

PLACEMENT

POUR CELLES ET CEUX QUI ONT DES CŒURS

emmanuel.courcelle@toulouse-inp.fr

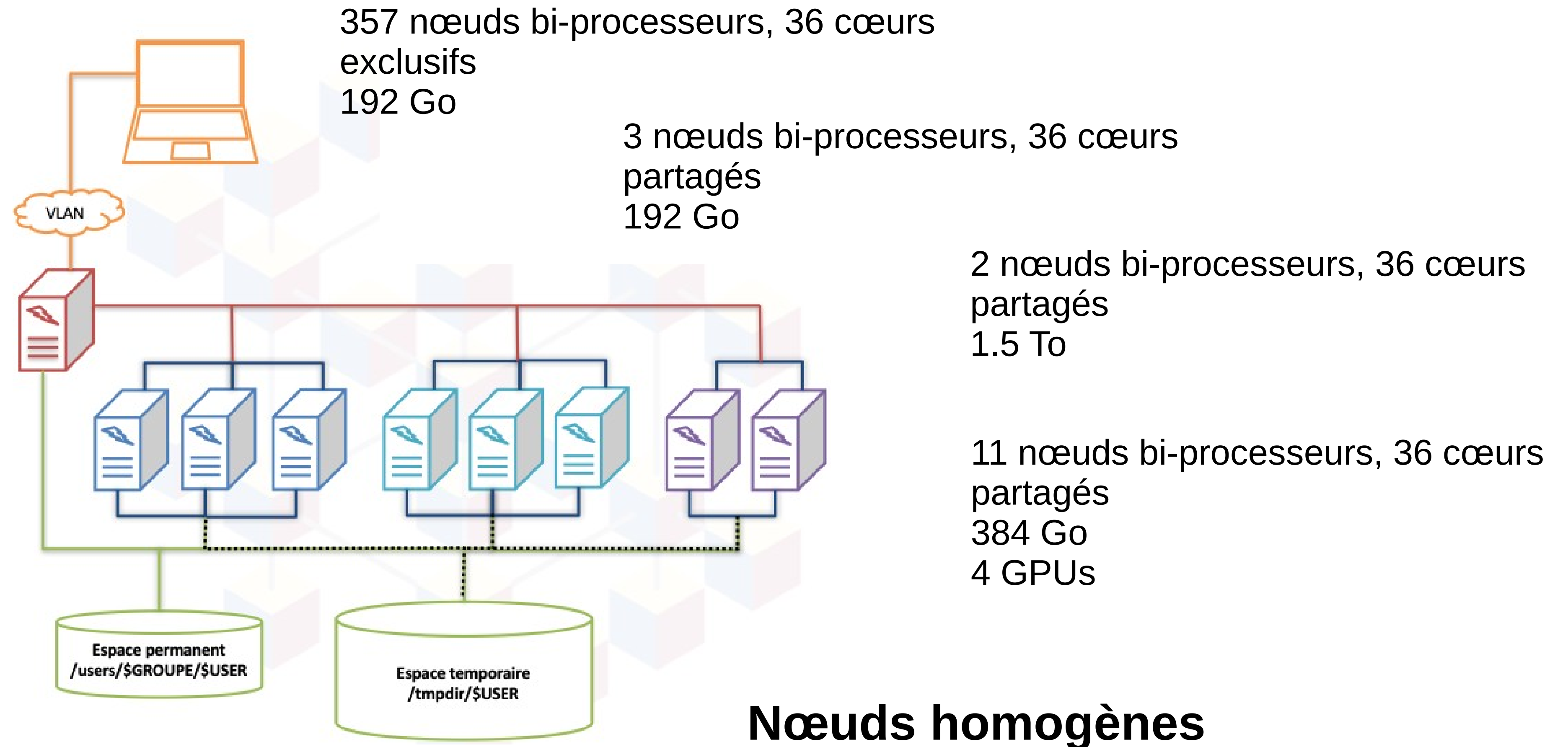
Café calcul – BBB - Jeudi 21 janvier 2021



