SG6 - HPC

# TD2-3/Lab-2 – Part 1: Deployment of an MPI application on a PC cluster

**Stéphane Vialle**

# Deployment of an MPI application on a PC cluster
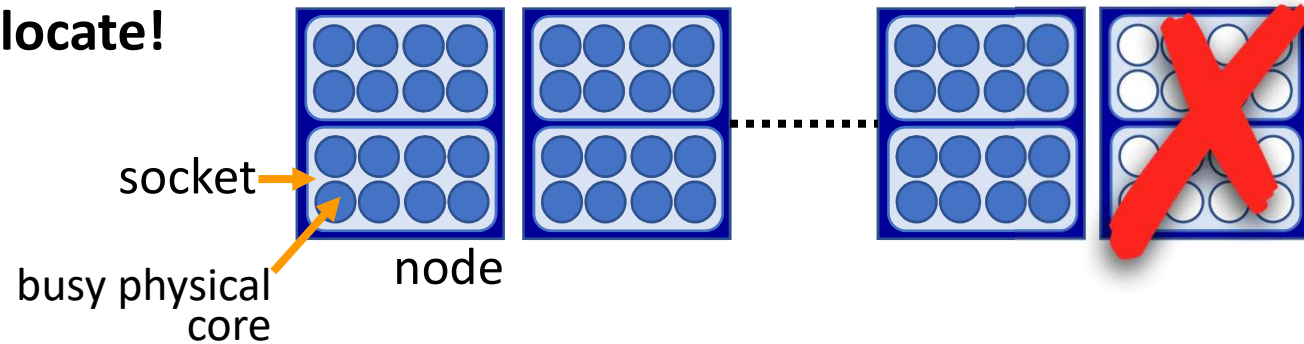
**Deployment rules & communication scheme**

1st deployment: using processes and threads
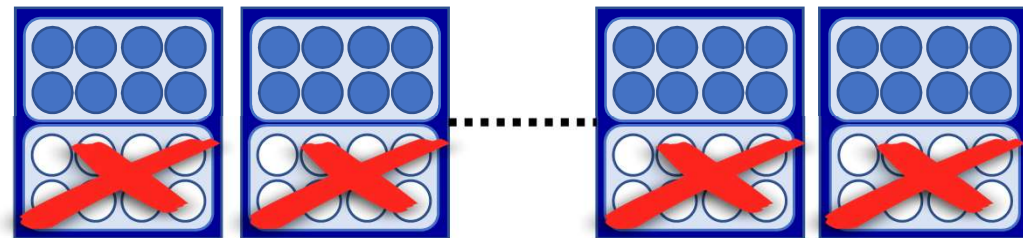
2nd deployment: using only processes

