



# Introduction to CROCO and Parallelisation

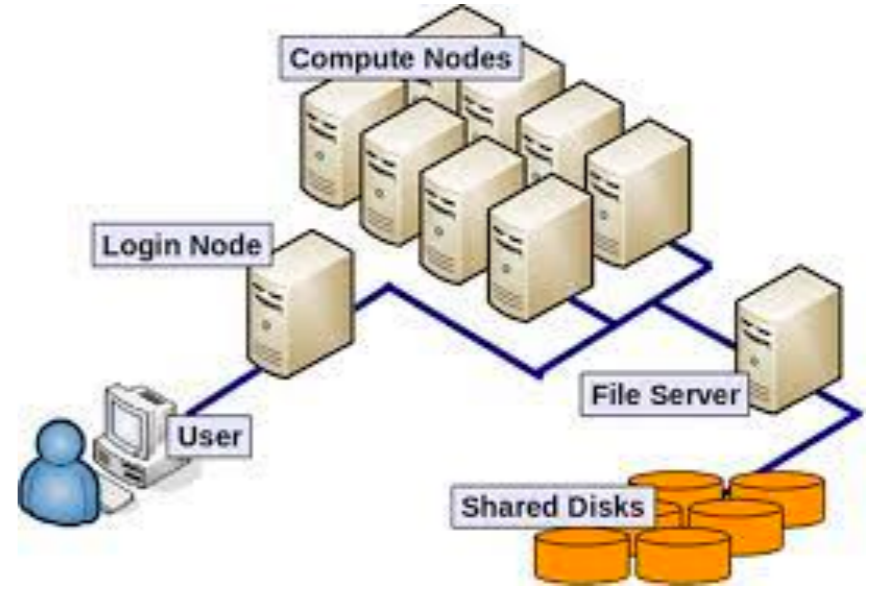
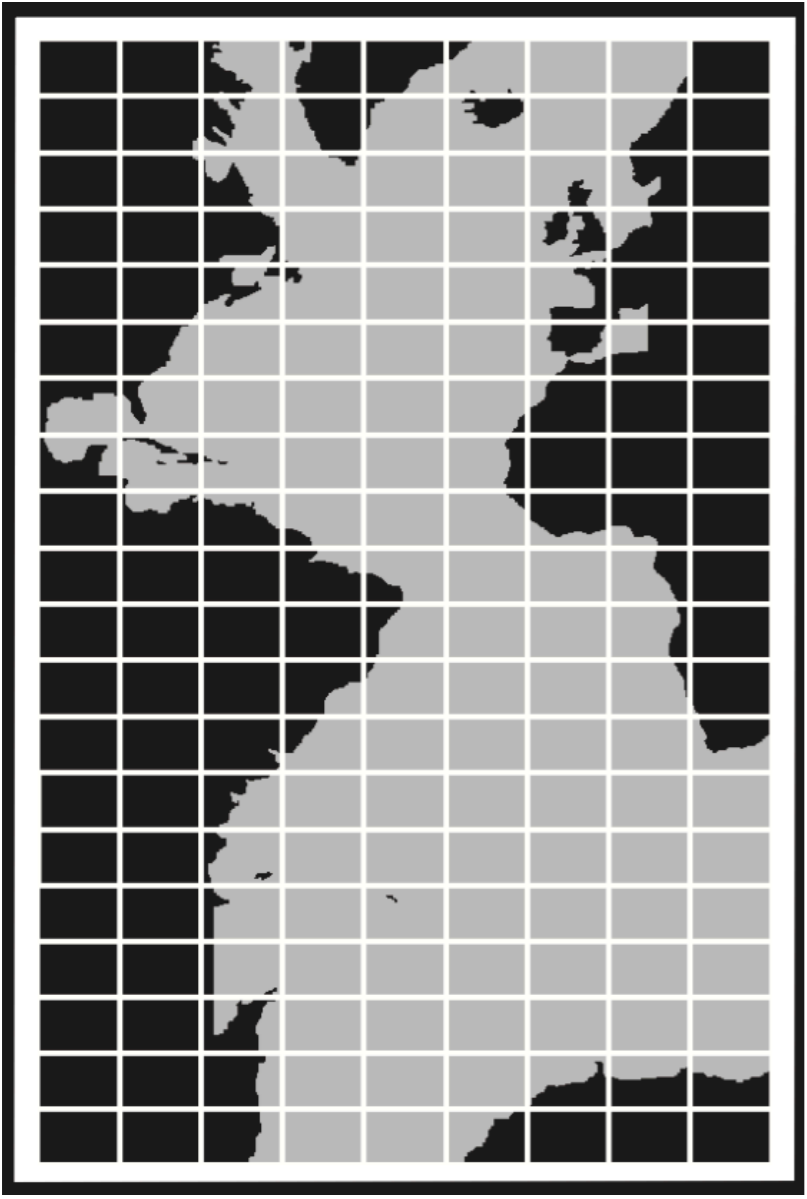
Rachid Benshila

**Concept et techniques**

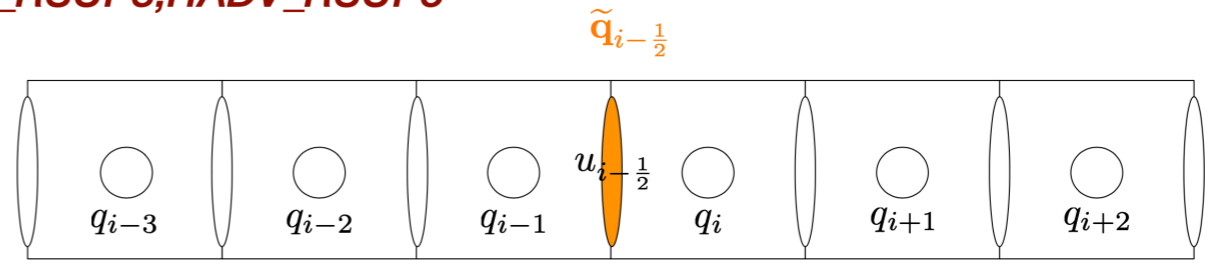
**Utilisation dans CROCO**

**Aspects connexes**

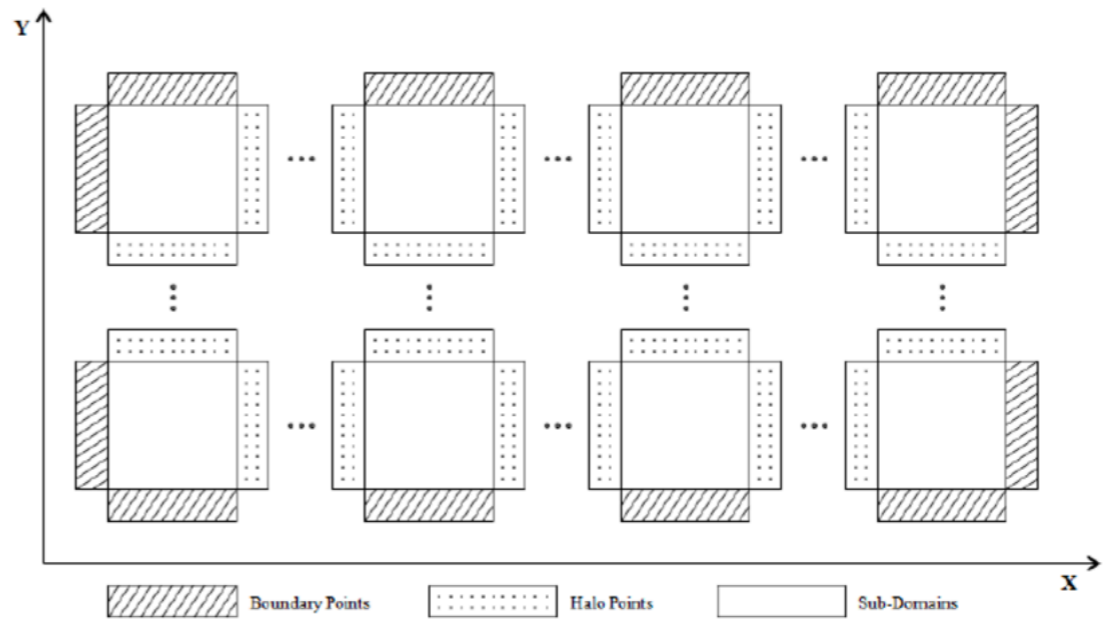
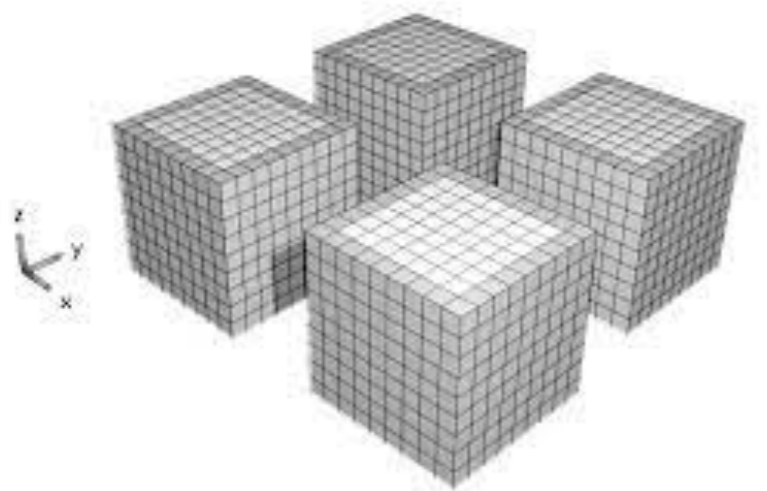
# Parallelisation : décomposition de domaine