CALMIP (UMS 3667)
Espace Clément Ader
www.calmip.univ-toulouse.fr

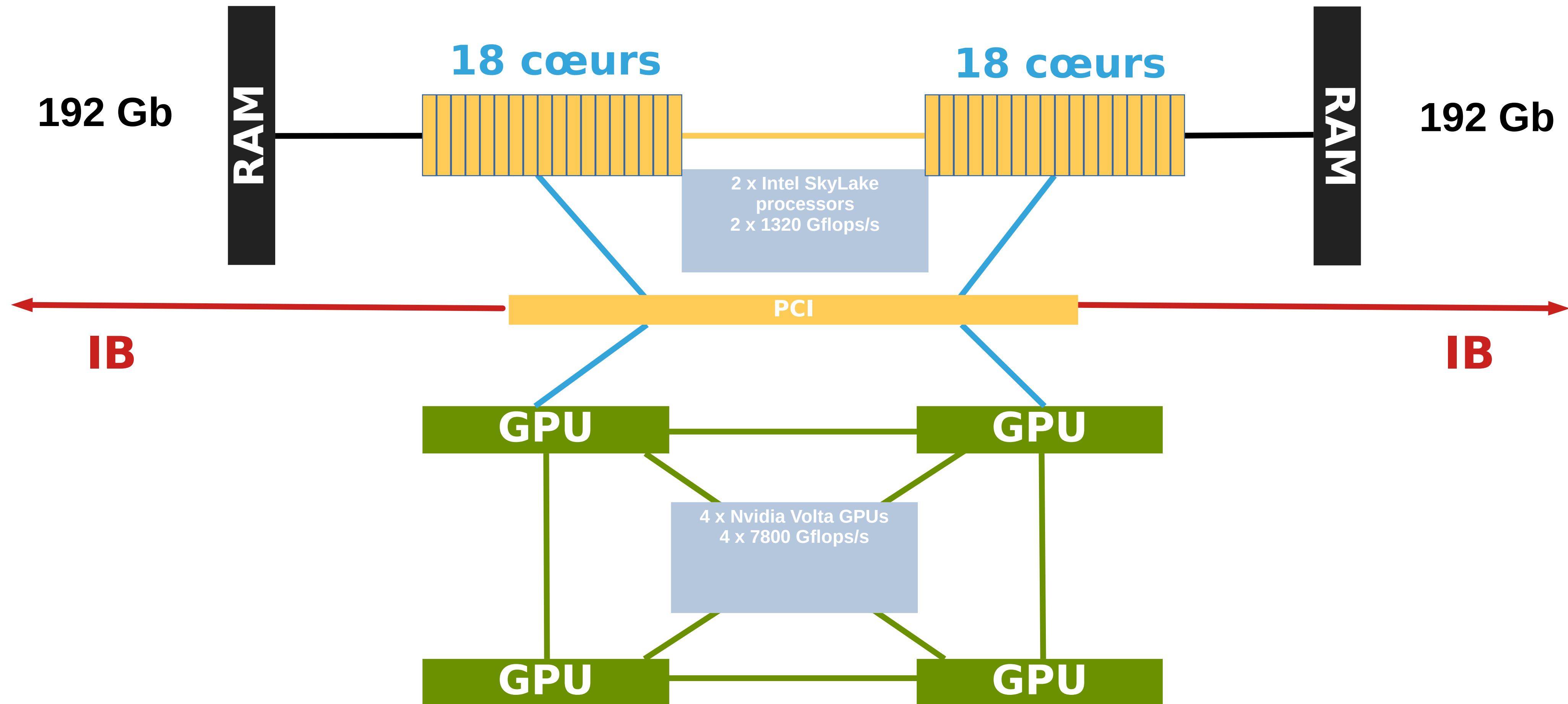


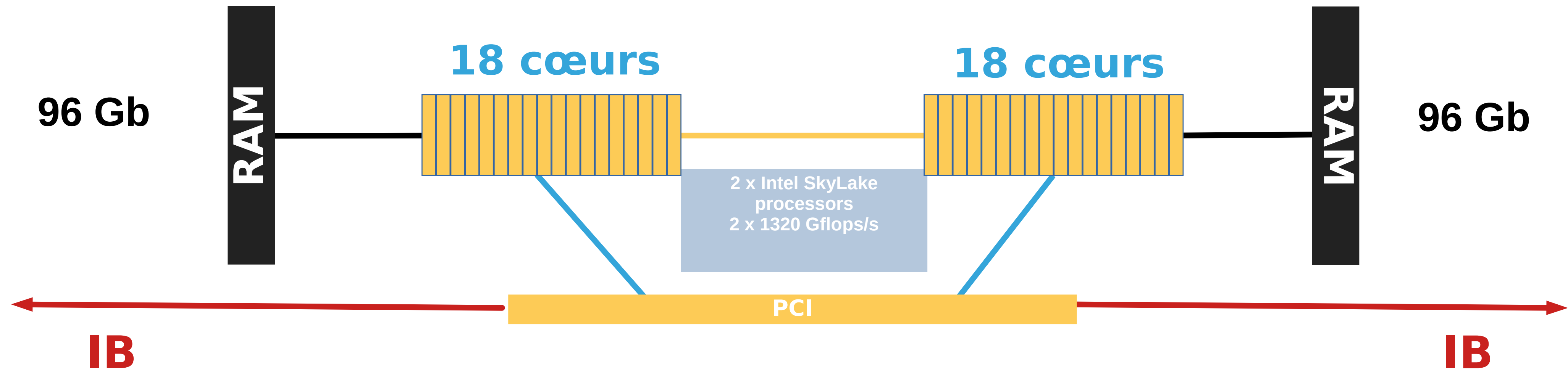
Olympe, un cluster de calcul



