

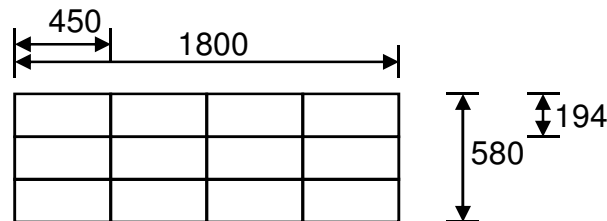
Topology aware Cartesian grid mapping with MPI

Christoph Niethammer
Rolf Rabenseifner
HLRS

Examples

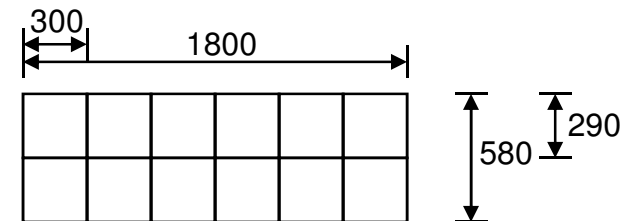
- Application topology awareness
 - 2-D example with 12 MPI processes and gridsize 1800x580

• `MPI_Dims_create` → 4x3 ☹️



Boundary of a subdomain = $2(450+194) = 1288$ ☹️

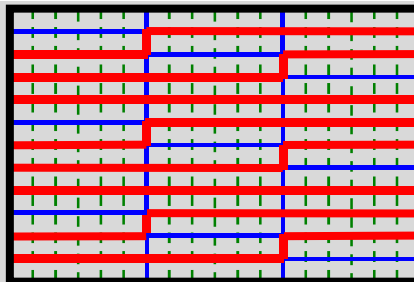
• `grid aware` → 6x2 😊 processes



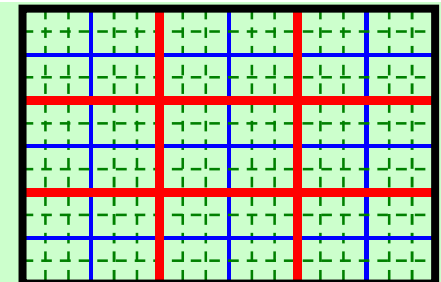
Boundary of a subdomain = $2(300+290) = 1180$ 😊

- Hardware topology awareness
 - 2-D example with 9 nodes x 4 CPUs x 6 cores → 12*18 Cartesian processes

Without re-numbering:
150 node-to-node ☹️
72 CPU-to-CPU
180 core-to-core



With new ranks:
60 node-to-node 😊
90 CPU-to-CPU
252 core-to-core



- `MPI_Dims_create` is limited, e.g., 3-D on 625 nodes x 2 CPUs x 12 cores
→ 25 x 25 x 24 ☹️ processes

→ We do

- Node level: $625 = 5 \times 25 \times 5$
- CPU level: $2 = 2 \times 1 \times 1$
- Core level: $12 = 3 \times 1 \times 4$

Result (product): $30 \times 25 \times 20$ 😊



Topology aware MPI process grid mapping 2018
Slide 0 Niethammer, Rabenseifner

Slides for the poster at EuroMPI 2018

Further info, see <https://fs.hlrs.de/projects/par/mpi/EuroMPI2018-Cartesian/>

Duplex accumulated ring bandwidth per node

(each message is counted twice, as outgoing and incoming)

