

Message Passing Interface (MPI-1)

Jalel Chergui
Isabelle Dupays
Denis Girou
Stéphane Requena
<Prénom.Nom@idris.fr>

1	Introduction	6
1.1	Définitions	6
1.2	Concepts de l'échange de messages	12
1.3	Historique	17
1.4	Bibliographie	18
2	Environnement	19
2.1	Description	19
2.2	Exemple	22
3	Communications point à point	23
3.1	Notions générales	23
3.2	Types de données de base	26
3.3	Autres possibilités	28
3.4	Exemple:anneau de communication	31
3.5	Construction et reconstruction de messages	38
4	Communications collectives	41
4.1	Notions générales	41

43	4.2 – Synchronisation globale:MPI_BARRIER()
44	4.3 – Diffusion générale:MPI_BCAST()
46	4.4 – Diffusion sélective:MPI_SCATTER()
49	4.5 – Collecte:MPI_GATHER()
51	4.6 – Collecte générale:MPI_ALLGATHER()
54	4.7 – Echanges croisés:MPI_ALLTOALL()
57	4.8 – Réductions réparties
66	4.9 – Compléments
67	5 – Optimisations
67	5.1 – Introduction
68	5.2 – Programme modèle
71	5.3 – Temps de communication
72	5.4 – Quelques définitions
76	5.5 – Que fournit MPI?
78	5.6 – Envoi synchrone bloquant
80	5.7 – Envoi synchrone non-bloquant

5.8 – Conseils 1 84

5.9 – Communications persistantes 88

5.10 – Conseils 2 95

6 – Types de données dérivés 96

6.1 – Introduction 96

6.2 – Types contigus 98

6.3 – Types avec un pas constant 99

6.4 – Descriptif des sous-programmes 102

6.5 – Exemples 103

6.6 – Types homogènes à pas variable 109

6.7 – Types hétérogènes 116

6.8 – Sous-programmes annexes 121

6.9 – Conclusion 123

7 – Topologies 124

7.1 – Introduction 124

7.2 – Topologies de processus 125

7.3 – Topologies cartésiennes	126
7.4 – Graphe de processus	141
8 – Communicateurs	148
8.1 – Introduction	148
8.2 – Communicateur par défaut	149
8.3 – Groupes et communicateurs	153
8.4 – Communicateur issu d'un groupe	156
8.5 – Communicateur issu d'un autre	163
8.6 – Subdiviser une topologie cartésienne	169
8.7 – Intra et intercommunicateurs	175
8.8 – Exemple récapitulatif	176
8.9 – Conclusion	183
9 – Évolution de MPI:MPI-2	184

1 — Introduction

1.1 — Définitions

❶ Le modèle de programmation séquentiel :

☞ le programme est exécuté par un et un seul processeur ;

☞ toutes les variables et constantes du programme sont allouées dans la mémoire centrale du processeur.

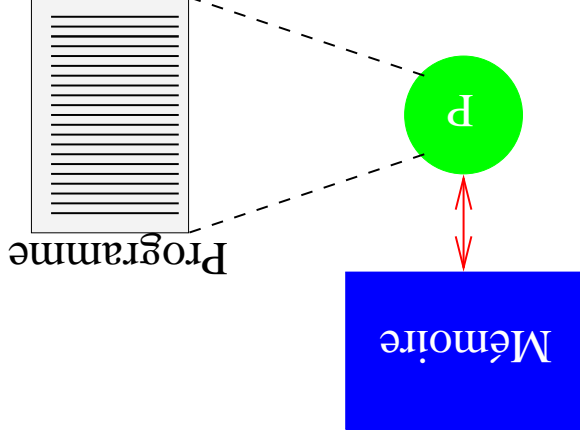


FIG. 1 – *Modèle séquentiel*

② Le modèle de programmation par échange de messages :

👉 le programme est écrit dans un langage classique (*Fortran*, *C* ou *C++*) ;

👉 chaque processus exécute éventuellement des parties différentes d'un

programme ;

👉 toutes les variables du programme sont privées et résident dans la mémoire

locale de chaque processus ;

👉 une donnée est échangée entre deux ou plusieurs processus via un appel, dans le programme, à des sous-programmes particuliers.

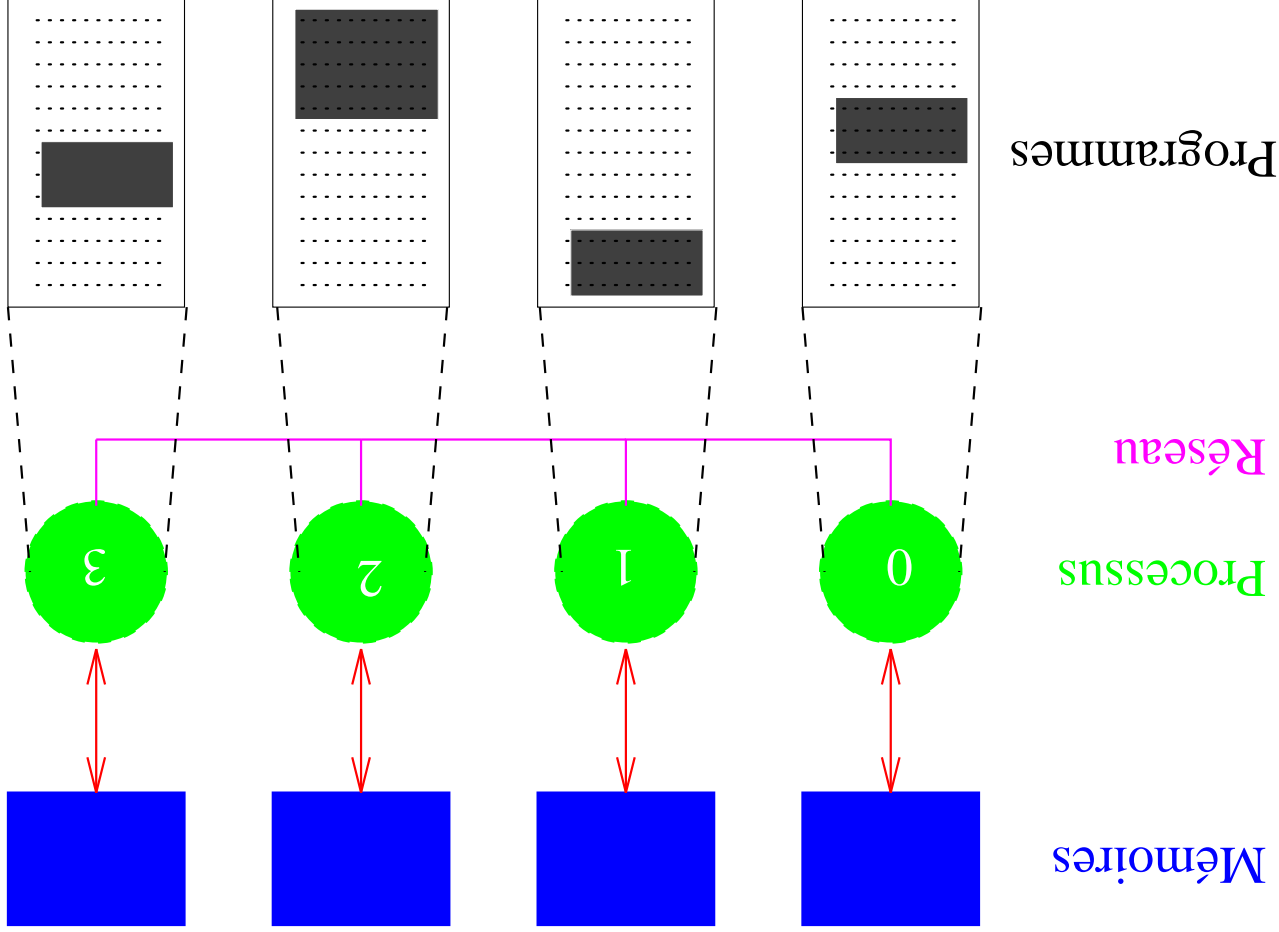


FIG. 2 – Modèle à échange de messages

③ Le modèle d'exécution *SPMD* :

👉 *Single Program, Multiple Data* ;

👉 le même programme s'exécute pour tous les processus ;

👉 toutes les machines supportent ce modèle de programmation et certaines ne

supportent que celui-là ;

👉 c'est un cas particulier du modèle plus général *MPMD*, qu'il peut en revanche

émuler.

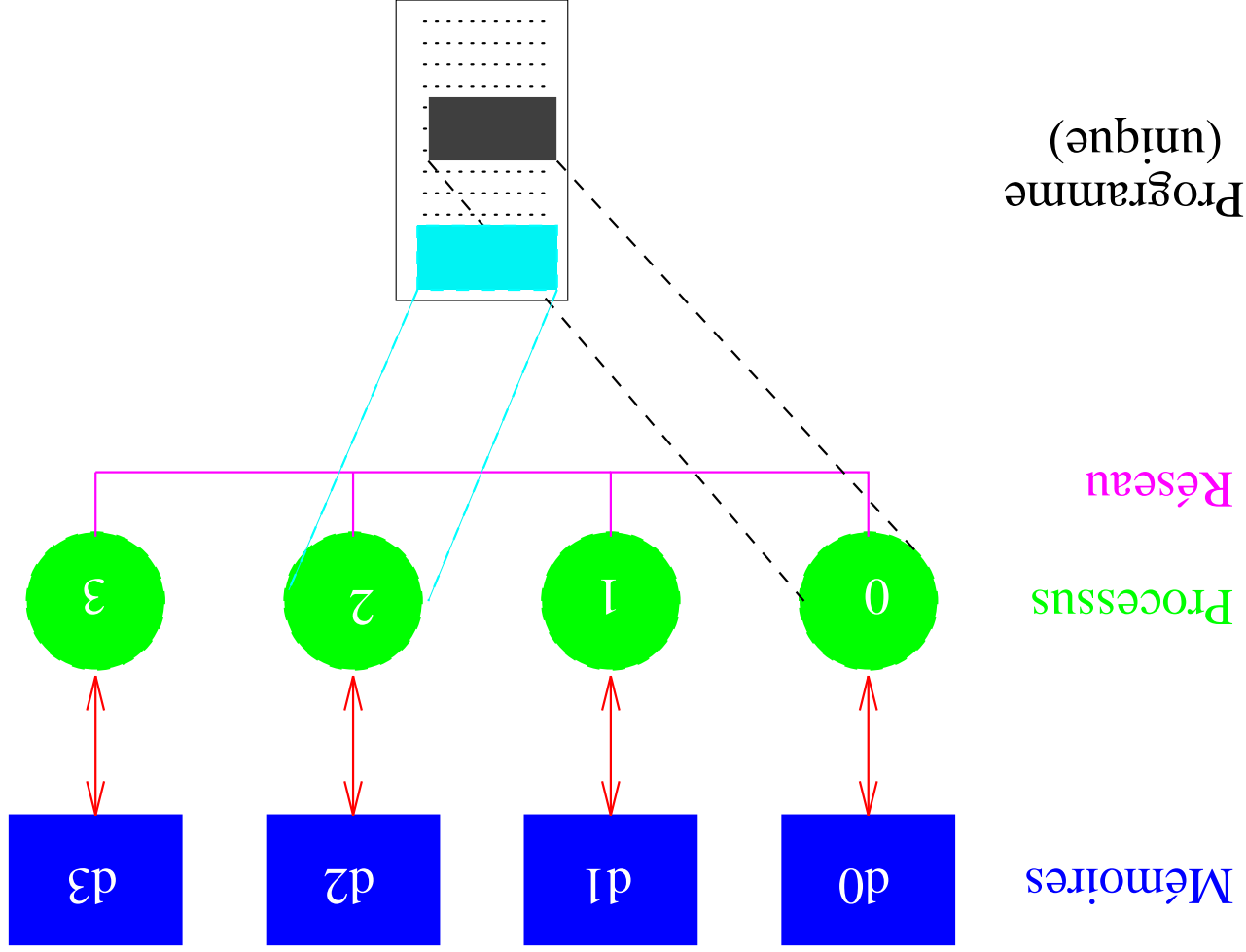


FIG. 3 – *Single Program, Multiple Data*

④ Exemple en Fortran d'émulation MPMD en SPMD

```
program spmd
  if (ProcessusMaitre) then
    call LeChef(Arguments)
  else
    call LesOuvriers(Arguments)
  endif
end program spmd
```

1.2 – Concepts de l'échange de messages

☞ Si un message est envoyé à un processus, celui-ci doit ensuite le recevoir

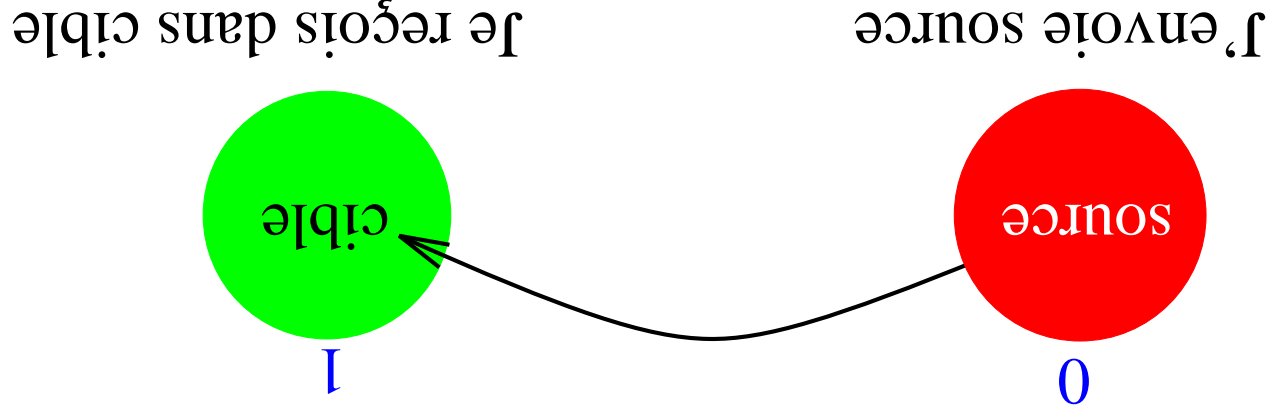


FIG. 4 – L'échange de messages

- ➡ Un message est constitué de paquets de données transitant du processus émetteur au(x) processus récepteur(s)
- ➡ En plus des données (variables scalaires, tableaux, etc.) à transmettre, un message doit contenir les informations suivantes :
 - ➡ l'identificateur du processus émetteur ;
 - ➡ le type de la donnée ;
 - ➡ sa longueur ;
 - ➡ l'identificateur du processus récepteur.