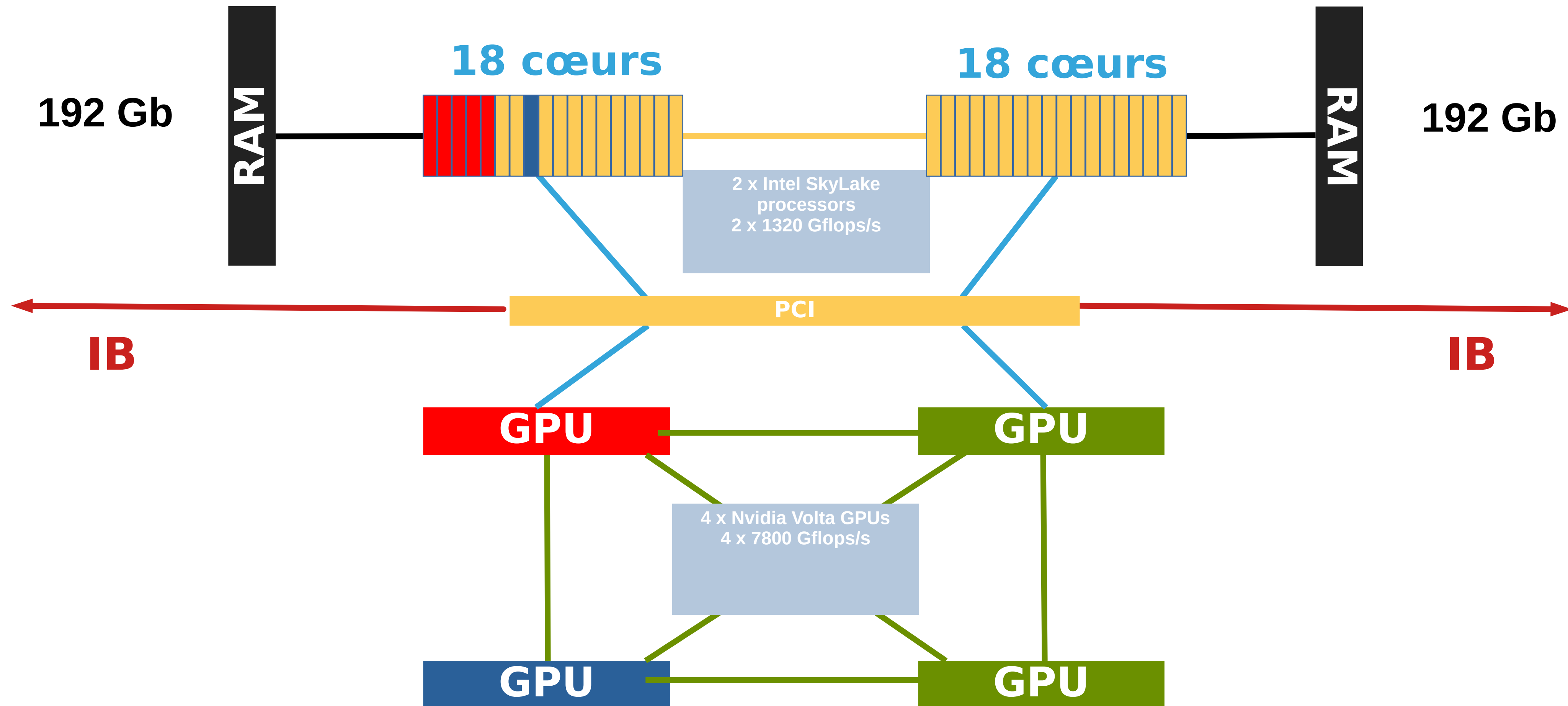
Plusieurs sortes de nœuds

Un nœud avec des gpus...

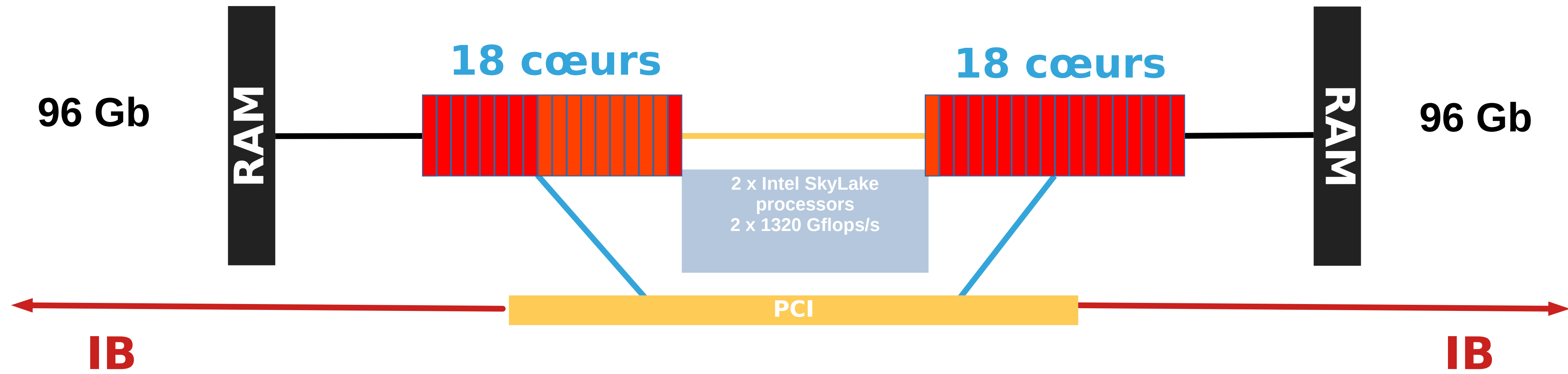




Des nœuds partagés...

