

# Topology aware Cartesian grid mapping with MPI

Christoph Niethammer    Rolf Rabenseifner

niethammer@hirs.de

rabenseifner@hirs.de

High Performance Computing Center (HLRS), University of Stuttgart, Germany

**HLRS, Stuttgart, Sep 25, 2018**

Slides for the poster at EuroMPI 2018

For further information, see <https://fs.hirs.de/projects/par/mpi/EuroMPI2018-Cartesian/>

2018

Höchstleistungsrechenzentrum Stuttgart



# The Problems of MPI\_Dims\_create + MPI\_Cart\_create

- The factorization of a given amount of MPI processes must be
  - Application topology aware [1]
  - Hardware topology aware
- Current definition of MPI\_Dims\_create is not prepared for this
- Extreme differences in latency and accumulated bandwidth between **inter**-node and **intra**-node communication
- The reordering by MPI\_Cart\_create:
  - Many implementations do nothing
  - A perfect reordering may require complex domain decomposition algorithms (e.g. Metis) [3]

Slides 1+2

Slides 3-5

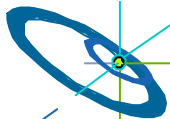
Slide 6

We propose a new and fast algorithm, which is application and hardware topology aware

Slides 7-17

[1, 3] see References on last slide

Slide 19

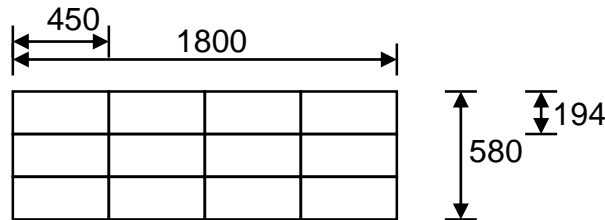


# Examples

- Application topology awareness

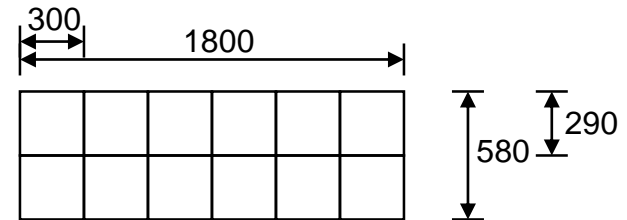
- 2-D example with 12 MPI processes and gridsize 1800x580

- **MPI\_Dims\_create** → 4x3



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

- **grid aware** → 6x2 processes

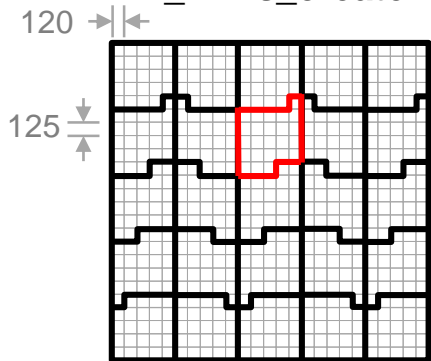


Boundary of a subdomain =  $2(300+290) = 1180$  😊

- Hardware topology awareness

- 2-D example with 25 nodes x 24 cores and gridsize 3000x3000

- **MPI\_Dims\_create** → 25 x 24

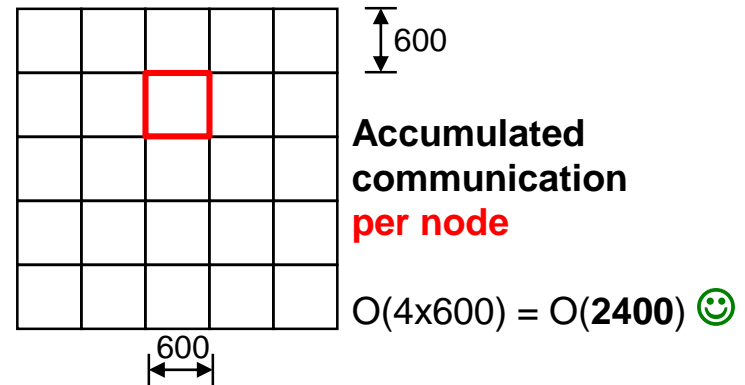


**Accumulated communication per node**

$O(10 \times 120 + 12 \times 125)$   
 $= O(2700)$  ☹️

- **Hardware aware**

→ (5 nodes x 6 cores) X (5 nodes x 4 cores)