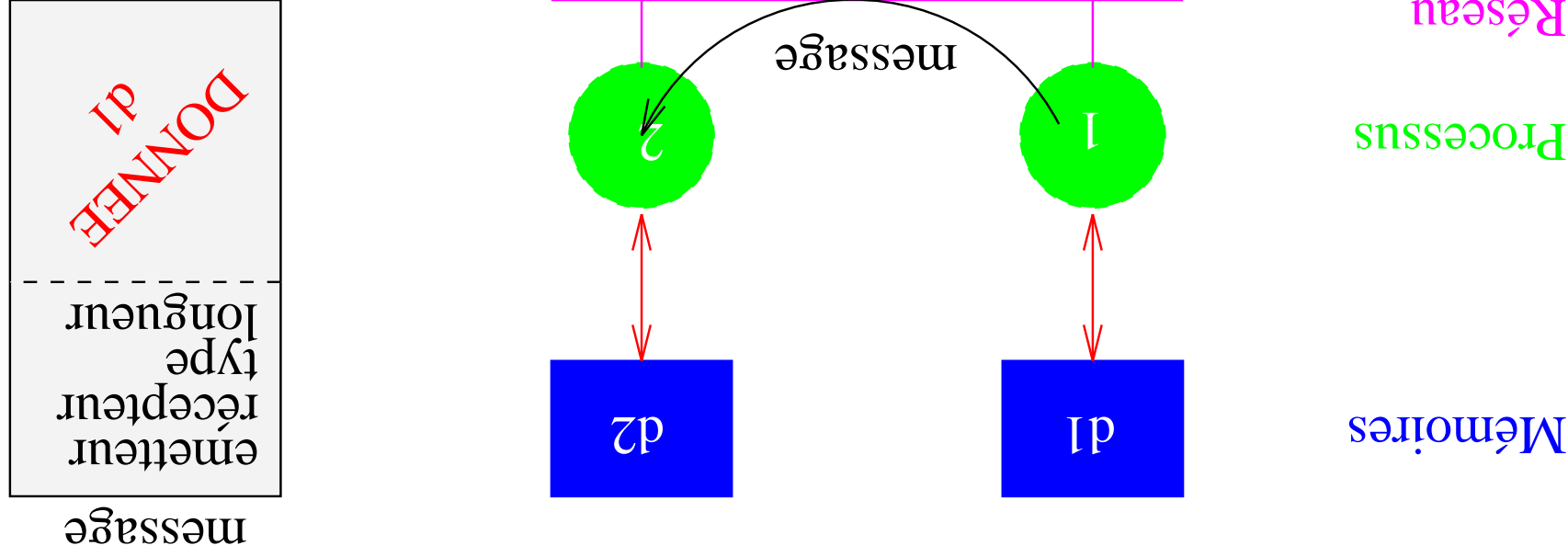


FIG. 5 – Constitution d'un message

☞ Les messages échangés sont interprétés et gérés par un environnement qui peut être comparé :

☞ à la téléphonie ;

☞ à la télécopie ;

☞ au courrier postal ;

☞ à une messagerie électronique ;

☞ etc.

☞ Le message est envoyé à une adresse déterminée

☞ Le processus récepteur doit pouvoir classer et interpréter les messages qui lui ont

été adressés

☞ L'environnement en question est *MPI-1 (Message Passing Interface)*. Une application *MPI* est un ensemble de processus autonomes exécutant chacun leur propre code et communiquant via des appels à des sous-programmes de la bibliothèque *MPI*

Ces sous-programmes peuvent être classés dans les grandes catégories suivantes :

- ① environnement ;
- ② communications point à point ;
- ③ communications collectives ;
- ④ types de données dérivés ;
- ⑤ topologies ;
- ⑥ groupes et communicateurs.

1.3 – Historique

- ➡ Novembre 92 (*Supercomputing '92*) : « formalisation » d'un groupe de travail créé en avril 92 et décision d'adopter les structures et les méthodes du groupe HPF (*High Performance Fortran*).
- ➡ Participants, américains (essentiellement) et européens, aussi bien constructeurs que représentants du monde académique.
- ➡ « Brouillon » de *MPI-1* présenté en novembre 93 (*Supercomputing '93*), finalisé en mars 1994.
- ➡ **Visé à la fois la portabilité et la garantie de bonnes performances.**
- ➡ « Brouillon » de *MPI-2* en novembre 96, suite à des travaux ayant commencés au printemps 95.
- ➡ « Standard » non élaboré par les organismes officiels de normalisation (ISO, ANSI, etc.).

1.4 – Bibliographie

☞ Les spécifications de la norme : *MPI: A Message Passing Interface Standard*, Mars 1994.

<ftp://ftp.irisa.fr/pub/netlib/mpl/drafts/draft-final.ps>

☞ Quelques ouvrages :

1. Marc Snir & al. *MPI: The Complete Reference*. Second edition. MIT Press, 1998. Volume 1, *The MPI core*. Volume 2, *MPI-2*.

2. William Gropp, Ewing Lusk et Anthony Skjellum. *Using MPI: Portable Parallel Programming with the Message Passing Interface*. Second edition. MIT Press

1999 ;

3. Peter S. Pacheco. *Parallel Programming with MPI*. Morgan Kaufman Ed., 1997.

☞ Une documentation complémentaire et une description des différentes

implémentations :

<http://www.idris.fr/su/Parallel1e1e/Page-generale.html>

<http://www.mcs.anl.gov/mpl/>

2 – Environnement

2.1 – Description

👉 Le sous-programme `MPI_INIT()` permet d'initialiser l'environnement nécessaire :

```
integer, intent(out) :: code
call MPI_INIT(code)
```

👉 Réciproquement, le sous-programme `MPI_FINALIZE()` désactive cet environnement :

```
integer, intent(out) :: code
call MPI_FINALIZE(code)
```

→ Toutes les opérations effectuées par *MPI* portent sur des communicateurs. Le communicateur par défaut est `MPI_COMM_WORLD` qui comprend tous les processus actifs.

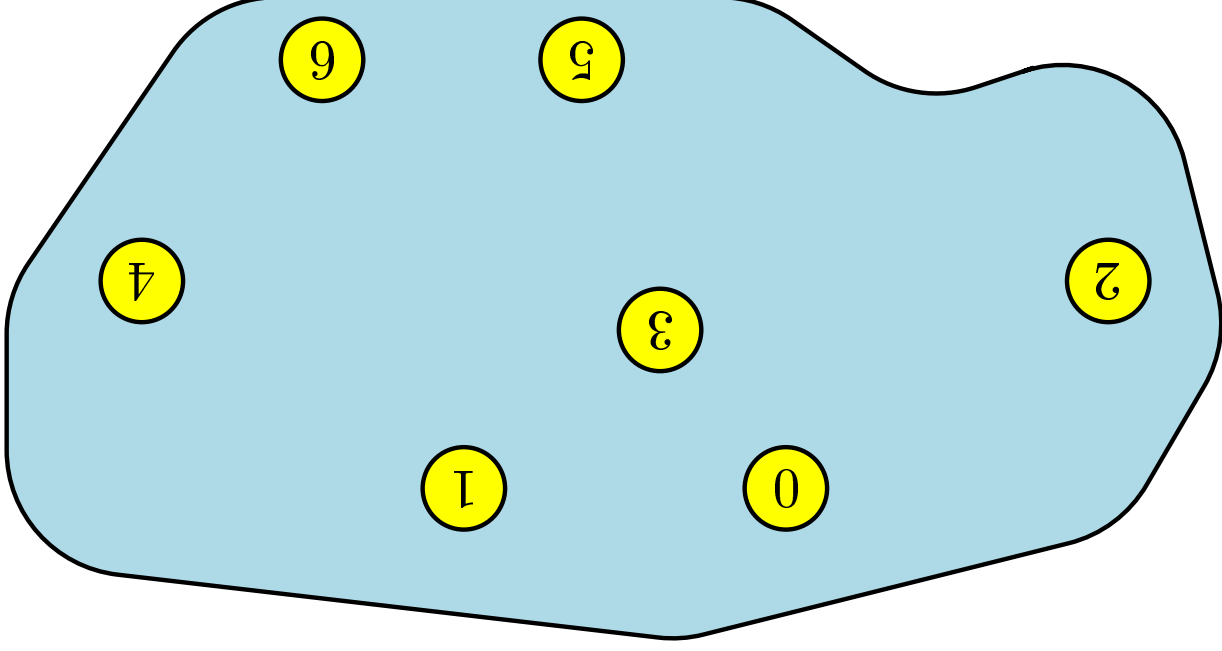


Fig. 6 – *Communicateur MPI_COMM_WORLD*

2 — Environnement : description

À tout instant, on peut connaître le nombre de processus gérés par un communicateur donné par le sous-programme `MPI_COMM_SIZE()` :

```
integer, intent(out) :: nb_procs, code
call MPI_COMM_SIZE(MPI_COMM_WORLD, nb_procs, code)
```

De même, le sous-programme `MPI_COMM_RANK()` permet d'obtenir le rang d'un processus (i.e. son numéro d'instance, qui est un nombre compris entre 0 et la valeur renvoyée par `MPI_COMM_SIZE()` - 1) :

```
integer, intent(out) :: rang, code
call MPI_COMM_RANK(MPI_COMM_WORLD, rang, code)
```

```

> mpiexec -np 7 qui-je-suis
je suis le processus 3 parmi 7
je suis le processus 0 parmi 7
je suis le processus 4 parmi 7
je suis le processus 1 parmi 7
je suis le processus 5 parmi 7
je suis le processus 2 parmi 7
je suis le processus 6 parmi 7
    
```

```

1 program qui-je-suis
2   implicit none
3   include 'mpif.h'
4   integer :: nb-procs, rang, code
5
6   call MPI_INIT(code)
7
8   call MPI_COMM_SIZE(MPI_COMM_WORLD, nb-procs, code)
9   call MPI_COMM_RANK(MPI_COMM_WORLD, rang, code)
10
11   print *, 'je suis le processus ', rang, ' parmi ', nb-procs
12
13   call MPI_FINALIZE(code)
14 end program qui-je-suis
    
```

2.2 – Exemple

3 – Communications point à point

3.1 – Notions générales

➔ Une communication dite **point à point** a lieu entre deux processus, l'un appelé processus **émetteur** et l'autre processus **destinataire**.

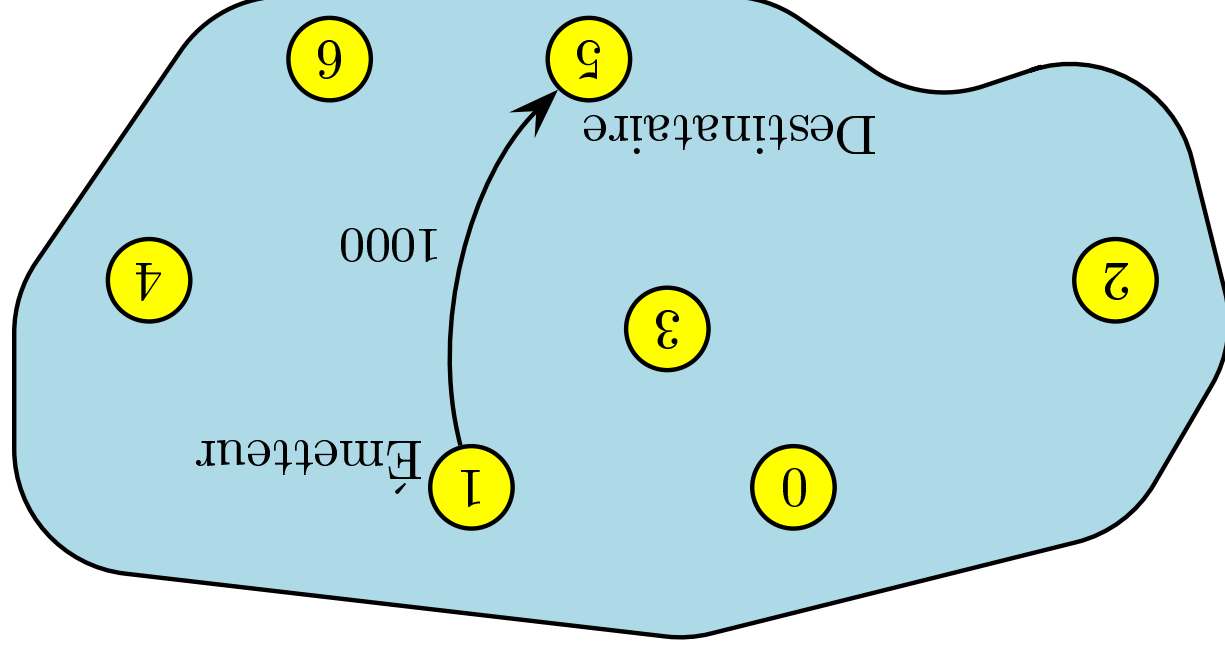


FIG. 7 – *Communication point à point*

➡ L'émetteur et le destinataire sont identifiés par leur **rang** dans le communicateur (i.e. leur numéro d'instance).

➡ Ce que l'on appelle l'**enveloppe d'un message** est constituée :

① du numéro du processus émetteur ;

② du numéro du processus récepteur ;

③ de l'étiquette (*tag*) du message ;

④ du nom du communicateur qui définira le contexte de communication de l'opération.

➡ Les données échangées sont **typées** (entiers, réels, etc. ou types dérivés personnels).

➡ Il existe dans chaque cas plusieurs **modes** de transfert, qui font appel à des protocoles différents. Ils seront vus au chapitre 5.

3 — Communications point à point : notions générales

```
1 program point-a-point
2 implicit none
3
4 include 'mpi.h'
5 integer, dimension(MPI_STATUS_SIZE) :: statut
6 integer, parameter :: etiquette=100
7 integer :: rang, valeur, code
8
9 call MPI_INIT(code)
10
11 call MPI_COMM_RANK(MPI_COMM_WORLD, rang, code)
12
13 if (rang == 1) then
14     valeur=1000
15     call MPI_SEND(valeur, 1, MPI_INTEGER, 5, etiquette, MPI_COMM_WORLD, code)
16 elseif (rang == 5) then
17     call MPI_RECV(valeur, 1, MPI_INTEGER, 1, etiquette, MPI_COMM_WORLD, statut, code)
18     print *, 'Moï, processus 5, j''ai reçu ', valeur, ' du processus 1.'
19 end if
20
21 call MPI_FINALIZE(code)
22
23 end program point-a-point
```

< mpiexec -np 7 point-a-point
Moï, processus 5, j'ai reçu 1000 du processus 1

Type MPI	Type Fortran
MPI_INTEGER	INTEGER
MPI_REAL	REAL
MPI_DOUBLE_PRECISION	DOUBLE PRECISION
MPI_COMPLEX	COMPLEX
MPI_LOGICAL	LOGICAL
MPI_CHARACTER	CHARACTER
MPI_PACKED	Types hétérogènes

Tab. 1 – Principaux types de données de base (Fortran)

3.2 – Types de données de base

TAB. 2 – Principaux types de données de base (C)

Type MPI	Type C
MPI_CHAR	signed char
MPI_SHORT	signed short
MPI_INT	signed int
MPI_LONG	signed long int
MPI_UNSIGNED_CHAR	unsigned char
MPI_UNSIGNED_SHORT	unsigned short
MPI_UNSIGNED_INT	unsigned int
MPI_UNSIGNED_LONG	unsigned long int
MPI_FLOAT	float
MPI_DOUBLE	double
MPI_LONG_DOUBLE	long double
MPI_PACKED	Types hétérogènes

3.3 – Autres possibilités

- ➡ À la réception d'un message, le rang du processus et l'étiquette peuvent être des « *jokers* », respectivement `MPI_ANY_SOURCE` et `MPI_ANY_TAG`.
- ➡ Une communication avec le processus « fictif » de rang `MPI_PROC_NULL` n'a aucun effet.
- ➡ On peut créer des structures de données plus complexes à l'aide de sous-programmes tels que `MPI_TYPE_CONTIGUOUS()`, `MPI_TYPE_VECTOR()`, `MPI_TYPE_INDEXED()` et `MPI_TYPE_STRUCT()` (voir le chapitre 6).

```

1 program sendrecv
2   implicit none
3   include 'mpi.h'
4   integer, dimension(MPI_STATUS_SIZE) :: statut
5   integer, parameter :: etiquette=100
6   integer :: rang, valeur, num_proc, code
7
8   call MPI_INIT(code)
9   call MPI_COMM_RANK(MPI_COMM_WORLD, rang, code)
10
11   ! On suppose avoir exactement 2 processus
12   num_proc=mod(rang+1,2)
13
14   call MPI_SENDRCV(rang+100,1,MPI_INTEGER,num_proc,etiquette,valeur,1,MPI_INTEGER, &
15   & num_proc,etiquette,MPI_COMM_WORLD,statut,code)
16
17   print *, 'Moï, processus ',rang,', j''ai reçu ',valeur,' du processus ',num_proc
18
19   call MPI_FINALIZE(code)
20 end program sendrecv

```

< mpiexec -np 2 sendrecv

Moï, processus 1, j'ai reçu 1000 du processus 0
 Moï, processus 0, j'ai reçu 1001 du processus 1

Il convient de noter que si le sous-programme `MPI_SEND` est implémenté de façon **bloquante** (voir le chapitre 5), ce code sera en situation de verrouillage, puisque chacun des deux processus attendra un ordre de réception qui ne viendra jamais.