# Do not waste resources!

CentraleSupélec

**Use ALL nodes you allocate!**
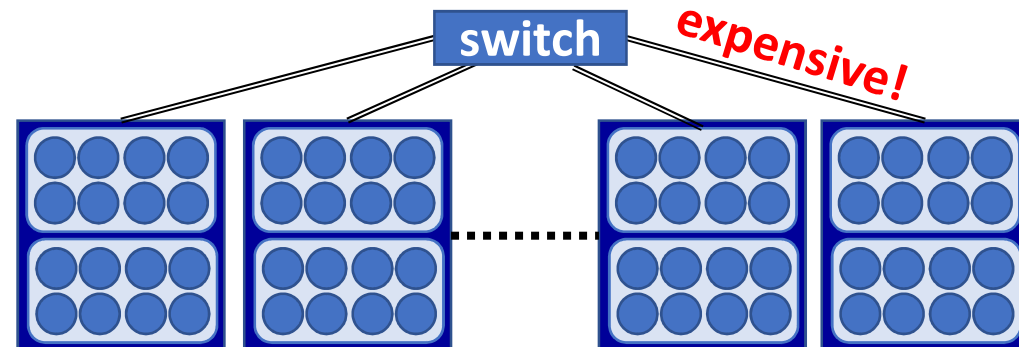
socket →

busy physical
core

node

**Use ALL physical cores of your nodes!**
(with processes or threads)

**Minimize the communication
cost across the interconnection
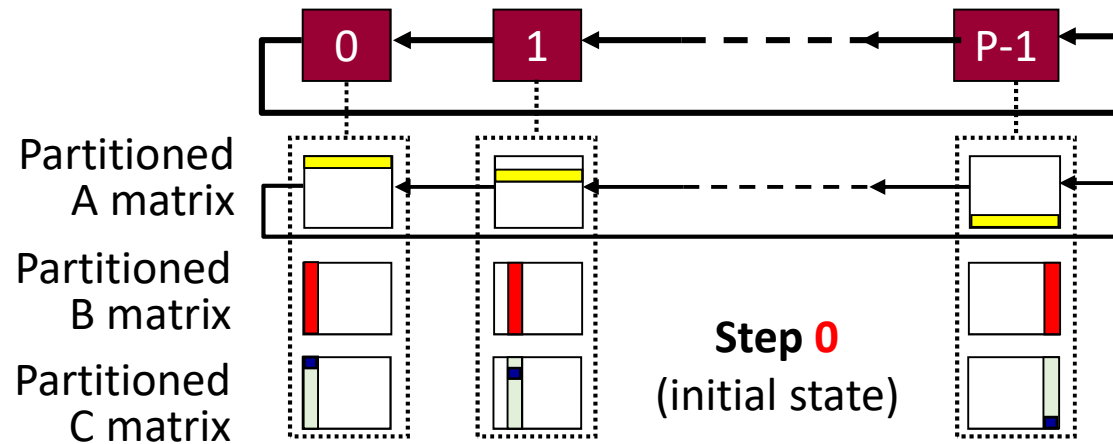network** (maximizing comm.
inside each node)

switch

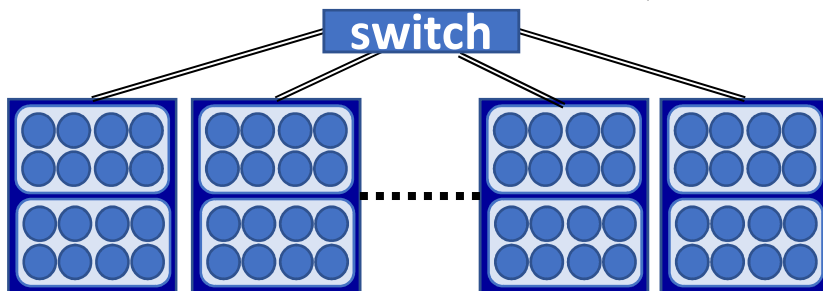expensive!

# Virtual ring of processes

**Distributed *Matrix Product* algorithm:**

- ring comm. scheme
- $P_i$ communicates only with $P_{i-1}$ and $P_{i+1}$



Partitioned A matrix

Partitioned B matrix

Partitioned C matrix

**Step 0** (initial state)

**Distributed & multithreaded implementation:**

- MPI + OpenMP
- OpenBLAS

process

threads

switch

**For a given nb of allocated nodes ($N_n$):**

→ **Find 2 relevant *mpirun* commands**

- Not wasting any resource
- Minimizing the comm. cost

# Deployment of an MPI application on a PC cluster
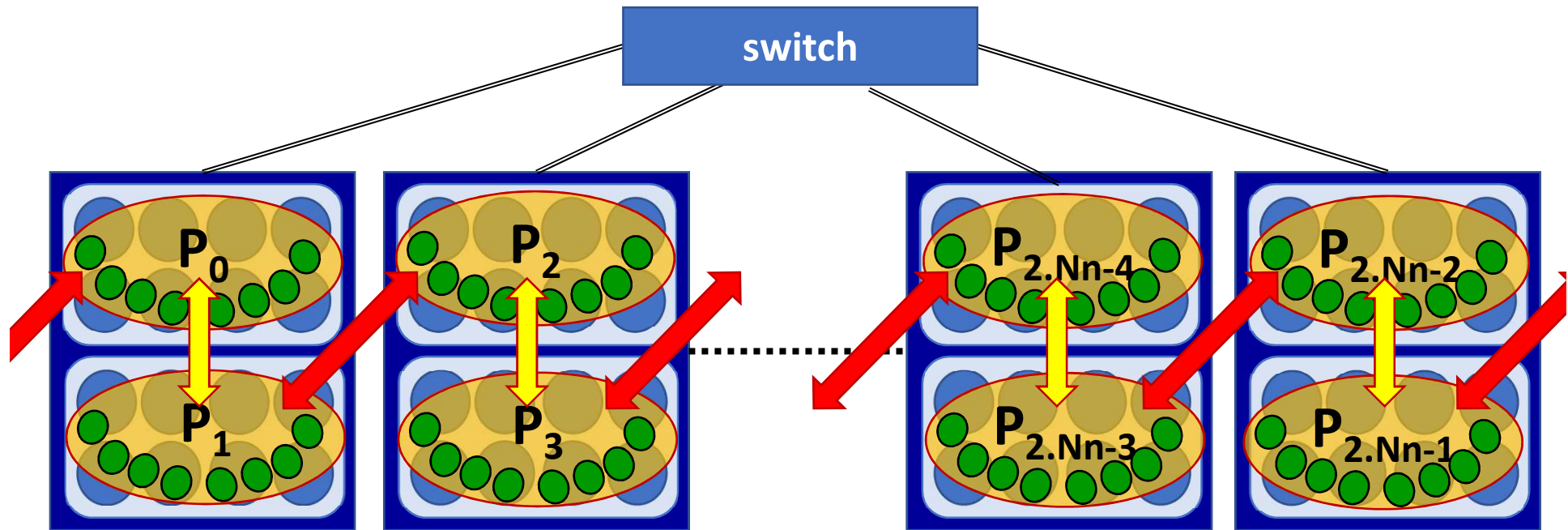
Deployment rules & communication scheme