**Accumulated communication per node**

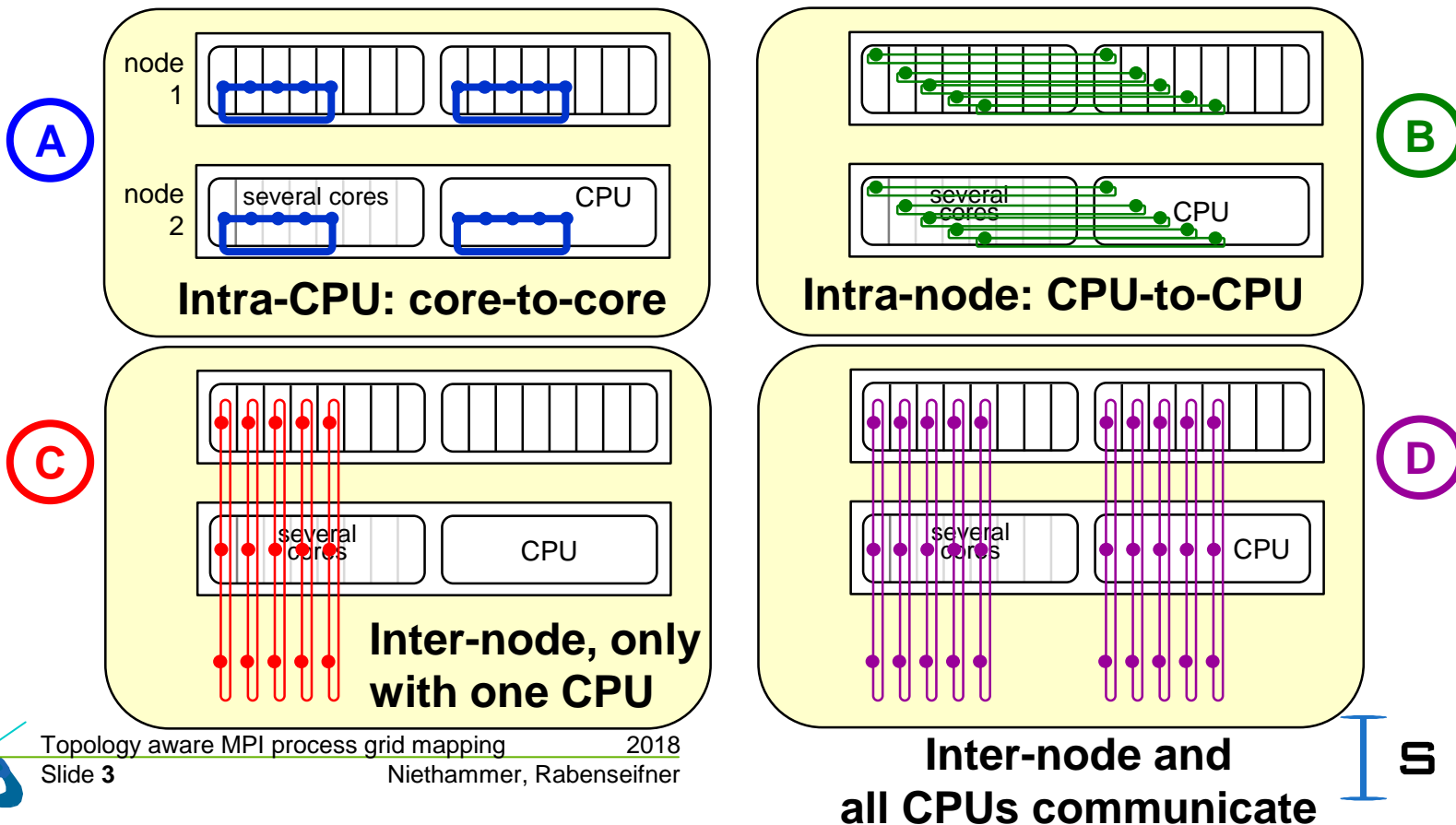
$O(4 \times 600) = O(2400)$  😊

# Ring Benchmarks for Inter- and Intra-node Communication

Benchmark `halo_irecv_send_multiplelinks_toggle.c`

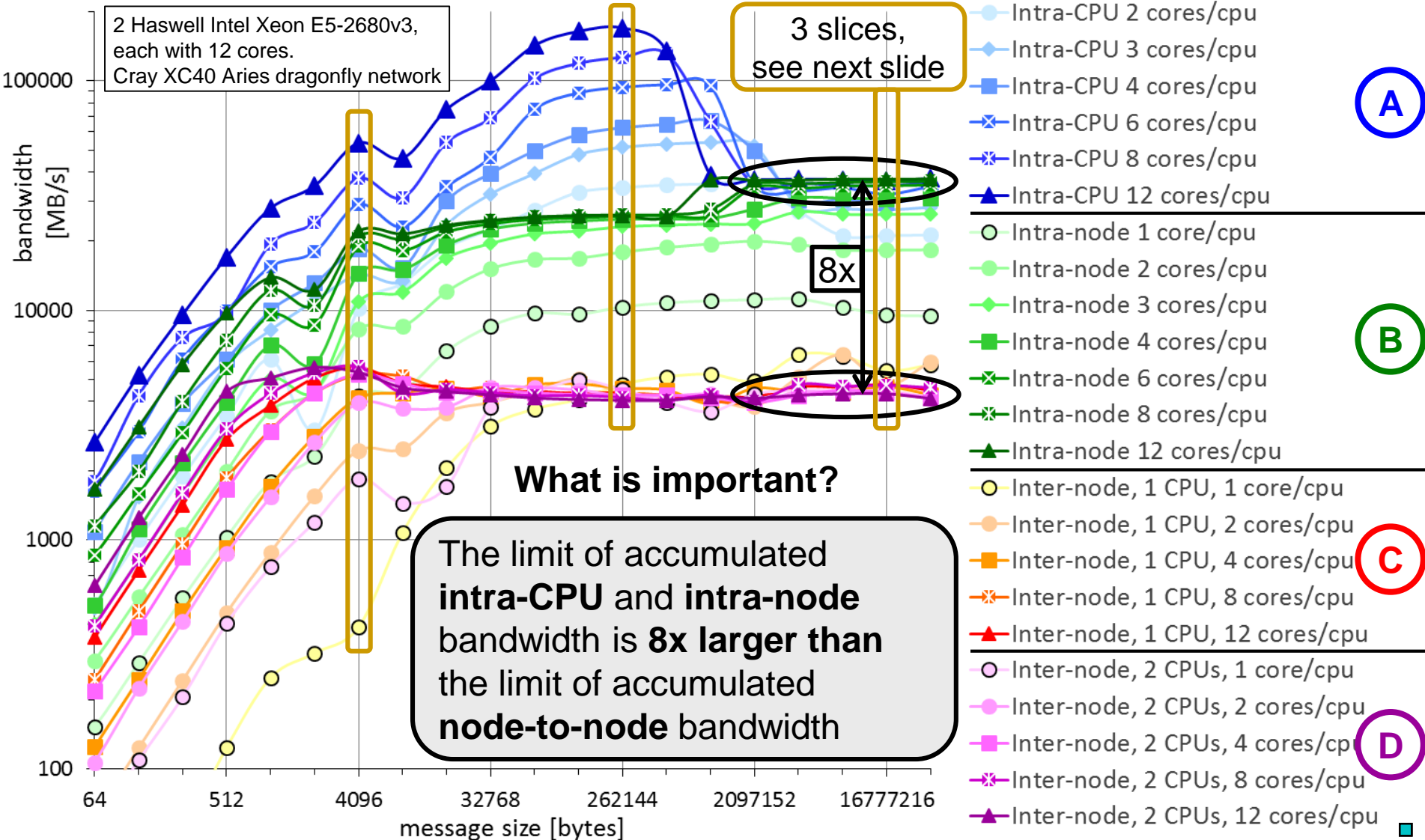
- Varying message size,
- number of **communication cores per CPU**, and
- four communication schemes (example with 5 **communicating cores per CPU**)

See HLRS online courses  
<http://www.hlrs.de/training/par-prog-ws/>  
→ Practical → MPI.tar.gz  
→ subdirectory MPI/course/C/1sided/