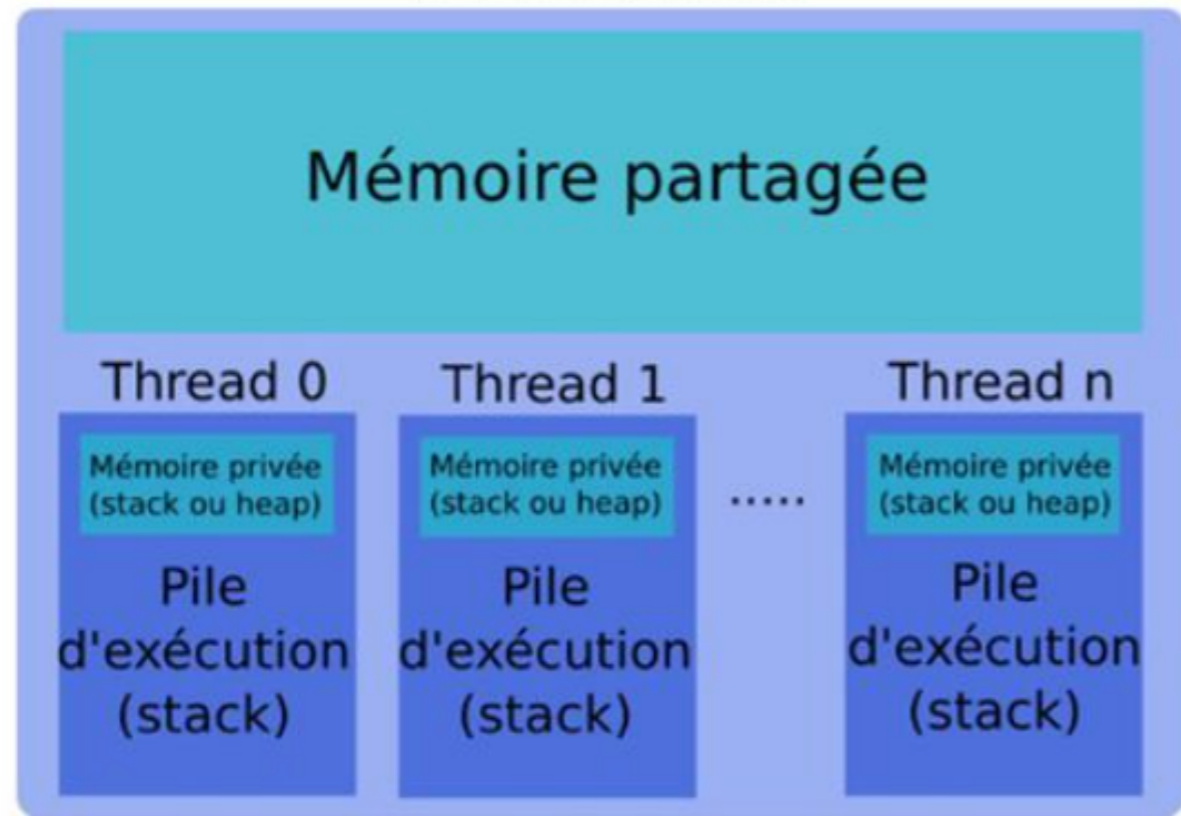
*HADV\_RSUP3, HADV\_RSUP5*



$$\partial_x(uq)|_{x=x_i} = \frac{1}{\Delta x_i} \left\{ u_{i+\frac{1}{2}} \tilde{q}_{i+\frac{1}{2}} - u_{i-\frac{1}{2}} \tilde{q}_{i-\frac{1}{2}} \right\}$$



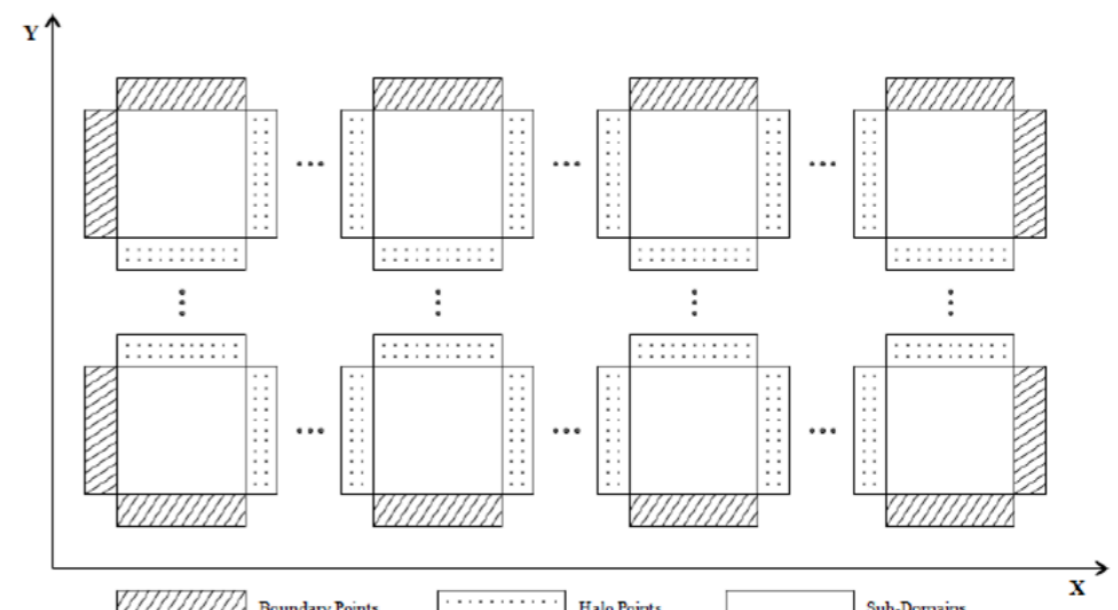
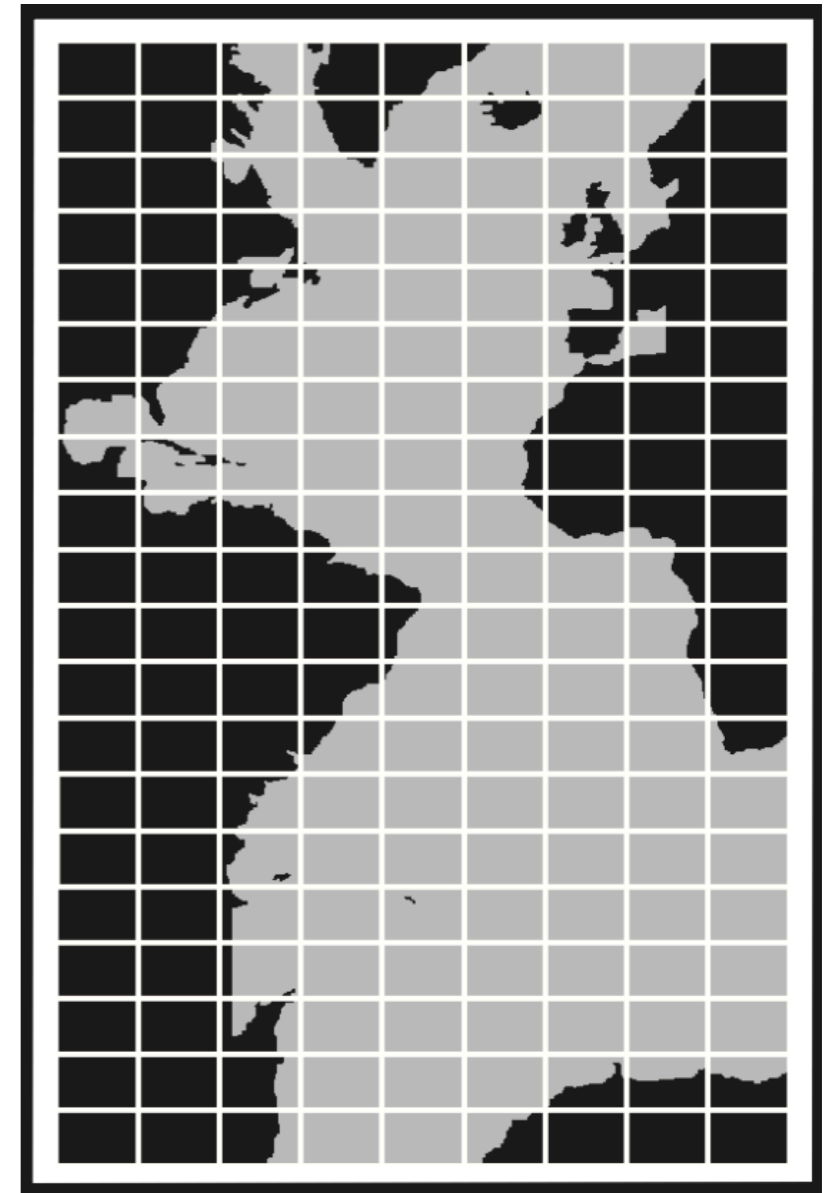
# Approche 1 : mémoire partagée



