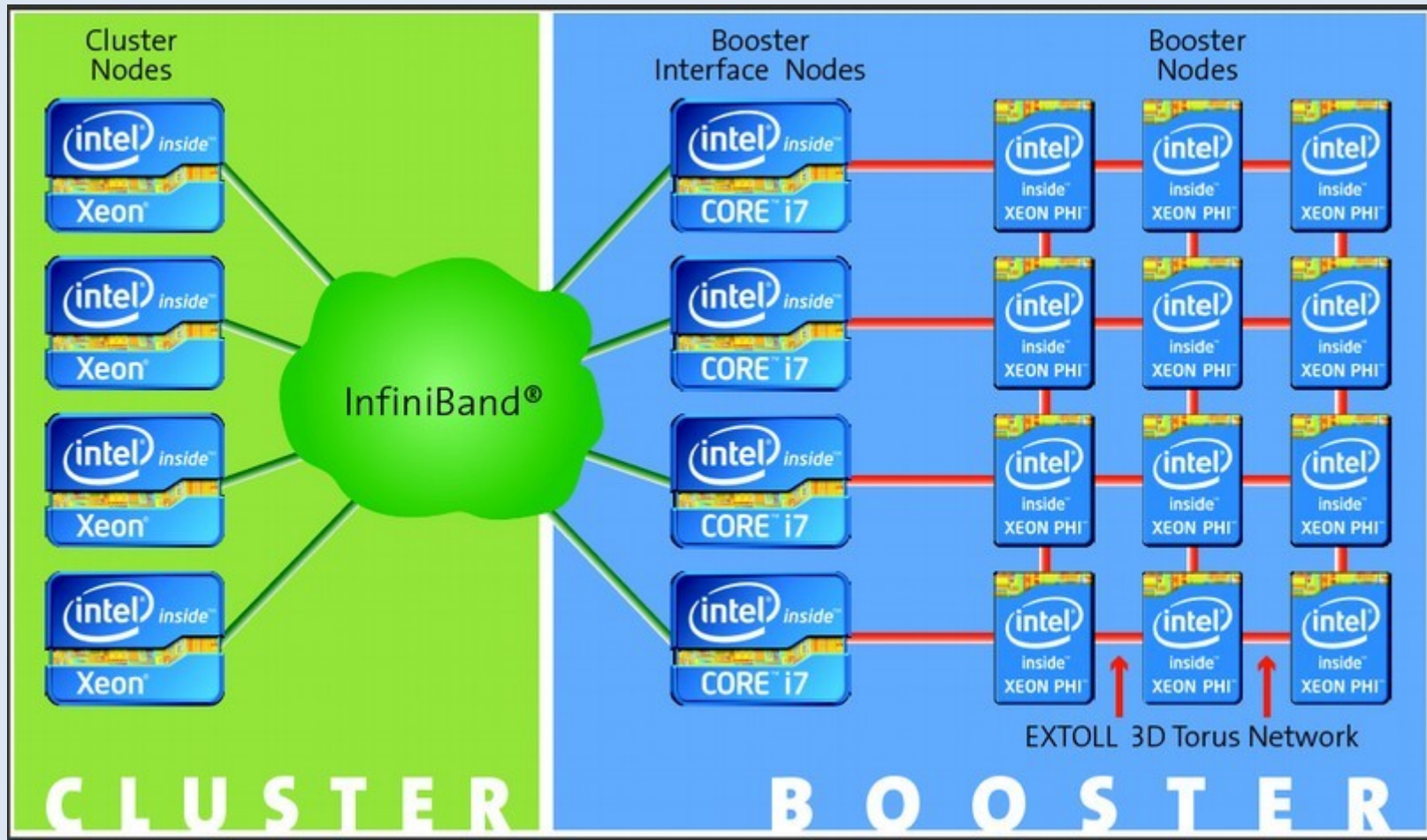


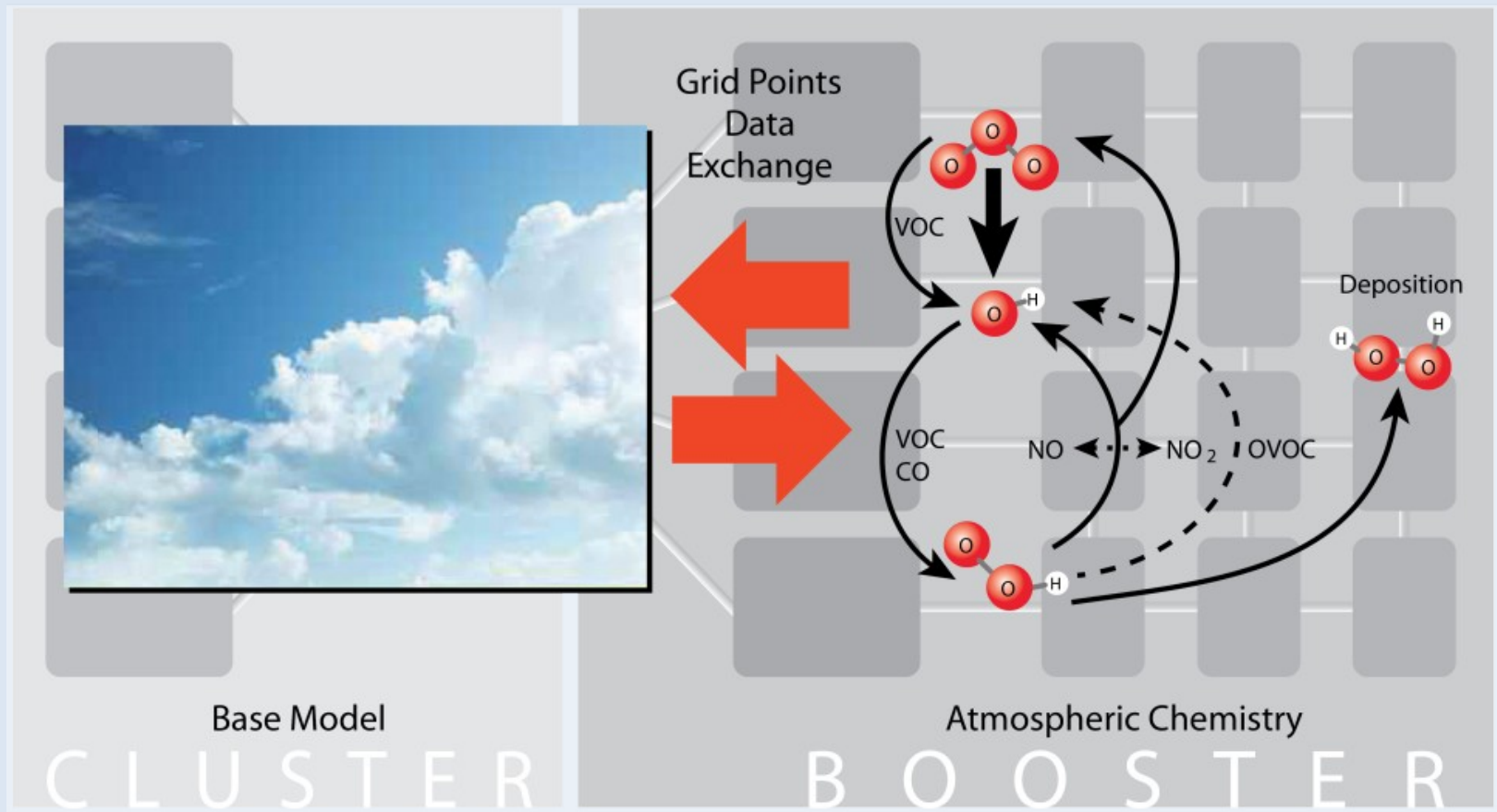




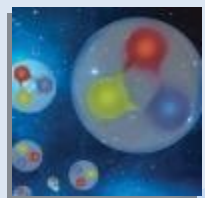
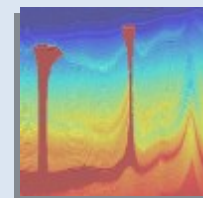
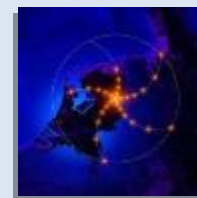
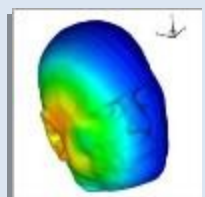
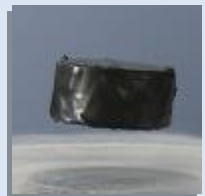
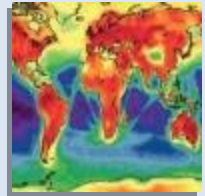
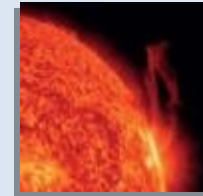
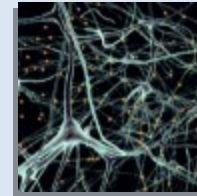