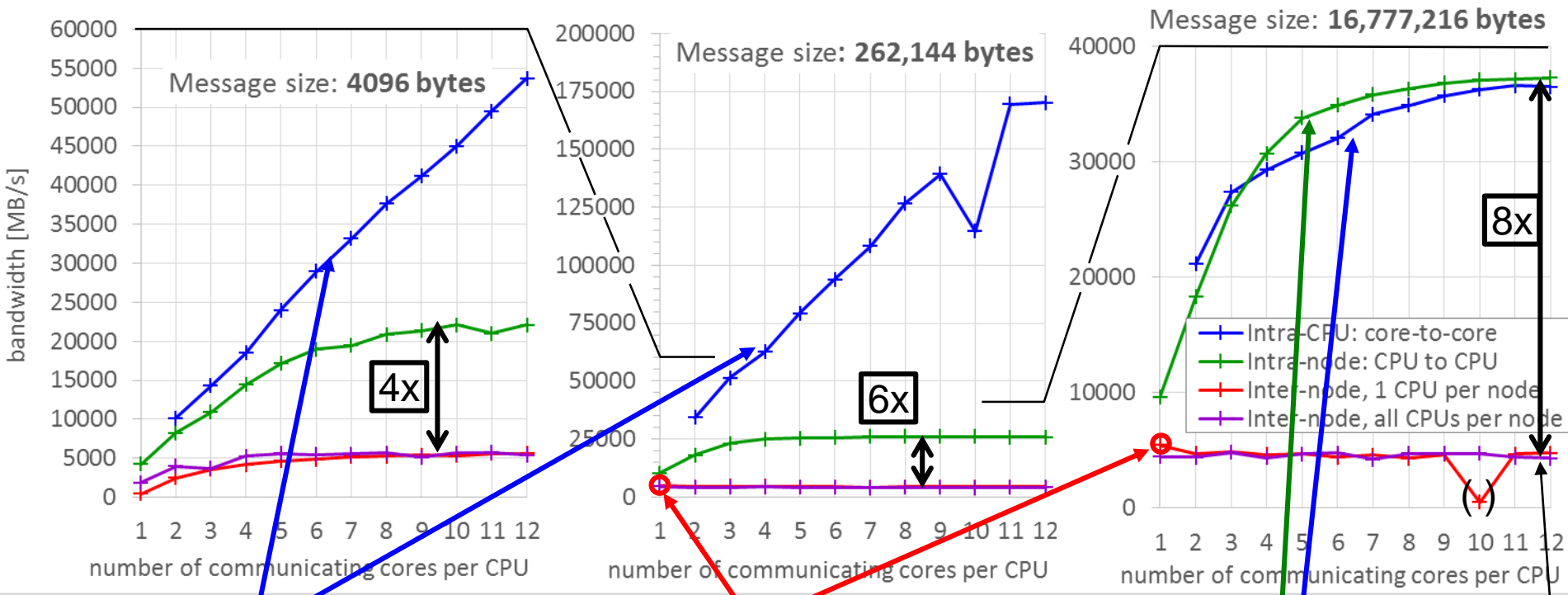


# Duplex accumulated ring bandwidth per node

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



# Duplex accumulated ring bandwidth per node – scaling vs. asymptotic behavior





Core-to-core:  
Linear scaling for small to medium size messages due to caches

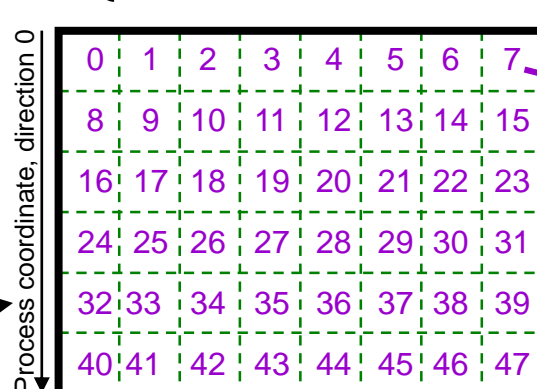
Node-to-node:  
One duplex link by **one core** already fully saturates the network

Core-to-core & CPU-to-CPU:  
**Long messages:**  
Same asymptotic limit through **memory bandwidth**

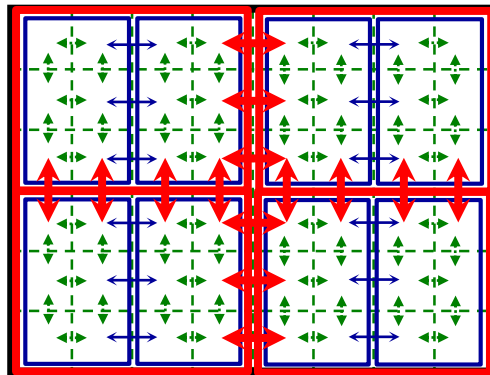
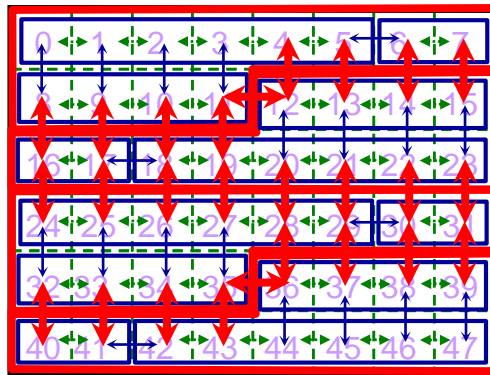
Result: The limit of accumulated **intra-CPU** and **intra-node** bandwidth is **8x larger** than the limit of accumulated **node-to-node** bandwidth

# Re-numbering on a cluster of SMPs (cores / CPUs / nodes)

- Example with 48 cores on:
  - 4 ccNUMA nodes
  - each node with 2 CPUs ,
  - each CPU with 6 cores 
- 2-dim application with 6000 x 8080 gridpoints
  - Minimal communication with 2-dim domain composition with 1000 x 1010 gridpoints/core (shape as quadratic as possible → minimal circumference → minimal halo communication)
  - virtual 2-dim process grid: 6 x 8
- How to locate the MPI processes on the hardware?
  - Using sequential ranks in MPI\_COMM\_WORLD
  - Optimized placement → Proposed algorithm in slides 7-15

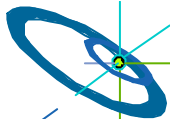


Order of the new ranks:  
Last coordinate is running contiguously  
→ Perfect basis for MPI\_Cart\_create() without reorder, i.e. with reorder==0 / .FALSE.



