



Programmation parallèle :

**M**essage **P**assing **I**nterface (MPI-1)

Geneviève Moguilny

[moguilny@ipgp.jussieu.fr](mailto:moguilny@ipgp.jussieu.fr)

Institut de Physique du Globe de Paris

# Architectures et modèles de programmation

D'après Flynn (1972),

**S** : Single, **I** : Instruction, **M** : Multiple, **D** : Data.

Séquentielles ou Parallèles

		Contrôle (Inst.)	
		S	I
Données	S	SISD	MISD
	M	SIMD	MIMD

Parmi les MIMD, **systèmes répartis** = ensemble de machines reliées par un réseau.

Programmation parallèle  $\Rightarrow$  écriture d'un (ou plusieurs) programme(s) qui s'exécutera(ont) simultanément sur plusieurs processeurs afin d'atteindre un but commun :

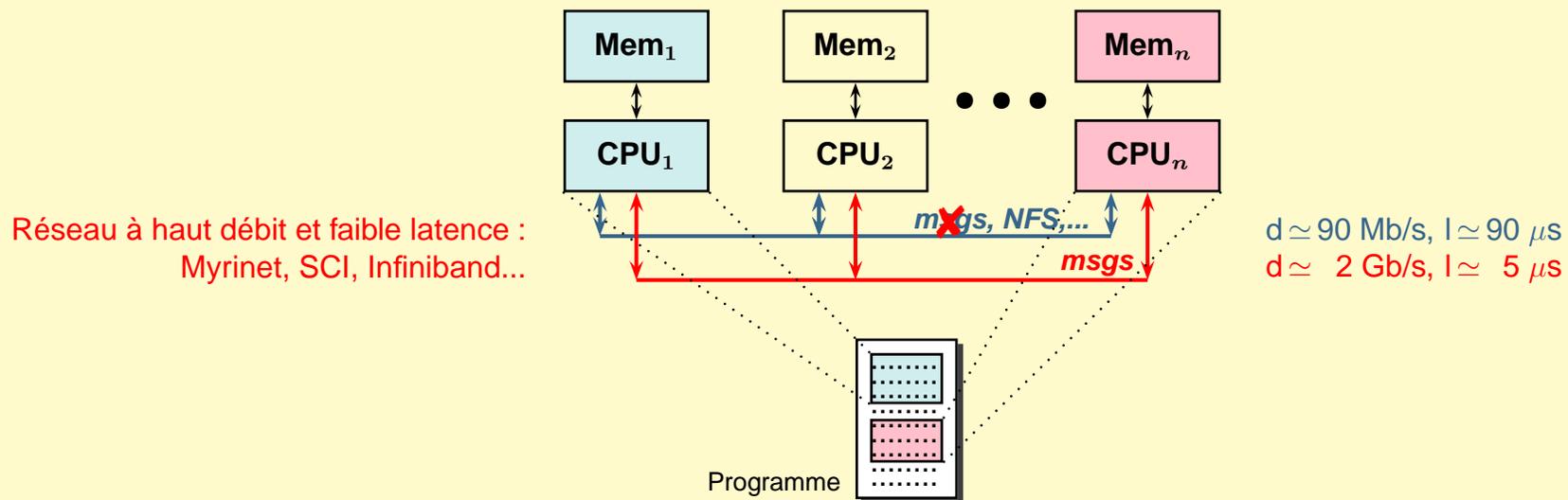
- chaque processeur exécutera une instance de programme,
- échange de données (message)  $\Rightarrow$   
appel à un sous-programme d'une **bibliothèque d'échange (ou de passage) de messages** comme PVM (*Parallel Virtual Machine*) ou MPI (*Message Passing Interface*).

Modèles de programmation sur les systèmes répartis à mémoire distribuée (**P** = Program) :

- MPMD pour un réseau hétérogène  $\Rightarrow$  PVM ou MPI-2,
- **SPMD** pour un réseau homogène  $\Rightarrow$  PVM (obsolète) ou MPI-[1 | 2].