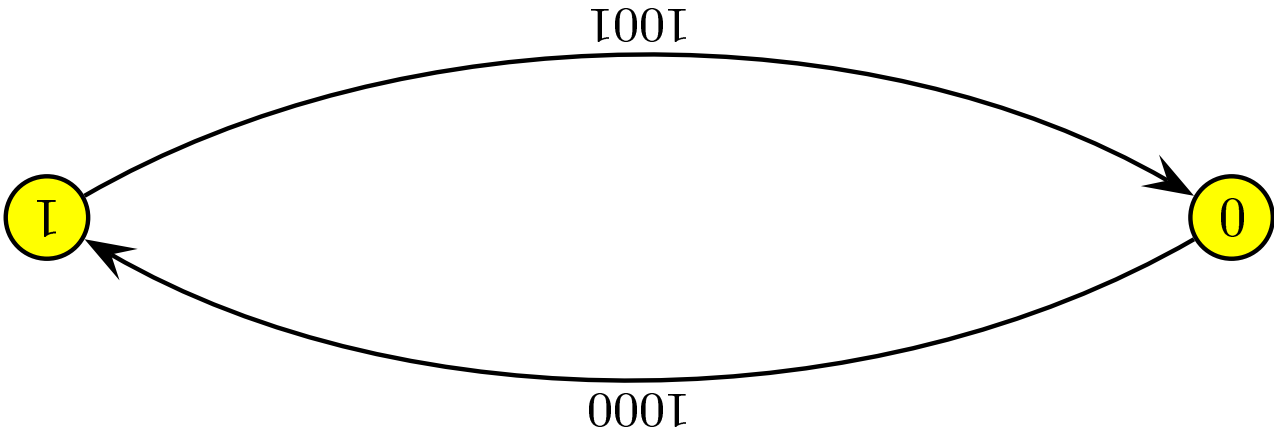


Fig. 8 – *Communication sendrecv entre les processus 0 et 1*

3.4 – Exemple : anneau de communication

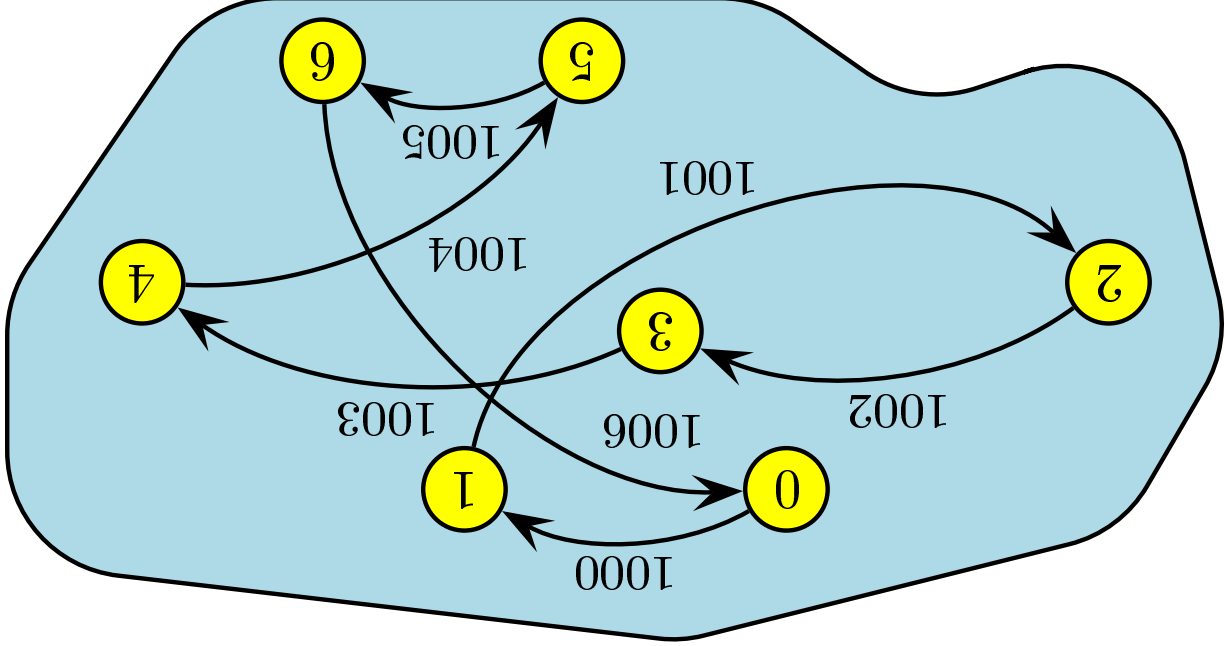


FIG. 9 – *Anneau de communication*

```

1  program anneau_1
2  implicit none
3  include 'mpi.h'
4  integer, dimension(MPI_STATUS_SIZE) :: statut
5  integer, parameter :: etiquette=100
6  integer :: nb_procs, rang, valeur, &
7  num_proc_precedent, num_proc_suivant, code
8
9  call MPI_INIT(code)
10 call MPI_COMM_SIZE(MPI_COMM_WORLD, nb_procs, code)
11 call MPI_COMM_RANK(MPI_COMM_WORLD, rang, code)
12
13 num_proc_suivant=mod(rang+1, nb_procs)
14 num_proc_precedent=mod(nb_procs+rang-1, nb_procs)
15
16 call MPI_SEND(rang+1000, 1, MPI_INTEGER, num_proc_suivant, etiquette, MPI_COMM_WORLD, code)
17 call MPI_RECV(valeur, 1, MPI_INTEGER, num_proc_precedent, etiquette, MPI_COMM_WORLD, &
18 statut, code)
19
20 print *, 'Moi, processus ', rang, ', j''ai reçu ', valeur, ' du processus ', &
21 num_proc_precedent
22
23 call MPI_FINALIZE(code)
24 end program anneau_1
    
```



```
> mpi_run -np 7 anneau_1
```

```
Moï, processus 1, j'ai reçu 1000 du processus 0
Moï, processus 3, j'ai reçu 1002 du processus 2
Moï, processus 5, j'ai reçu 1004 du processus 4
Moï, processus 0, j'ai reçu 1006 du processus 6
Moï, processus 2, j'ai reçu 1001 du processus 1
Moï, processus 4, j'ai reçu 1003 du processus 3
Moï, processus 6, j'ai reçu 1005 du processus 5
```

Cet exemple peut se programmer autrement, en utilisant une seule opération pour faire un envoi et une réception par processus, puisqu'à chaque fois les processus récepteur et émetteur sont définis de façon similaire.

```

1  program anneau_2
2  implicit none
3  include 'mpi.h'
4  integer, dimension(MPI_STATUS_SIZE) :: statut
5  integer, parameter :: etiquette=100
6  integer :: nb_procs, rang, valeur, &
7  num_proc_precedent, num_proc_suivant, code
8
9  call MPI_INIT(code)
10 call MPI_COMM_SIZE(MPI_COMM_WORLD, nb_procs, code)
11 call MPI_COMM_RANK(MPI_COMM_WORLD, rang, code)
12 valeur=rang+1000
13 num_proc_suivant=mod(rang+1, nb_procs)
14 num_proc_precedent=mod(nb_procs+rang-1, nb_procs)
15 call MPI_SENDRECV_REPLACE(valeur, 1, MPI_INTEGER, num_proc_suivant, etiquette, &
16 num_proc_precedent, etiquette, MPI_COMM_WORLD, statut, code)
17 print *, 'Mot, processus, rang, j, ai reçu, valeur, du processus', num_proc
18 call MPI_FINALIZE(code)
19 end program anneau_2

```

Mais en fait on ne fait pas ici de communications en **anneau**, puisque celles-ci se font toutes simultanément.. Il faut procéder différemment si on veut imposer véritablement l'organisation suivante, où les processus représentent du point de vue logique une topologie en **anneau** :

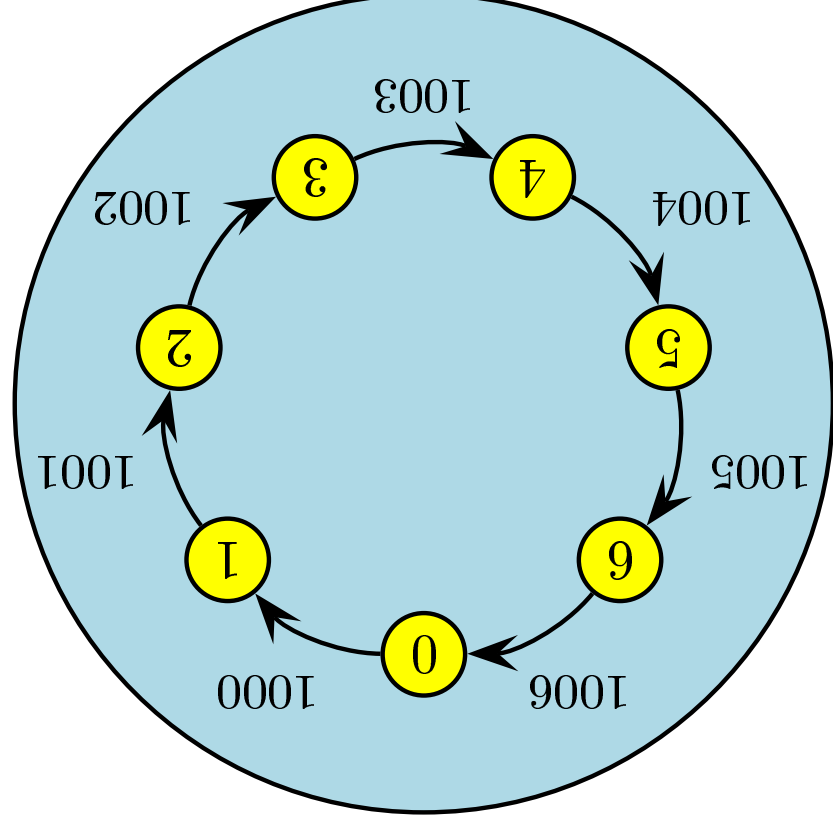


FIG. 10 – *Anneau de communication*

Dans ce cas, il faut donc programmer autrement pour **ordonner** les écritures. Il est alors nécessaire de gérer autrement l'ordre dans lequel vont se faire les communications :

```

17 if (rang == 0) then
18   call MPI_SEND(rang+100, 1, MPI_INTEGER, num_proc_suivant, etiquette, &
19               MPI_COMM_WORLD, code)
20   call MPI_RECV(valeur, 1, MPI_INTEGER, num_proc_precedent, etiquette, &
21               MPI_COMM_WORLD, statut, code)
22 else
23   call MPI_RECV(valeur, 1, MPI_INTEGER, num_proc_precedent, etiquette, &
24               MPI_COMM_WORLD, statut, code)
25   call MPI_SEND(rang+100, 1, MPI_INTEGER, num_proc_suivant, etiquette, &
26               MPI_COMM_WORLD, code)
27 end if

```

```
> mpiexec -np 7 anneau_3
Moï, processus 1, j'ai reçu 1000 du processus 0
Moï, processus 2, j'ai reçu 1001 du processus 1
Moï, processus 3, j'ai reçu 1002 du processus 2
Moï, processus 4, j'ai reçu 1003 du processus 3
Moï, processus 5, j'ai reçu 1004 du processus 4
Moï, processus 6, j'ai reçu 1005 du processus 5
Moï, processus 0, j'ai reçu 1006 du processus 6
```

3.5 – Construction et reconstruction de messages

Dans le cas où le processus récepteur ne connaît pas le nombre d'éléments à recevoir, il est possible de n'envoyer qu'un seul message regroupant ces informations. Ceci évite d'abord un message contenant le nombre d'éléments puis un autre contenant les éléments eux-mêmes. Pour ce faire, il faut construire en émission le contenu du message en passant par une zone tampon, via le sous-programme `MPI_PACK`. En réception, le sous-programme `MPI_UNPACK` permet de façon réciproque de reconstruire le message.

```

1  program pack_unpack
2  implicit none
3
4  include 'mpif.h'
5  integer, dimension(MPI_STATUS_SIZE) :: statut
6  integer, parameter :: etiquette=999, &
7  longueur_msg_max=100 i En octets
8  integer :: rang,nb_valeurs,position,code
9  real(kind=8), allocatable :: valeurs
10 real(kind=8), dimension(longueur_msg_max) :: message
11
12 call MPI_INIT(code)
13 call MPI_COMM_RANK(MPI_COMM_WORLD,rang,code)

```

```

14 nb_valuers=rang*3 + 2
15 allocate(valuers(nb_valuers))
16 call random_number(valuers)
17
18 position=1
19 call MPI_PACK (nb_valuers,1,MPI_INTEGER,message,longueur_msg_max, &
20 position,MPI_COMM_WORLD,code) ; position est mis à jour en sortie
21 call MPI_PACK (valuers,nb_valuers,MPI_DOUBLE_PRECISION,message,longueur_msg_max, &
22 position,MPI_COMM_WORLD,code) ; position est mis à jour en sortie
23 call MPI_SEND (message,position,MPI_PACKED,mod(rang+1,2),etiquette,MPI_COMM_WORLD, &
24 code)
25 deallocate(valuers)
26
27 call MPI_RECV (message,longueur_msg_max,MPI_PACKED,mod(rang+1,2),etiquette, &
28 position=1,MPI_COMM_WORLD,statut,code)
29
30 position=1
31 call MPI_UNPACK (message,longueur_msg_max,position, & i position mis à jour en sortie
32 nb_valuers,1,MPI_INTEGER,MPI_COMM_WORLD,code)
33 allocate(valuers(nb_valuers))
34 call MPI_UNPACK (message,longueur_msg_max,position, & i position mis à jour en sortie
35 valuers,nb_valuers,MPI_DOUBLE_PRECISION,MPI_COMM_WORLD,code)
36
37 print *, 'Mot, processus ',rang,', j''ai reçu ',nb_valuers,' valeurs du processus ', &
38 mod(rang+1,2), ' : ',valuers(1:nb_valuers)
39 deallocate(valuers)
40
41 call MPI_FINALIZE(code)
42 end program pack_unpack

```

```
> mpiexec -np 2 pack_unpack  
Moï, processus 1, j'ai reçu 2 valeurs du processus 0 : 0.58, 0.95  
Moï, processus 0, j'ai reçu 5 valeurs du processus 1 : 0.58, 0.95, 0.78, 0.29, 0.45
```


4 – Communications collectives

4.1 – Notions générales

Les communications **collectives** permettent de faire en une seule opération une série de communications point à point.

Une communication collective concerne toujours les processus du

communicateur indiqué.

Pour chacun des processus, l'appel se termine lorsque la participation de celui-ci à l'opération collective est achevée, au sens des communications point-à-point (donc quand la zone mémoire concernée peut être modifiée).

Il est inutile d'ajouter une synchronisation globale (barrière) après une opération

collective.

La gestion des **étiquettes** dans ces communications est transparente et à la charge du système. Elles ne sont donc jamais définies explicitement lors de l'appel à ces sous-programmes. Cela a entre autres pour avantage que les communications collectives n'interfèrent jamais avec les communications point à point.

