**Légende :**
- Accès
- Connexion frontale / nœuds
- Interconnexion rapide
- Accès disques
- ..... 40 Go/s

- Nœud de login et de transfert
- Nœuds de calcul CPU
- Nœuds de calcul GPU
- Nœuds grande mémoire

VPN

Espace permanent
/users/$GROUPE/$USER

Espace temporaire
/tmpdir/$USER

## Un exemple : Olympe (CALMIP)
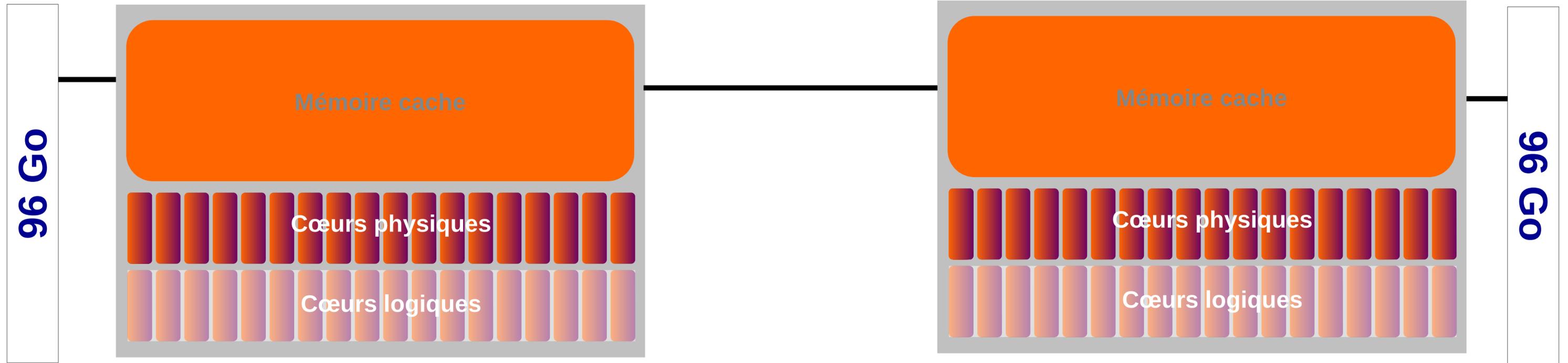
3 frontales

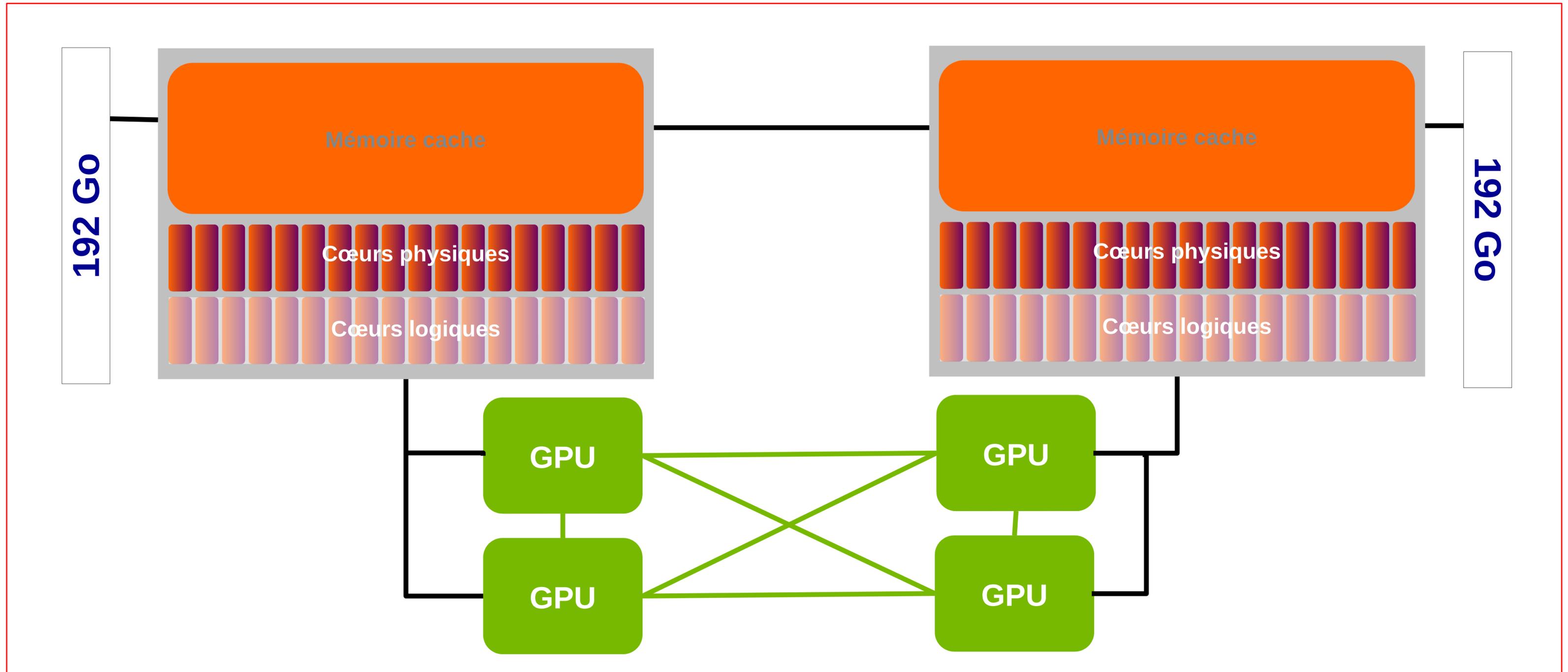360 nœuds bi-processeurs
    2 x 18 cœurs
    192 Go