# Application SPMD

- Un programme écrit dans un langage classique (Fortran, C...).
- Chaque processus reçoit la même occurrence du programme mais par le biais d'instructions conditionnelles, exécute le même code ou une partie différente du code, sur les mêmes données ou des données différentes.
- Les variables du programme sont privées et résident dans la mémoire du processeur allouée au processus.
- Les données sont échangées entre les processus par des appels à des sous-programmes d'une librairie d'échange de messages, comme **MPI**.



# La librairie de passage de messages MPI

## Historique :

- 1992 : Nécessité de créer des applications portables avec de bonnes performances ⇒ création d'un groupe de travail (principalement américano-européen) afin d'adopter les méthodes HPF (*High Performance Fortran*).
- 1994 : Version 1.0 de MPI.
- 1997 : Définition d'une norme pour MPI-2 (contrôle dynamique des tâches, E/S parallèles...) disponible aujourd'hui.

## Principales implémentations du domaine public :

- LAM MPI (*Indiana University*),
- MPICH (*Argonne National Laboratory*);

et sur celles-ci,

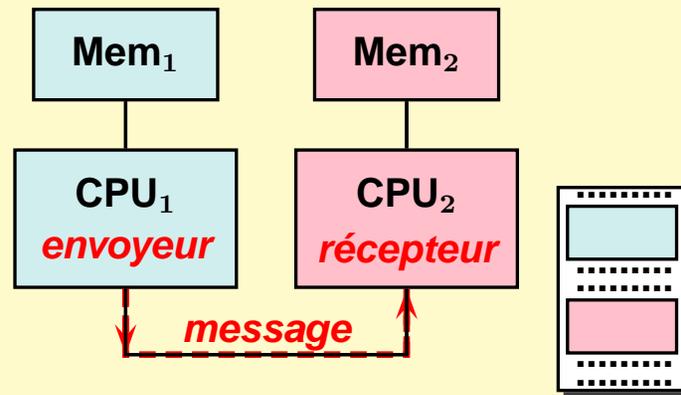
- librairies spécifiques à un réseau (MPI-GM : implémentation de MPICH sur Myrinet).

# Concept d'échange de message

**Application MPI** = ensemble de processus autonomes exécutant leur propre code, et communiquant via des appels à des sous-programmes de la bibliothèque MPI.

**Communication**  $\Rightarrow$  **Échange de messages**

**Message** = identificateur(processus émetteur) + type des données (simple ou dérivé) + longueur + identificateur(processus récepteur) + **donnée**



*Communication point à point*

# Principales catégories de l'API MPI

## ① Gestion de l'environnement

- activation / désactivation de l'environnement (**MPI\_INIT**, **MPI\_FINALIZE**),
- identification des processus (**MPI\_COMM\_SIZE**, **MPI\_COMM\_RANK**).

## ② **Communications** (envoi et réception de messages)

- **point à point**,
- **collectives**.

## ③ **Gestion des groupes et communicateurs**

**Communicateur** : ensemble de processus sur lesquels portent les opérations MPI, groupe de processeurs + contexte de communication (géré par MPI).

Création d'un communicateur à partir :

- d'un groupe de processus préalablement créé par le programmeur,
- d'un autre communicateur (**MPI\_COMM\_WORLD** = l'ensemble des processus).