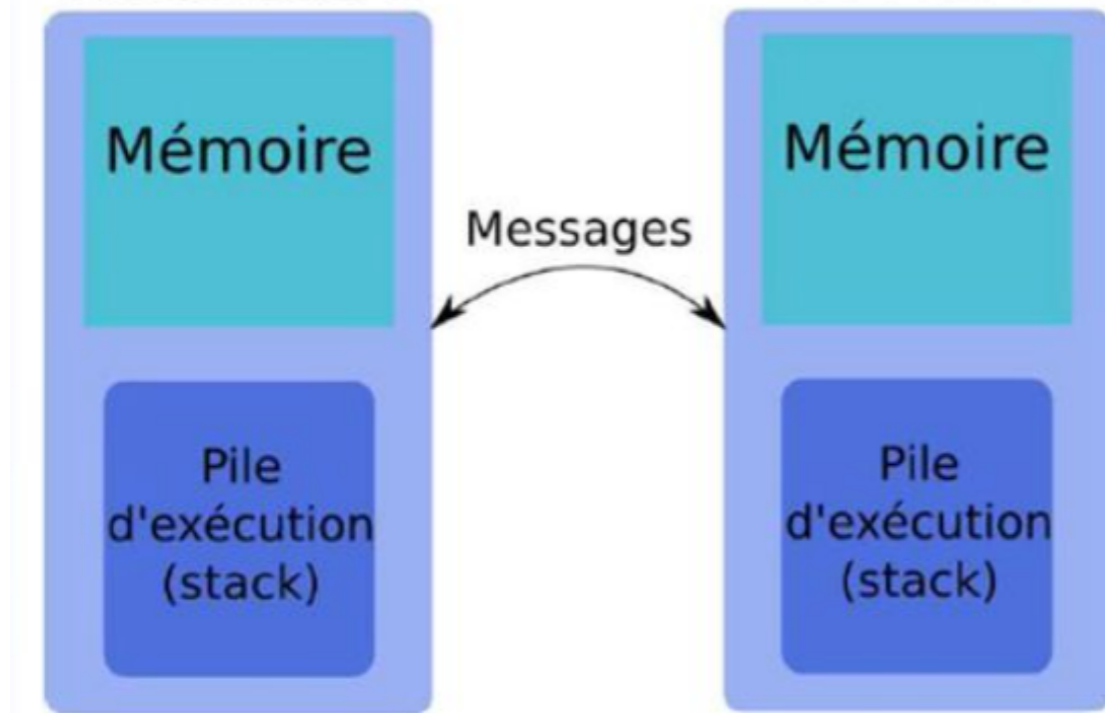
=> les coeurs de calcul ont accès à une mémoire commune

=> échanges par copie mémoire

**Standard OpenMP**  
**(Open Multi-Processing)**



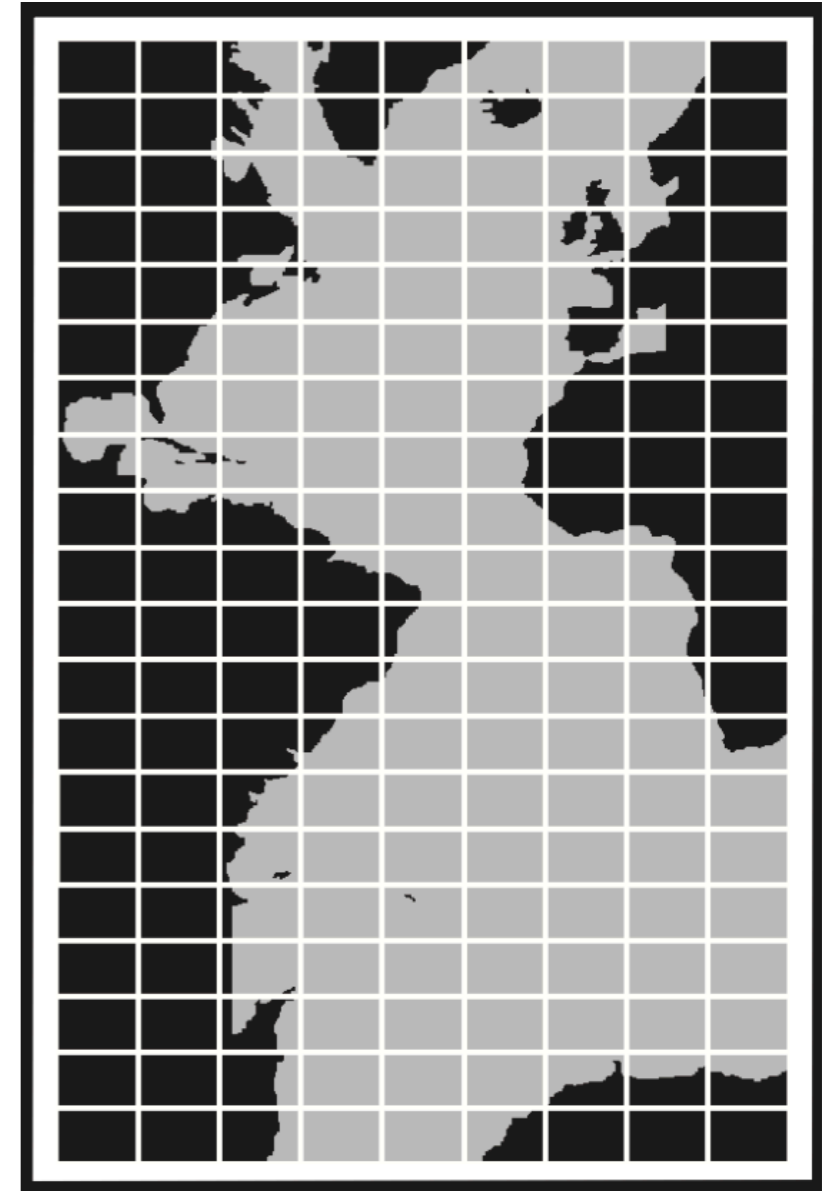
## Approche 2 : mémoire distribuée



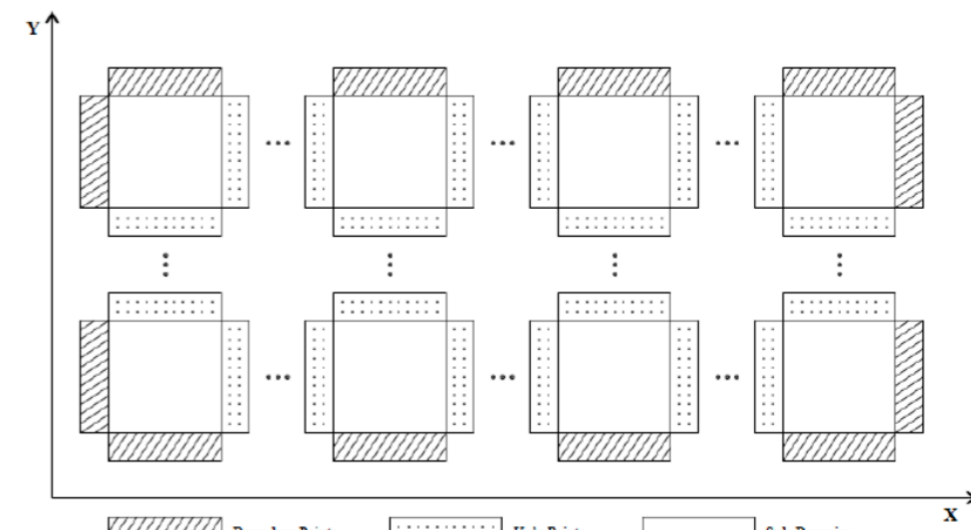
=> les coeurs de calcul n'ont pas accès à une mémoire commune