2 nœuds grande mémoire
    2 x 18 cœurs
    1500 Go

12 nœuds GPU
    2 x 18 cœurs CPU
    384 Go
    4 GPUs connectés :
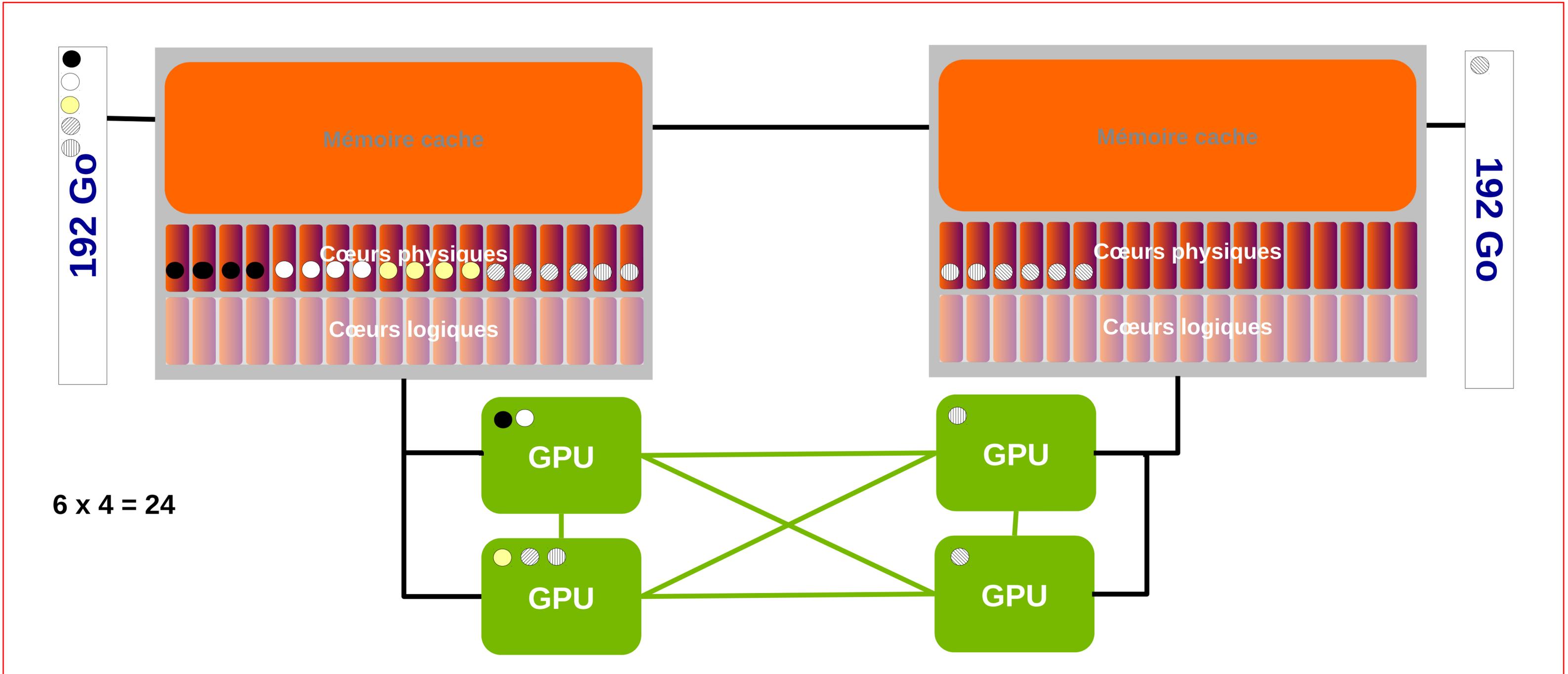        Avec les CPUS en PCI (2 x 2)
        Entre eux en nvLink