# Tout petit exemple

- Source : HelloMPI.f90

```
1: program HelloMPI
2:
3: implicit none
4: include 'mpif.h'
5: integer :: nb_procs, rang, ierr
6:
7: call MPI_INIT(ierr)
8:
9: call MPI_COMM_SIZE(MPI_COMM_WORLD, nb_procs, ierr)
10: call MPI_COMM_RANK(MPI_COMM_WORLD, rang, ierr)
11: print *, 'Je suis le processus ', rang, ' parmi ', nb_procs
12:
13: call MPI_FINALIZE(ierr)
14:
15: end program HelloMPI
```

## Tout petit exemple (suite)

```
1: program HelloMPI
2:
3: implicit none
4: include 'mpif.h'
5: integer :: nb_procs, rang, ierr
6:
7: call MPI_INIT(ierr)
8:
9: call MPI_COMM_SIZE(MPI_COMM_WORLD, nb_procs, ierr)
10: call MPI_COMM_RANK(MPI_COMM_WORLD, rang, ierr)
11: print *, 'Je suis le processus ', rang, ' parmi ', nb_procs
12:
13: call MPI_FINALIZE(ierr)
14:
15: end program HelloMPI
```

- **Compilation et Édition de liens :**

```
mpif90 HelloMPI.f90 -o HelloMPI
```



```
ifort -c -Iincludes_path HelloMPI.f90
```

```
ifort -Llibs_path HelloMPI.o -static-libcxa -o HelloMPI -lmpichf90 -lmpich
```

- **Lancement de l'exécution :**

```
mpirun -np 4 -machinefile mf HelloMPI
```

Lancement de l'exécution de `HelloMPI` (par `ssh` ou `rsh`) sur 4 processeurs (`np = number of processors`) des machines parmi celles citées dans le fichier `mf`.

- **Résultat obtenu :**

```
Je suis le processus      0  parmi      4
Je suis le processus      3  parmi      4
Je suis le processus      1  parmi      4
Je suis le processus      2  parmi      4
```

# Les communications point à point

Communications entre 2 processus identifiés par leur rang dans le communicateur.

## Modes de transfert :

*standard* : MPI effectue ou non une copie temporaire suivant l'implémentation.

*Synchronous* : envoi terminé quand réception achevée.

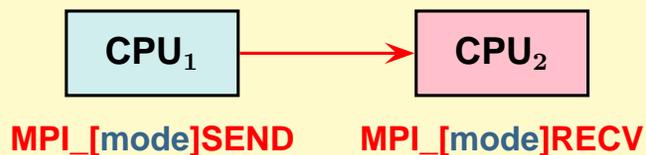
*Buffered* : copie temporaire dans un buffer créé par le programmeur, envoi terminé lorsque copie achevée.

*Ready* : envoi effectué lorsque réception initiée (clients/serveurs).

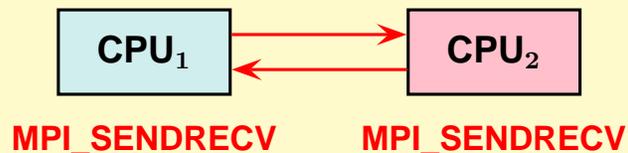
Envoi bloquant (envoi et réception couplés) ou non, réception Immédiate ou non.

Deux types d'échange :

- 1 1 expéditeur et 1 récepteur distincts (MPI\_SEND, MPI\_SSEND, MPI\_IBSEND ...) :



- l'expéditeur reçoit aussi et inversement (MPI\_SENDRECV, MPI\_SENDRECV\_REPLACE) :