=> échanges par message réseau

=> en pratique MPI gère aussi la mémoire partagée



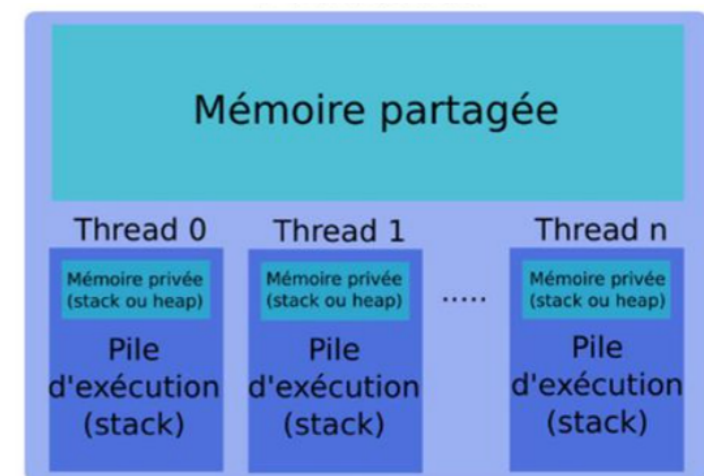
**Standard MPI**  
**(Message Passing Interface)**



# Implementation dans CROCO : OPENMP

- étape 1 : 2 fichiers à éditer
  - **param.h**  
spécifier la décomposition en x et y => NPP=4
  - **cppdefs.h** :  
activer OpenMP => **#define OPENMP**
- étape 2 : compiler  
**./jobcomp**
- étape 3 : exécuter
  - **export OMP\_NUM\_THREADS=4**  
spécifie le nombre de cores à l'environnement
  - **./croco**

```
! Domain subdivision parameters
! =====
!
! NPP          Maximum allowed number of parallel threads;
! NSUB_X,NSUB_E Number of SHARED memory subdomains in XI- and
!                                     ETA-directions;
! NNODES       Total number of MPI processes (nodes);
! NP_XI,NP_ETA Number of MPI subdomains in XI- and ETA-directions;
!
integer NSUB_X, NSUB_E, NPP
#ifdef MPI
integer NP_XI, NP_ETA, NNODES
parameter (NP_XI=1, NP_ETA=4, NNODES=NP_XI*NP_ETA)
parameter (NPP=1)
parameter (NSUB_X=1, NSUB_E=1)
#elif defined OPENMP
parameter (NPP=4)
# ifdef AUTOTILING
common/distrib/NSUB_X, NSUB_E
# else
parameter (NSUB_X=1, NSUB_E=NPP)
# endif
#else
parameter (NPP=1)
```



# Implementation dans CROCO : MPI

## - étape 1 : fichiers à éditer

### - **param.h**

spécifier la décomposition  
en x et y => NP\_XI, NP\_ETA

### - **cppdefs.h** :

activer MPI => *#define MPI*

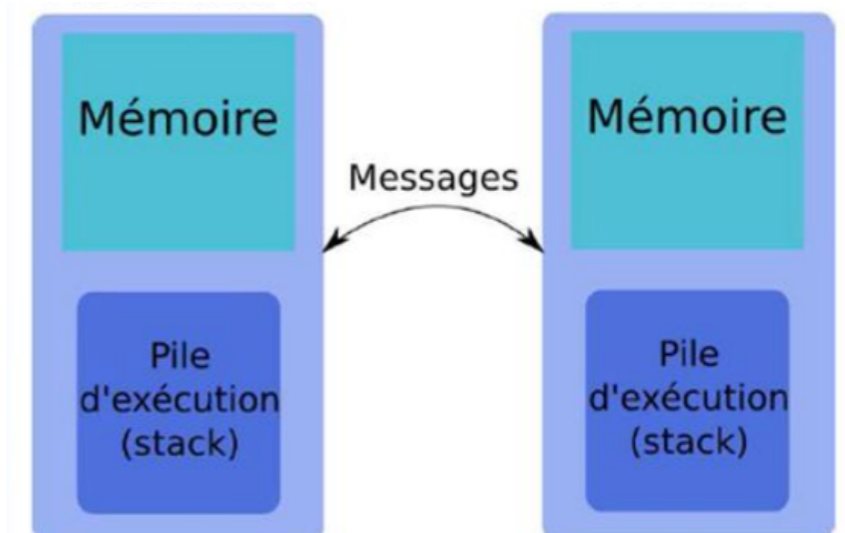
## - étape 2 : compiler

**./jobcomp**

## - étape 3 : exécuter

- **mpirun -n 4 ./croco**  
(ou mpiexec ou autre)

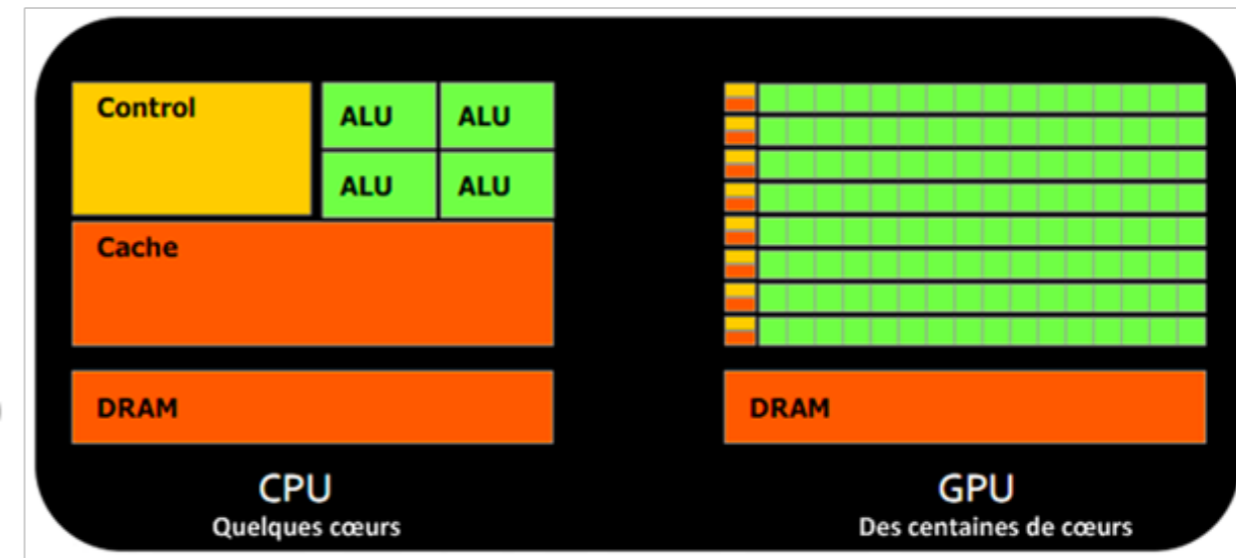
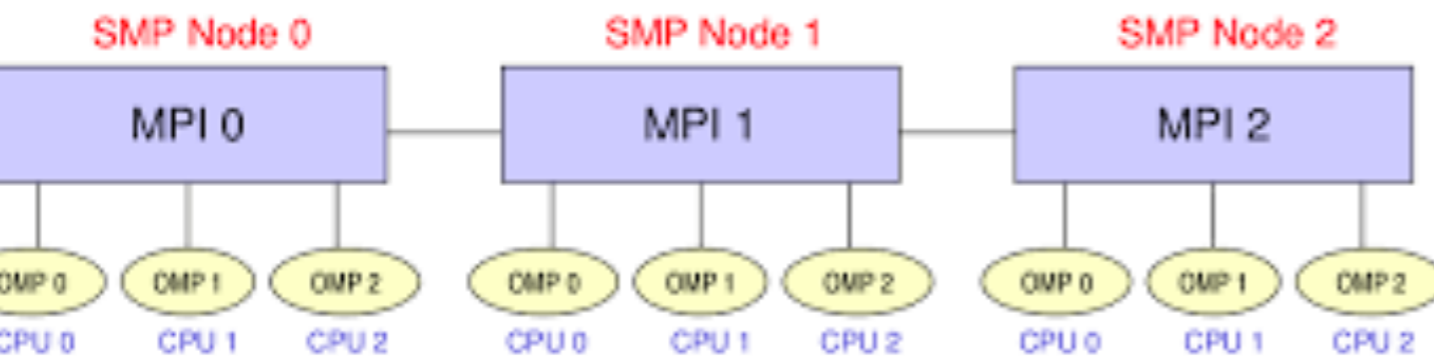
```
! Domain subdivision parameters
! =====
!
! NPP          Maximum allowed number of parallel threads;
! NSUB_X,NSUB_E Number of SHARED memory subdomains in XI- and
!                                     ETA-directions;
! NNODES       Total number of MPI processes (nodes);
! NP_XI,NP_ETA Number of MPI subdomains in XI- and ETA-directions;
!
integer NSUB_X, NSUB_E, NPP
#ifdef MPI
integer NP_XI, NP_ETA, NNODES
parameter (NP_XI=1, NP_ETA=4, NNODES=NP_XI*NP_ETA)
parameter (NPP=1)
parameter (NSUB_X=1, NSUB_E=1)
#elif defined OPENMP
parameter (NPP=4)
# ifdef AUTOTILING
common/distrib/NSUB_X, NSUB_E
# else
parameter (NSUB_X=1, NSUB_E=NPP)
# endif
#else
parameter (NPP=1)
```