# Proposed Mapping Algorithm

- To keep the algorithm small (and fast):
  - Use multi-level Cartesian subdomains
- Based on the benchmark results:
  - The first and major optimization goal is minimizing the **inter**-node communication volume
  - Using the algorithm from [4] for the multi-dimensional factorization of the number of nodes, but with a modified optimization goal that is application topology aware
- Using the same principles for each further hardware level to minimize
  - intra-node (i.e., CPU-to-CPU) communication
  - intra-CPU (i.e., core-to-core) communication





# The Optimization Algorithm – First level

**Given:**  $d$ -dimensional Cartesian grid with a total of  $T = \prod_{i=1}^d t_i$  elements

Level 1 (= outer level = node level) on  $N$  nodes:

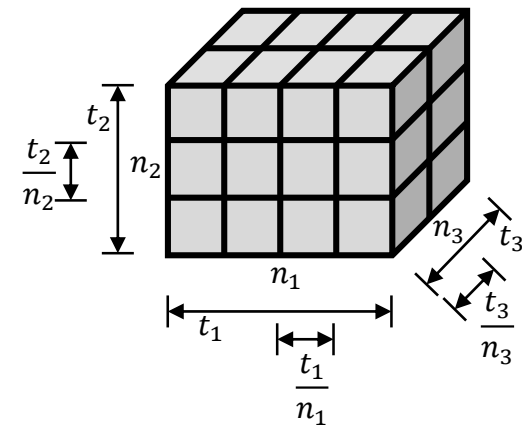
- Factorization of  $N$  into factors  $(n_i)_{i=1,d}$  with  $N = \prod_{i=1}^d n_i$
- Communication cost in both directions of each dimension (example for  $d = 3$ ):

$$2 \frac{t_2 t_3}{n_2 n_3}, \quad 2 \frac{t_1 t_3}{n_1 n_3}, \quad 2 \frac{t_1 t_2}{n_1 n_2}$$

- Minimizing the communication costs  $c$

$$c^{(level=1)} = 2 \sum_{i=1}^d \prod_{\substack{j=1 \\ j \neq i}}^d \frac{t_j}{n_j} = 2 \frac{T}{N} \sum_{i=1}^d \frac{n_i}{t_i}$$

The total grid divided into subdomains, one on each node



**Summary of Level 1:** One must search factors  $(n_i)_{i=1,d}$

- that factorize  $N$  with  $N = \prod_{i=1}^d n_i$
- and minimize the term  $\sum_{i=1}^d \frac{n_i}{t_i}$

The subdomain on each node is now divided into processors

# Second Level Optimization

Level 2:

- Each node has  $P$  processors (or cores)
- Factorization of  $P$  into factors  $p_i$  with  $P = \prod_{i=1}^d p_i$
- Communication costs in both directions of each dimension:

$$2 \frac{t_2 t_3}{n_2 p_2 n_3 p_3}, \quad 2 \frac{t_1 t_3}{n_1 p_1 n_3 p_3}, \quad 2 \frac{t_1 t_2}{n_1 p_1 n_2 p_2}$$

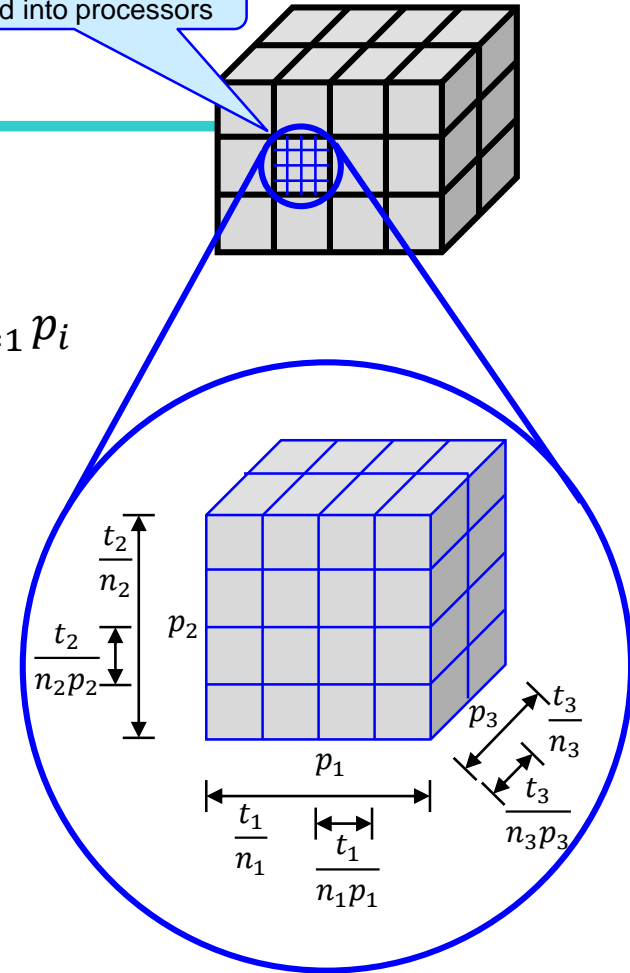
- Minimizing the communication costs  $c$ ,

$$c^{(level=2)} = 2 \sum_{i=1}^d \prod_{\substack{j=1 \\ j \neq i}}^d \frac{t_j}{n_j p_j} = 2 \frac{T}{NP} \sum_{i=1}^d \frac{n_i p_i}{t_i}$$