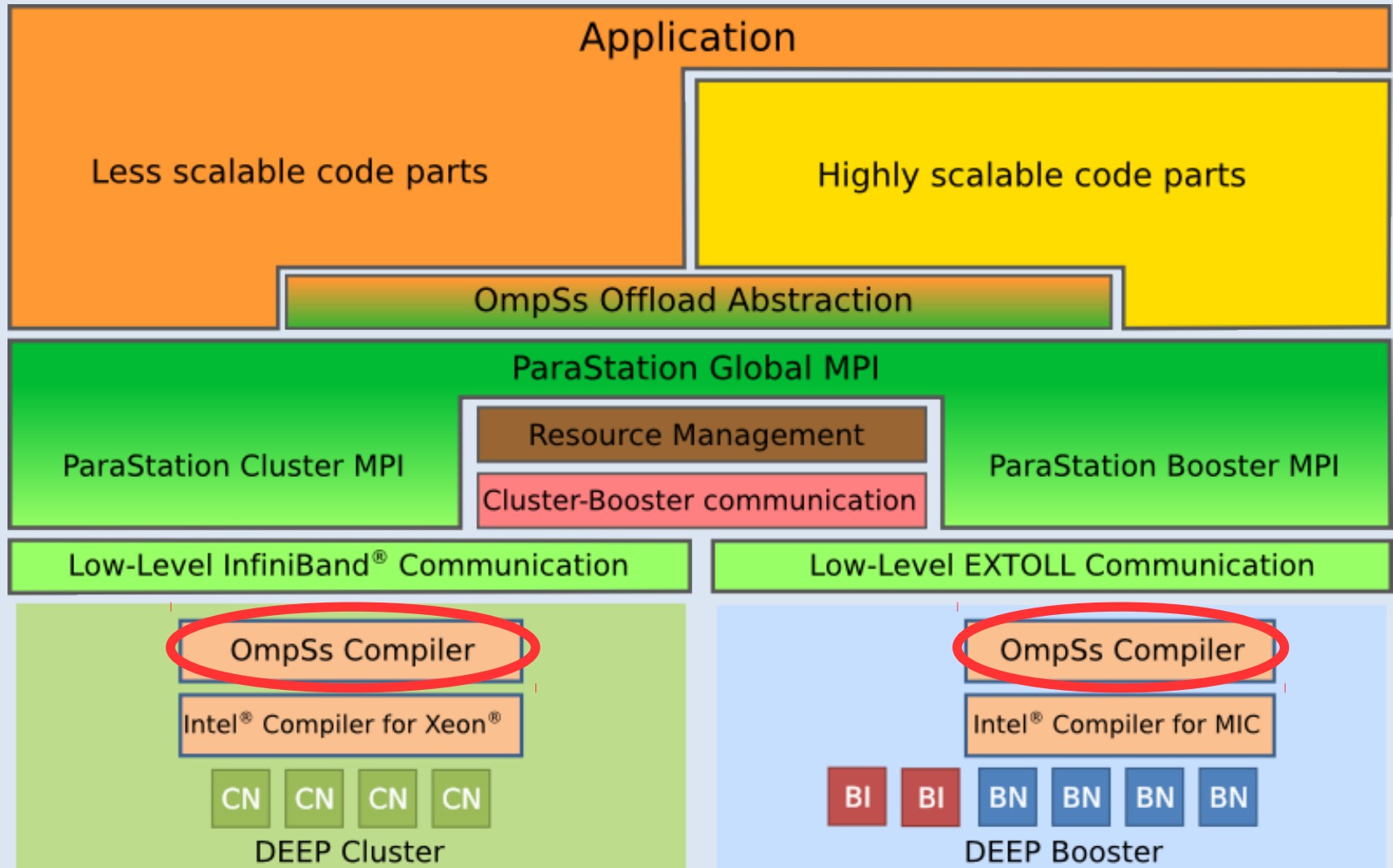




- Brain simulation (EPFL)
- Space weather simulation (KULeuven)
- Climate simulation (CYI)
- Computational fluid engineering (CERFACS)
- High  $T_c$  superconductivity (CINECA)
- Seismic imaging (CGG)
- Human exposure to electromagnetic fields (INRIA)
- Geoscience (BADW-LRZ)
- Radio astronomy (Astron)
- Oil exploration (BSC)
- Lattice QCD (UREG)



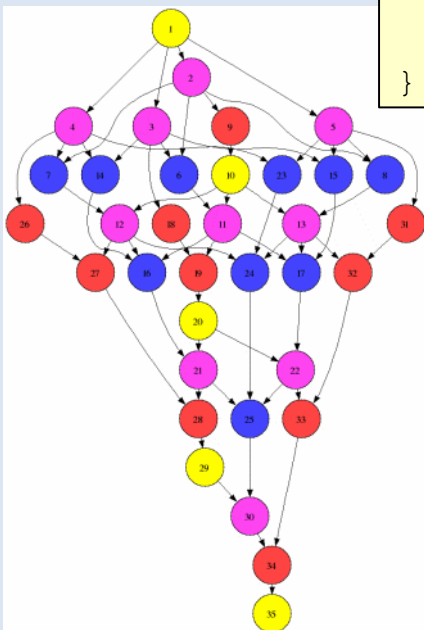
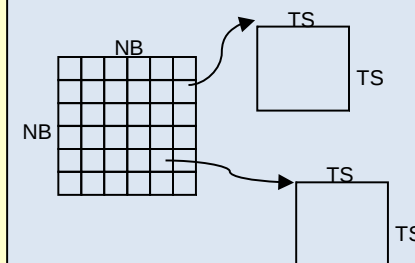
# Programming Paradigm



# OmpSs: tasks, dependencies, heterogeneity



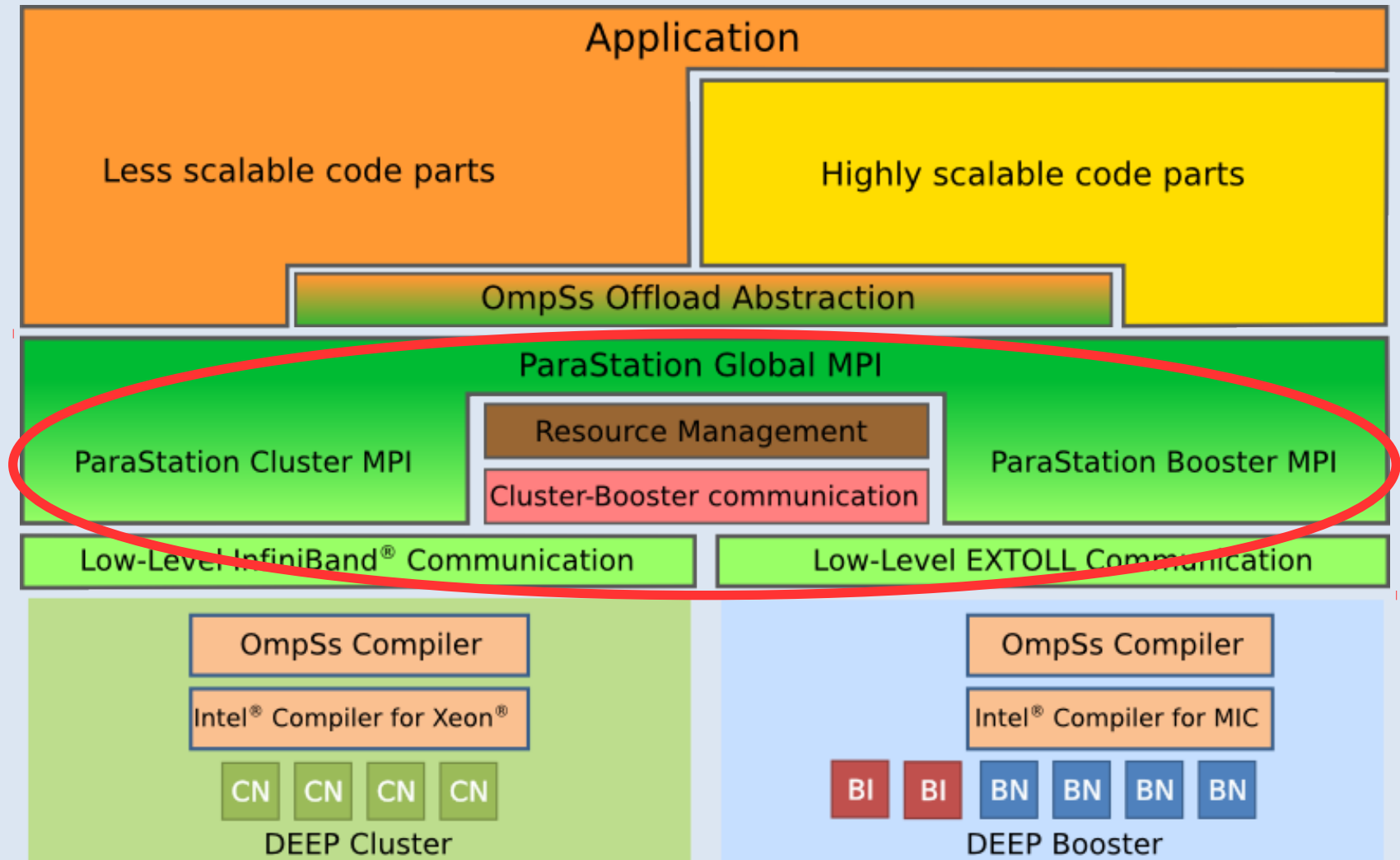
```
void Cholesky( float *A[NT] ) {
  int i, j, k;
  for (k=0; k<NT; k++) {
    spotrf (A[k][k]) ;
    for (i=k+1; i<NT; i++)
      strsm (A[k][k], A[k][i]);
    for (i=k+1; i<NT; i++) {
      for (j=k+1; j<i; j++)
        sgemm( A[k][i], A[k][j], A[j][i]);
      ssyrk (A[k][i], A[i][i]);
    }
  }
}
```

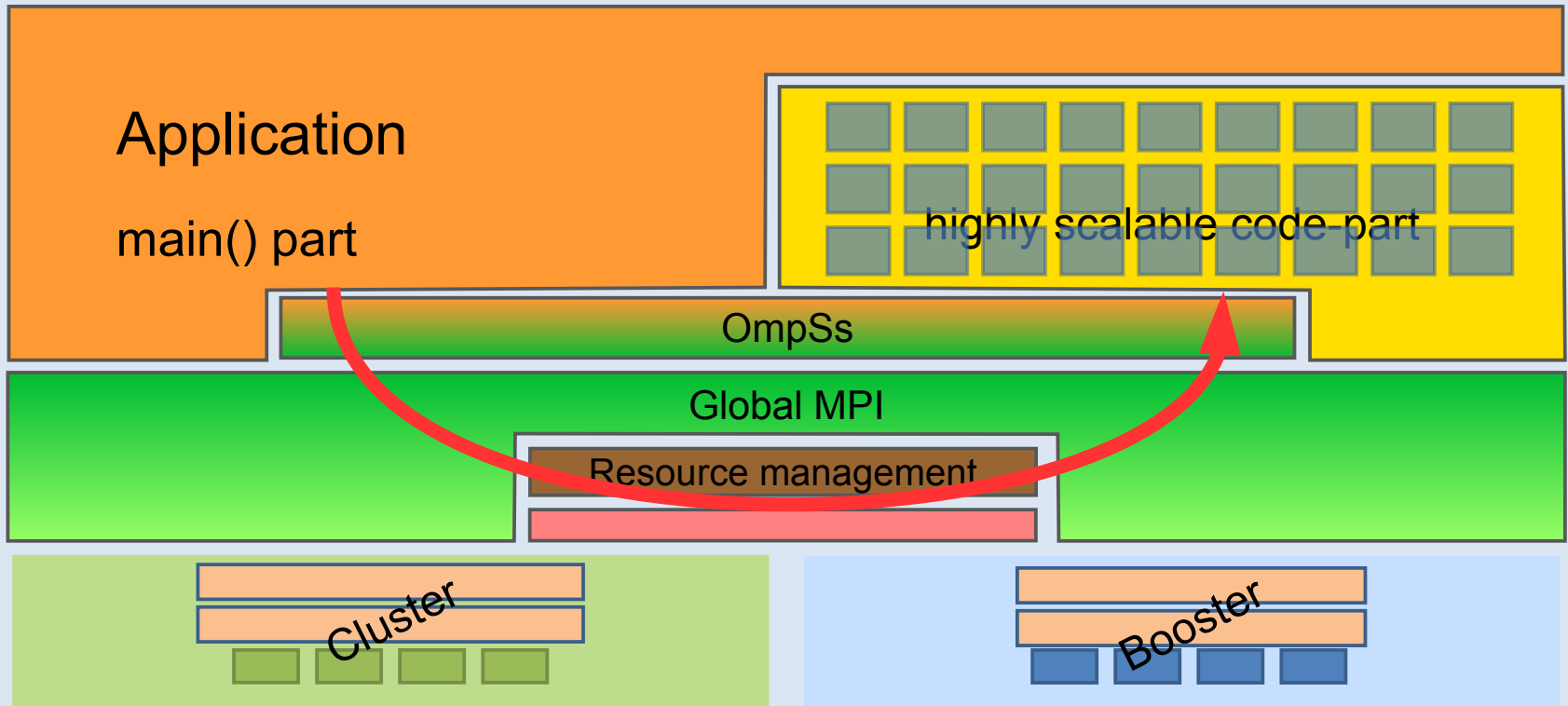


```
#pragma omp task inout ([TS][TS]A)
void spotrf (float *A); ●
#pragma omp task input ([TS][TS]T) inout ([TS][TS]B)
void strsm (float *T, float *B); ●
#pragma omp task input ([TS][TS]A,[TS][TS]B) inout ([TS][TS]C)
void sgemm (float *A, float *B, float *C); ●
#pragma omp task input ([TS][TS]A) inout ([TS][TS]C)
void ssyrk (float *A, float *C); ●
```