# Résumé et perspectives

---

- 2 paradigmes disponibles MPI et OpenMP
- code à recompiler !!
- MPI à privilégier (plus utilisé)
- Direction ETA à privilégier



pas de version hybride MPI/OpenMP

version GPU en développement

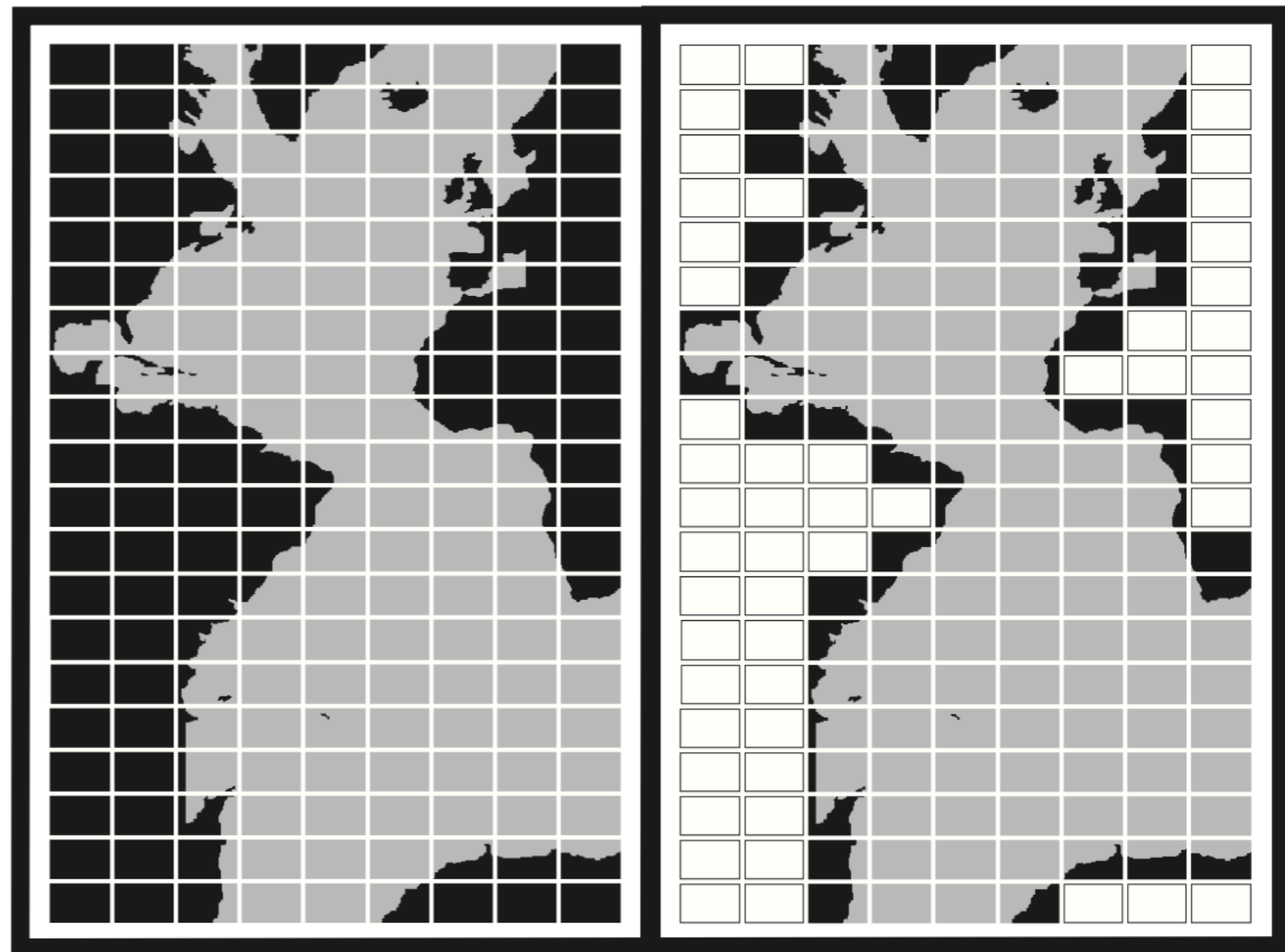


# Parallelisation : mais aussi ....

---

Quelque subtilités :

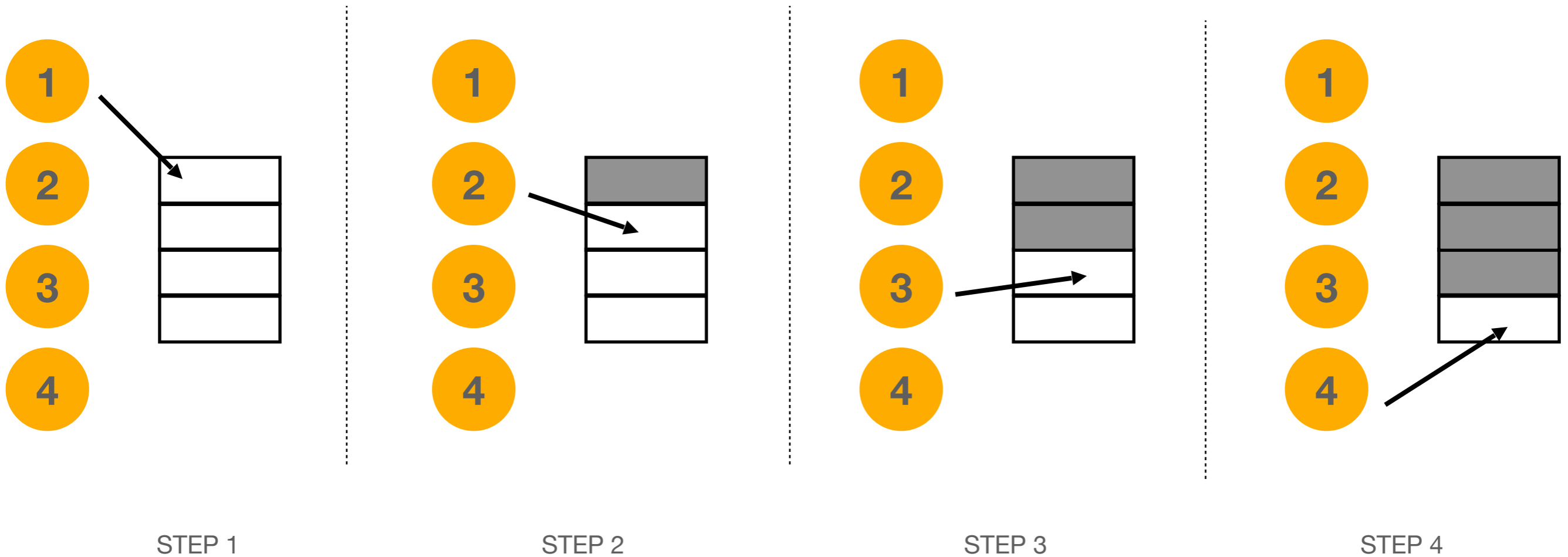
- le cas des fichiers de sortie (MPI)
- le cas des points terre