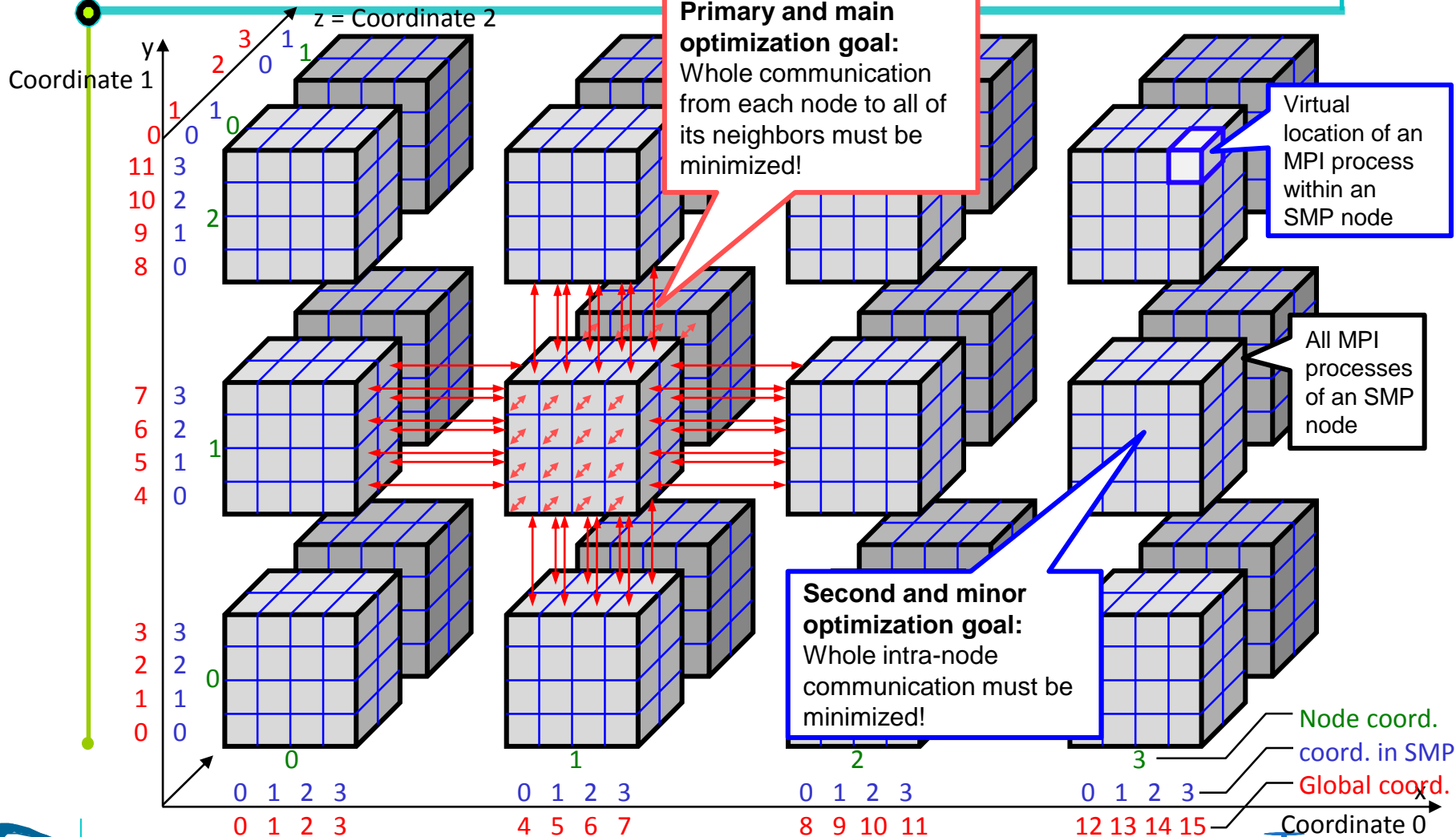
## Summary of Level 2:

One must search factors  $(p_i)_{i=1,d}$

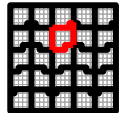
- that factorize  $P$  with  $P = \prod_{i=1}^d p_i$
- and minimize the term  $\sum_{i=1}^d \frac{n_i p_i}{t_i}$



# Hierarchical Cartesian Domain Decomposition



# Example

- **Given:**  $d=3$  dimensions,  $N=625$  nodes,  $P=24$  cores, and all  $t_i$  are identical
- $\rightarrow$  in all formulas, the  $1/t_i$  can be ignored
- Level 1:
  - Search  $(n_i)_{i=1,3}$  that  $n_1 n_2 n_3 = 625$  and  $\sum_{i=1}^3 n_i$  minimal
  - Result:  $(n_i)_{i=1,3} = (25,5,5)$  with  $\sum_{i=1}^3 n_i = 35$
- Level 2:
  - Search  $(p_i)_{i=1,3}$  that  $p_1 p_2 p_3 = 24$  and  $\sum_{i=1}^3 n_i p_i$  minimal
  - Result:  $(p_i)_{i=1,3} = (1,6,4)$  with  $\sum_{i=1}^3 n_i p_i = 25 + 30 + 20 = 75$
- **Optimized result:**
  - with *(nodes x cores)* in each dimension:  $(25 \times 1) \times (5 \times 6) \times (5 \times 4)$
  - Total core numbers as used for MPI\_Cart\_create:  $25 \times 30 \times 20$
- Comparison with existing MPI\_Dims\_create:
  - MPI\_Dims\_create would result in  $25 \times 25 \times 24$
  - $\rightarrow$  complex mapping to the hardware is needed, see also slide 2 
  - or no reordering is done  $\rightarrow (25 \times 1) \times (25 \times 1) \times (1 \times 24)$
  - with  $\sum_{i=1}^3 n_i = 51$  😞 (instead of 35 😊),  $\rightarrow$  **46% more inter-node communication!**

# Generalized multi-level optimization

**Given:**  $d$ -dimensional Cartesian grid with a total of  $T = \prod_{i=1}^d t_i$  elements

Number of hardware levels  $L$

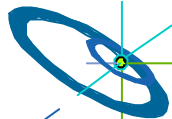
and for each level the number of processors  $N^{(l)}$ ,  $l = 1, L$

Communication costs on each level  $l$ :

$$c^{(l)} = 2 \frac{T}{\prod_{k=1}^l N^{(k)}} \sum_{i=1}^d \frac{\prod_{k=1}^l n_i^{(k)}}{t_i} \quad \text{for a factorization } N^{(l)} = \prod_{i=1}^d n_i^{(l)}$$

Search:

- On each level  $l = 1..L$ , one must search a factorization  $(n_i^{(l)})_{i=1,d}$  of  $N^{(l)}$
  - with  $N^{(l)} = \prod_{i=1}^d n_i^{(l)}$
  - and minimal sum  $\sum_{i=1}^d \frac{\prod_{k=1}^l n_i^{(k)}}{t_i}$ ,
- i.e., with minimal sum  $\sum_{i=1}^d a_i^{(l)} n_i^{(l)}$  with  $a_i^{(l)} = \frac{\prod_{k=1}^{l-1} n_i^{(k)}}{t_i}$