Il y a trois types de sous-programmes :

① celle qui assure les synchronisations globales : `MPI_BARRIER()`.

② celles qui ne font que transférer des données :

diffusion globale de données : `MPI_BCAST()` ;

diffusion sélective de données : `MPI_SCATTER()` ;

collecte de données réparties : `MPI_GATHER()` ;

collecte par tous les processus de données réparties : `MPI_ALLGATHER()` ;

diffusion sélective, par tous les processus, de données réparties :

`MPI_ALLTOALL()`.

③ celles qui, en plus de la gestion des communications, effectuent des opérations

sur les données transférées :

opérations de réduction, qu'elles soient d'un type prédéfini (somme, produit, maximum, minimum, etc.) ou d'un type personnel : `MPI_REDUCE()` ;

opérations de réduction avec diffusion du résultat (il s'agit en fait d'un `MPI_REDUCE()` suivi d'un `MPI_BCAST()`) : `MPI_ALLREDUCE()`.

4.2 – Synchronisation globale : MPI_BARRIER()

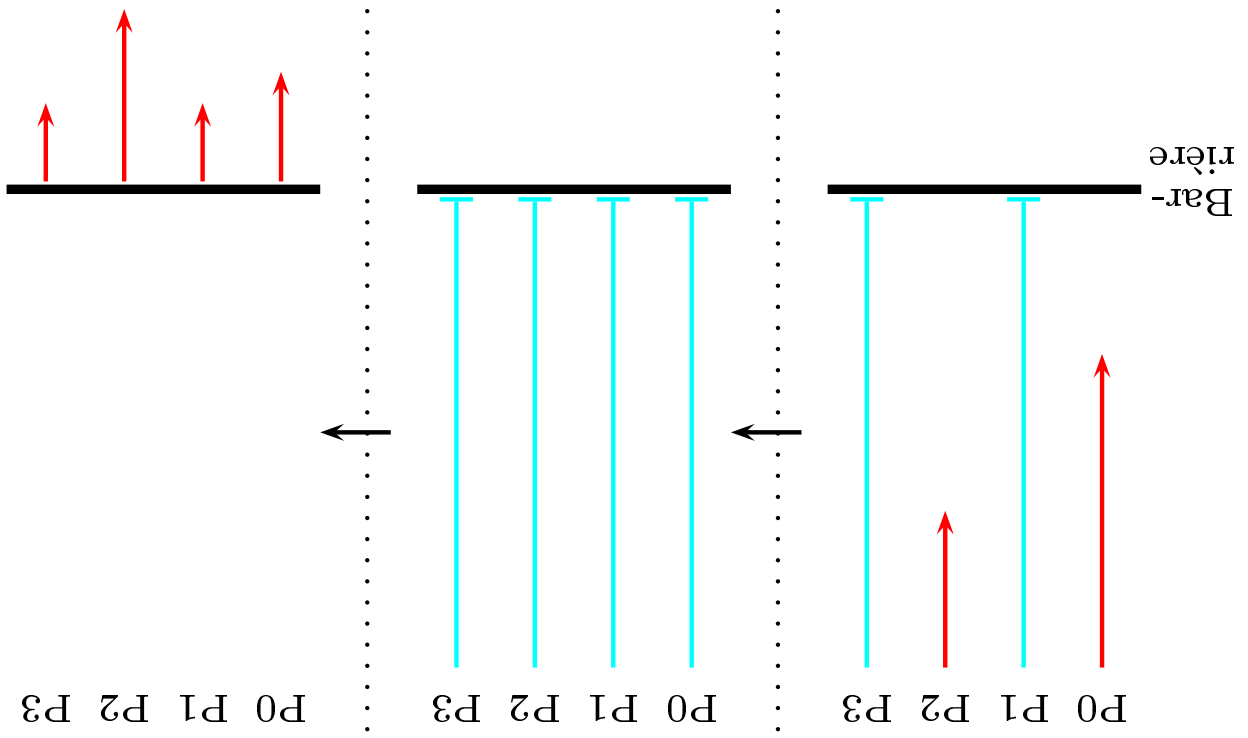


FIG. 11 – Synchronisation globale : MPI_BARRIER()

```
integer, intent(out) :: code
call MPI_BARRIER(MPI_COMM_WORLD, code)
```

4.3 – Diffusion générale : MPI_BCAST()

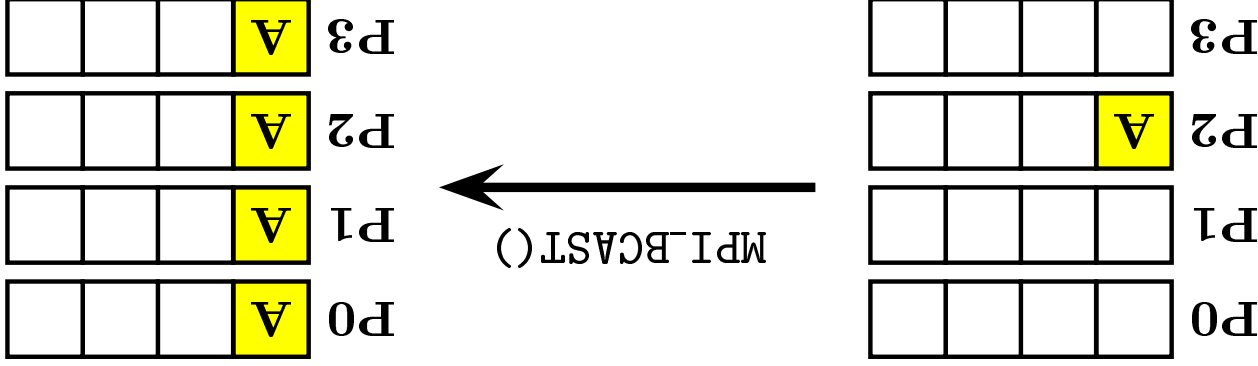
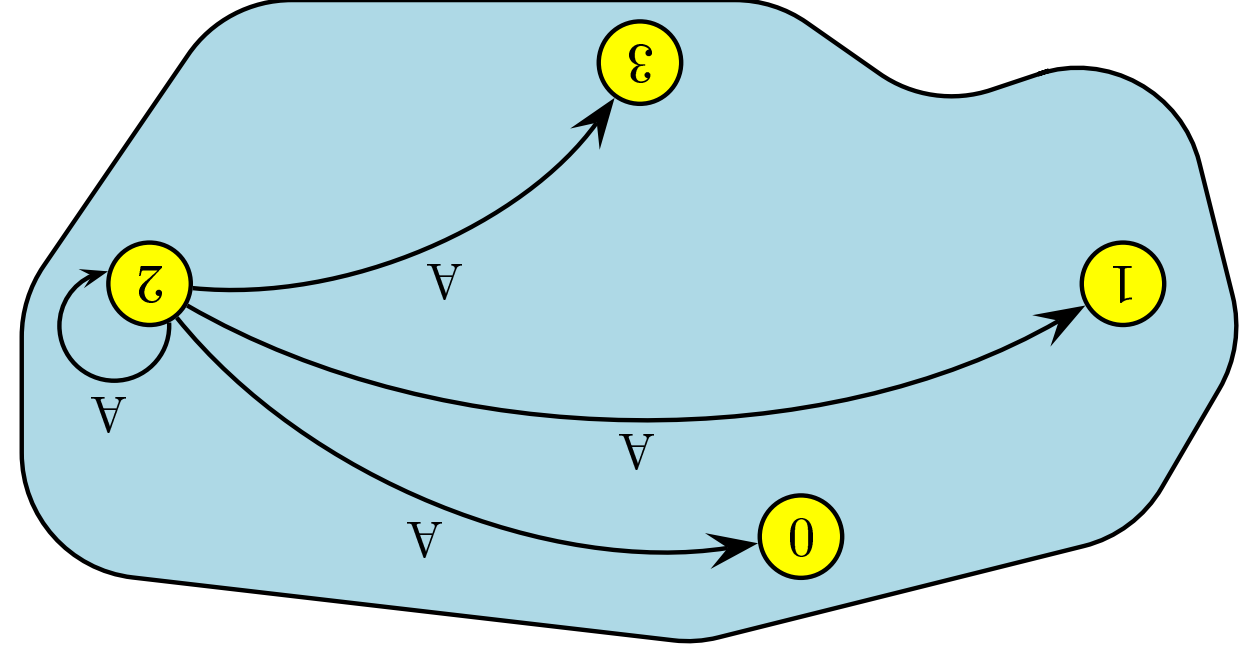


FIG. 12 – Diffusion générale : MPI_BCAST()

> mpiexec -np 4 bcast

Moï, processus 2, j'ai reçu 1002 du processus 2
 Moï, processus 0, j'ai reçu 1002 du processus 2
 Moï, processus 1, j'ai reçu 1002 du processus 2
 Moï, processus 3, j'ai reçu 1002 du processus 2

```

1  program bcast
2  implicit none
3
4  include 'mpi.h'
5  integer :: rang, valeur, code
6
7  call MPI_INIT(code)
8  call MPI_COMM_RANK(MPI_COMM_WORLD, rang, code)
9
10 if (rang == 2) valeur=rang+1000
11
12 call MPI_BCAST(valeur, 1, MPI_INTEGER, 2, MPI_COMM_WORLD, code)
13
14 print *, 'Moï, processus ', rang, ', j'ai reçu ', valeur, ' du processus 2'
15
16 call MPI_FINALIZE(code)
17
18 end program bcast
    
```

4.4 – Diffusion sélective : MPI_SCATTER()

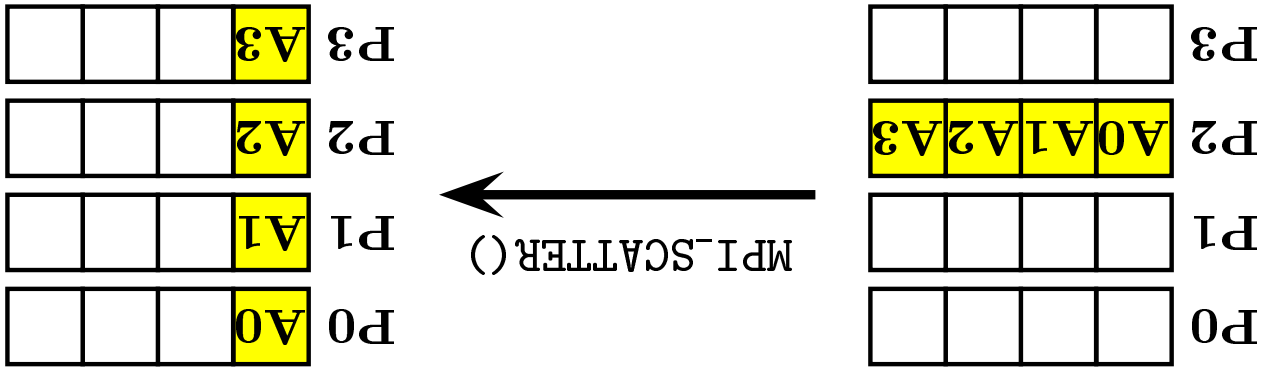
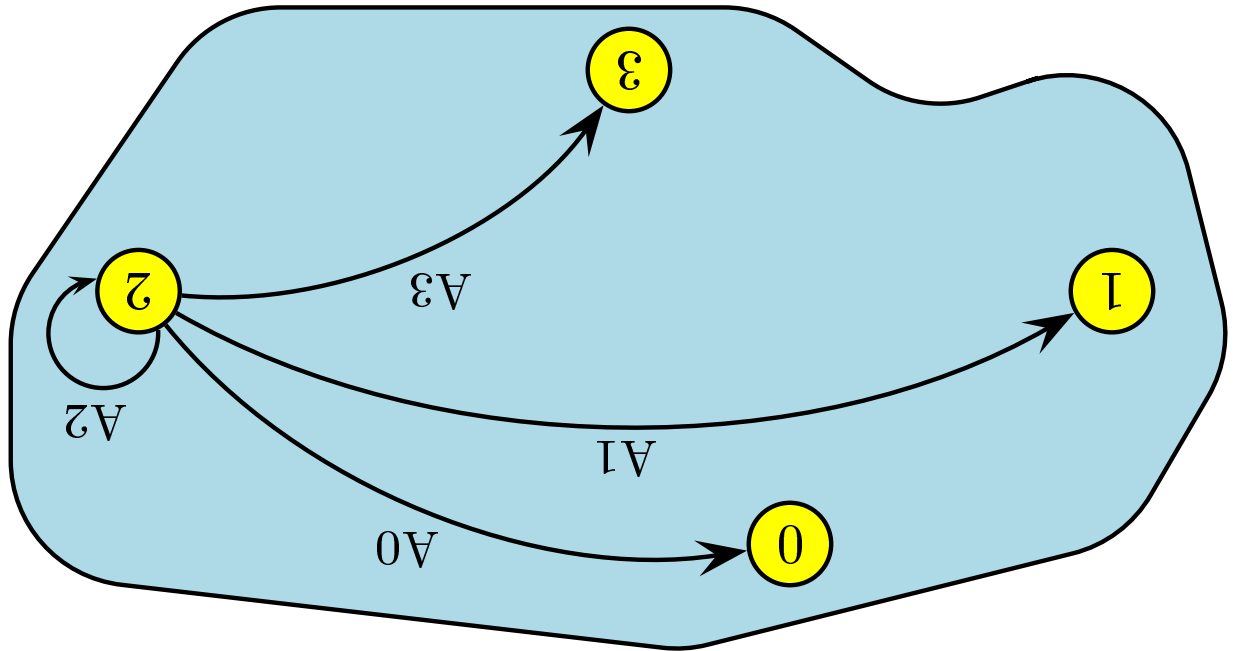


FIG. 13 – Diffusion sélective : MPI_SCATTER()

```

1  program scatter
2  implicit none
3
4  include 'mpif.h'
5  integer, parameter :: nb_valuers=128
6  integer :: nb_procs, rang, longueur_tranche, i, code
7  real, allocatable, dimension(:) :: valuers, donnees
8
9  call MPI_INIT(code)
10
11  call MPI_COMM_SIZE(MPI_COMM_WORLD, nb_procs, code)
12  call MPI_COMM_RANK(MPI_COMM_WORLD, rang, code)
13
14  longueur_tranche=nb_valuers/nb_procs
15  allocate(donnees(longueur_tranche))
16
17  if (rang == 2) then
18      allocate(valuers(nb_valuers))
19      valuers(:)=(/(1000.+i, i=1, nb_valuers)/)
20  end if
21
22  call MPI_SCATTER(valuers, longueur_tranche, MPI_REAL,
23                  donnees, longueur_tranche, MPI_REAL,
24                  &
25                  print *, 'Moi, processus ', rang, ', j'ai reçu ', donnees(1), ' à ', &
26                  donnees(longueur_tranche), ' du processus 2',
27
28  call MPI_FINALIZE(code)
29  end program scatter

```

```
> mpiexec -np 4 scatter
```

Moï, pr.	0,	j'ai reçu	1001.	à	1032.	du pr.	2
Moï, pr.	1,	j'ai reçu	1033.	à	1064.	du pr.	2
Moï, pr.	3,	j'ai reçu	1097.	à	1128.	du pr.	2
Moï, pr.	2,	j'ai reçu	1065.	à	1096.	du pr.	2

4.5 – Collecte : MPI_GATHER()