# Écriture de fichiers MPI 1/4 : par défaut

---

`mpirun -np 4 ./croco. (NP_ETA=4)`



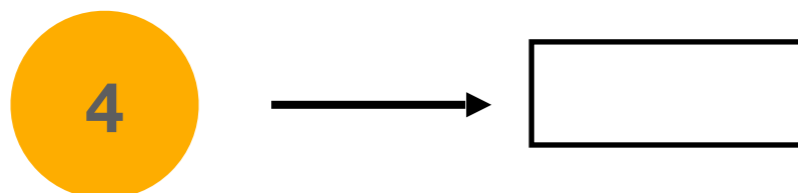
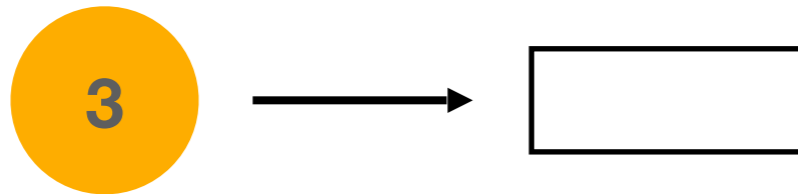
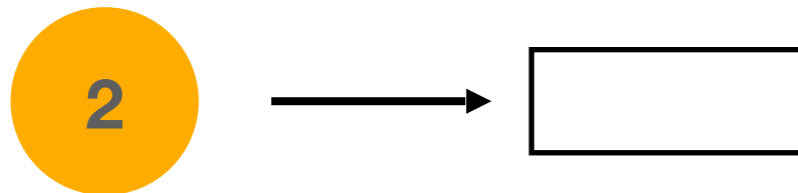
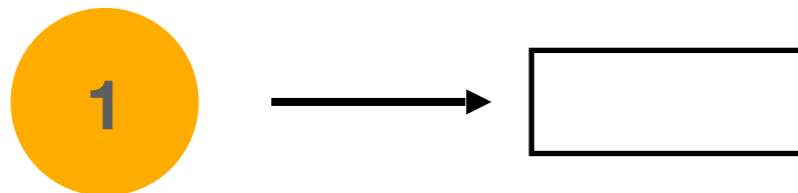
**Très peu efficace !!!!!!!**

# Ecriture de fichiers MPI 2/4: fichiers parallèles

---

## #define PARALLEL\_FILES

```
mpirun -np 4 ./croco. (NP_ETA=4)
```



Rapide mais bcp de fichiers à l'arrivée.

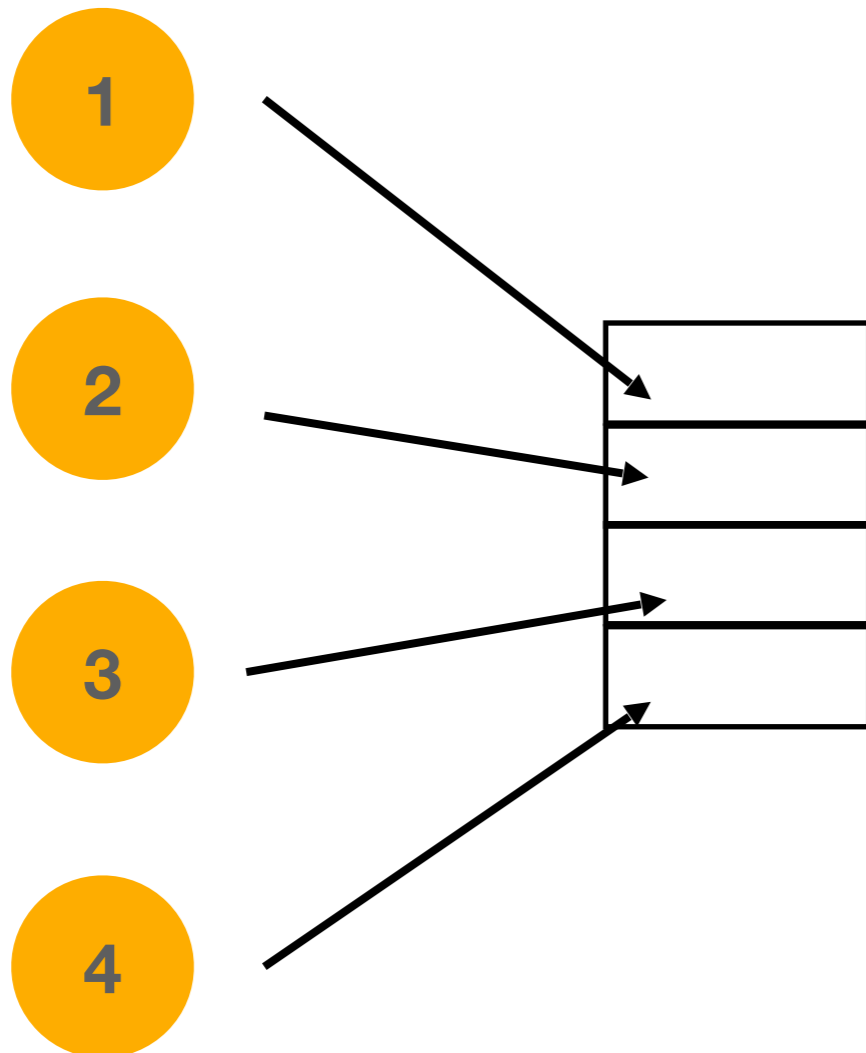
Besoin de les recombinaer (cf utilitaire ncjoin)