Decouple how we write (think sequential) from how it is executed

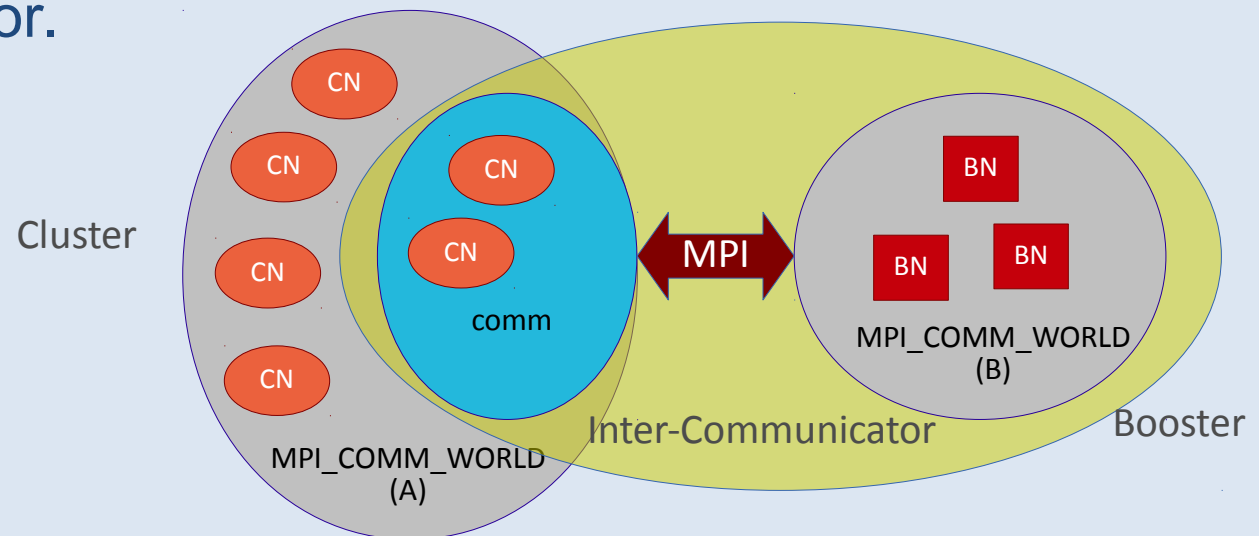


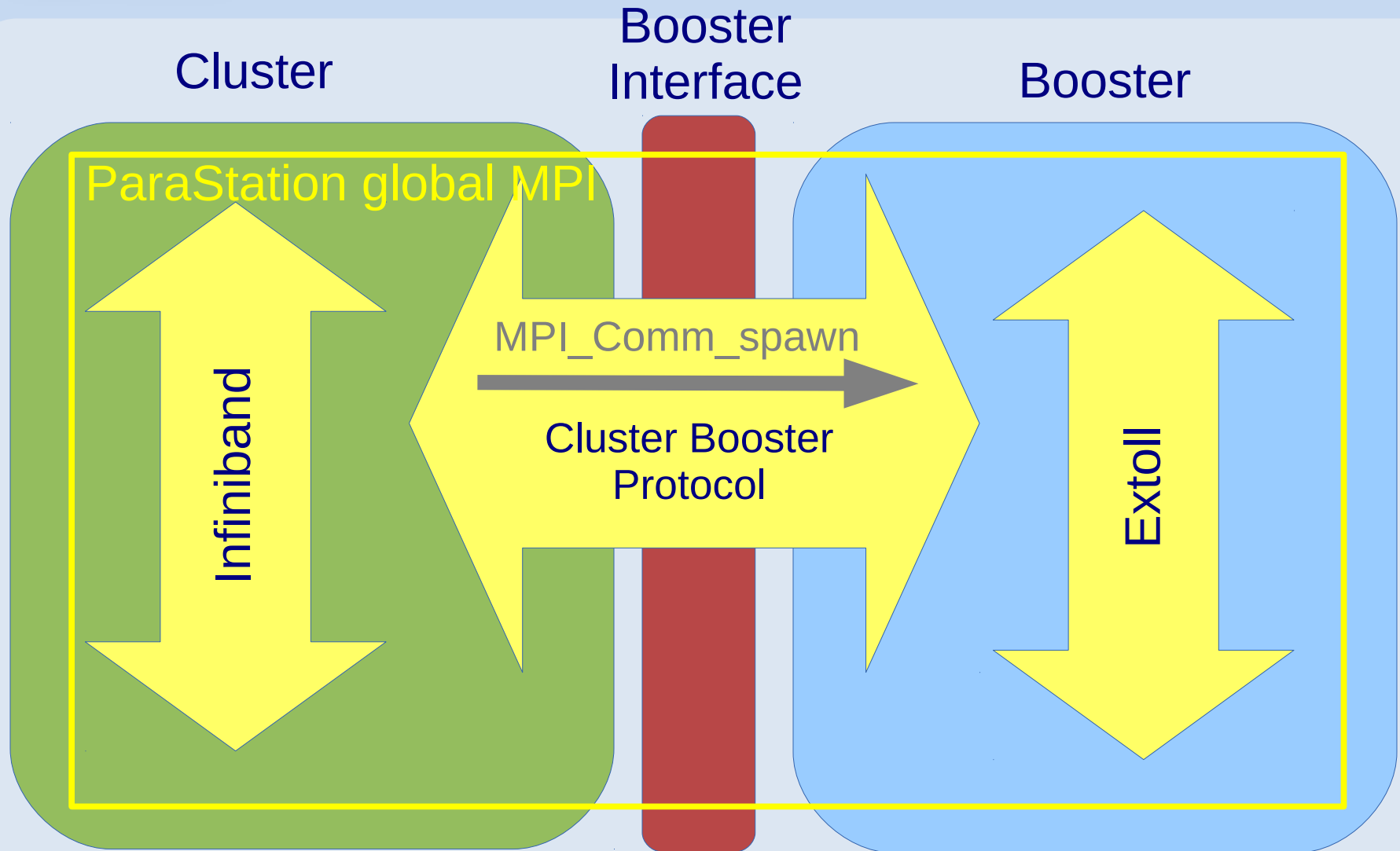


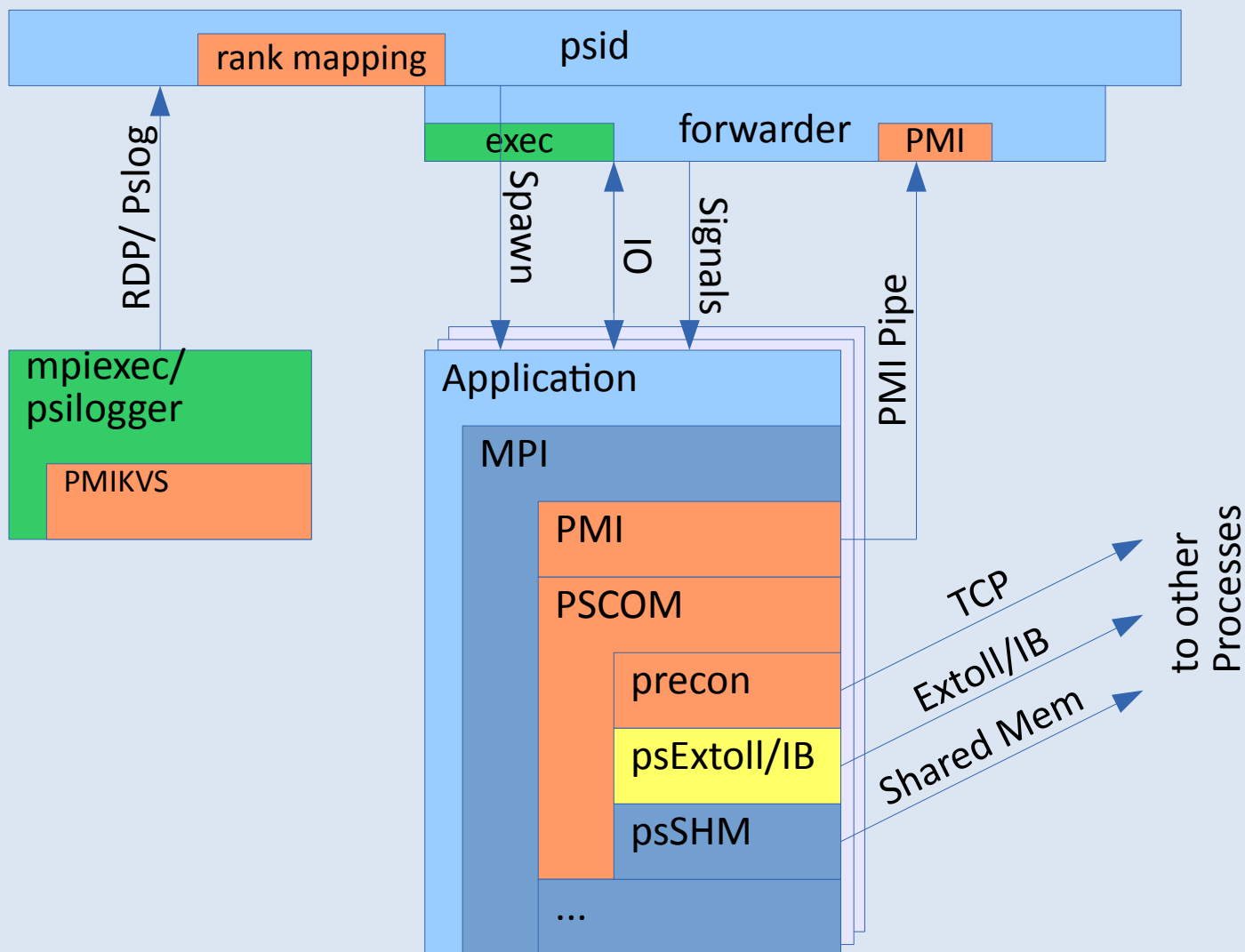


- Application's main()-part runs on Cluster-nodes (CN) only
- Actual spawn done via global MPI
- OmpSs acts as an abstraction layer
- Spawn is a collective operation of Cluster-processes
- Highly scalable code-parts (HSCP) utilize multiple Booster-nodes (BN)

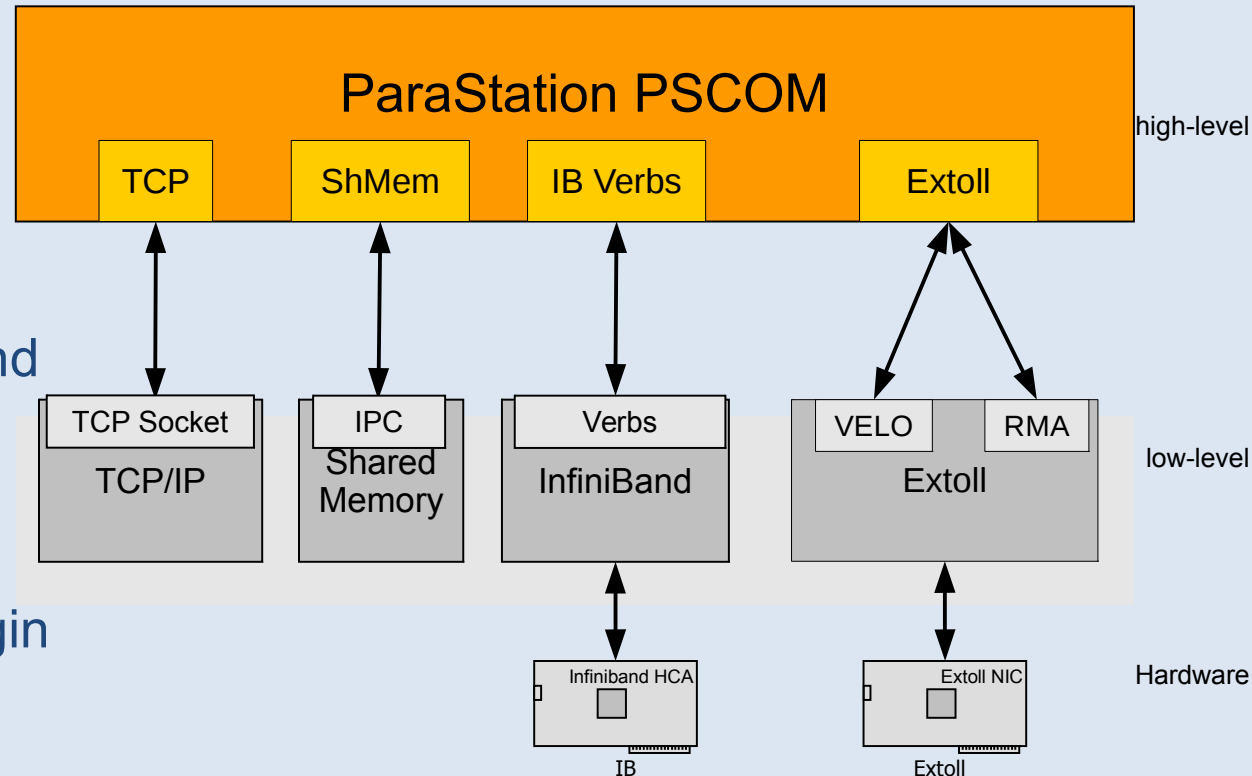
- The inter-communicator contains all parents on the one side and all children on the other side.