Les nœuds sont **attribués** à un **utilisateur exclusif** durant un **temps limité**
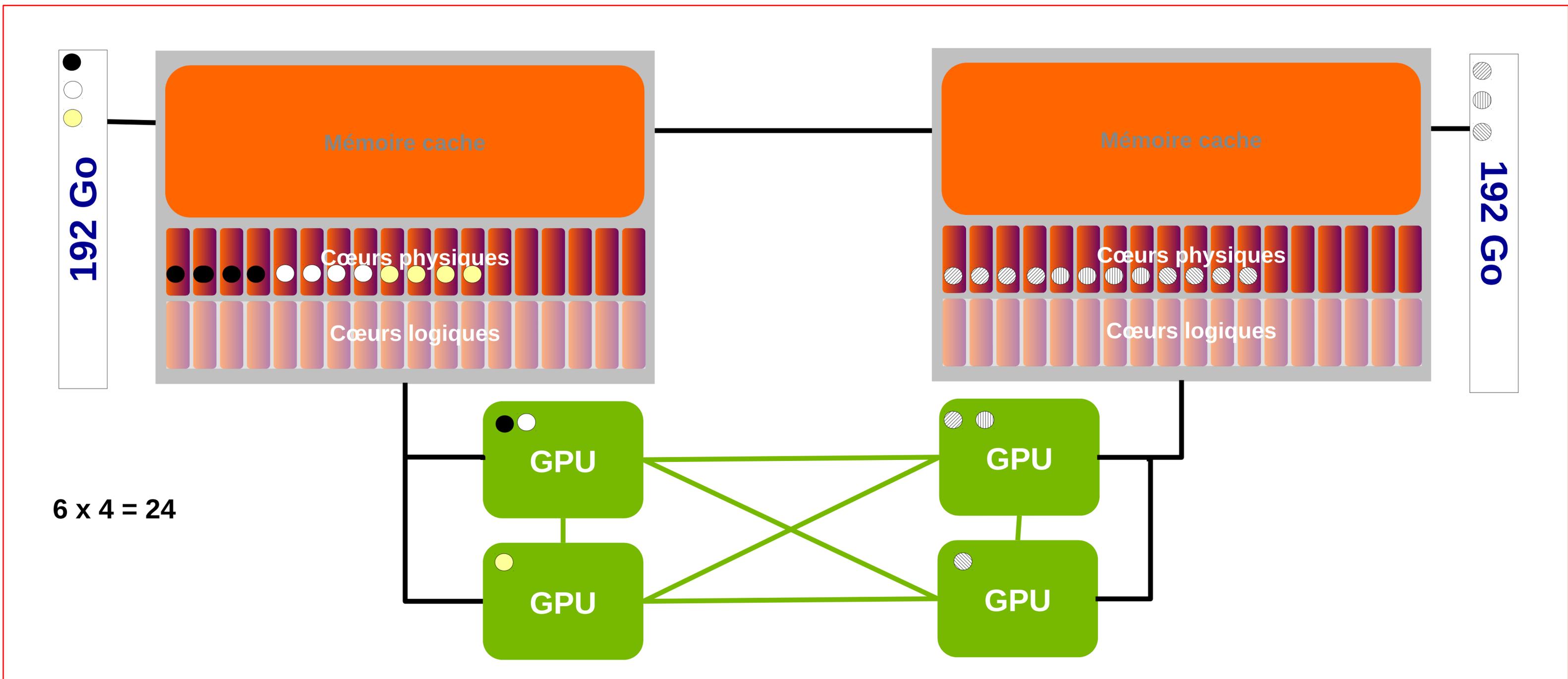
# UN MEILLEUR PLACEMENT