# Écriture de fichiers MPI 3/4 : écriture parallèle

---

**#define KEY NC4PAR**

`mpirun -np 4 ./croco. (NP_ETA=4)`

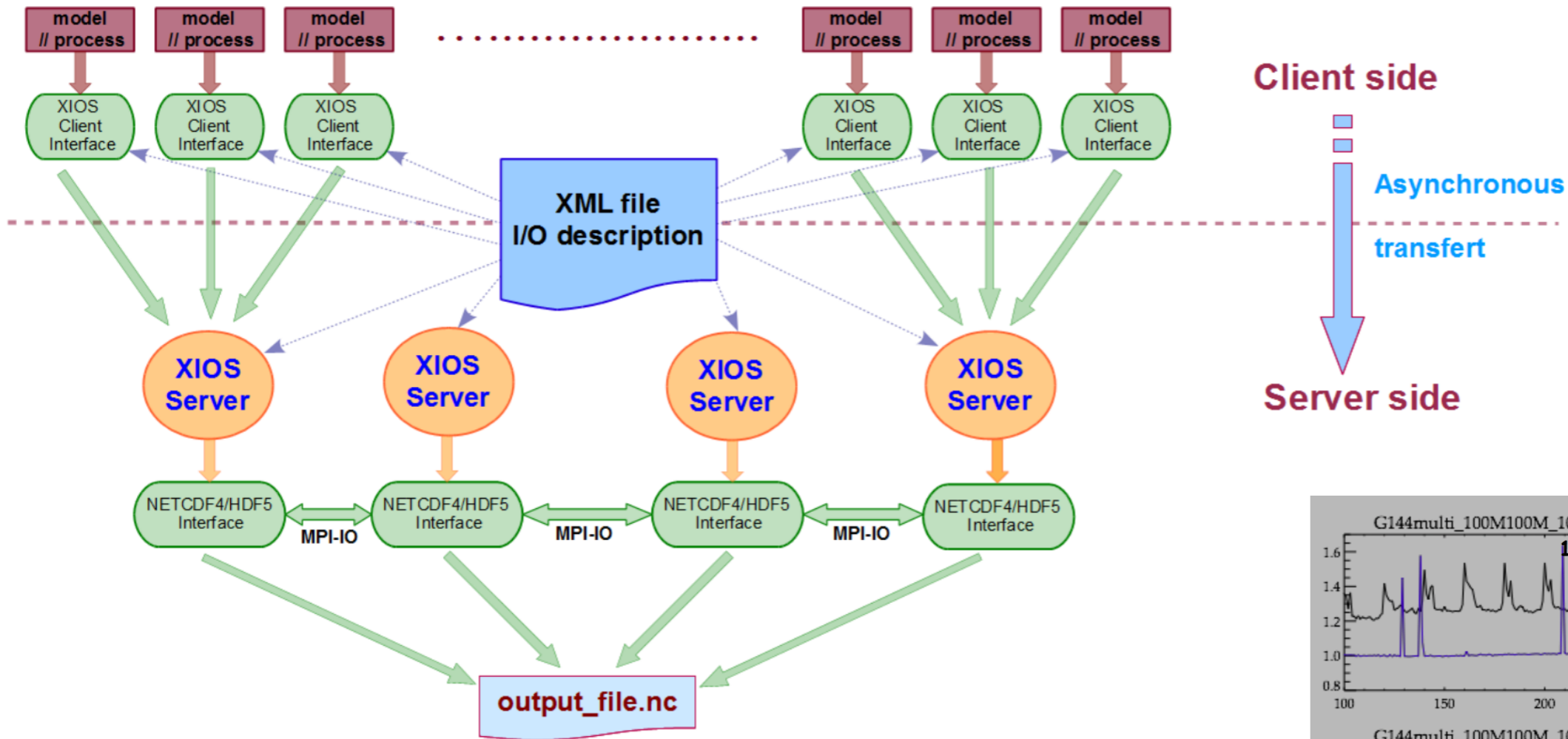


Rapide avec un unique fichier à l'arrivée !!!

Nécessite la librairie NetCDF4 installée avec support parallèle

# Ecriture de fichiers MPI 4/4 : XIOS

## XIOS

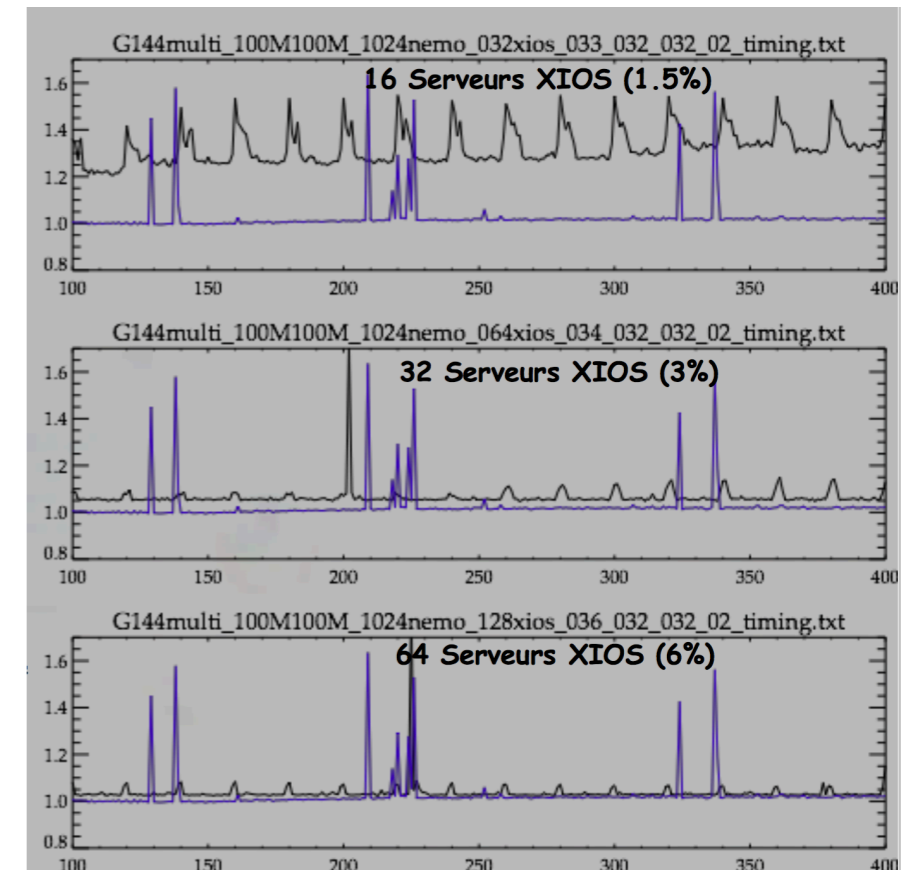


## Strategy for outputs

XIOS : external server developed at IPSL

<http://forge.ipsl.jussieu.fr/iomserver>

**Exemple**  
(S. Masson, from NEMO ...)



# Ecriture de fichiers MPI 4/4 : XIOS

---

## XIOS general

- Originally, a library dedicated to Input/Output management of large climate coupled models (e.g. CMIP simulations for IPCC with NEMO and other code)
- Written and managed at (LSCE-IPSL) by Y. Meurdesoif et al.
- XIOS creates output NetCDF files
- Implemented in other codes (ROMS, MARS3D, CROCO) by non-xios-expert developers despite of a light existing documentation.
- All documentation at <http://forge.ipsl.jussieu.fr/ioserver> with tutorials, user guide
- Installation of XIOS could be not an easy task to do on a new machine, be sure it is already well installed with the right netcdf4 library !
- In the next croco version, XIOS version  $\geq 2$

# Ecriture de fichiers MPI 4/4 : XIOS

---

## XIOS why and when ?

- I/O becomes a bottleneck in parallel computing with using a large amount of processors

e.g. Atlantic model at **1km** resolution :

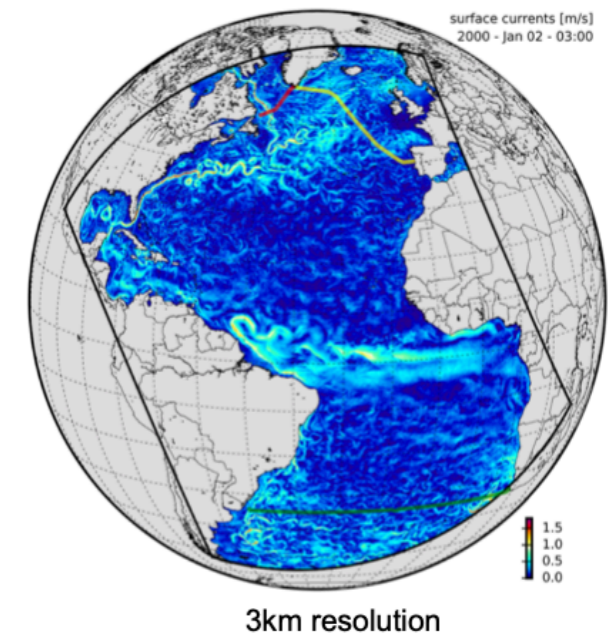
10000 x 14000 x 200 grid points ; using up to ~50000 procs

=> Very difficult or impossible to manage such amount of output datas with classical netcdf library.

- Only an external configuration file is needed to configure the outputs (no need to compile each time)
  - create new files
  - create new variables from referenced variables
  - use time filter (instantaneous, average, cumulate, ...)

1. Efficiency in production of data on supercomputer parallel file system

2. Flexibility and “simplicity” in management of I/O and data definition



Remark : It is may be not so “ simple ” for beginners because you need to understand how to modify the configuration file written in xml