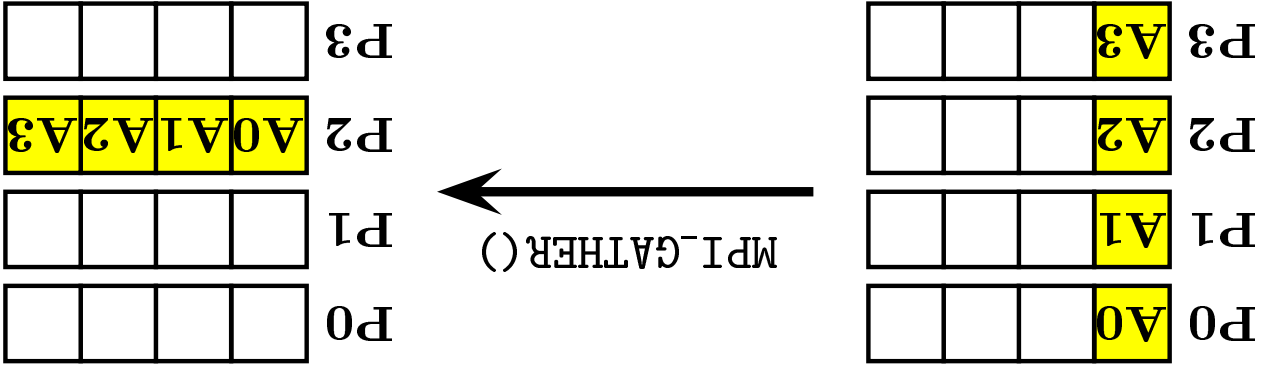
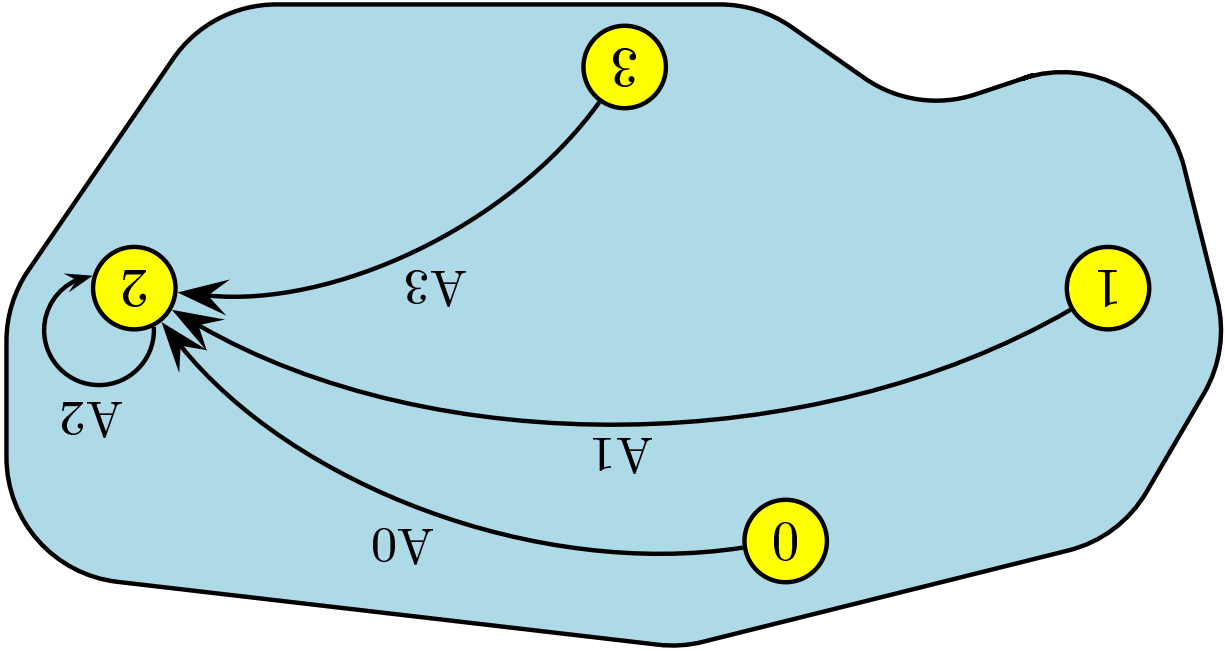


FIG. 14 – Collecte : MPI_GATHER()

```
> mpiexec -np 4 gather  
Moï, pr. 2, j'ai regu 1001. ... 1033. ... 1128.
```

```
1  program gather
2  implicit none
3  include 'mpi.h'
4  integer, parameter :: nb_valuers=128
5  integer :: nb_procs, rang, longueur_tranche, i, code
6  real, dimension(nb_valuers) :: donnees
7  real, allocatable, dimension(:) :: valuers
8
9  call MPI_INIT(code)
10 call MPI_COMM_SIZE(MPI_COMM_WORLD, nb_procs, code)
11 call MPI_COMM_RANK(MPI_COMM_WORLD, rang, code)
12 longueur_tranche=nb_valuers/nb_procs
13 allocate(valuers(longueur_tranche))
14
15 valuers(:)=(/(1000.+rang*longueur_tranche+i, i=1, longueur_tranche)/)
16
17 call MPI_GATHER(valuers, longueur_tranche, MPI_REAL, MPI_REAL, MPI_COMM_WORLD, code), &
18 if (rang == 2) print *, 'Moï, processus 2, j'ai regu ', donnees(1), ' ... ', &
19 donnees(longueur_tranche+1), ' ... ', donnees(nb_valuers)
20
21 call MPI_FINALIZE(code)
22
23 end program gather
24
```

4.6 – Collecte générale : MPI_ALLGATHER()

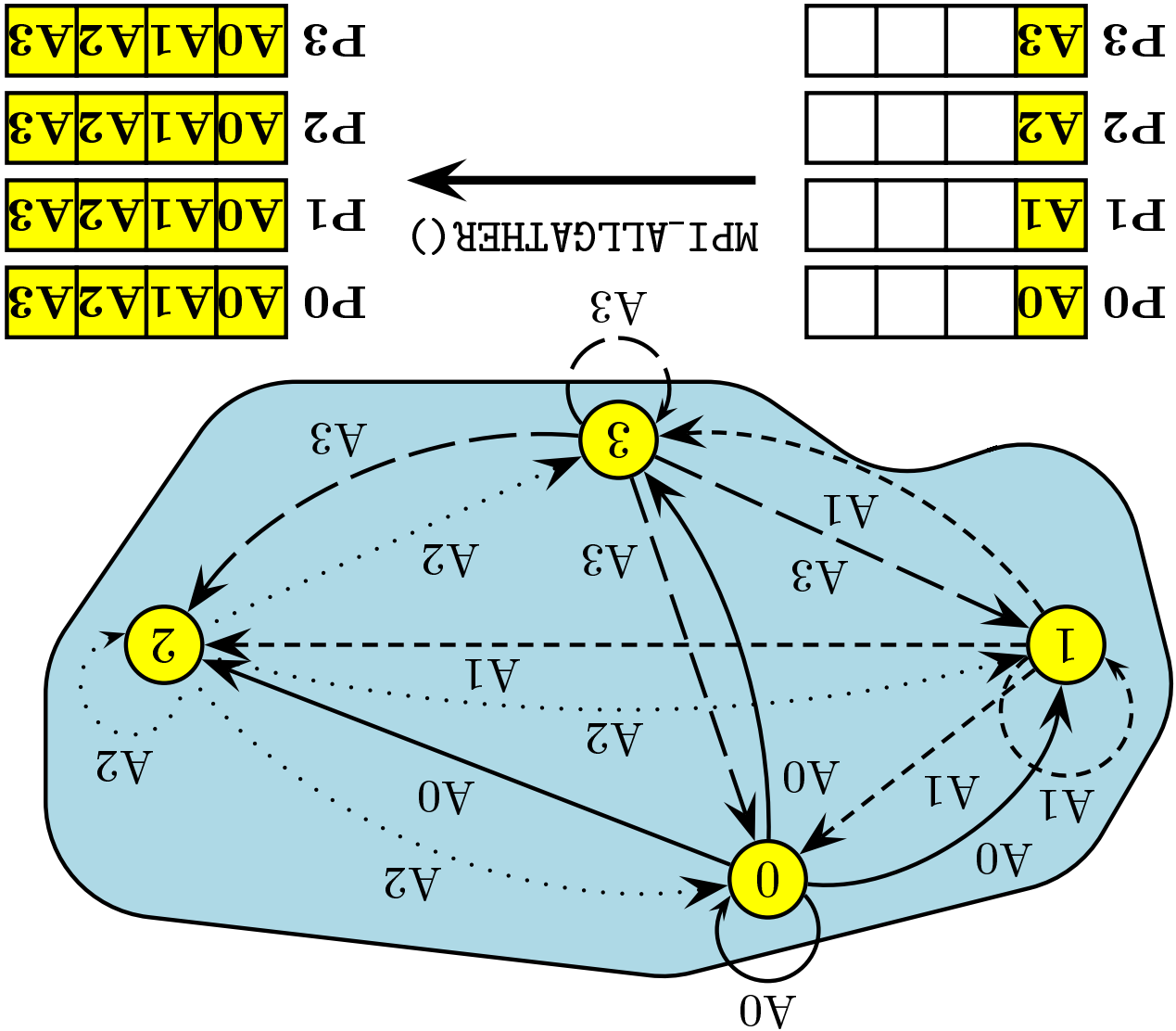


FIG. 15 – Collecte générale : MPI_ALLGATHER()



```
1 program allgather
2 implicit none
```

```
3 include 'mpi.h'
```

```
4 integer, parameter
```

```
5 :: nb_values=128
```

```
6 integer :: nb_procs, rang, longueur_tranche, i, code
```

```
7 real, dimension(nb_values)
```

```
8 real, allocatable, dimension(:) :: valeurs
```

```
9 call MPI_INIT(code)
```

```
10 call MPI_COMM_SIZE(MPI_COMM_WORLD, nb_procs, code)
```

```
11 call MPI_COMM_RANK(MPI_COMM_WORLD, rang, code)
```

```
12 longueur_tranche=nb_values/nb_procs
13 allocate(valeurs(longueur_tranche))
```

```
14 valeurs(:)=(/(1000.+rang*longueur_tranche+i, i=1, longueur_tranche)/)
```

```
15 call MPI_ALLGATHER(valeurs, longueur_tranche, MPI_REAL, MPI_REAL, MPI_COMM_WORLD, code)
16 print *, 'Moi, processus ', rang, ', j'ai reçu ', donnees(1), ' ... ', &
```

```
17 donnees(longueur_tranche+1), ' ... ', donnees(nb_values)
18 end program allgather
```

```
28
```

```
26 call MPI_FINALIZE(code)
```

```
27
```

```
28
```

```
> mpirun -np 4 allgather
```

Moï, pr.	1	, j'ai regu	1001.	...	1033.	...	1128.
Moï, pr.	3	, j'ai regu	1001.	...	1033.	...	1128.
Moï, pr.	2	, j'ai regu	1001.	...	1033.	...	1128.
Moï, pr.	0	, j'ai regu	1001.	...	1033.	...	1128.

4.7 - Échanges croisés : MPI_ALLTOALL()

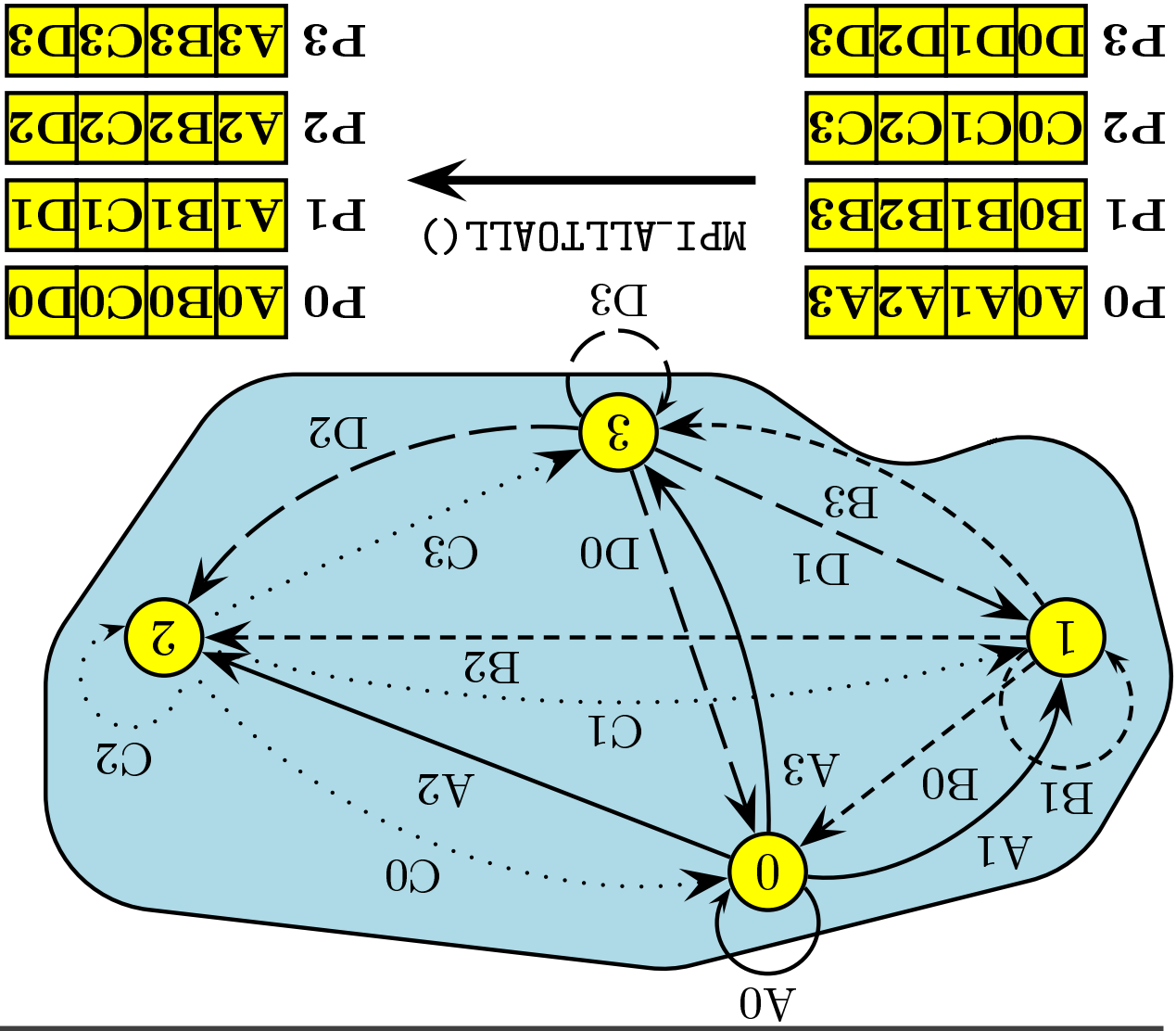


FIG. 16 - Échanges croisés : MPI_ALLTOALL()

P3	D0	D1	D2	D3
P2	C0	C1	C2	C3
P1	B0	B1	B2	B3
P0	A0	A1	A2	A3

MPI_ALLTOALL()

P3	A3	B3	C3	D3
P2	A2	B2	C2	D2
P1	A1	B1	C1	D1
P0	A0	B0	C0	D0

```

1 program alltoall
2 implicit none
3
4 include 'mpit.h'
5 integer, parameter :: nb_values=128
6 integer :: nb_procs, rang, longueur_tranche, i, code
7 real, dimension(nb_values) :: valeurs, donnees
8
9 call MPI_INIT(code)
10 call MPI_COMM_SIZE(MPI_COMM_WORLD, nb_procs, code)
11 call MPI_COMM_RANK(MPI_COMM_WORLD, rang, code)
12
13 valeurs(:)=(/ (1000.+rang*nb_values+i, i=1, nb_values) /)
14 longueur_tranche=nb_values/nb_procs
15
16 call MPI_ALLTOALL(values, longueur_tranche, MPI_REAL, MPI_REAL,
17                 donnees, longueur_tranche, code)
18
19 print *, 'Mot, processus, rang, j, ai regu, donnees(1), ...', &
20         ' donnees(longueur_tranche+1), ...', donnees(nb_values)
21
22 call MPI_FINALIZE(code)
23
24 end program alltoall

```

```
> mpiexec -np 4 alltoall1
```

Moï, processus 0, j'ai reçu	1001.	...	1129.	...	1416.
Moï, processus 2, j'ai reçu	1065.	...	1193.	...	1480.
Moï, processus 1, j'ai reçu	1033.	...	1161.	...	1448.
Moï, processus 3, j'ai reçu	1097.	...	1225.	...	1512.