  - Returned by `MPI_Comm_spawn` for the parents
  - Returned by `MPI_Get_parent` by the children
- Rank numbers are the same as in the the corresponding intra-communicator.



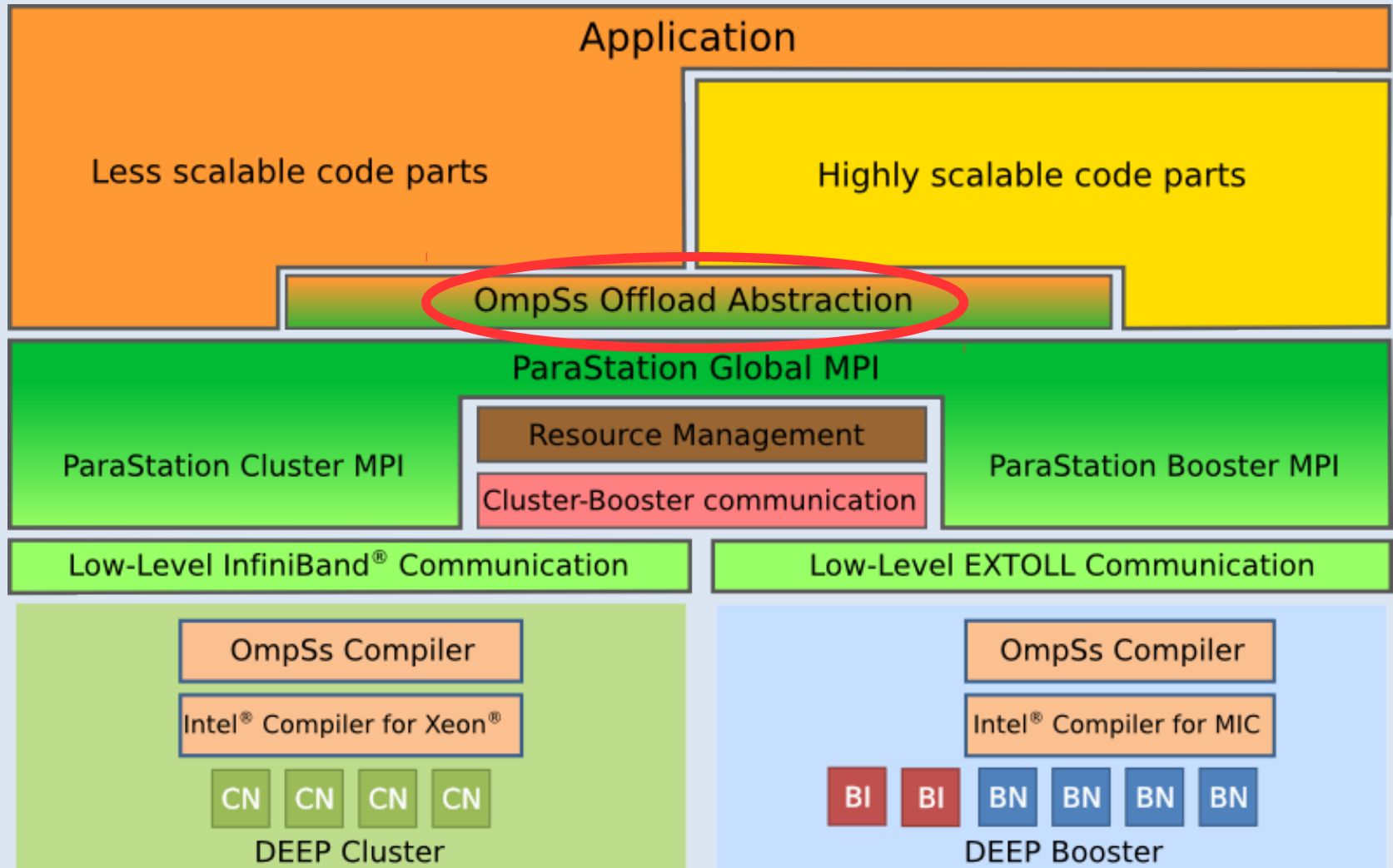




- Unified comm layer
- pscom plugins
  - Modular
  - Flexible to extend
  - Verbs plugin
  - VELO/RMA plugin
  - Easy enabling of Cluster-Booster protocol







Source Code

```
int main(int argc, char *argv[]){
    /*...*/
    for(int i=0; i<3; i++){
        #pragma target device (comm:size*rank+i) copy_deps
        #pragma omp task input(...) output(...)
        foo_mpi(i, ...);}
}
```