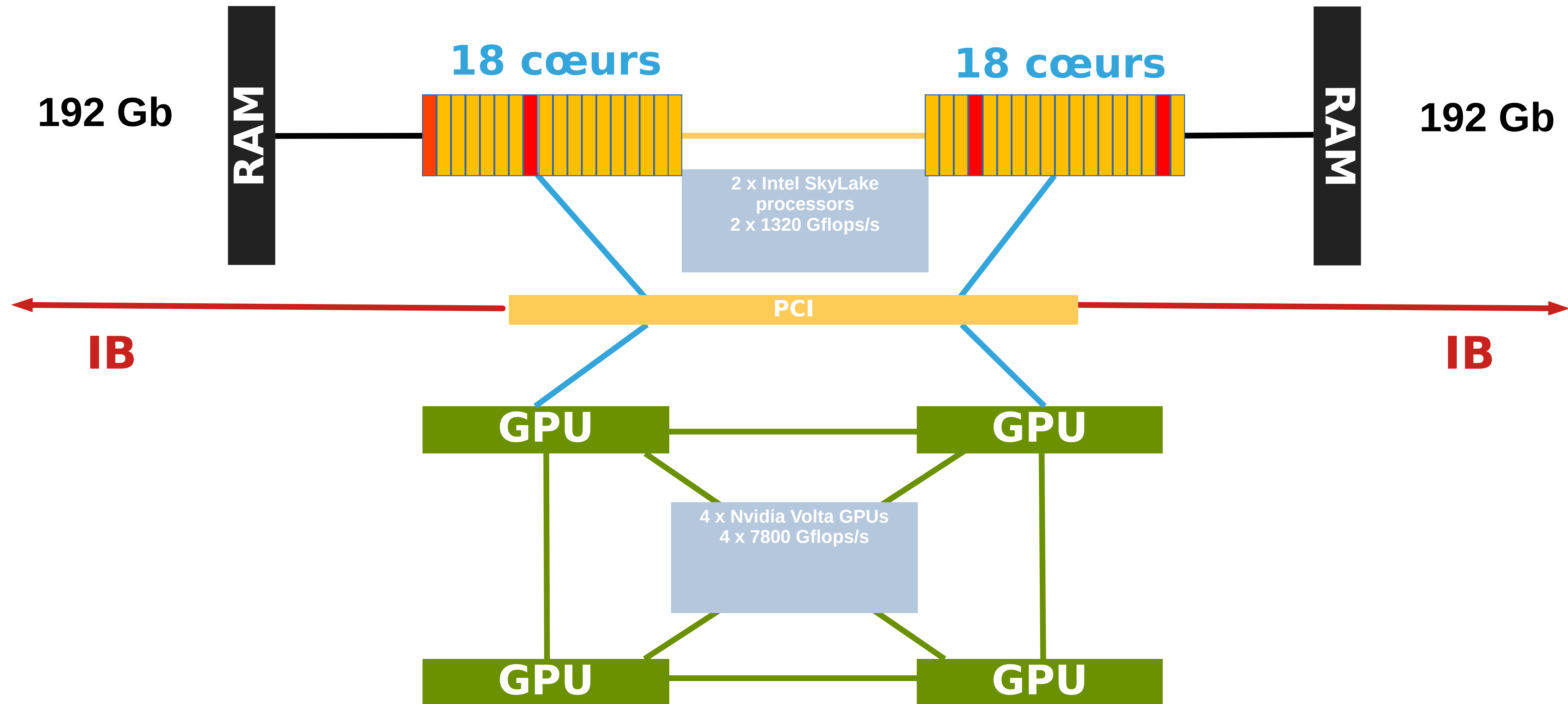


Deux utilisateurs se partagent le nœud



Un utilisateur a lancé un job mpi, 36 processus

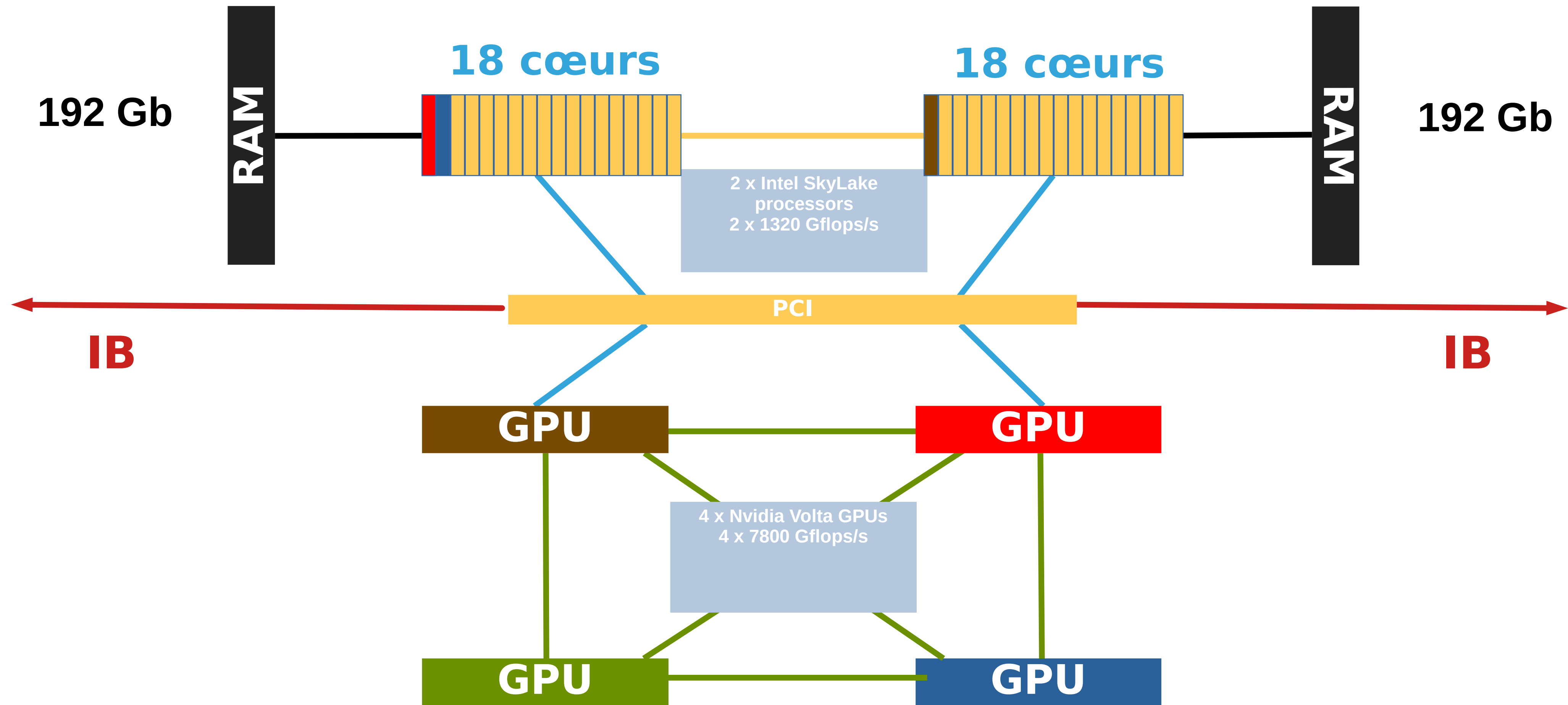
Ça ne se passe pas toujours bien



Placement permet de le détecter

Un utilisateur a réservé 4 GPUS
mais ne s'en sert pas !