4.8 – Réductions réparties

→ Une **réduction** est une opération appliquée à un ensemble d'éléments pour en obtenir une seule valeur. Des exemples typiques sont la somme des éléments d'un vecteur `SUM(A(:))` ou la recherche de l'élément de valeur maximum dans un vecteur `MAX(V(:))`.

→ *MPI* propose des sous-programmes de haut-niveau pour opérer des réductions sur des données réparties sur un ensemble de processus, avec récupération du résultat sur un seul processus (`MPI_REDUCE()`) ou bien sur tous (`MPI_ALLREDUCE()`), qui est en fait seulement un `MPI_REDUCE()` suivi d'un `MPI_BCAST()`.

→ Si plusieurs éléments sont concernés par processus, la fonction de réduction est appliquée à chacun d'entre eux.

→ Le sous-programme `MPI_SCAN()` permet en plus d'effectuer des réductions partielles en considérant, pour chaque processus, les processus précédents du groupe.

→ Les sous-programmes `MPI_OP_CREATE()` et `MPI_OP_FREE()` permettent de définir des opérations de réduction personnelles.

TAB. 3 – Principales opérations de réduction prédéfinies (il existe aussi d'autres opérations logiques)

Nom	Opération
MPI_SUM	Somme des éléments
MPI_PROD	Produit des éléments
MPI_MAX	Recherche du maximum
MPI_MIN	Recherche du minimum
MPI_MAXLOC	Recherche de l'indice du maximum
MPI_MINLOC	Recherche de l'indice du minimum
MPI_LAND	ET logique
MPI_LOR	OU logique
MPI_LXOR	OU exclusif logique

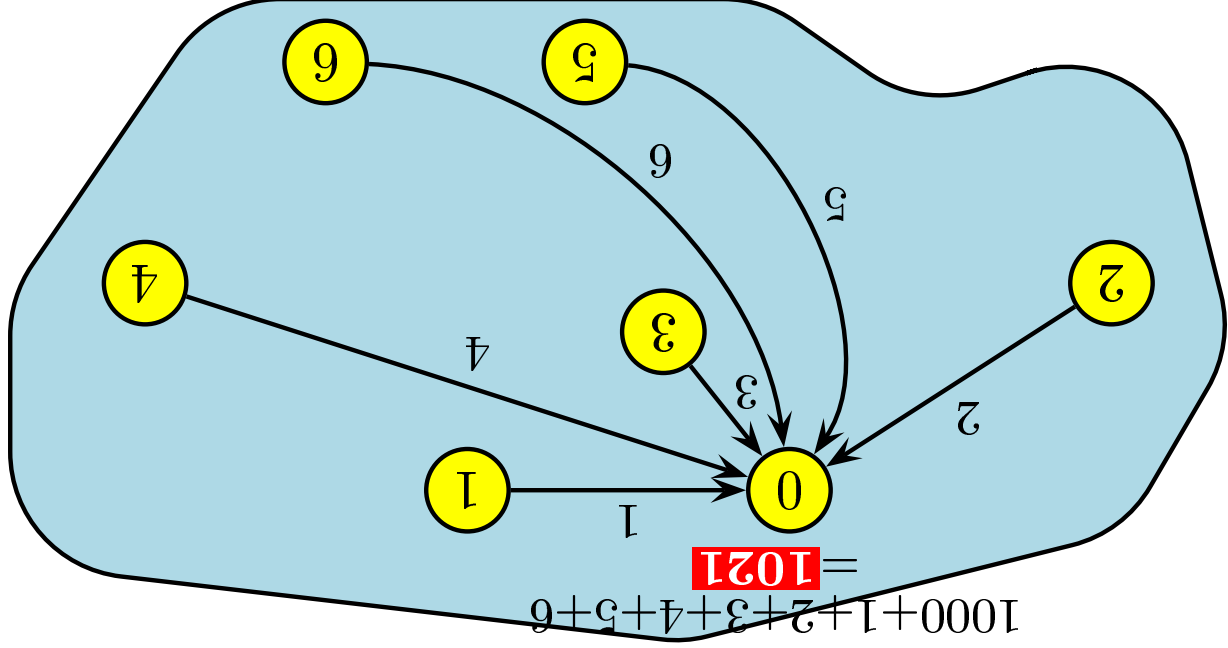


FIG. 17 – Réduction répartie (somme)

Moi, processus 0, j'ai pour v. de la somme globale 1021

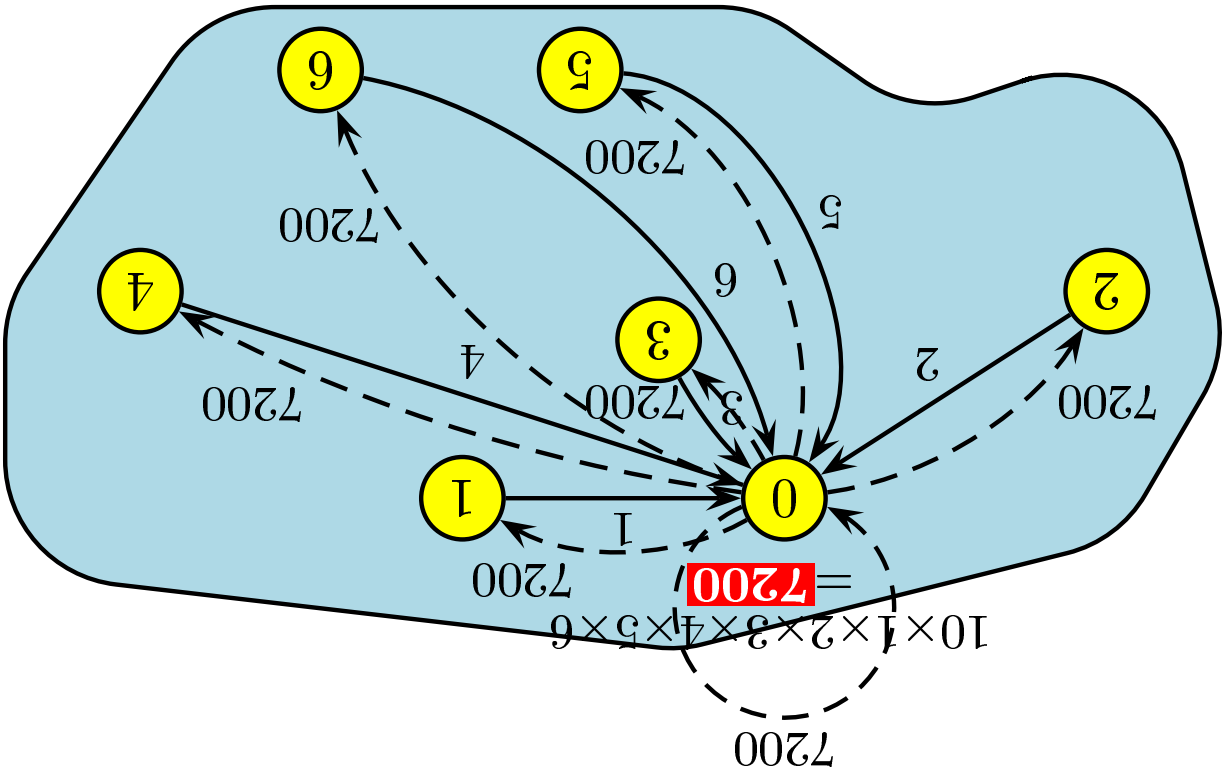
< mpirun -np 7 reduce

```

1 program reduce
2   implicit none
3   include 'mpif.h'
4   integer :: nb_procs, rang, valeur, somme, code
5
6   call MPI_INIT(code)
7   call MPI_COMM_SIZE(MPI_COMM_WORLD, nb_procs, code)
8   call MPI_COMM_RANK(MPI_COMM_WORLD, rang, code)
9
10  if (rang == 0) then
11    valeur=1000
12  else
13    valeur=rang
14  endif
15
16  call MPI_REDUCE(valeur, somme, 1, MPI_INTEGER, MPI_SUM, 0, MPI_COMM_WORLD, code)
17
18  if (rang == 0) then
19    print *, 'Moi, processus 0, j'ai pour ', valeur de la somme globale ', somme
20  endif
21
22  call MPI_FINALIZE(code)
23  end program reduce

```

FIG. 18 – Réduction répartie avec diffusion du résultat (produit)



```

1 program allreduce
2
3 implicit none
4 include 'mpi.h'
5 integer :: nb-procs, rang, valeur, produit, code
6
7 call MPI_INIT(code)
8 call MPI_COMM_SIZE(MPI_COMM_WORLD, nb-procs, code)
9 call MPI_COMM_RANK(MPI_COMM_WORLD, rang, code)
10
11 if (rang == 0) then
12     valeur=10
13 else
14     valeur=rang
15 endif
16
17 call MPI_ALLREDUCE(valeur, produit, 1, MPI_INTEGER, MPI_PROD, MPI_COMM_WORLD, code)
18
19 print *, 'Moi, processus ', rang, ', j''ai reçu ', la valeur du produit global ', produit
20
21 call MPI_FINALIZE(code)
22
23 end program allreduce

```

```
> mpiexec -np 7 allreduce
```

```
Moï, pr. 6, j'ai reçu val. du produit global 7200  
Moï, pr. 2, j'ai reçu val. du produit global 7200  
Moï, pr. 0, j'ai reçu val. du produit global 7200  
Moï, pr. 4, j'ai reçu val. du produit global 7200  
Moï, pr. 5, j'ai reçu val. du produit global 7200  
Moï, pr. 3, j'ai reçu val. du produit global 7200  
Moï, pr. 1, j'ai reçu val. du produit global 7200
```

```

1 program ma_reduction
2 implicit none
3 include 'mpi.h'
4 integer :: rang, code, i, mon_operation
5 integer, parameter :: n=4
6 complex, dimension(n) :: a, resultat
7 external mon_produit
8
9 call MPI_INIT(code)
10 call MPI_COMM_RANK(MPI_COMM_WORLD, rang, code)
11
12 i Initialisation du vecteur A sur chaque processus
13 a(:) = (/ (cmplx(rang+1, rang+1+1), i=1, n) /)
14
15 i Creation de l'opération commutative mon_operation
16 call MPI_OP_CREATE(mon_produit, .true., mon_operation, code)
17
18 i Collecte sur le processus 0 du produit global
19 call MPI_REDUCE(a, resultat, n, MPI_COMPLEX, mon_operation, 0, MPI_COMM_WORLD, code)
20
21 i Affichage du resultat
22 if (rang == 0) then
23     print *, 'Valeur du produit', resultat
24 end if
25
26 call MPI_FINALIZE(code)
27 end program ma_reduction
    
```


4 — Communications collectives : réductions réparties

```

1  i Définition du produit terme à terme de deux vecteurs de nombres complexes
2
3  integer fonction mon_produit(vecteur1, vecteur2, longueur, type_donnee) result(itnutilise)
4      implicit none
5
6      complex, dimension(longueur) :: vecteur1, vecteur2
7      integer :: longueur, type_donnee, i
8
9      do i=1, longueur
10         vecteur2(i) = cmplx(real(vecteur1(i))*real(vecteur2(i)) - &
11             aimag(vecteur1(i))*aimag(vecteur2(i)), &
12             real(vecteur1(i))*aimag(vecteur2(i)) + &
13             aimag(vecteur1(i))*real(vecteur2(i)))
14     end do
15
16     itnutilise=0
17
18 end fonction mon_produit
    
```

< mpiexec -np 5 ma_reduction

Valeur du produit (155., -2010.), (-1390., -8195.), (-7215., -23420.), (-22000., -54765.)

4.9 – Compléments

Les sous-programmes `MPI_SCATTERV()`, `MPI_GATHERV()` et `MPI_ALLGATHERV()` étendent `MPI_SCATTER()`, `MPI_GATHER()` et `MPI_ALLGATHER()` au cas où le nombre d'éléments à diffuser ou collecter est différent suivant les processus.

5 – Optimisations

5.1 – Introduction

☞ L'optimisation doit être un souci essentiel lorsque la part des communications par rapport aux calculs devient assez importante

☞ L'optimisation des communications peut s'accomplir à différents niveaux dont les principaux sont :

- ① recouvrir les communications par des calculs ;
- ② éviter si possible la copie du message dans un espace mémoire temporaire (*buffering*) ;
- ③ minimiser les surcoûts induits par des appels répétitifs aux sous-programmes de communication.

5.2 – Programme modèle

```

1 program Optimiser
2 implicit none
3 include 'mpif.h'
4
5 integer, parameter
6 integer, parameter
7 real, dimension(na,na)
8 real, dimension(nb,nb)
9 real, dimension(na)
10 real, dimension(nb)
11 real, dimension(m,m)
12 integer
13 real(kind=8)
14 integer, dimension(MPI_STATUS_SIZE)
15
16 *** Initialisation MPI
17 call MPI_INIT(code)
18 call MPI_COMM_SIZE(MPI_COMM_WORLD, nb_procs, code)
19 call MPI_COMM_RANK(MPI_COMM_WORLD, rang, code)
20
21 *** Initialisation des tableaux
22 call random_number(a)
23 call random_number(b)
24 call random_number(c)

```

```

25 temps_debut = MPI_WTIME ()
26 if (rang == 0) then
27     *** J'envoie un gros message
28     call MPI_SEND (c,m*m,MPI_REAL,1,etiquette,MPI_COMM_WORLD,code)
29     *** Ce calcul (factorisation LU avec LAPACK) modifie le contenu du tableau A
30     call sgetrf(na, a, na, pivota, code)
31     *** Ce calcul modifie le contenu du tableau C
32     c(1:nb,1:nb) = matmul(a(1:nb,1:nb),b)
33     elseif (rang == 1) then
34         *** Je calcule
35         call sgetrf(na, a, na, pivota, code)
36         *** Je reçois le gros message
37         call MPI_RECV (c,m*m,MPI_REAL,0,etiquette,MPI_COMM_WORLD,statut,code)
38         *** Ce calcul dépend du message précédent
39         a(:, :) = transpose(c(1:na,1:na)) + a(:, :)
40         *** Ce calcul est indépendant du message
41         call sgetrf(nb, nb, b, nb, pivota, code)
42     end if
43     temps_fin = (MPI_WTIME () - temps_debut)
44     *** Temps de restitution maximum
45     call MPI_REDUCE (temps_fin, temps_fin_max, 1, MPI_DOUBLE_PRECISION, MPI_MAX, 0, &
46         MPI_COMM_WORLD, code)
47     if (rang == 0) print(,"Temps : ",t6.3, " secondes"), temps_fin_max
48     call MPI_FINALIZE (code)
49     end program Aoptimiser
50
51

```

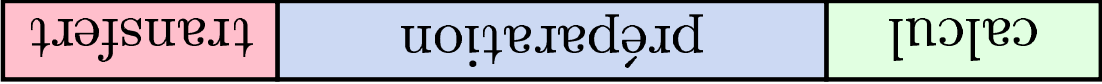
```
> mpiexec -np 2 Optimiser  
Temps : 0.7 secondes
```

Hors communications, le maximum des temps de calcul par processus est de **0,15 secondes**. Ce qui veut dire que les communications prennent environ 78% du temps global !