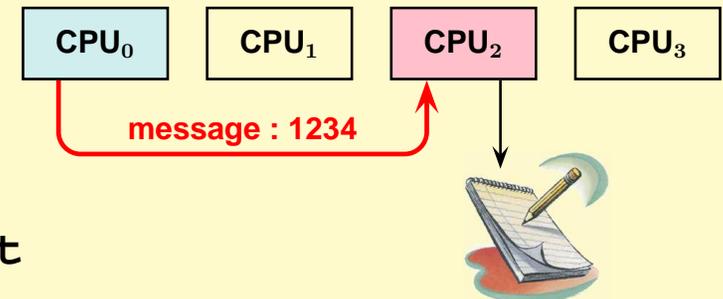
# Exemple de communication point à point simple

```
1: !! point_a_point.f90 : Exemple d'utilisation de MPI_SEND et MPI_RECV
2: !! Auteur : Denis GIROU (CNRS/IDRIS - France) <Denis.Girou@idris.fr> (1996)
3: program point_a_point
4:   implicit none
5:   include 'mpif.h'
6:   integer, dimension(MPI_STATUS_SIZE) :: statut
7:   integer, parameter                :: etiquette=100
8:   integer                            :: rang,valeur,ierr
9:   call MPI_INIT(ierr)
10:  call MPI_COMM_RANK(MPI_COMM_WORLD,rang,ierr)
11:
12:  if (rang == 0) then
13:    valeur=1234
14:    call MPI_SEND(valeur,1,MPI_INTEGER,2,etiquette,MPI_COMM_WORLD,ierr)
15:  elseif (rang == 2) then
16:    call MPI_RECV(valeur,1,MPI_INTEGER,0,etiquette,MPI_COMM_WORLD,statut,ierr)
17:    print *, 'Moi, processus ', rang, ' j'ai reçu ',valeur,' du processus 0.'
18:  end if
19:
20:  call MPI_FINALIZE(ierr)
21: end program point_a_point
```

```
mpif90 point_a_point.f90 -o point_a_point
```

```
mpirun -np 4 -machinefile mf point_a_point
```

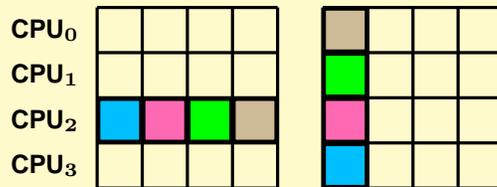
```
Moi, processus          2  j'ai reçu          1234  du processus 0.
```



# Les communications collectives

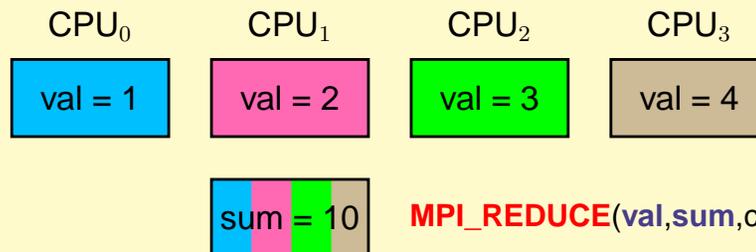
Ensemble de communications point à point à l'intérieur d'un groupe, fait en une opération.

- Synchronisations (BARRIER)
- Mouvement de données
  - Diffusions générales (BCAST) ou sélectives (SCATTER)



**MPI\_SCATTER**(valeurs,ltranche,MPI\_REAL, &  
donnees,ltranche,MPI\_REAL,2,MPI\_COMM\_WORLD,code)

- Collectes réparties (GATHER) ou générales (ALLGATHER)
- Échanges croisés (ALLTOALL)
- Réductions (REDUCE avec SUM, PROD, MAX, MIN, MAXLOC, MINLOC, IAND, IOR, IXOR)



**MPI\_REDUCE**(val,sum,count,MPI\_INTEGER,**MPI\_SUM**,0,MPI\_COMM\_WORLD,code)

# Création de communicateurs : les topologies

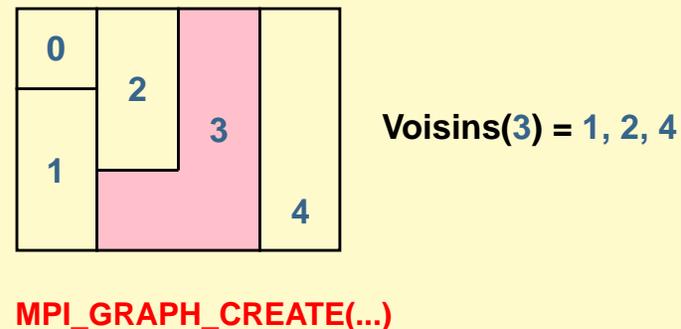
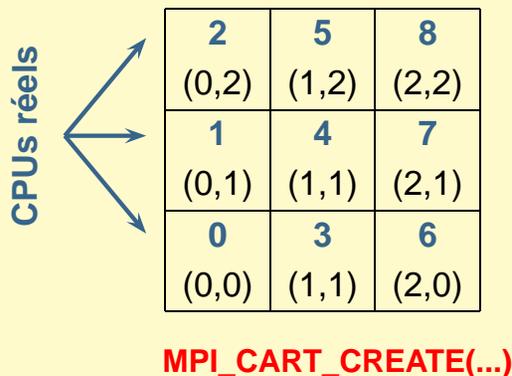
Pour les applications qui utilisent une **décomposition de domaine**,  
il est intéressant de définir une grille virtuelle des processus afin d'en faciliter la manipulation :