**1st deployment: using processes and threads**

2nd deployment: using only processes

# Deployment strategy

1st deployment: using processes and threads

mpirun –np  *XX=2×N_n* –machinefile machines.txt

   –map-by **ppr:1:socket**

   –rank-by **socket**

   –bind-to **socket**

   ./MatrixProduct –klc *YY* –k 1 **–nt 8**

**comms: only 50% are expensive!**

# TO DO (1)

**Questions:**

1.  **Measure performances (Gflops) on 4, 8 and 16 nodes, with –k 1 –klc 16**

    → **Use OAR « *batch mode* » with « *myrun* » shell script:**

    - **Ex: oarsub -p "cluster='kyle'" -l nodes=4 '*./myrun 8 16*'**
    - after unzipping the archive, don't forget:

      *dos2linux myrun*    and    *chmod 700 myrun*

2.  **Compare to previous measurements on Kyle cluster** (S. Vialle – 27/12/2019)**:**

    | Nb of nodes | 1 | 2 | 4 | 8 | 16 | 32 |
    |---|---|---|---|---|---|---|
    | Gflops | 369 | 529 | 819 | 1103 | 1359 | 1387 |

3.  **Draw performance curves**

4.  **Analyse the performance curves**
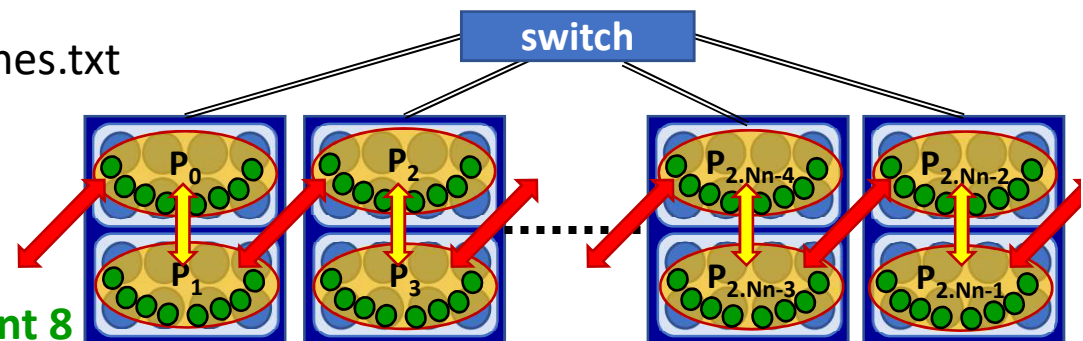
mpirun –np  $XX=2 \times N_n$ –machinefile machines.txt