Compiler

OmpSs Compiler

Application  
Binaries

Cluster  
Executable

Booster  
Executable

DEEP Runtime

Cluster MPI

ParaStation Global MPI

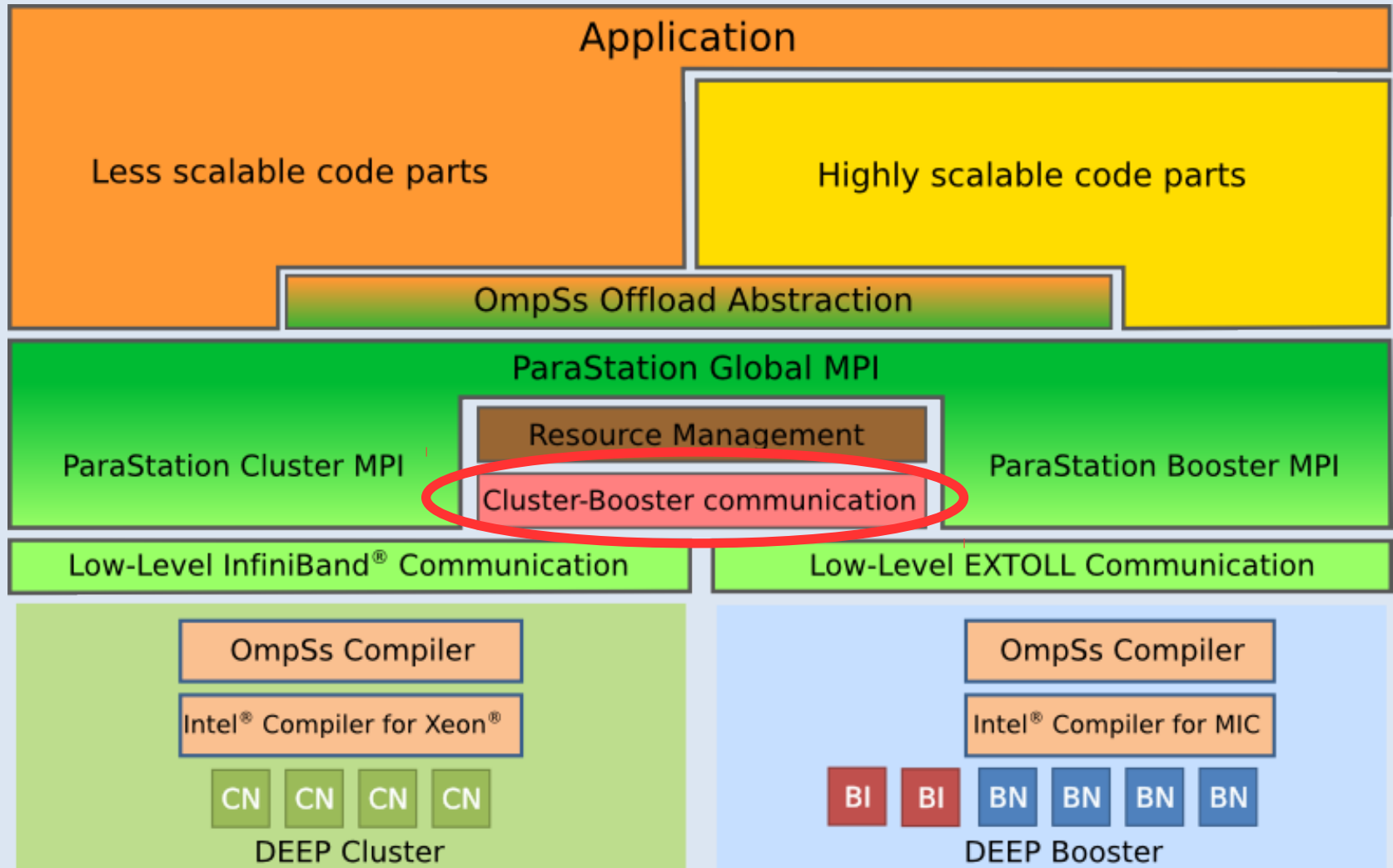
**DEEP Runtime**

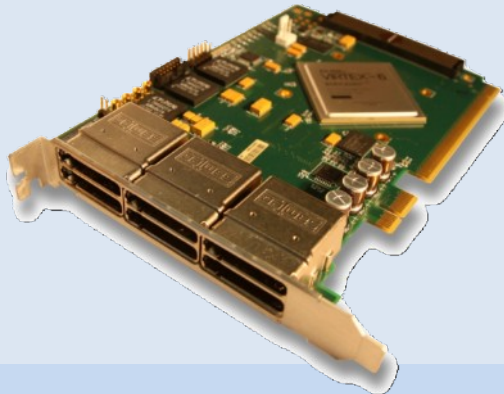
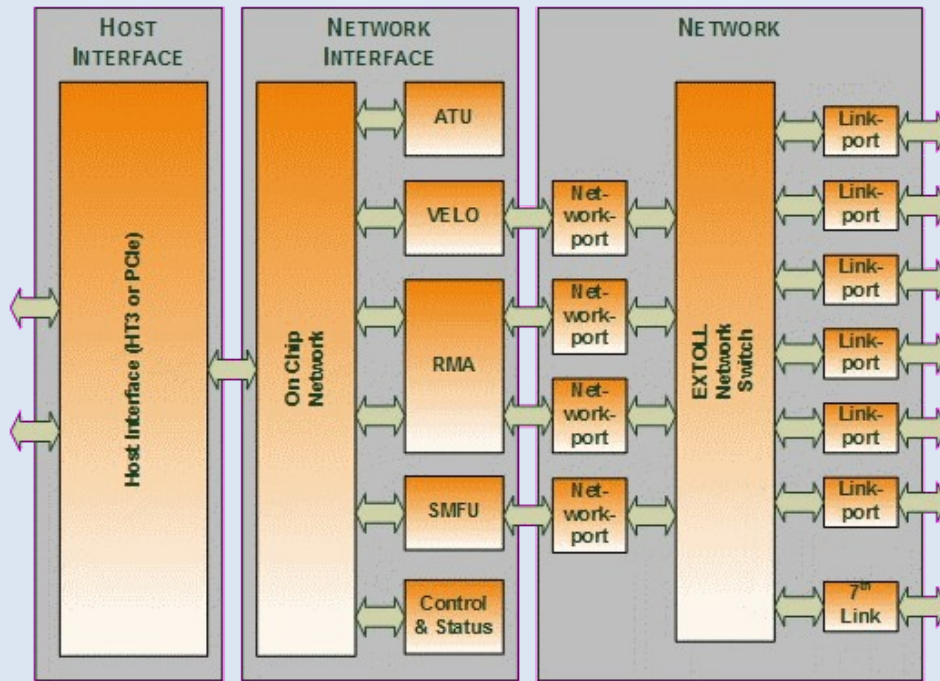
Booster MPI

OmpSs Runtime

CLUSTER

BOOSTER



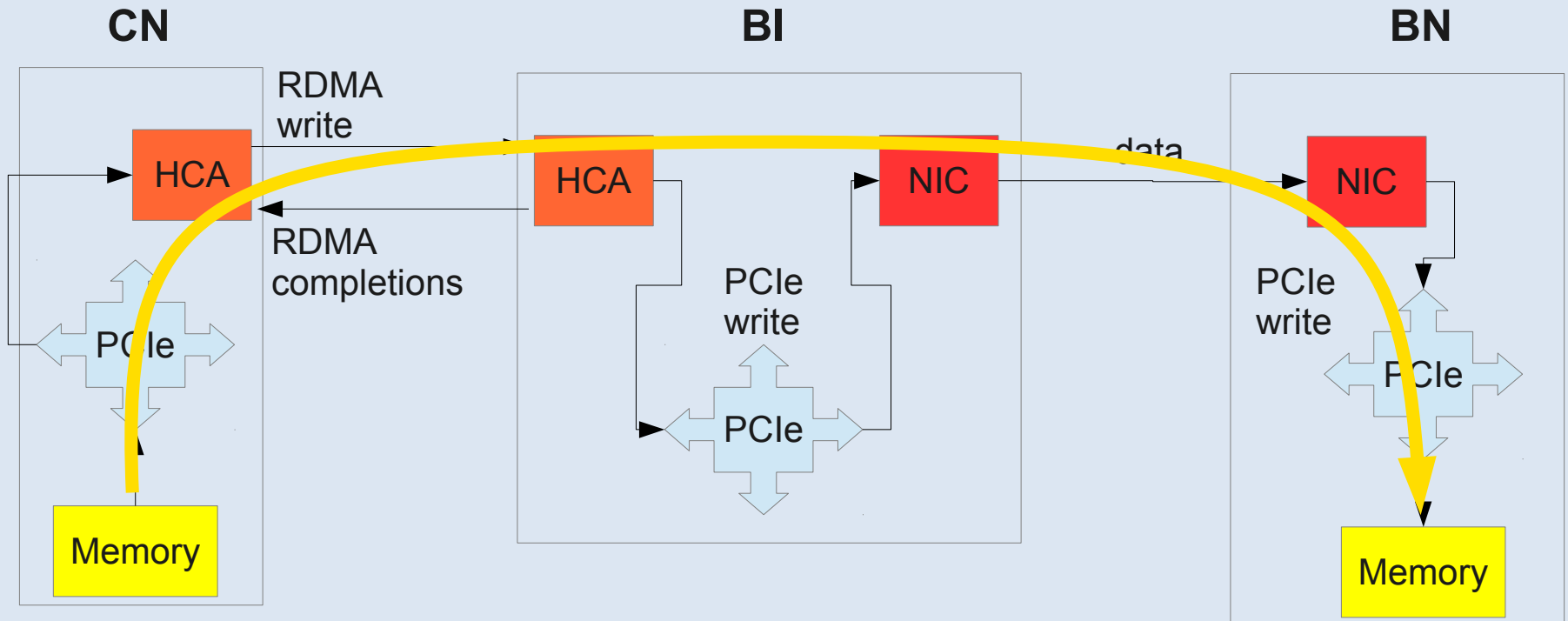


## Relevant features for DEEP

- Low latency, high bandwidth
- RMA engine for remote memory access, bulk data transfer
- VELO communication engine (zero-copy MPI)
- SMFU engine for bridging to InfiniBand
- 6 links for 3D torus topology
- 7<sup>th</sup> link for general devices
- Built-in PCIe root-port
- RAS features: CRC/ECC protection, link level retransmission
- Many status & control registers
- Access from host, via I2C bus or over EXTOLL