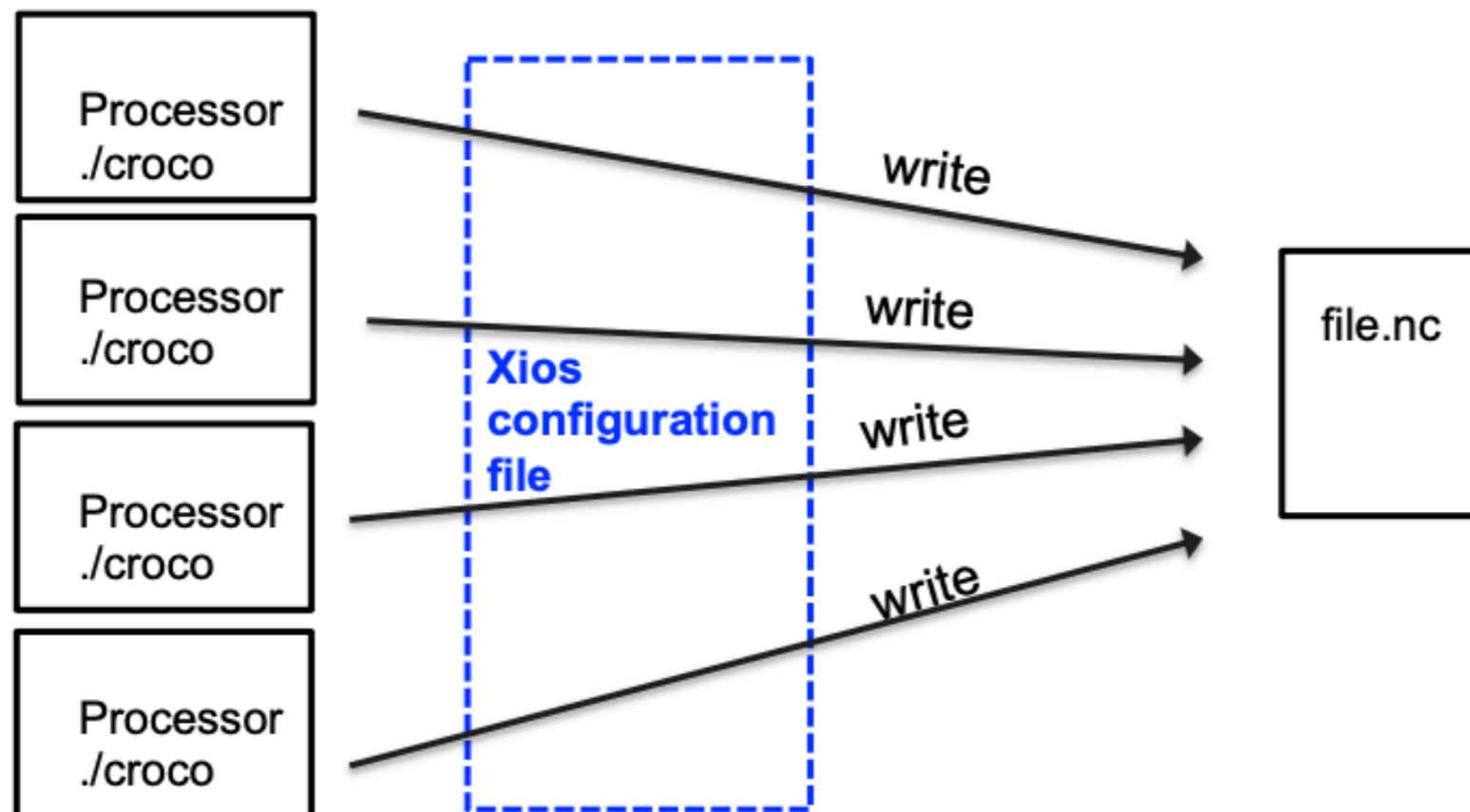
# Ecriture de fichiers MPI 4/4 : XIOS

---

## XIOS : attached mode

Using xios in **attached mode** :

each croco executable **compute** and **write** (like a classical library)



Ergonomy AND efficient parallel writing BUT writing overhead

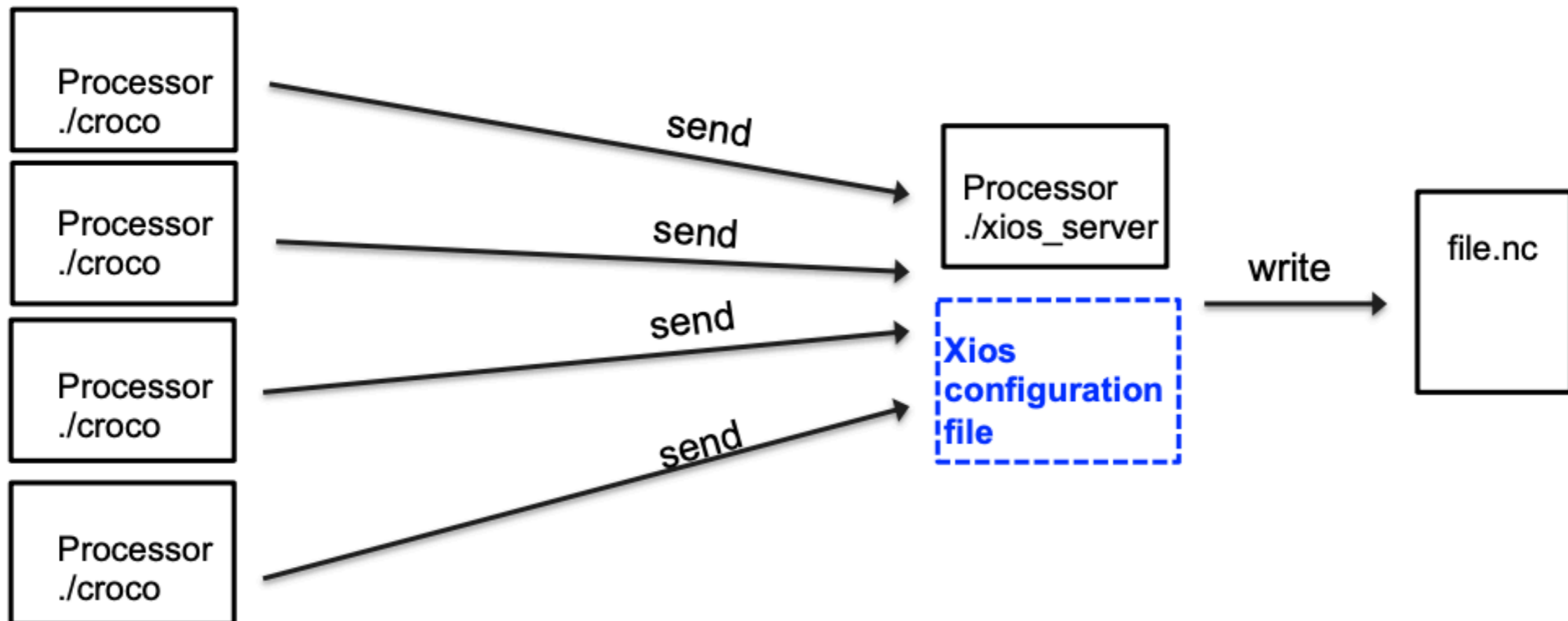


# Ecriture de fichiers MPI 4/4 : XIOS

---

## XIOS : detached mode (server mode)

each croco executable compute and send field to the server



- croco executables for computing only
- only xios server writes output
- Flexibility AND efficient parallel writing AND (almost) no overhead

# Ecriture de fichiers MPI 4/4 : XIOS

---

## XIOS : in practice

- In `cppdefs.h` add ccp keys : `#define XIOS`
- Add the XIOS library path in `jobcomp`
- Compile once : `./jobcomp`
- Edit/modify xios configuration file : `iodef.xml`
- To run :
  - in attached mode : as usual
  - in detached mode : like a coupled model ...  
`mpirun -np 10 ./croco -np 2 ./xios.exe`

# Le cas des points terre

## 1. Préprocessing

Dans croco/MPI\_NOLAND :

- lire le README
- compiler: edit makefile + make
- éditer la namelist :
  - nom du fichier grille
  - nombre de procs max
- exécuter : `./mpp_optimize`
- visualiser :  
`./mpp_plot.py croco_grd.nc benguela-008x005_033`
- re-lire le README ...

## 2. Avant de compiler CROCO

- `cppdefs.h` : `#define MPI_NOLAND`
- `param.h` : insérer les valeurs pour `NP_XI`, `NP_ETA` and `NPP` données par le preprocessing  
( $NPP \leq NP\_XI \times NP\_ETA$ )
- exécution habituelle (`mpirun -np etc`)

WARNING : grid file as to be called `croco_grd.nc` (or to be changed in `MPI_Setup.F`)

```
.....
namproc
.....
jprocx = maximum number of proc

&NAMPROC
  jprocx=220
.....
namfile of filename
.....
NAMELIST /namfile/ cbathy
cbathy = name of the bathy/mask file(nc)
covdta = Root for the overdata file name .
Complete name will be {covdta}.{NP_XI}x{NP_ETA}_{NPP}
&NAMFILE
  cbathy='croco_grd.nc'
  covdta = 'benguela'
/
```

```
-bash
(base) sdb-benshila:MPP_PREP rblod$ ./mpp_optimiz
Number of pts      :    1936
Number of sea pts  :    1411

optimum choice
=====
--> Number of CPUs : NNODES =          33

NP_XI =           8  NP_ETA =           5
Lm =           6  Mm =           9

number of sea CPUs          33
number of land CPUs         7
average overhead            0.69463869463869454
minimum overhead            0.138461545
maximum overhead            1.000000000

nb of overhead p. < 10 %      0
nb of overhead p. 10 < nb < 30 %  3
nb de overhead p 30 < nb < 50 %  5
number of integration points   4290
number of additional pts       2398
% sup                          2.26744175

(base) sdb-benshila:MPP_PREP rblod$
```