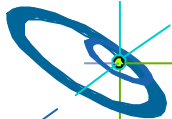
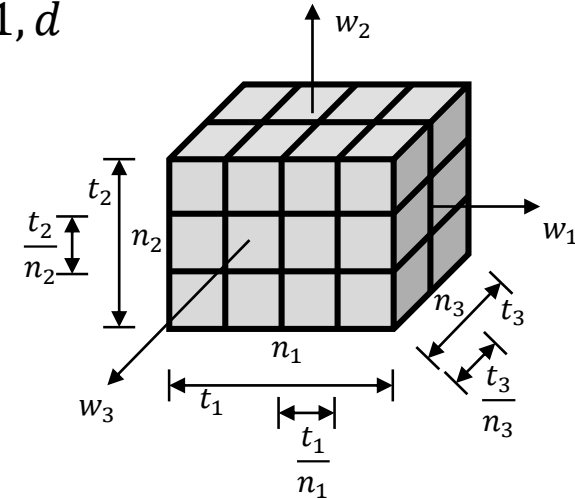


# Generalization with weighted communication in each direction, e.g. by different halo width

- If the communication cost in each direction  $i = 1, d$  is multiplied with a halo width  $w_i$ , e.g.,  $2 \frac{t_2 t_3}{n_2 n_3} w_1$
- On each  $l = 1..L$ , the sum to be minimized is

$$\sum_{i=1}^d \frac{\prod_{k=1}^l n_i^{(k)}}{(t_i/w_i)}$$

- i.e.,  $\sum_{i=1}^d a_i^{(l)} n_i^{(l)}$  with  $a_i^{(l)} = \frac{\prod_{k=1}^{l-1} n_i^{(k)}}{(t_i/w_i)}$



# Using weighted MPI\_Dims\_create for application topology awareness

- Task: In general, on each hardware topology level  $l$  with  $l = 1, L$  and for given  $N^{(l)}$ , find factorizations  $(n_i^{(l)})_{i=1,d}$  with  $N^{(l)} = \prod_{i=1}^d n_i^{(l)}$  and the following sum is minimal:  $\sum_{i=1}^d a_i^{(l)} n_i^{(l)}$  with  $a_i^{(l)} = \frac{\prod_{k=1}^{l-1} n_i^{(k)}}{(t_i/w_i)}$
- Algorithm on each level  $l = 1..L$

New MPI\_Dims\_weighted\_create

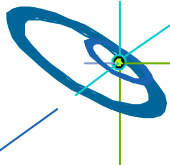
- Sorting indexes  $i = 1..d$  so that the  $(a_{i'})_{i'=1..d}$  to be non-decreasing
- Calculate all divisors of  $N^{(l)}$  i.e. with  $n_i \geq n_{i+1}$
- Loop over all non-increasing factorizations and find optimum according to
  - **1st criterion: a factorization is better if  $\sum_{i=1}^d a_i^{(l)} n_i^{(l)}$  is smaller**
  - **2nd criterion: if the  $\sum_{i=1}^d a_i^{(l)} n_i^{(l)}$  is the same then a factorization is better if  $\Delta = n_1^{(l)} - n_d^{(l)}$  is smaller**
  - **3rd criterion: if  $\sum_{i=1}^d n_i$  and  $\Delta = n_1^{(l)} - n_d^{(l)}$  are the same then a factorization is better if  $n_1^{(l)}$  is smaller**
- Revert the index mapping  $(n_{i'})_{i'=1..d} \rightarrow (n_i^{(l)})_{i=1..d}$

Additional tricks:

- Calculate divisors only upto sqrt(N), calculate the rest by reciprocal values.
- Loop over divisors from highest to smallest, recursively over  $i = 1, d$ ,
- Start value for  $n_{i+1}$  is next real divisor equal or smaller then  $n_i$

Re-mapping of all process ranks according to  $(n_i^{(l)})_{i=1..d, l=1..L}$

and creation of the new Cartesian communicator according to algorithm in [2]



# Proposed Interfaces

e.g., with  
 $25 \times 25 \times 24 = 15000$  processes  
on **625** ccNUMA nodes with  
**2 CPUs/node** and **12 cores/CPU**

```
MPI_Cart_ml_create_from_types ( MPI_Comm comm_old,  
    int ntype_levels,    int type_levels[ntype_levels],  
    int ndims,           double dim_weights[ndims],  
    int periods[ndims], MPI_Info info,  
/*OUT*/ int dims[ndims], MPI_Comm *comm_cart );
```

e.g.,  
{ MPI\_COMM\_TYPE\_SHARED,  
 MPI\_COMM\_TYPE\_NUMA }

e.g., 3 dimensions with a data  
grid with 1000 x 1100 x 950  
elements → dim\_weights[] =  
{ 1.0/1000, 1.0/1100, 1.0/950 }

Rank mapping is based on:

- Node level: **625** = 5 x 25 x 5
- CPU level: **2** = 2 x 1 x 1
- Core level: **12** = 3 x 1 x 4

Result (product): **30 x 25 x 20**

The Cartesian communicator reflects this result: **30 x 25 x 20**

Next steps:

```
MPI_Comm_rank ( comm_cart, &my_rank );  
MPI_Cart_coords ( comm_cart, my_rank, ndims, coords)
```

```
MPI_Cart_ml_create_from_comms (int nlevels,  
    MPI_Comm level_comms[nlevels],  
    int ndims, double dim_weights[ndims], int periods[ndims], MPI_Info info,  
/*OUT*/ int dims[ndims], MPI_Comm *comm_cart );
```

e.g., level\_comms[0] is comm\_old, level\_comms[1  
and 2] are the result recursively called MPI\_  
Comm\_split\_type with the type\_levels from above.

Same as above

Same as above

```
MPI_Dims_weighted_create ( int nnodes, int ndims, double dim_weights[ndims],  
/*OUT*/ int dims[ndims] );
```

```
MPI_Dims_ml_create ( int nnodes, int ndims, double dim_weights[ndims],  
    int nlevels, int sizes[nlevels], /*OUT*/ int dims_ml[ndims][nlevels] );
```

Same as above

Substitute for / enhancement to existing MPI-1

MPI\_Dims\_create ( size\_of\_comm\_old, ndims, dims );  
MPI\_Cart\_create ( comm\_old, ndims, dims, periods,  
reorder, \*comm\_cart );



# Conclusion and Outlook

## Conclusions

- We developed a new algorithm to minimize the total communication time in a cluster of ccNUMA nodes with multi-core CPUs.
- It is needed, due to the significant bandwidth differences between inter- and intra-node communication.
- It can be implemented based on the algorithm in [4], but with a modified optimization goal, and repeated for each hardware level.