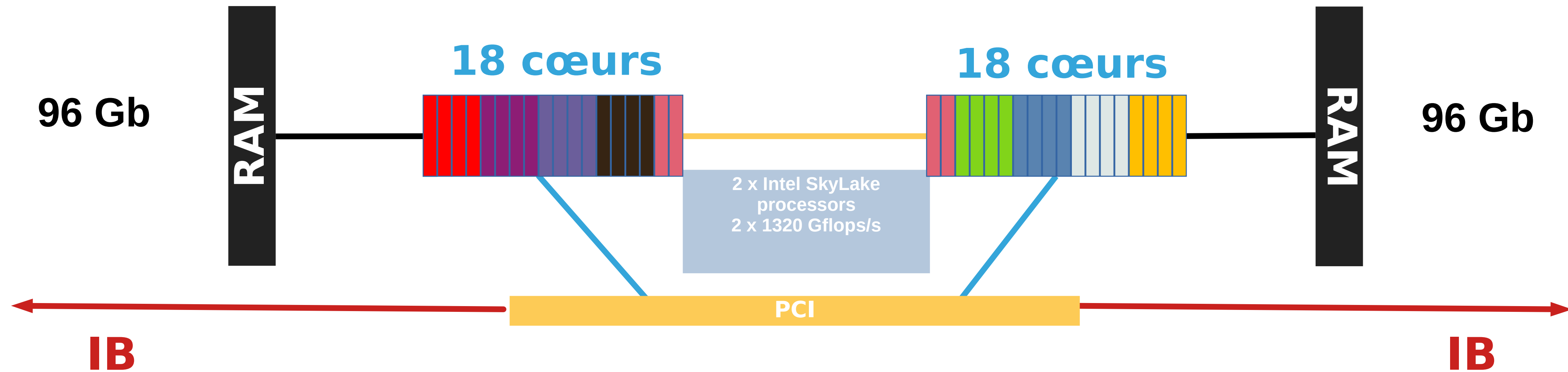
Ça ne se passe pas toujours bien



Placement permet de le détecter

Trois utilisateurs se partagent le nœud
Les transferts de données ne sont pas optimaux

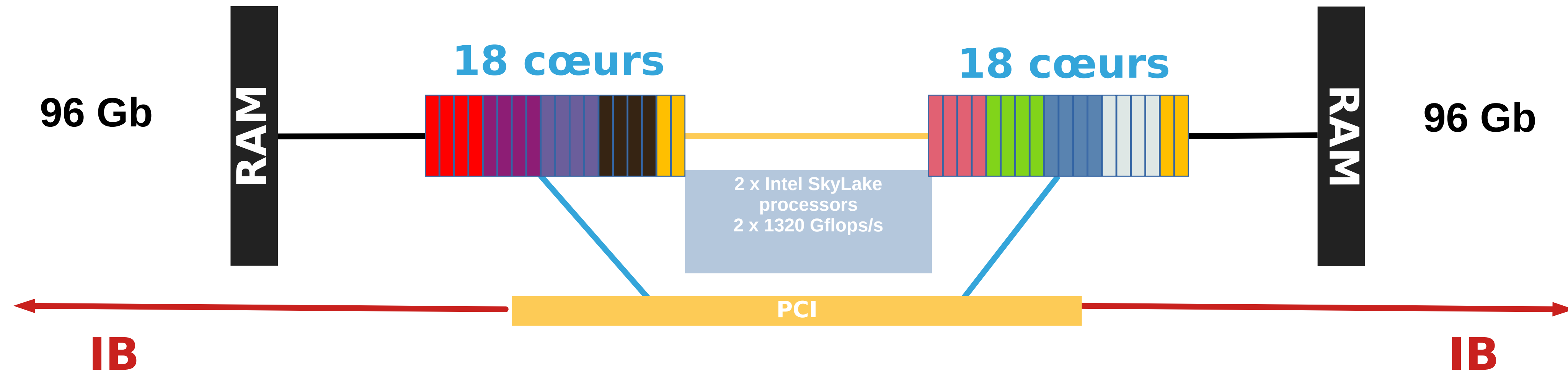
Ça ne se passe pas toujours bien



Le nœud est dédié (mode exclusif)
Un utilisateur a lancé un job hybride en mode dépeuplé
8 processus utilisant 4 threads
Le cinquième est **à cheval sur deux sockets**
Ce n'est pas optimal

Placement permet de l'empêcher
Placement permet de le détecter

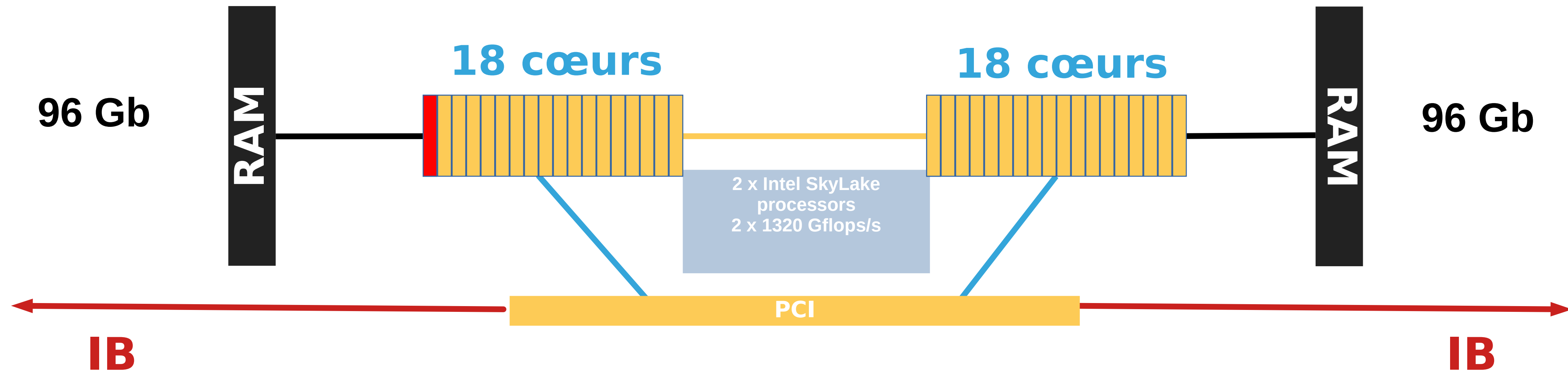
On aimerait bien avoir ça !



Le nœud est dédié (mode exclusif)
Un utilisateur a lancé un job hybride en mode dépeuplé
8 processus utilisant 4 threads
Cette fois les deux sockets sont équilibrés !

Placement permet de le contrôler aisément
Placement permet de le vérifier

Ça se passe parfois très mal



Le nœud est dédié (mode exclusif)
Un utilisateur a lancé un job mpi
36 processus mono-threadés
Ils utilisent tous le cœur **zéro**
La performance est catastrophique

Placement permet de l'empêcher
Placement permet de le détecter

Quel est mon objectif ?