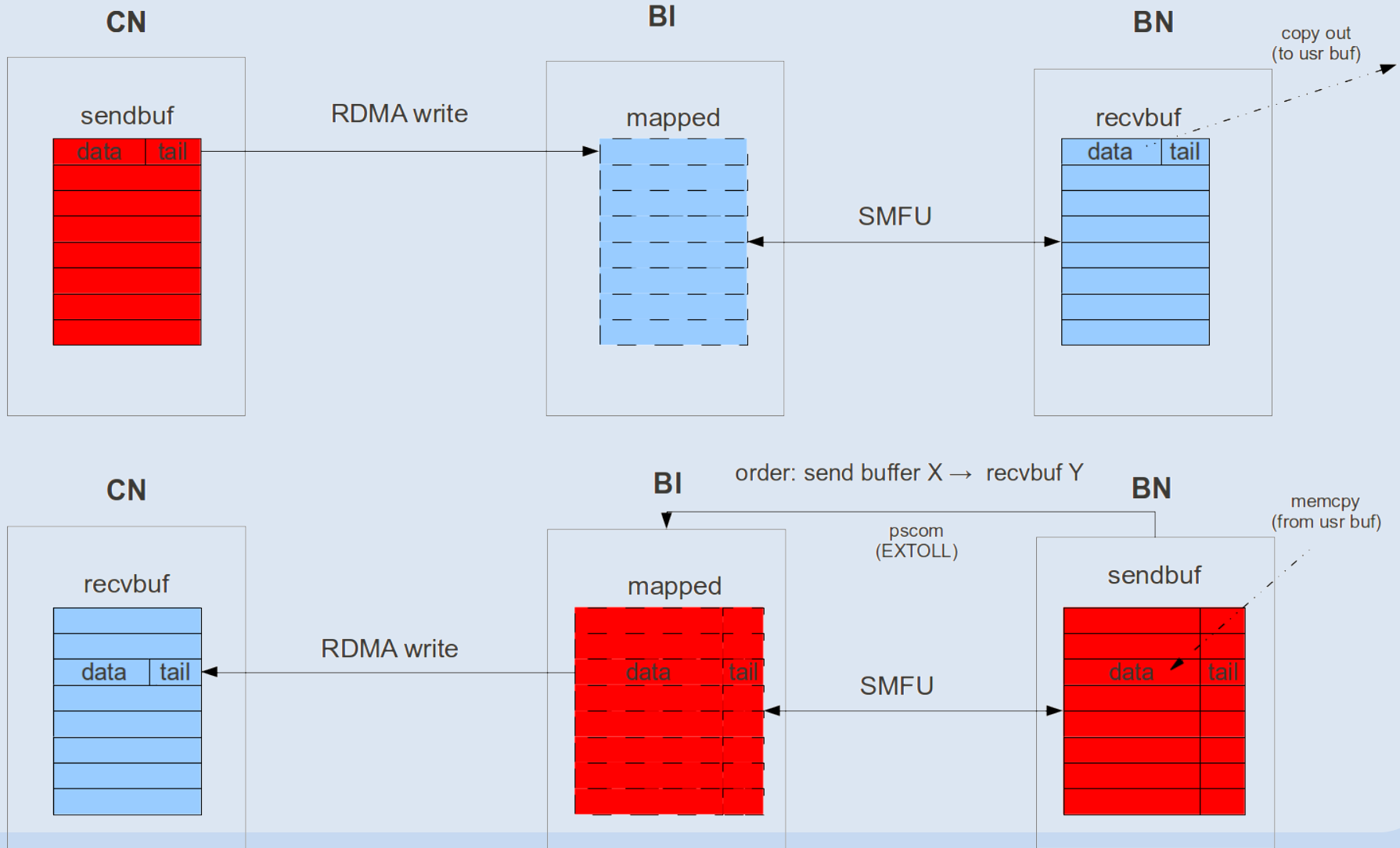


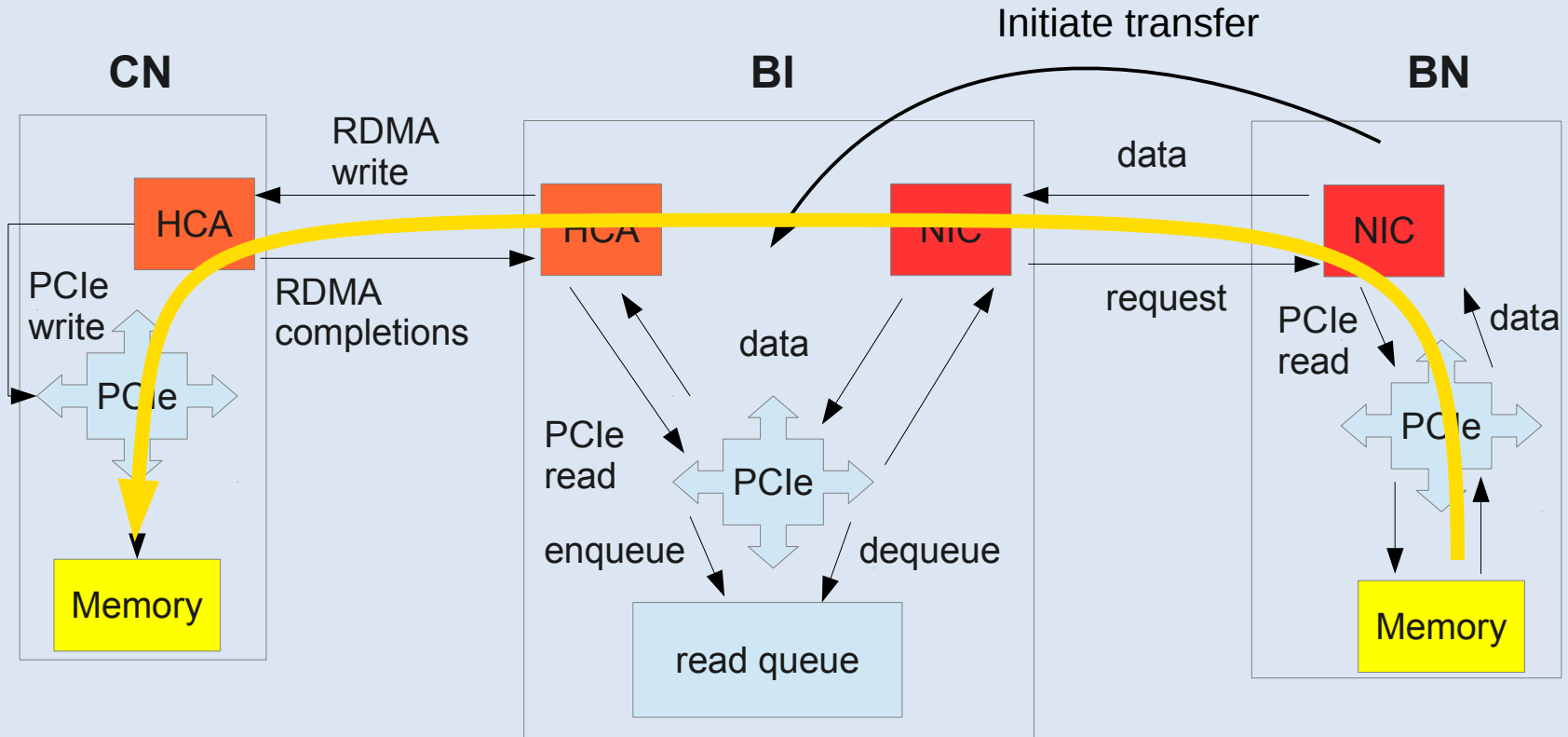
- High-bandwidth, low latency communication between Cluster Nodes and Booster Nodes needed
- Booster Interfaces bridge between InfiniBand and EXTOLL
- A Booster Interface can act in two different modes:
  - Bridge: receive packets, actively forward to end-point
  - Set-up and enable RDMA via SMFU directly between memory locations on CN and BN



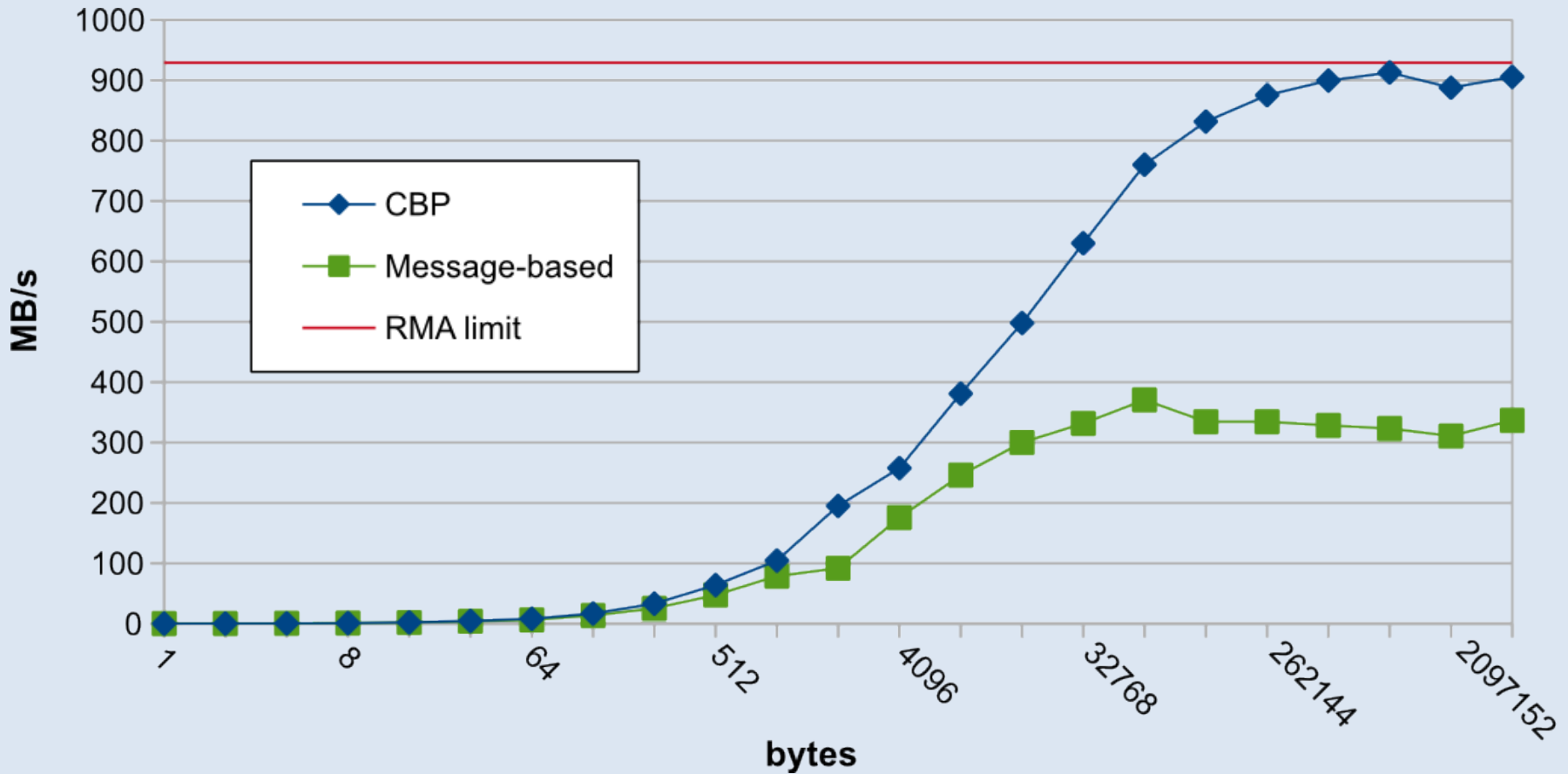
The RDMA write operation *initiated on the Cluster Node* is “forwarded” via PCIe and EXTOLL to the Booster Node