5.3 – Temps de communication

Que comprend le temps que l'on mesure avec `MPI_WTIME()` ?



↳ Latence : temps d'initialisation des paramètres réseaux.

↳ Surcoût : temps de préparation du message ; caractéristique liée à l'implémentation

MPI et au mode de transfert.

↳ Ce que l'on mesure en général est la somme des deux.

5.4 – Quelques définitions

- ❶ *Recopie temporaire d'un message*. C'est la copie du message dans une mémoire tampon locale (*buffer*) avant son envoi. Cette opération peut parfois être prise en charge par le système *MPI*. C'est le cas par exemple pour les sous-programmes `MPI_SEND()` et `MPI_RECV()` sur *T3E*.

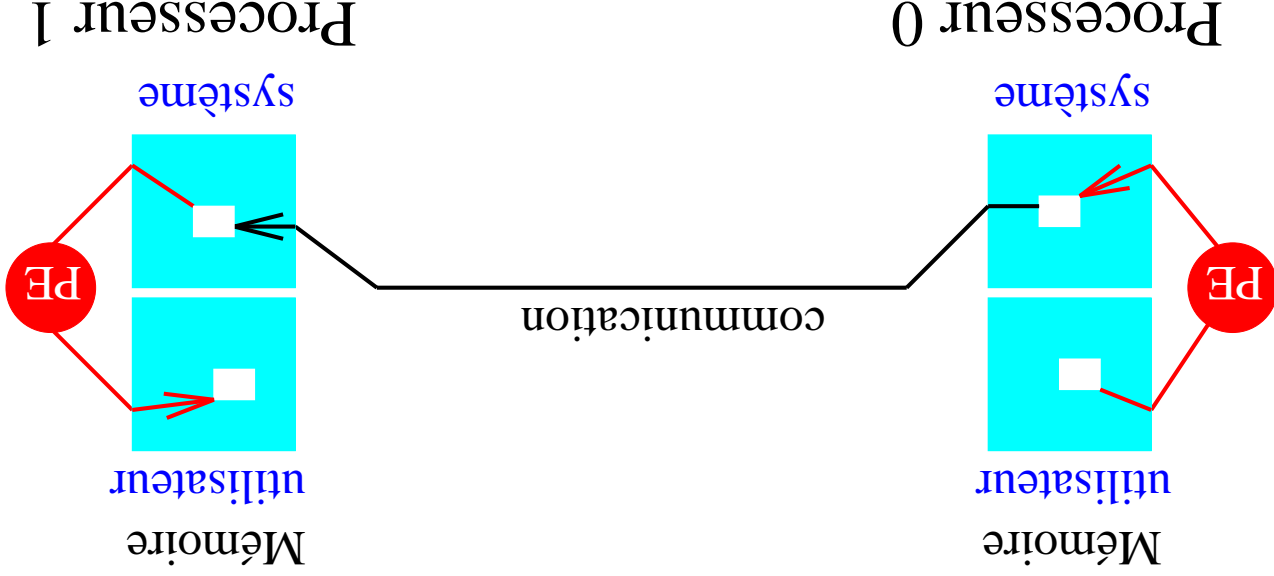


FIG. 19 – *Recopie temporaire d'un message*

② *Envoi non bloquant avec recopie temporaire, non couplé avec la réception.* L'appel à un sous-programme de ce type retourne au programme appelant même quand la réception n'a pas été postée. La recopie temporaire des messages est l'un des moyens d'implémenter un envoi non bloquant afin de découpler l'envoi de la réception.

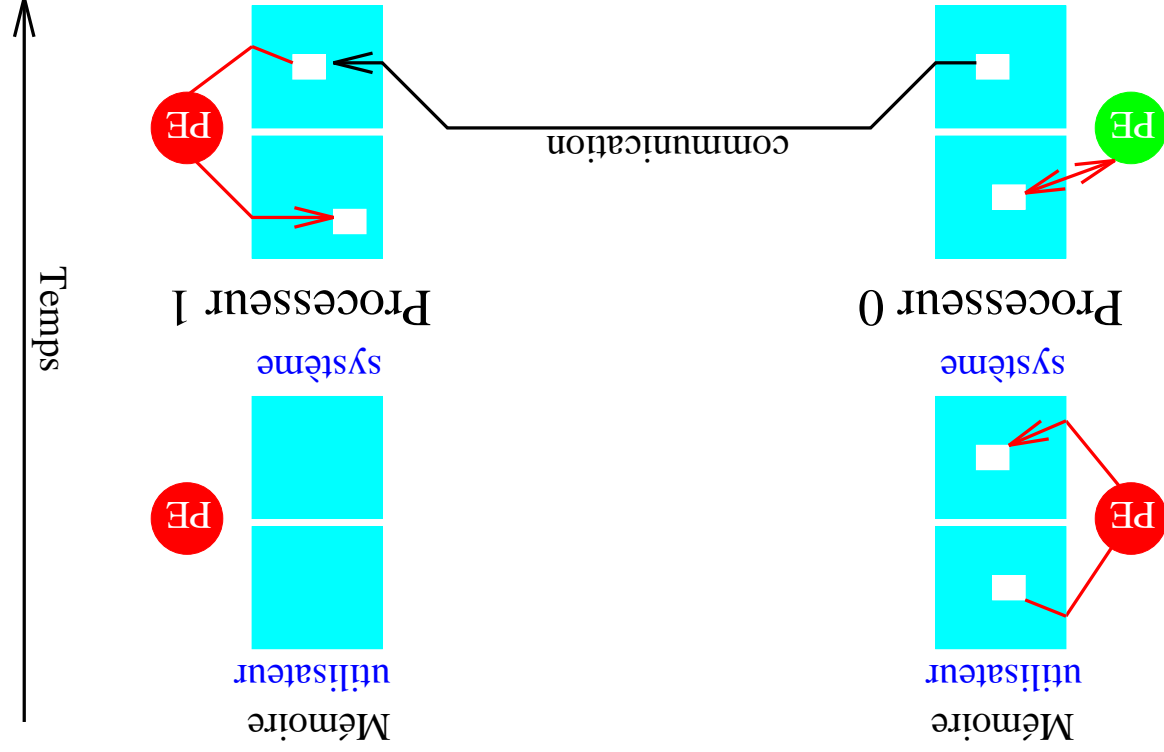


FIG. 20 – *Envoi non bloquant avec recopie temporaire du message*

③ *Envoi bloquant sans copie temporaire, couple avec la réception.* Le message ne quitte le processus émetteur que lorsque le processus récepteur est prêt à le recevoir.

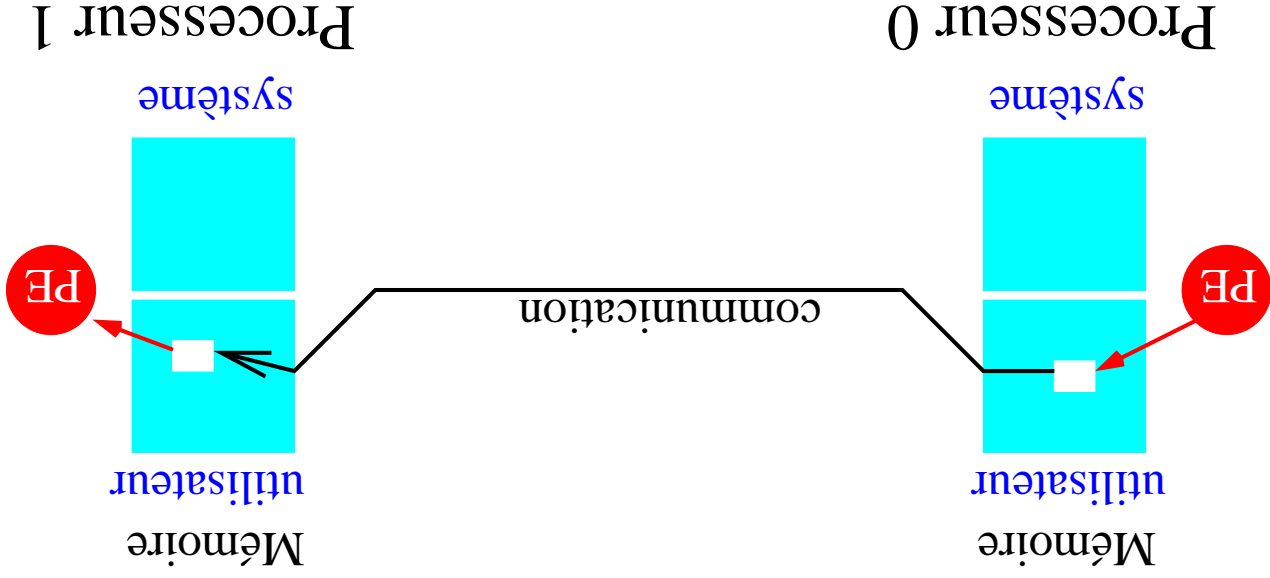


FIG. 21 – *Envoi bloquant couple avec la réception*

④ *Envoi non bloquant sans recopie temporaire, couple avec la réception.* L'appel à ce type de sous-programmes retourne immédiatement au programme appelant bien que l'envoi effectif du message reste couple avec la réception. Il est donc à la charge du programmeur de s'assurer que le message est bien arrivé à sa destination finale avant de pouvoir modifier les données envoyées.

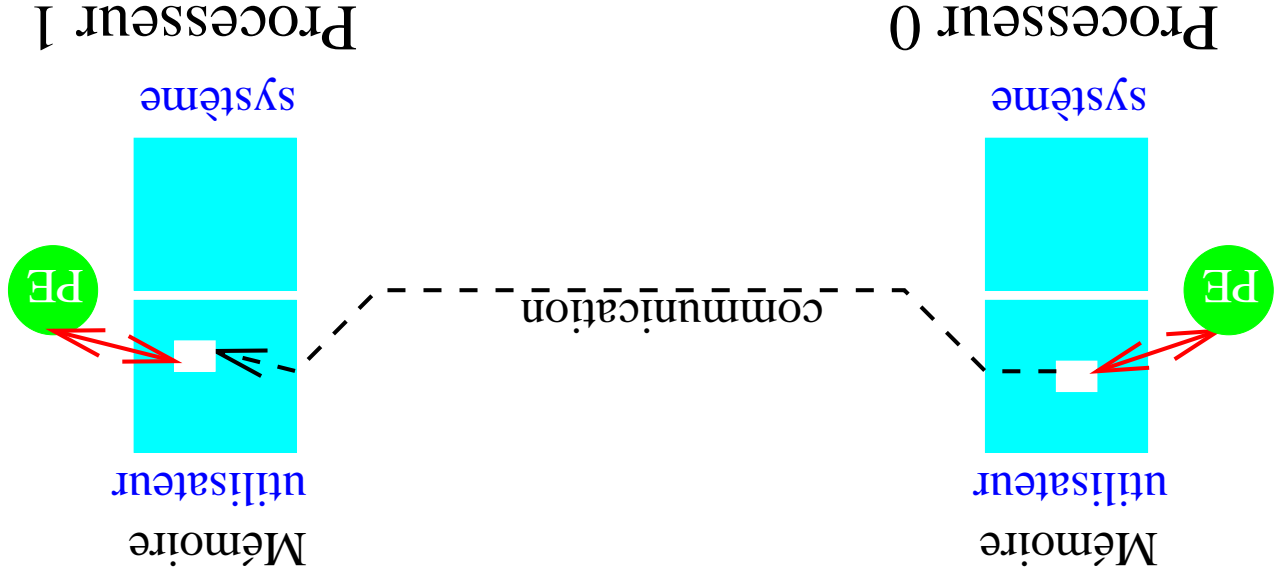


FIG. 22 – *Envoi non bloquant couple avec la réception*

5.5 – Que fournit MPI ?

👉 Avec *MPI* l'envoi d'un message peut se faire suivant différents modes :

- ① ***standard*** : il est à la charge de *MPI* d'effectuer ou non une copie temporaire du message. Si c'est le cas, l'envoi se termine lorsque la copie temporaire est achevée (l'envoi est ainsi découplé de la réception). Dans le cas contraire, l'envoi se termine quand la réception du message est achevée.
- ② ***synchronous*** : l'envoi du message ne se termine que si la réception a été postée et la lecture du message terminée. C'est un envoi couplé avec la réception.
- ③ ***buffered*** : il est à la charge du programmeur d'effectuer une copie temporaire du message. L'envoi du message se termine lorsque la copie temporaire est achevée. L'envoi est ainsi découplé de la réception.
- ④ ***ready*** : l'envoi du message ne peut commencer que si la réception a été postée auparavant (ce mode est intéressant pour les applications clients-serveurs).

À titre indicatif voici les différents cas envisagés par la norme sachant que les implémentations peuvent être différentes :

modes	bloquant	non-bloquant
envoi <i>standard</i>	MPI_SEND () ^a	MPI_ISEND ()
envoi <i>synchronous</i>	MPI_SSEND ()	MPI_ISSEND ()
envoi <i>buffered</i>	MPI_BSEND ()	MPI_IBSEND ()
envoi <i>ready</i>	MPI_RSEND ()	MPI_IRSEND ()
réception	MPI_RECV ()	MPI_IRECV ()

Remarques : pour une implémentation *MPI* donnée, un envoi standard peut-être bloquant avec recopie temporaire ou synchrone avec la réception, selon la taille du message à envoyer.

a. En réalité non bloquant dans certaines implémentations.

5.6 – Envoi synchrone bloquant

Ce mode d'envoi (`MPI_SEND()`) de messages permet d'éviter la recopie temporaire des messages et, par conséquent, les surcoûts que cela peut engendrer. L'efficacité de ce mode de communication sur *T3E*, particulièrement notable pour les applications *bien équilibrées*, provient surtout du fait qu'il est bâti sur des sous-programmes permettant la copie directe de mémoire à mémoire.

Dans le programme modèle, il suffit de remplacer `MPI_SEND()` par `MPI_SEND()` pour gagner un **facteur 2** !



< mpiexec -np 2 Optimiser
0.36 secondes Temps :

```

22 temps_debut = MPI_WTIME ()
23 if (rang == 0) then
24     *** J'envoie un gros message
25     call MPI_SEND (c,m*m,MPI_REAL,1,etiquette,MPI_COMM_WORLD,code)
26     *** Je calcule : factorisation LU avec LAPACK
27     call sgetrf(na, na, a, na, pivota, code)
28     *** Ce calcul modifie le contenu du tableau C
29     c(1:nb,1:nb) = matmul(a(1:nb,1:nb),b)
30     elseif (rang == 1) then
31         *** Je calcule
32         call sgetrf(na, na, a, na, pivota, code)
33         *** Je recois le gros message
34         call MPI_RECV (c,m*m,MPI_REAL,0,etiquette,MPI_COMM_WORLD,statut,code)
35         *** Ce calcul dépend du message précédent
36         a(:,:) = transpose(c(1:na,1:na)) + a(:,:)
37         *** Ce calcul est indépendant du message
38         call sgetrf(nb, nb, b, nb, pivota, code)
39     end if
40     temps_fin = (MPI_WTIME () - temps_debut)
    
```

5.7 – Envoi synchrone non-bloquant

L'utilisation des sous-programmes `MPI_ISSEND()` et `MPI_IRecv()` conjointement aux sous-programmes de synchronisation précédents permet principalement de recouvrir les communications par des calculs.