nombre de sous-domaines = nombre de processus,  
travail sur les sous-domaines puis finalisation par calculs aux frontières.

⇒ création de nouveaux **communicateurs** sur ces nouvelles topologies.

Ces topologies peuvent être :

- **cartésiennes** : grille périodique ou non, identification des processus par leurs coordonnées,
- de type **graphe** pour les topologies plus complexes.



# Développement d'une application MPI

## Pourquoi MPI ?

- Permet de changer l'ordre de grandeur du problème à traiter par une augmentation possible quasi-illimitée de la mémoire disponible, avec du matériel et du logiciel de base.

## Optimisation

- Choix du mode de communication : éviter les recopies inutiles
- Contention des messages
- Équilibrage de la charge en faisant attention aux goulets d'étranglement (E/S)
- Recouvrement des temps de calcul et de communications
- Utilisation des communications persistantes

## Aide au développement

- Débogueurs graphiques supportant MPI : TotalView, DDT
- Visualisation des communications : Vampir/VampirTrace, Blade

**Application bien parallélisée avec MPI  $\Rightarrow$   
temps d'exécution inversement proportionnel au nombre de processeurs utilisés**

# Documentation MPI

- Programmation parallèle et répartie :

[http://www.gulus.org/spip/IMG/pdf/midi\\_linux\\_mpi.pdf](http://www.gulus.org/spip/IMG/pdf/midi_linux_mpi.pdf)

- Cours MPI :

[http://www.idris.fr/data/cours/parallel/mpi/choix\\_doc.html](http://www.idris.fr/data/cours/parallel/mpi/choix_doc.html)

[www.obspm.fr/~websio/documentations/parallel/archi\\_mpi3.pdf](http://www.obspm.fr/~websio/documentations/parallel/archi_mpi3.pdf)

[http://www.cines.fr/textes/support\\_para.html](http://www.cines.fr/textes/support_para.html)

- Implémentations MPI du domaine public :

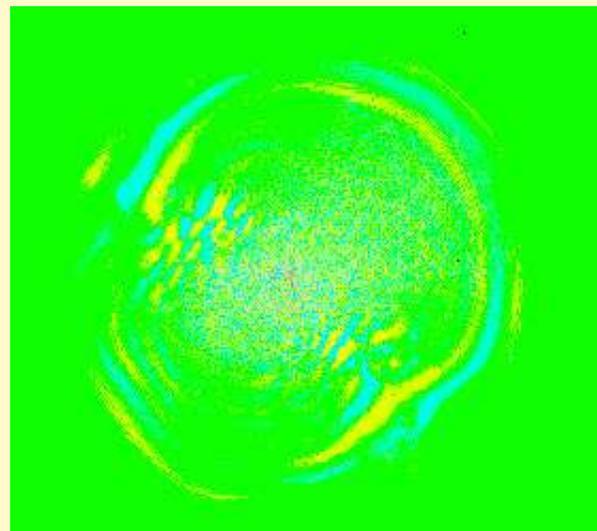
- LAM : <http://www.lam-mpi.org>

- MPICH : <http://www.mcs.anl.gov/mpi/mpich>

## Exemple d'application MPI tournant sur EGEE : SPECFEM3D

Résolution de problèmes de propagation d'ondes sismiques dans des milieux complexes et à haute fréquence à l'échelle régionale, par la méthode des éléments spectraux.