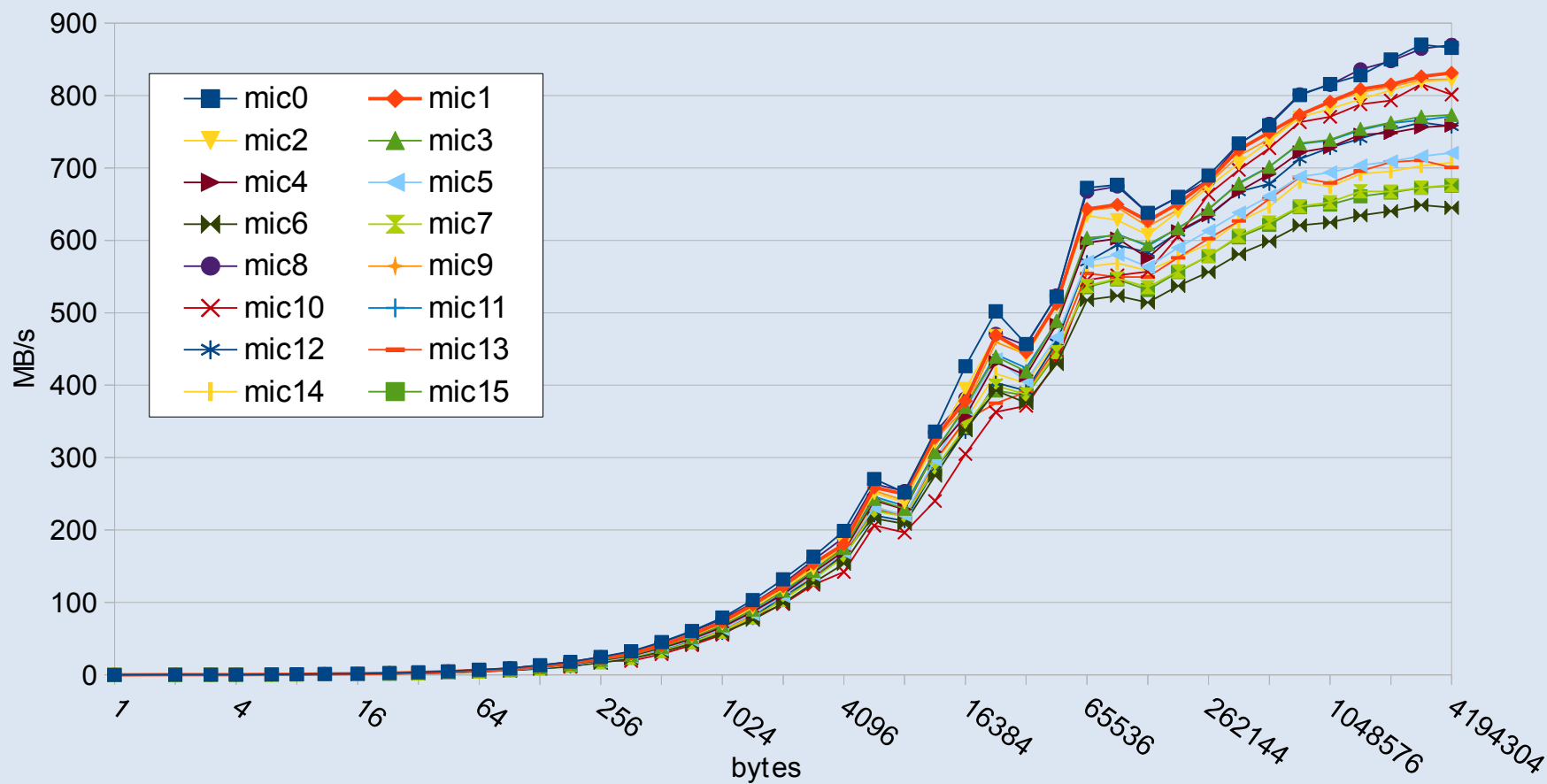


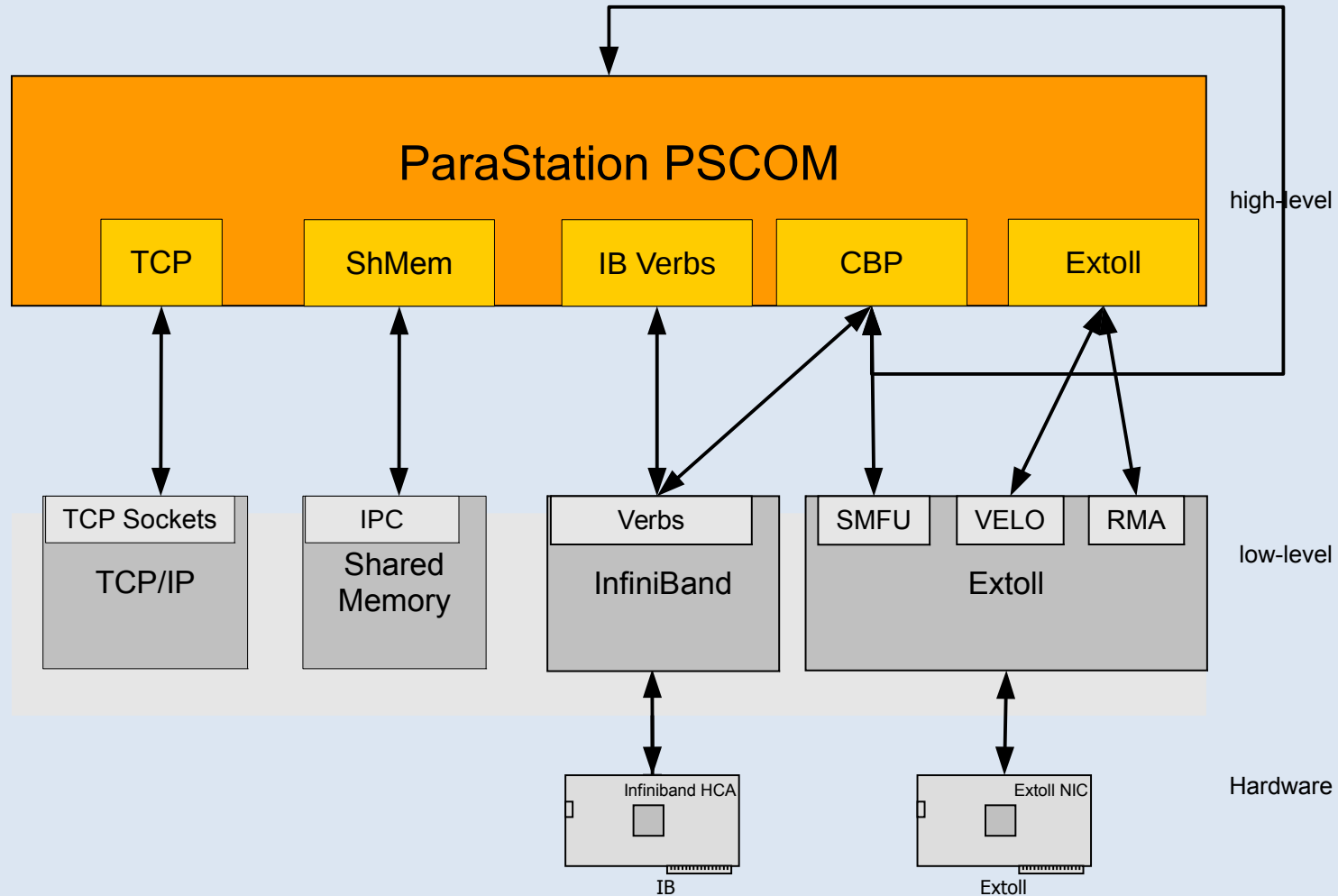




- The *initiating Booster Node* asks the Booster Interface's HCA to do an RDMA write operation to the Cluster Node's Memory
- The resulting PCIe read operation is “forwarded” by Extoll