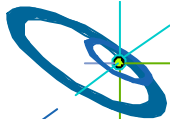
## Outlook

- We plan to provide a portable implementation, and
- compare it with existing solutions with `MPI_Dims_create` + `MPI_Cart_create`.
- We plan to propose an appropriate interface for the next MPI standard, because MPI libraries may internally have faster access to the hardware topology information for a given communicator.



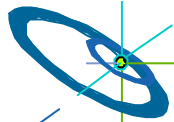
# References

- [1] Pavan Balaji et al. 2009-2012. Topology awareness in MPI Dims create. <https://github.com/mpi-forum/mpi-forumhistoric/issues/195> Accessed 2018-07-19.
- [2] Bill Gropp. 2018. Using Node Information to Implement MPI Cartesian Topologies. In *Proceedings of the 25nd European MPI Users' Group Meeting (EuroMPI '18)*, September 23–26, 2018, Barcelona, Spain. ACM, New York, NY, USA, 9 pages.
- [3] T. Hoefler and M. Snir. 2011. Generic Topology Mapping Strategies for Large-scale Parallel Architectures. In *Proceedings of the 2011 ACM International Conference on Supercomputing (ICS'11)*. ACM, 75–85.
- [4] Jesper Larsson Träff and Felix Donatus Lübbe. 2015. Specification Guideline Violations by MPI Dims Create. In *Proceedings of the 22nd European MPI Users' Group Meeting (EuroMPI '15)*. ACM, New York, NY, USA, Article 19, 2 pages.



# Appendix

- Additional material that is not part of the poster

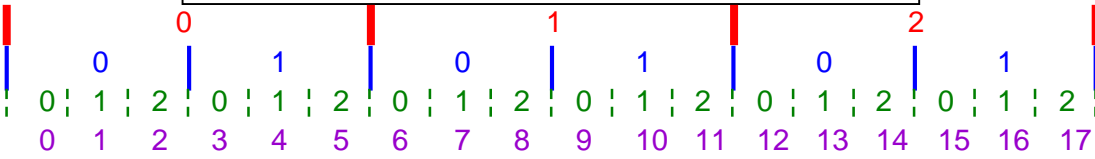


2-D example on hierarchical hardware with 9 nodes x 4 CPUs x 6 cores resulting in 12\*18 Cartesian processes

# Renumbering in a 2-dim example with 3 levels (nodes / CPUs / cores)

$$\text{dim1} = \text{inner\_d1} * \text{mid\_d1} * \text{outer\_d1} = 3 * 2 * 3 = 18$$

oc0 = outer coordinate  
mc0 = middle coordinate  
ic0 = inner coordinate  
c0 = process coordinate



direction 1 →  
oc1 = outer coordinate in 0..(outer\_d1-1)  
mc1 = middle coordinate in 0..(mid\_d1-1)  
ic1 = inner coordinate in 0..(inner\_d1-1)  
c1 = process coordinate

**Outer** = several ccNUMA nodes  
**Middle** = CPUs of a ccNUMA node  
**Inner** = cores of a CPU

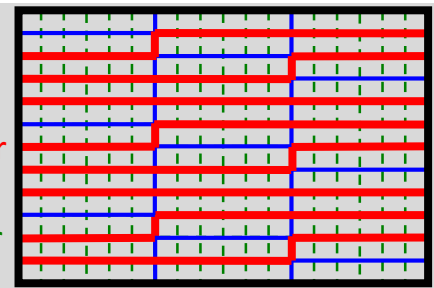
|           |           |           |           |           |           |     |    |    |     |    |    |      |     |
|-----------|-----------|-----------|-----------|-----------|-----------|-----|----|----|-----|----|----|------|-----|
| 0         | 1         | 2         | 6         | 7         | 8         | 24  | 25 | 26 | 30  | 31 | 32 | 48   | 54  |
| →0→1→2    | →3→4→5    | →6→7→8    | →9→10→11  | →12→13→14 | →15→16→17 |     |    |    |     |    |    |      |     |
| 3         | 4         | 5         | 9         | 10        | 11        | 27  | 28 | 29 | 33  | 34 | 35 |      |     |
| →18→19→20 | →21→22→23 | →24→25→26 | →27→28→29 | →30→31→32 | →33→34→35 |     |    |    |     |    |    |      |     |
| 12        | 13        | 14        | 18        | 19        | 20        | 36  | 37 | 38 | 42  | 43 | 44 | 60   | 66  |
| →36→37→38 | →         | →         | →         | →         | →         | →   | →  | →  | →   | →  | →  | →    | →   |
| 15        | 16        | 17        | 21        | 22        | 23        | 39  | 40 | 41 | 45  | 46 | 47 |      |     |
| →54→      | →         | →         | →         | →         | →         | →   | →  | →  | →   | →  | →  | →    | →   |
| 72        | 73        | 74        | 78        |           |           | 96  |    |    | 102 |    |    | 120  | 126 |
| →72       |           |           |           |           |           |     |    |    |     |    |    |      |     |
| 75        | 76        | 77        |           |           |           |     |    |    |     |    |    |      |     |
| →90       |           |           |           |           |           |     |    |    |     |    |    |      |     |
| 84        |           |           | 90        |           |           | 108 |    |    | 114 |    |    | 132  | 138 |
| →108      |           |           |           |           |           |     |    |    |     |    |    |      |     |
| 87        |           |           |           |           |           |     |    |    |     |    |    |      |     |
| →126      |           |           |           |           |           |     |    |    |     |    |    |      |     |
| 144       |           |           | 150       |           |           | 168 |    |    | 174 |    |    | 192  | 198 |
| →144      |           |           |           |           |           |     |    |    |     |    |    |      |     |
| 147       |           |           |           |           |           |     |    |    |     |    |    |      |     |
| →162      |           |           |           |           |           |     |    |    |     |    |    |      |     |
| 156       |           |           | 162       |           |           | 180 |    |    | 186 |    |    | 204  | 210 |
| →180      |           |           |           |           |           |     |    |    |     |    |    |      |     |
| 159       |           |           |           |           |           |     |    |    |     |    |    |      |     |
| →198      |           |           |           |           |           |     |    |    |     |    |    |      |     |
|           |           |           |           |           |           |     |    |    |     |    |    | 215  | 215 |
|           |           |           |           |           |           |     |    |    |     |    |    | →215 |     |