Application écrite initialement par D. Komatitsch (Université de Pau), utilisée dans plus de 75 laboratoires dans le monde, essentiellement pour des études de tremblements de terre, tomographie sismique, ultrasons dans les cristaux, topographie, bassins et vallées sédimentaires, ondes d'interface, et cristaux anisotropes.



Informations supplémentaires : <http://www.univ-pau.fr/~dkomatil>

Positionnement **EGEE** : NA4, VO ESR, catégorie *Physique de la Terre solide*

## SPECFEM3D : spécificités techniques

- $\simeq$  20 000 lignes de Fortran 90 / MPI.
- Peut tourner sur 4 CPUs, a tourné sur 1944 CPUs sur le *EarthSimulator* (Japon).  
Sur EGEE, a tourné sur 64 CPUs au Nikhef (Pays-Bas), et sur 16 CPUs à SCAI (Allemagne), au LAL, au CPPM, et à la CGG.
- **Contraintes :**
  - Optimisation de la mémoire  $\Rightarrow$  recompilation et mise à jour de fichiers d'entrée sur SE nécessaires à chaque changement des paramètres d'entrée.
  - Beaucoup de sorties provisoires ( $\Rightarrow$  écriture dans les `/tmp` propres à chaque nœud), + sorties écrites dans un répertoire partagé à récupérer ( $\Rightarrow$  **/home montés NFS**).
  - Nécessite le lancement successif de **deux mpirun** qui doivent utiliser les mêmes nœuds, alloués dans le même ordre.

## SPECFEM3D sur EGEE : Fem4.jdl

```
1: # Fem4.jdl
2:   NodeNumber = 4;
3:   VirtualOrganisation = "esr";
4:   Executable = "Fem";
5:   Arguments = "xmeshfem3D4 xspecfem3D4 4";
6:   StdOutput = "Fem.out";
7:   StdError = "Fem.err";
8:   JobType = "mpich";
9:   LRMS_type = "pbs";
10:  Type = "Job";
11:  InputSandbox = { "Fem", "xmeshfem3D4", "xspecfem3D4" };
12:  OutputSandbox = { "output_mesher.txt", "output_solver.txt", "xmeshfem3D4.out",
13:                   "xspecfem3D4.out", "Fem.out", "Fem.err", "Par_file" };
14: # Requirements = (other.GlueCEInfoTotalCPUs >= 4);
15:
```

# SPECFEM3D sur EGEE : script Fem (1/2)

```
1: #!/bin/sh
2: # copy of scai:/home/lcg2/mpirun_g95 modified - 11/03/2005
3: #
4: EXE1=$1
5: EXE2=$2
6: CPU_NEEDED=$3
7: TMP_DIR=/tmp/DATABASES_MPI_DIMITRI4
8: institut=`echo $HOSTNAME | cut -f2-5 -d.`
9: # Different path for mpirun on LAL and CGG
10: if [ $institut == "lal.in2p3.fr" -o $institut == "private.egee.fr.cgg" ] ; then
11:     MPIHOME=/usr/bin
12: else
13:     MPIHOME=/opt/mpich/bin
14: fi
15: if [ -f "$PWD/.BrokerInfo" ] ; then
16: TEST_LSF=`edg-brokerinfo getCE | cut -d/ -f2 | grep lsf`
17: else
18: TEST_LSF=`ps -ef | grep sbatchd | grep -v grep`
19: fi
20: if [ "x$TEST_LSF" = "x" ] ; then
21:     HOST_NODEFILE=$PBS_NODEFILE
22: else
23:     echo "LSF Hosts: $LSB_HOSTS"
24:     HOST_NODEFILE=`pwd`/lsf_nodefile.$$
25: fi
26: for host in `cat ${HOST_NODEFILE}`
27: do
28:     ssh $host rm -rf $TMP_DIR
29:     ssh $host mkdir $TMP_DIR
30:     ssh $host chmod 775 $TMP_DIR
31: done
32: LocalFn="DATA.tgz"
33: echo "trying edg-rm..."
34: cmd="edg-rm --vo esr copyFile lfn:DATA${CPU_NEEDED}.tgz file://`pwd`/$LocalFn"
35: echo ">>> $cmd"
36: $cmd
37:
38: if [ ! -f $LocalFn ] ; then
39:     echo "edg-rm failed, trying lcg-cp..."
40:     export LCG_GFAL_INFOSYS="mu3.matrix.sara.nl:2170"
41:     cmd="lcg-cp --vo esr lfn:DATA${CPU_NEEDED}.tgz file://`pwd`/$LocalFn"
42:     echo ">>> $cmd"
43:     $cmd
44: fi
```