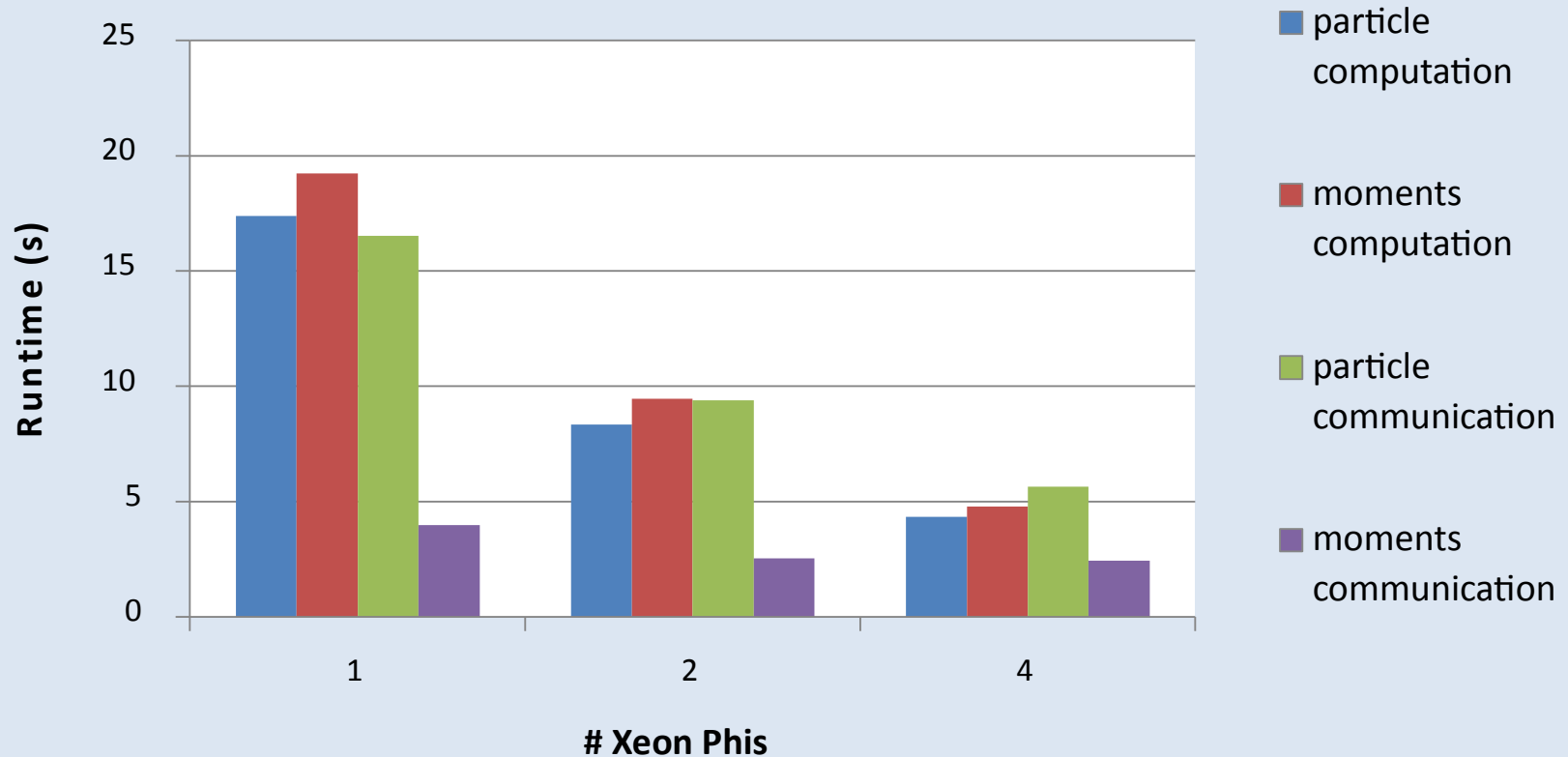






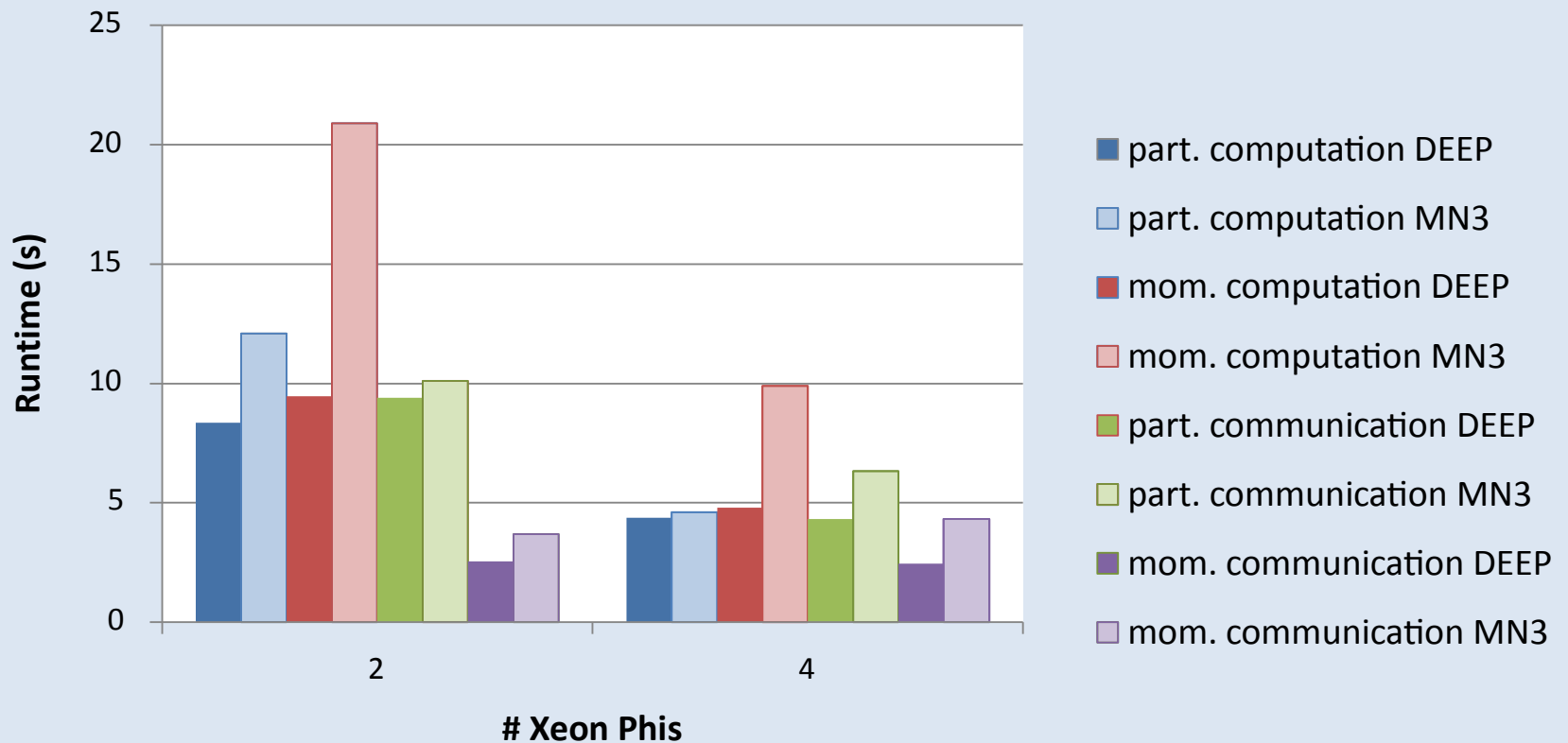
# (Few) Results

## Runtime of different phases on the DEEP Booster



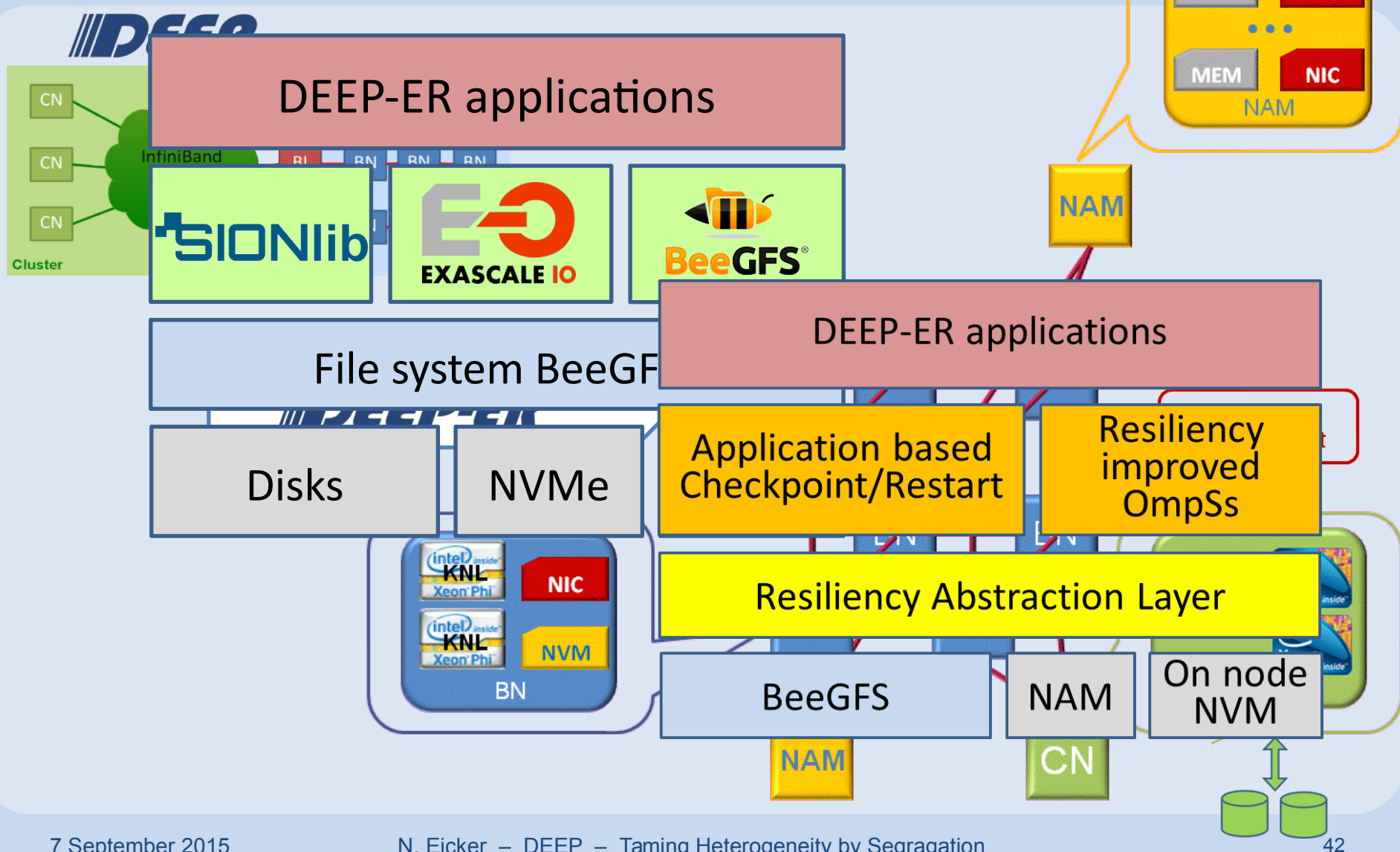


## Phases comparison DEEP BNs and MN3 XeonPhis



- 14 Partners
  - 4 PRACE hosts
  - 4 Research Centers
  - 4 Industry Partners
  - 4 Universities
  - Coordinator JSC
- 7 Countries
- Duration: 3 years
- Budget: 10.0 M€
- EU funding: 6.4 M€





- DEEP and DEEP-ER explore new ways to use and manage heterogeneity
  - Cluster Booster Architecture and it's implementation
  - Programming Model
- DEEP's 384-KNC Booster is up and running – mostly
- As second 32-KNC Booster utilizing immersion-cooling is on the way
- Both are application driven – co-design is important
- Try to hide the details via OmpSs' abstraction layer
- DEEP-ER extends the concept to I/O, Resiliency and innovative memory technologies like NVMe and NAM
- Jülich plans to extend its new Cluster JURECA by a 10 PF/s Booster in 2016 (based on KNL)
- More info: <http://www.deep-project.eu>  
<http://www.deep-er.eu>