● **Processus A (4 threads)**  ○ **Processus B (4 threads)**  ◉ **Processus C (4 threads)**  ◍ **Processus D (4 threads)**  ◫ **Processus E (4 threads)**  ◍ **Processus F (4 threads)**

**Contrôler** le placement (cœurs) lors de l'exécution d'une application

Vérifier **visuellement** que les instructions de placement seront pertinentes

**Vérifier** que le placement est correct lorsque l'application est **en exécution**

# 1/ **Variables d'environnement** (1 processus, multithreadé)

```
export GOMP_CPU_AFFINITY="0 1 18 19"
export KMP_AFFINITY="granularity=fine,explicit,proclist=[0,1,18,19]"
```

# 2/ **numactl**

*Pour un code mpi/openmp , une commande différente pour chaque rang !*

Rang 0 → `numactl --physcpubind=0-3 my_exe`

Rang 1 → `numactl --physcpubind=18-21 my_exe`

# 3/ **srun (slurm)**

*Des switches de ligne de commande*

```
srun --cpu_bind=mask_cpu:0xf,0x3c0000 my_exe
```

**Exemples de la démo :**

**+ 8 %** lorsque le placement est correct

**+ 15 %** lorsqu'on utilise les GPU

# Galerie de photos

```
manu@olympelogin1.bullx: ~/tmpdir/DEMO/JCAD2018-placement/3

[manu@olympelogin1 3]$ placement --jobid=63481
jobid 63481
host  olympevolta10
             0000000000000000000 000000000000000000
             0000000000011111111 1122222222222333333
             0123456789012345 890123456789012345
     PID   TID                                      %CPU %MEM
A  72529  72529 A................. .................. 99.7  0.0
   GPU 0
USE             ****************. 97%
MEMORY          ................. 2%
POWER           *****............. 26%
PROCESSES       A
USED MEMORY     ■

   GPU 1
USE             ................. 0%
MEMORY          ................. 0%
POWER           **............... 13%
PROCESSES
USED MEMORY

   GPU 2
USE                              ................. 0%
MEMORY                           ................. 0%
POWER                            **............... 13%
PROCESSES
USED MEMORY

   GPU 3
USE                              ................. 0%
MEMORY                           ................. 0%
POWER                            **............... 13%
PROCESSES
USED MEMORY


[manu@olympelogin1 3]$
```

# Galerie de photos

# Galerie de photos

# Galerie de photos



```
manu@olympelogin1.bullx: ~/tmpdir/DEMO/JCAD2018-placement/3    —  +  x

[manu@olympelogin1 3]$ placement --jobid=63474
jobid 63474
host  olympevolta6
                    00000000000000000 00000000000000000
                    00000000011111111 11222222222233333
                    01234567890123456 78901234567890123
      PID    TID                                          %CPU %MEM
A  64019  64019  .A................. .................   96.2  1.2
B  64020  64020  .....B............. .................   96.0  1.0
C  64021  64021  ..........C........ .................   96.0  1.0
D  64022  64022  ...............D... .................   96.0  0.9
E  64023  64023  .................. ...E.............   96.0  1.0
F  64024  64024  .................. ........F.........   96.1  0.9
G  64025  64025  .................. .............G.....   96.1  0.9
H  64026  64026  .................. .................H....   96.2  1.0
   GPU 0
USE          ****************** 99%
MEMORY       *********......... 48%
POWER        *************..... 70%
PROCESSES    AB.
USED MEMORY  ■▖


   GPU 1
USE          ****************** 98%
MEMORY       *********.......... 47%
POWER        ****.............. 23%
PROCESSES    CD.
USED MEMORY  ■▖


   GPU 2
USE                            ****************** 99%
MEMORY                         ******........... 35%
POWER                          ****.............. 24%
PROCESSES                      EF.
USED MEMORY                    ■▖


   GPU 3
USE                            ****************** 99%
MEMORY                         *********......... 48%
POWER                          *********......... 50%
PROCESSES                      GH.
USED MEMORY                    ■▖



[manu@olympelogin1 3]$ ▯
```

# Galerie de photos

# Placement fonctionne :

**Avec** slurm

*Avec d'autres gestionnaires de batch ?*

→ **Contribuez !** 🙂

**Sans** gestionnaire de batch

**Avec** des processeurs intel (2 threads hardware / cœur)

*Avec d'autres processeurs ?*

→ **Contribuez !** 🙂

**https://github.com/calmip/placement**