*Définition du path de mpirun*

*Création des répertoires temporaires*

*Rapatriement des données d'un SE*

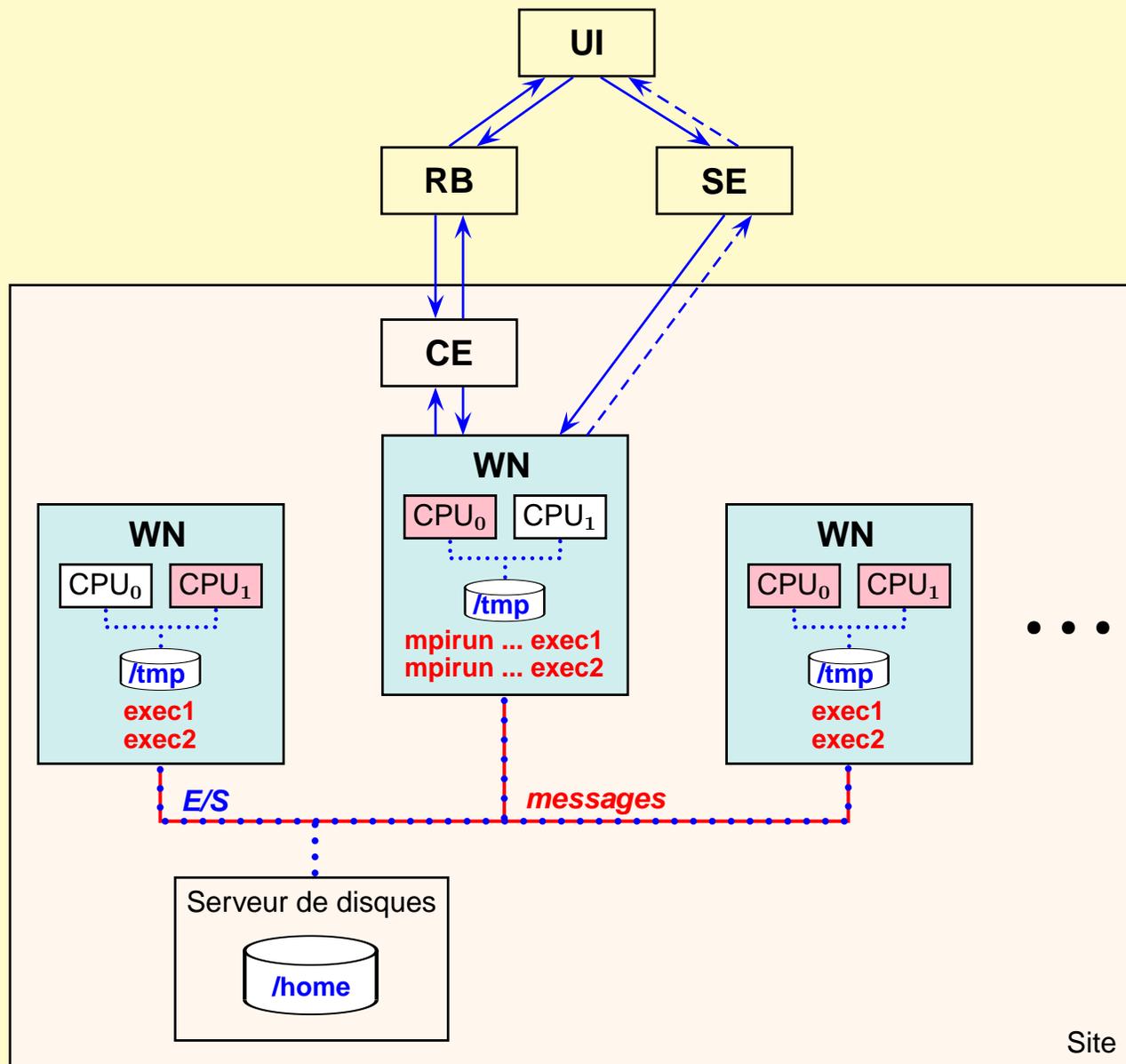
## SPECFEM3D sur EGEE : script Fem (2/2)

```
45: if [ ! -f $LocalFn ] ; then
46:     echo "$LocalFn not found."
47:     exit
48: fi
49: #
50: echo "*****"
51: tar xfz $LocalFn
52: rm $LocalFn
53: rm -rf OUTPUT_FILES
54: mkdir OUTPUT_FILES
55: cp DATA/Par_file .
56:
57: chmod 755 $EXE1
58: chmod 755 $EXE2
59: ls -l
60: time $MPIHOME/mpirun -np $CPU_NEEDED -machinefile $HOST_NODEFILE \
61:     `pwd`/$EXE1 > $EXE1.out
62: cp OUTPUT_FILES/output_mesher.txt .
63: time $MPIHOME/mpirun -np $CPU_NEEDED -machinefile $HOST_NODEFILE \
64:     `pwd`/$EXE2 > $EXE2.out
65: echo "Size of OUTPUT_FILES :"
66: du -sk OUTPUT_FILES
67: for host in `cat ${HOST_NODEFILE}`
68: do
69:     ssh $host echo "Size of $TMP_DIR on $host :"
70:     ssh $host du -sk $TMP_DIR
71:     ssh $host rm -rf $TMP_DIR
72: done
73: cp OUTPUT_FILES/output_solver.txt .
74: tar cfz OUTPUT_FILES${CPU_NEEDED}.tgz OUTPUT_FILES