Utiliser le plus de cœurs possibles avec des processus à 4 threads

```
# placement --ascii 9 4
PLACEMENT_ERROR_FOUND
PLACEMENT_ERROR - One task is straddling two sockets !
Please lower the number of tasks/node, max is 8
```

...je peux le forcer... :

```
# placement --ascii --mode=compact 9 4
P AAAABBBBCCCCDDDEE EFFFFFFGGGGHHHHIII
```

...voire même :

```
# placement --ascii --mode=compact --hyper 9 4
P AABBCDDEEFFGGHHII .....
L AABBCDDEEFFGGHHII .....
```

Quel est mon objectif ?

*Travailler en dépeuplé, 16 processus/4 threads (64 cœurs utilisés)
Équilibrer la charge des deux sockets :*

```
# placement --ascii 16 4
  S0----- S1-----
P AABBCDDEEFFGGHH.. IIJJKKLLMMNNOOPP..
L AABBCDDEEFFGGHH.. IIJJKKLLMMNNOOPP..
```

*Travailler en dépeuplé, 8 processus/4 threads (32 cœurs utilisés)
Équilibrer la charge des deux sockets :*

```
# placement --ascii 8 4
  S0----- S1-----
P AAAABBBBCCCCDDDD.. EEEEEFFFFGGGGHHHH..
```

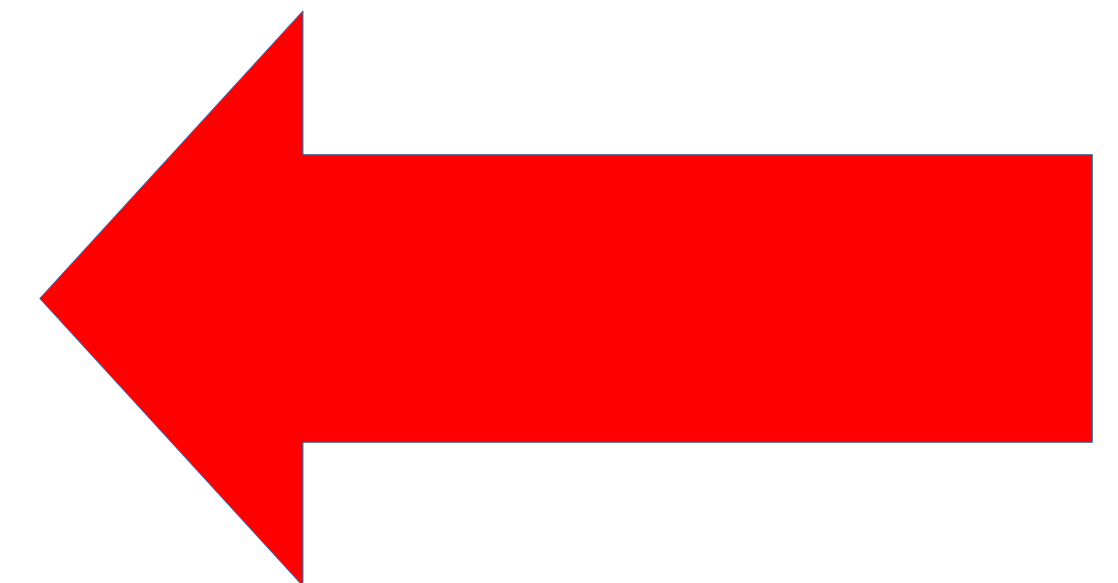
Quel est mon objectif ?

*Travailler en dépeuplé, 16 processus/4 threads (64 cœurs utilisés)
Équilibrer la charge des deux sockets :*

```
# placement --ascii 16 4
  S0----- S1-----
P AABBCDDEEFFGGHH.. IIJJKKLLMMNN00PP..
L AABBCDDEEFFGGHH.. IIJJKKLLMMNN00PP..
```

*Travailler en dépeuplé, 8 processus/4 threads (32 cœurs utilisés)
Équilibrer la charge des deux sockets :*

```
# placement --ascii 8 4
  S0----- S1-----
P AAAABBBBCCCCDDDD.. EEEFFFFFFGGGGHHHH..
```



Quel est mon lanceur ?

Je travaille avec srun, je place de manière explicite

```
# srun --cpu_bind=mask_cpu:0xf,0xf0,0xf00,0xf000,0x3c0000, \
    0x3c00000,0x3c000000,0x3c0000000 mon_appli.exe
```

...c'est plus cool de l'écrire comme ça :

```
# srun $(placement --srun 8 4) mon_appli.exe
```

Quel est mon lanceur ?

Je travaille avec mpirun, j'utilise numactl :

```
# numactl --physcpubind=0-3,4-7,8-11,12-15,18-21,22-25,\  
26-29,30-33 mpirun -np 8 mon_appli.exe
```

...c'est plus cool de l'écrire comme ça :

```
# numactl $(placement --numactl 8 4) mpirun -np 8 mon_appli.exe
```


Quel est mon lanceur ?

Je travaille avec mpihydra.exe (intelmpi) :

```
# export I_MPI_PIN_DOMAIN=[f,f0,f00,f000,3c0000,3c00000,3c000000,3c0000000]
# mpiexec.hydra -np 8 mon_appli.exe
```

...c'est plus cool de l'écrire comme ça :

```
# eval $(placement --intel_pin_domain 8 4)
# mpiexec.hydra -np 8 mon_appli.exe
```

Vérifier que ça fonctionne !