Old ranks in MPI\_COMM\_WORLD. The ranks must be sequential in the hardware, i.e., first through the CPUs, then through the ccNUMA nodes, and then through the cluster

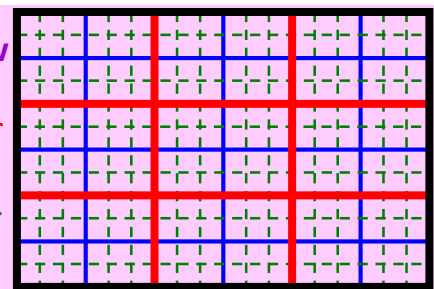
New ranks in optimized communicator.  
→ New global communicator with minimal node-to-node & optimal intra-node communication

Number of communication links:

Without re-numbering:  
150 outer  
72 mid  
180 inner



With new ranks:  
60 outer  
90 mid  
252 inner



2-D example with 12\*18 processes on hierarchical hardware

Order of the new ranks: last coordinate is running contiguously  
→ Perfect basis for MPI\_Cart\_create() with reorder=0 / .FALSE.



# Renumbering in a 2-dim example with 3 levels (nodes / CPUs / cores) – the code

This algorithm requires sequential ranking in MPI\_COMM\_WORLD

```
Product = number of cores/CPU number of CPUs/node number of nodes
/*Input: */ inner_d0=...; mid_d0=...; outer_d0=...;
             inner_d1=...; mid_d1=...; outer_d1=...;

dim0=inner_d0*mid_d0*outer_d0;          dim1=inner_d1*mid_d1*outer_d1;
idim=inner_d0*inner_d1; mdim=mid_d0*mid_d1; odim=outer_d0*outer_d1;
whole_size=dim0*dim1 /* or =idim*mdim*odim */;
ranks= malloc(whole_size*sizeof(int));
for (oc0=0; oc0<outer_d0; oc0++) /*any sequence of the nested loops works*/
  for (mc0=0; mc0<mid_d0; mc0++)
    for (ic0=0; ic0<inner_d0; ic0++)
      for (ocl1=0; ocl1<outer_d1; ocl1++)
        for (mcl1=0; mcl1<mid_d1; mcl1++)
          for (icl1=0; icl1<inner_d1; icl1++)
            { old_rank = icl1 + inner_d1*ic0 + (mcl1 + mid_d1 *mc0)*idim
              + (ocl1 + outer_d1*oc0)*idim*mdim;
              c0 = ic0 + inner_d0*mc0 + inner_d0*mid_d0*oc0;
              c1 = icl1 + inner_d1*mcl1 + inner_d1*mid_d1*ocl1;
              new_rank = c1 + dim1*c0;
              ranks[new_rank] = old_rank;
            }
/* Establishing new_comm with the new ranking in a array "ranks": */
MPI_Comm_group(MPI_COMM_WORLD, &world_group);
MPI_Group_incl(world_group, world_size, ranks, &new_group); free(ranks);
MPI_Comm_create(MPI_COMM_WORLD, new_group, &new_comm);
dims[0] = dim0; dims[1] = dim1; /* final output */
MPI_Cart_create(new_comm, 2, dims, periods, 0 /*=false*/, &comm_cart);
```

For an alternative with MPI\_Comm\_split, see MPI-3.1, Sec. 7.5.8, page 313, lines 7-13.

# MPI\_Dims\_create optimizing $\sum_{i=1}^d n_i$

- Based on Jesper Träff's algorithm *tuwdims.c* developed for [4]
  - Jesper's optimization criterion:
    - Minimizing  $\Delta = n_1 - n_d$
- New optimization criterion (three criteria)
  - 1<sup>st</sup> criterion: a factorization is better if  $\sum_{i=1}^d n_i$  is **smaller**
  - 2<sup>nd</sup> criterion: if the  $\sum_{i=1}^d n_i$  is the same then a factorization is better if  $\Delta = n_1 - n_d$  is smaller
  - 3<sup>rd</sup> criterion: if  $\sum_{i=1}^d n_i$  and  $\Delta = n_1 - n_d$  are the same then a factorization is better if  $n_1$  is smaller
- Both algorithms have nearly same execution time:  $\sim O(10\mu\text{s}/\text{call})$
- For  $n=2..10,000,000$  and  $\text{ndims}=2..10$ 
  - 4630x different factorization: with **same**  $\Sigma$  and  $\Delta$  and **smaller dims**[0]
    - e.g.,  $\text{ndims}=3$ :  $N=360 = (\text{old}) 10 \times 6 \times 6 = (\text{new}) 9 \times 8 \times 5$  (both  $\Sigma=22$  and  $\Delta=4$ )
  - 6066x different factorization: with **better**  $\Sigma$  and **mostly worse**  $\Delta$ 
    - e.g.,  $\text{ndims}=3$ :  $N=35200 = (\text{old}) 40 \times 40 \times 22$  ( $\Sigma=102, \Delta=18$ ) = (new)  $44 \times 32 \times 25$  ( $\Sigma=101, \Delta=19$ )
    - $N=37044 = (\text{old}) 42 \times 42 \times 21$  ( $\Sigma=105, \Delta=21$ ) = (new)  $49 \times 28 \times 27$  ( $\Sigma=104, \Delta=22$ )

Next below 10,000 are:

- $\text{ndims}=3$ :
  - $22 \times 14 \times 12 = 21 \times 16 \times 11$
  - $21 \times 16 \times 15 = 20 \times 18 \times 14$
  - $26 \times 16 \times 15 = 24 \times 20 \times 13$
- $\text{ndims}=4$ :
  - $10 \times F \times 6 \times 6 = 9 \times 8 \times F \times 5$  with  $F=6,7,8,9$

All other larger than 50,000

[4] Jesper Larsson Träff and Felix Donatus Lübbe. 2015. Specification Guideline Violations by MPI Dims Create. In *Proceedings of the 22Nd European MPI Users' Group Meeting (EuroMPI '15)*. ACM, New York, NY, USA, Article 19, 2 pages.