    –map-by **ppr:1:socket**

    –rank-by **socket**

    –bind-to **socket**

    ./MatrixProduct –klc *YY=16* –k 1 **–nt 8**

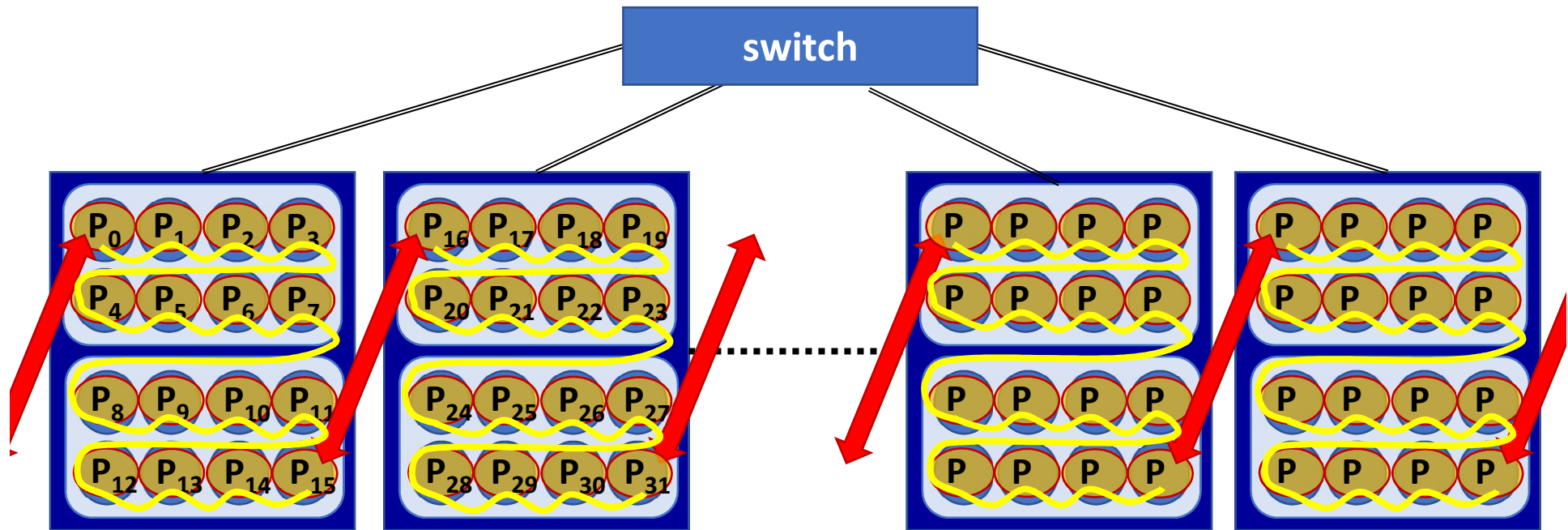# Deployment of an MPI application on a PC cluster

Deployment rules & communication scheme

1$^{st}$ deployment: using processes and threads

**2$^{nd}$ deployment: using only processes**

# Deployment strategy



switch

P0 P1 P2 P3
P4 P5 P6 P7
P8 P9 P10 P11
P12 P13 P14 P15

P16 P17 P18 P19
P20 P21 P22 P23
P24 P25 P26 P27
P28 P29 P30 P31

mpirun –np  *XX*=**??** –machinefile machines.txt

 –map-by **ppr:?:???**

 –rank-by **?**

 –bind-to **?**

 ./MatrixProduct –klc *YY* –k 1 **–nt ??**

**Comms: only 1/16 are expensive!**

**But 8x more steps.**

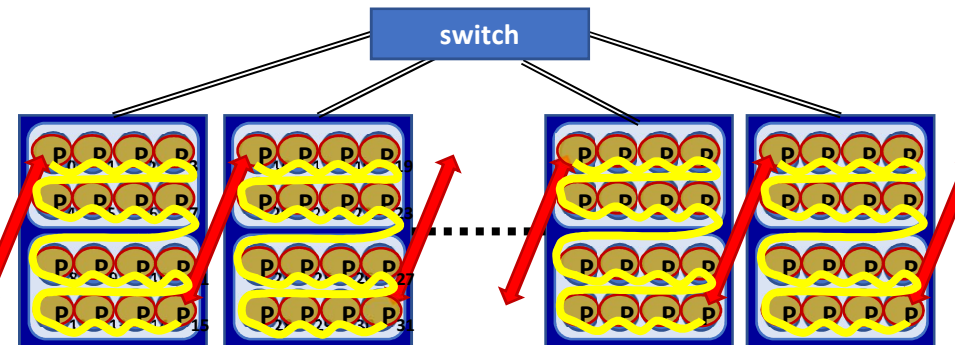**→ Same total volume on the interconnect (see course slides)**

# TO DO (2)

**Questions:**

1. **Measure performances (Gflops) on 4, 8, 16 and 32 nodes, with –k 1 –klc 16**
   - → Use OAR « *batch mode* » with « *myrun* » shell script
   - → **MODIFY *myrun* shell script and adapt oarsub command**

2. **Compare to previous measurements on Kyle cluster:**

| Nb of nodes | 1 | 2 | 4 | 8 | 16 | 32 |
|---|---|---|---|---|---|---|
| Gflops | 367 | 562 | 862 | 1290 | … | … |

3. **Draw performance curves (superpose with 1ˢᵗ deployement curves)**

4. **Analyse the performance curves**

mpirun –np  *XX=??* –machinefile machines.txt

    –map-by **ppr:?:???**

    –rank-by **?**

    –bind-to **?**

    ./MatrixProduct –klc *YY=16* –k 1 **–nt ??**

switch

# TD2-3/Lab-2 – Part 1:
# Deployment of an MPI application on a PC cluster

**End**