75: lcg-del --vo esr -s grid11.lal.in2p3.fr `lcg-lg --vo esr \
76:     lfn:OUTPUT_FILES${CPU_NEEDED}.tgz`
77: lcg-cr --vo esr -d grid11.lal.in2p3.fr \
78:     file://`pwd`/OUTPUT_FILES${CPU_NEEDED}.tgz \
79:     -l OUTPUT_FILES${CPU_NEEDED}.tgz
80: exit
```

} Lancement des *mpirun*

} Stockage des sorties sur SE

# SPECFEM3D sur EGEE : exécution



## Perspectives

- SPECFEM3D : des collaborations avec DEISA et Grid5000 vont être gérées par l'INRIA Bordeaux.
- Souhaits à court terme :
  - Installations homogènes de MPI (*path* standard ou tag MPIHOME).
  - Nœud  $\neq$  CPU
    - ⇒ remplacement de Torque par un autre LRMS comme **SGE** ou **OAR** ?
  - Queues spécifiques MPI contenant des CPUs de puissance équivalente.
- Dans le futur...
  - Réservation de CPUs.
  - MPI inter-sites (**MPICH-G2**, **MPICH-V**).
  - Disposer de réseaux performants (Myrinet, Infiniband...).
  - MPI sur des RC disposant de plusieurs centaines de nœuds (CCIN2P3 ?).
- Mais globalement, **MPI/Fortran90 sur EGEE marche !!**