Lorsque le programme est en train de tourner taper :

```
# placement --host olympevolta9 --memory  
# placement --jobid=12345 --memory  
# placement --checkme --memory
```

Les processus sont bien déployés

Les sockets sont équilibrés

Tous les cœurs alloués sont utilisés

La mémoire n'est pas saturée

Les bancs de mémoire sont affectés au plus proche

Les GPUs sont utilisés

```
[manu@olympeloin1 2020.2-mpi]$ placement --jobid 523614 --memory  
#jobid 523614  
#host olympevolta9  
000000000000000000000000 000000000000000000000000  
00000000001111111111 112222222222223333333333  
012345678901234567 890123456789012345  
PID TID %CPU %MEM SESS  
A 38416 38416 A ..... 98.9 0.8 38407  
A 38582 38582 A ..... 89.5 0.8  
A 38583 38583 A ..... 89.5 0.8  
B 38417 38417 B ..... 99.1 0.7  
B 38584 38584 B ..... 89.5 0.7  
B 38585 38585 B ..... 89.5 0.7  
C 38418 38418 C ..... 99.1 0.6  
C 38590 38590 C ..... 89.5 0.6  
C 38591 38591 C ..... 89.6 0.6  
D 38419 38419 D ..... 99.1 0.7  
D 38605 38605 D ..... 0.0 0.7  
D 38592 38592 D ..... 89.6 0.7  
D 38593 38593 D ..... 89.5 0.7  
E 38420 38420 E ..... 99.1 0.6  
E 38596 38596 E ..... 89.6 0.6  
E 38597 38597 E ..... 89.6 0.6  
F 38421 38421 F ..... 99.1 0.7  
F 38594 38594 F ..... 89.5 0.7  
F 38595 38595 F ..... 89.5 0.7  
G 38422 38422 G ..... 99.1 0.6  
G 38588 38588 G ..... 89.6 0.6  
G 38589 38589 G ..... 89.6 0.6  
H 38423 38423 H ..... 99.1 0.7  
H 38586 38586 H ..... 89.5 0.7  
H 38587 38587 H ..... 89.5 0.7  
  
DISTRIBUTION of the MEMORY among the sockets  
A 3.4Gb ***** 99% 0%  
B 2.7Gb ***** 99% 0%  
C 2.6Gb ***** 99% 0%  
D 2.7Gb ***** 99% 0%  
E 2.6Gb ***** 0% 99%  
F 2.8Gb ***** 0% 99%  
G 2.6Gb ***** 0% 99%  
H 2.7Gb ***** 0% 99%  
  
CPU  
USE ***** 5%  
MEMORY ***** 7%  
POWER ***** 20%  
PROCESSES AB  
USED MEMORY ■  
  
GPU 1  
USE ***** 24%  
MEMORY ***** 7%  
POWER ***** 20%  
PROCESSES CD  
USED MEMORY ■  
  
GPU 2  
USE ***** 24%  
MEMORY ***** 7%  
POWER ***** 24%  
PROCESSES EF  
USED MEMORY ■  
  
GPU 3  
USE ***** 23%  
MEMORY ***** 7%  
POWER ***** 20%  
PROCESSES GH  
USED MEMORY ■
```

Vérifier que ça fonctionne !

Lorsque le programme est en train de tourner

```
# placement --checkme
```

Le processus tourne sur le gpu le plus proche de lui

Chaque GPU calcule pour un processus

Chaque GPU « voit » les quatre processus

```
[manu@olympelogein1 bench_2v1]$ placement --checkme
jobid 523629
host olympevolta9
000000000000000000 000000000000000000
000000000001111111 112222222222333333
012345678901234567 890123456789012345
PID TID %CPU %MEM SESS
A 48323 48323 A..... 91.5 0.3 48314
B 48324 48324 B..... 91.7 0.3
C 48325 48325 C..... 91.6 0.3
D 48326 48326 D..... 91.4 0.3

GPU
USE ***** 79%
MEMOR ** ..... 13%
POWER ***** 53%
PROCESSES ABCD
USED MEMORY █

GPU 1
USE ***** 96%
MEMORY *** ..... 16%
POWER ***** 82%
PROCESSES ABCD
USED MEMORY █

GPU 2
USE ***** 95%
MEMORY *** ..... 16%
POWER ***** 80%
PROCESSES ABCD
USED MEMORY █

GPU 3
USE ***** 96%
MEMORY *** ..... 16%
POWER ***** 84%
PROCESSES ABCD
USED MEMORY █
```

Prérequis :

Avoir accès à python 3 ($> 3.5.3$)
éventuellement via un module

Accéder aux nœuds de calcul en ssh lorsque le job est lancé

Pas besoin d'être root

```
./install.sh
```

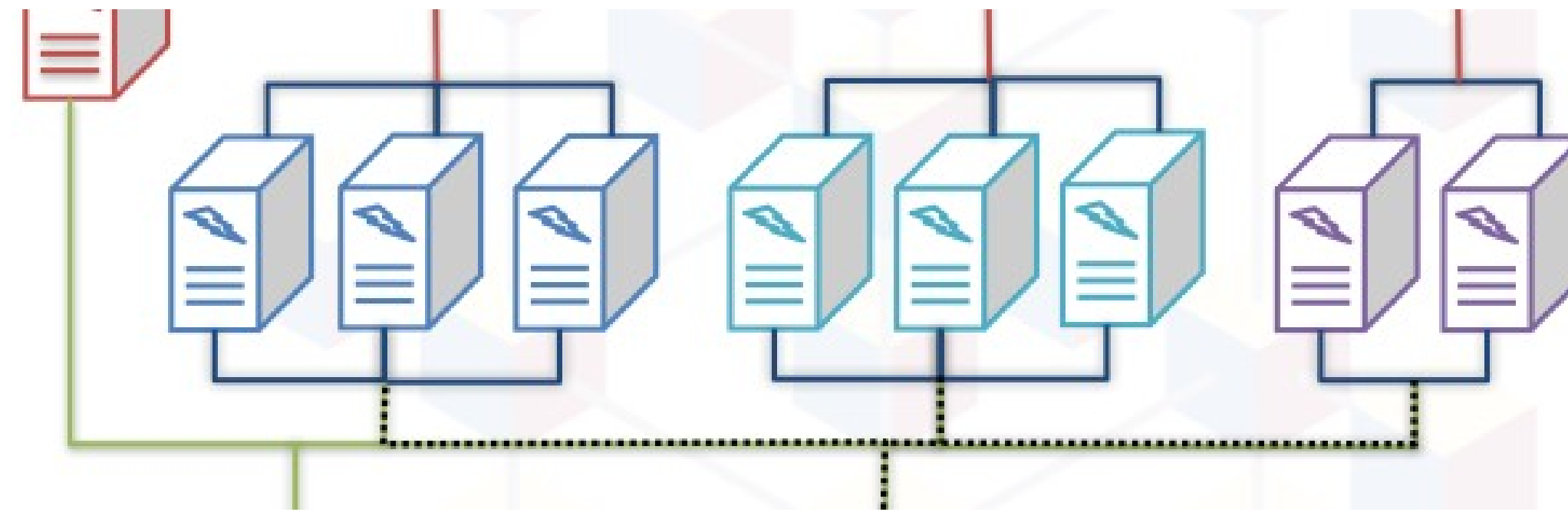
Configuration – voir après

```
cp bin/placement-dist bin/placement
```

et

Éditer **bin/placement** (*appeler la bonne version de python*)

Placement.conf



Décrire les trois **types** de nœuds

Correspondance entre **noms** des nœuds et **types** des nœuds

Configurer placement...

1/ Les types de nœuds :

```
[sequana]
SOCKETS_PER_NODE: 2
CORES_PER_SOCKET: 18
HYPERTHREADING: True
THREADS_PER_CORE: 2
MEM_PER_SOCKET: 98304
IS_SHARED: False
```

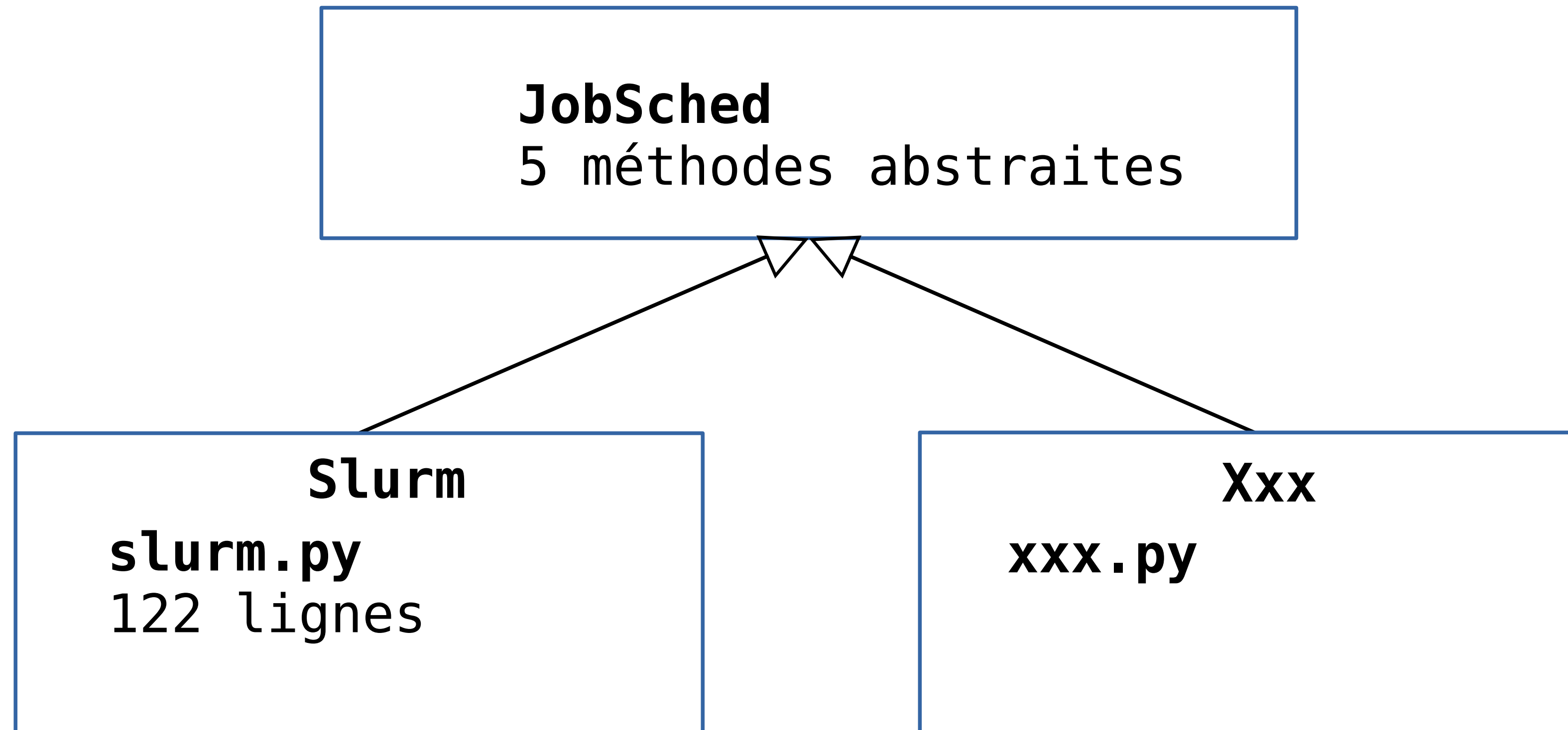
```
[sequana_gpu]
SOCKETS_PER_NODE: 2
CORES_PER_SOCKET: 18
HYPERTHREADING: True
THREADS_PER_CORE: 2
MEM_PER_SOCKET: 192000
IS_SHARED: True
GPUS: 0-1,2-3
```

2/ La correspondance noms – type des nœuds :

```
# The cluster compute nodes:  
olympcomp[0-356]: sequana
```

```
# Compute nodes with gpu  
olympvolta[0-11]: sequana_gpu
```


1°)



2°)

Fichier **front.py**
classe **FrontNode** :

Ligne 45 remplacer **Slurm()** par **Xxx()**

<https://github.com/calmip/placement/>

emmanuel.courcelle@toulouse-inp.fr

...merci pour ce café !