Le programme modèle, une fois modifié, donne des gains en performances atteignant environ un **facteur 3** par rapport à la version initiale !


```

1 program Optimiser
2 implicit none
3 include 'mpi.h'
4
5 integer, parameter :: na=256,nb=200
6 integer, parameter :: m=2048,etiquette=1111
7 real, dimension(na,na) :: a
8 real, dimension(nb,nb) :: b
9 real, dimension(na) :: pivota
10 real, dimension(nb) :: pivotb
11 real, dimension(m,m) :: c
12 integer :: nb_procs,rang,code,info,requete0,requete1
13 real(kind=8) :: temps_debut,temps_fin,temps_fin_max
14 integer, dimension(MPI_STATUS_SIZE) :: statut
15
16 *** Initialisation MPI
17 call MPI_INIT(code)
18 call MPI_COMM_SIZE(MPI_COMM_WORLD,nb_procs,code)
19 call MPI_COMM_RANK(MPI_COMM_WORLD,rang,code)
20
21 *** Initialisation des tableaux
22 call random_number(a)
23 call random_number(b)
24 call random_number(c)

```

calcul	communication
--------	---------------

< mpirun -np 2 Optimiser
0.23 secondes Temps :

```

25 temps_debut = MPI_WTIME ()
26 if (rang == 0) then
27     *** J'envoie un gros message
28     call MPI_ISSEND (c,m*m,MPI_REAL,1,etiquette,MPI_COMM_WORLD,requete0,code)
29     *** Je calcule : factorisation LU avec LAPACK
30     call sgetrf(na, na, a, na, pivota, code)
31     *** Ce calcul modifie le contenu du tableau C
32     call MPI_WAIT (requete0,statut,code)
33     c(1:nb,1:nb) = matmul(a(1:nb,1:nb),b)
34     elseif (rang == 1) then
35     *** Je calcule
36     call sgetrf(na, na, a, na, pivota, code)
37     *** Je recois le gros message
38     call MPI_IRECV (c,m*m,MPI_REAL,0,etiquette,MPI_COMM_WORLD,requete1,code)
39     *** Ce calcul recouvre la réception précédente
40     call sgetrf(nb, nb, b, nb, pivota, code)
41     *** Ce calcul dépend du message précédent
42     call MPI_WAIT (requete1,statut,code)
43     a(:, :) = transpose(c(1:na,1:na)) + a(:, :)
44     end if
45     temps_fin = (MPI_WTIME () - temps_debut)
    
```

En général, dans le cas d'un envoi (`MPI_XSEND()`) ou d'une réception (`MPI_RECV()`) non bloquant, il existe toute une palette de sous-programmes qui permettent :

- ☞ de synchroniser un processus (ex. `MPI_WAIT()`) jusqu'à terminaison de la requête ;
- ☞ ou de vérifier (ex. `MPI_TEST()`) si une requête est bien terminée ;
- ☞ ou encore de contrôler avant réception (ex. `MPI_PROBE()` ou `MPI_IPROBE()`) si un message particulier est bien arrivé.

5.8 – Conseils 1

→ Éviter si possible la recopie temporaire des messages en utilisant le sous-programme `MPI_SEND()`.

→ Recouvrir les communications par des calculs tout en évitant la recopie temporaire des messages en utilisant les sous-programmes non bloquants `MPI_ISSEND()` et `MPI_IRECV()`.

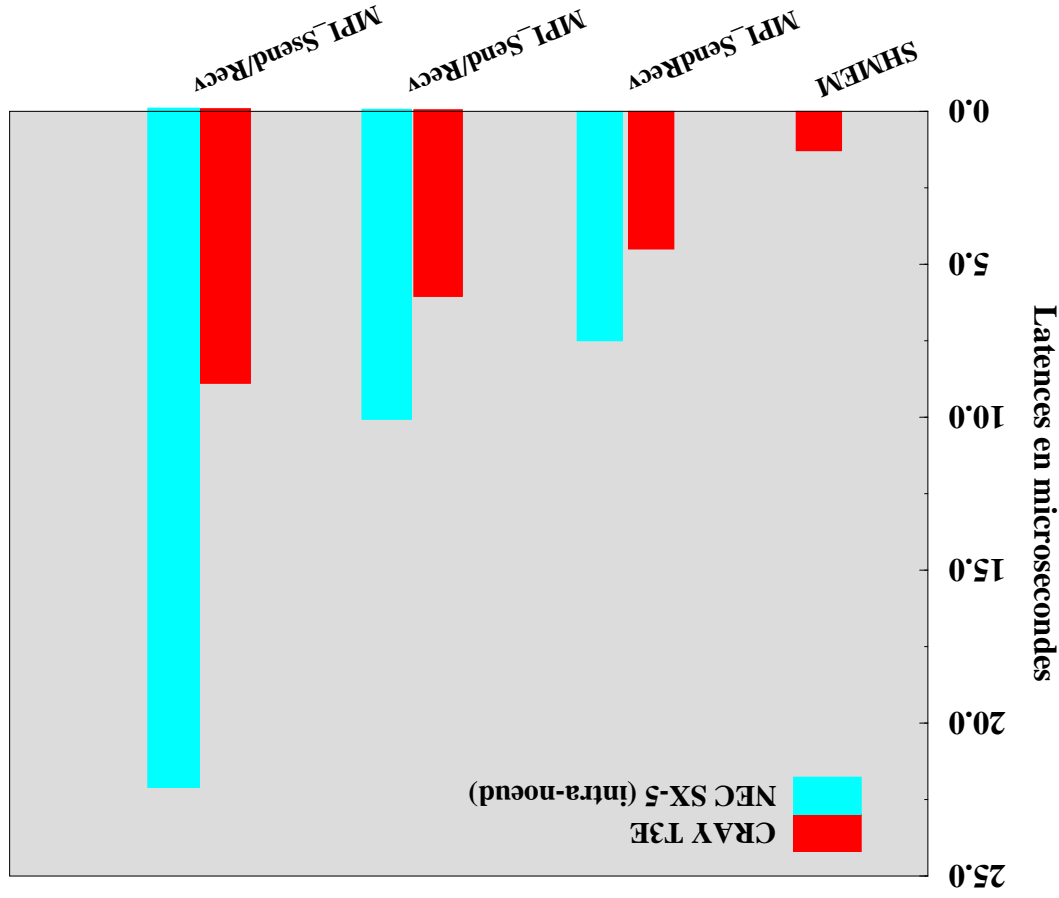
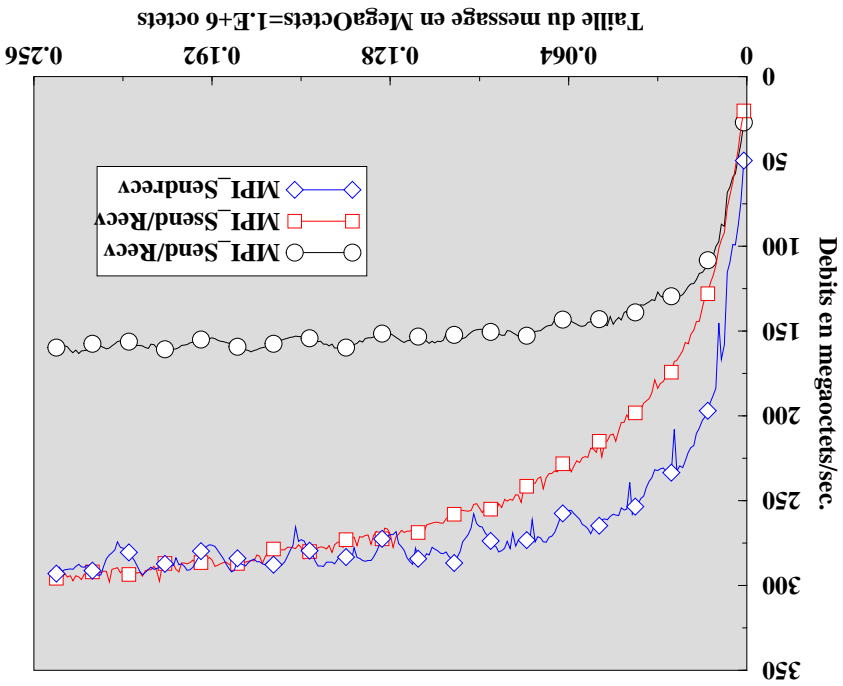


FIG. 23 – Latences des communications sur CRAY T3E-600 et NEC SX-5

FIG. 24 – Débits (messages courts)

(a) CRAY T3E-600



(b) NEC SX-5 (intra-nœud)

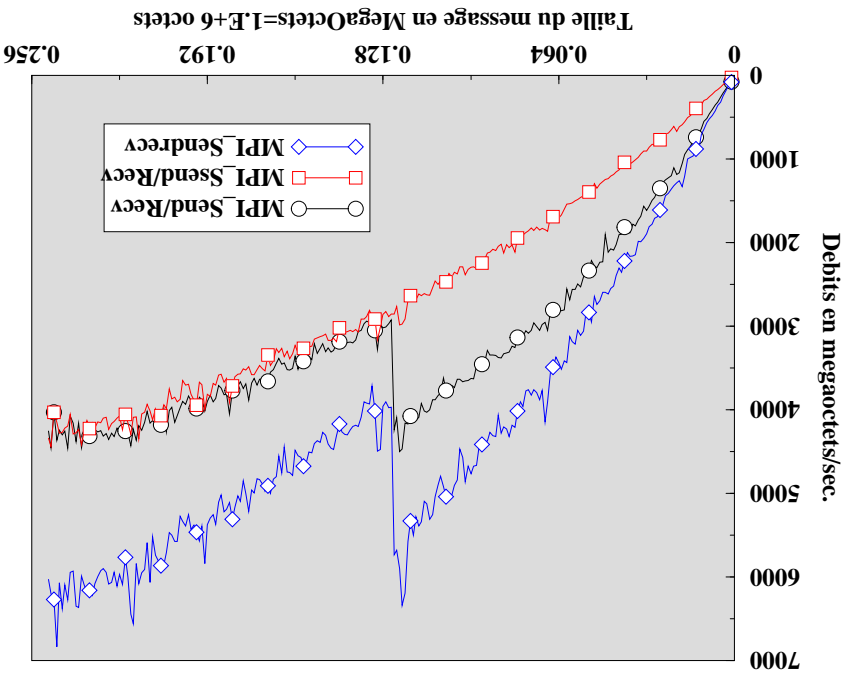
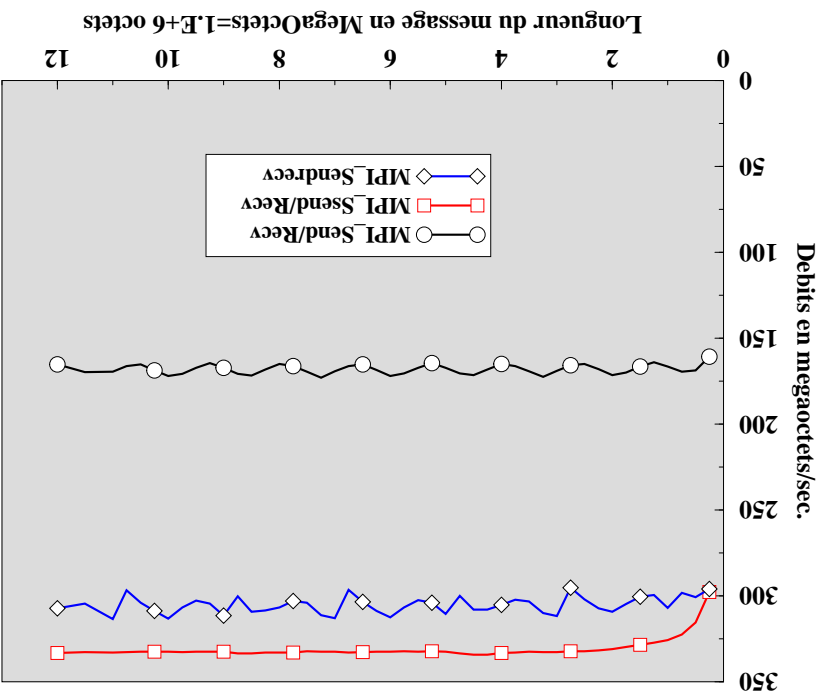
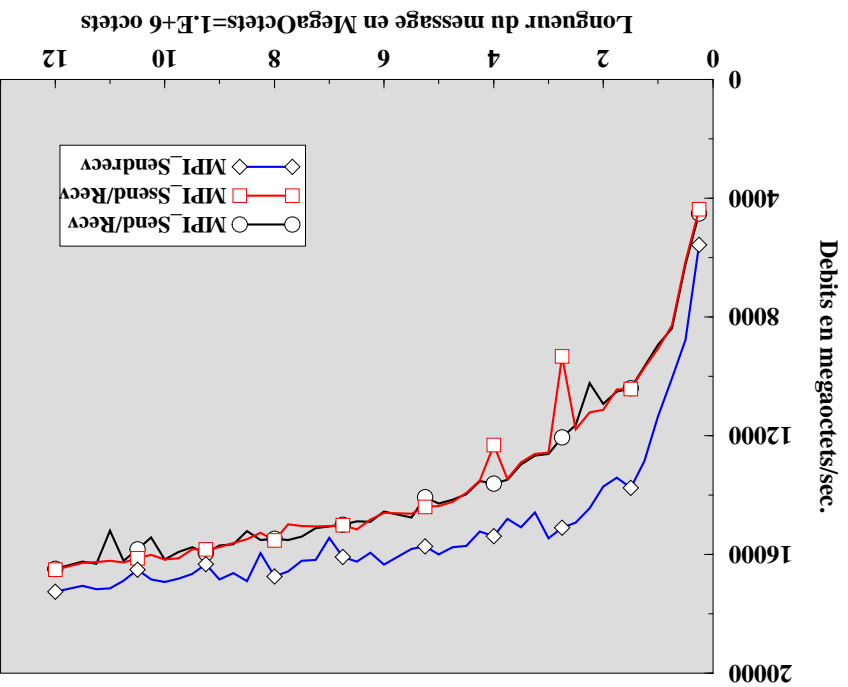


FIG. 25 – Débits (messages longs)

(a) CRAY T3E-600



(b) NEC SX-5 (intra-nœud)



5.9 – Communications persistantes

Dans un programme, il arrive parfois que l'on soit contraint de boucler un certain nombre de fois sur un envoi et une réception de message où la valeur des données manipulées change mais pas leurs adresses en mémoire ni leurs nombres ni leurs types. En outre, l'appel à un sous-programme de communication à chaque itération peut être très pénalisant à la longue d'où l'intérêt des communications persistantes.

Elles consistent à :

- ❶ créer un schéma persistant de communication une fois pour toute (à l'extérieur de la boucle) ;

- ❷ activer réellement la requête d'envoi ou de réception dans la boucle ;
- ❸ libérer, si nécessaire, la requête en fin de boucle.

MPI_SEND_INIT()	envoi <i>standard</i>
MPI_SSEND_INIT()	envoi <i>synchronous</i>
MPI_BSEND_INIT()	envoi <i>buffered</i>
MPI_RECV_INIT()	réception <i>standard</i>

Reprenons le programme modèle...

```

23 if (rang == 0) then
24   do k = 1, 1000
25     *** J'envoie un gros message
26     call MPI_ISSEND(c,m*m,MPI_REAL,1,etiquette,MPI_COMM_WORLD,requete0,code)
27     *** Je calcule : factorisation LU avec LAPACK
28     call sgetrf(na, na, a, na, pivota, code)
29     call MPI_WAIT(requete0,statut,code)
30     *** Ce calcul modifie le contenu du tableau C
31     c(1:nb,1:nb) = matmul(a(1:nb,1:nb),b)
32   end do
33   elseif (rang == 1) then
34     do k = 1, 1000
35       *** Je calcule
36       call sgetrf(na, na, a, na, pivota, code)
37       *** Je reçois le gros message
38       call MPI_IRECV(c,m*m,MPI_REAL,0,etiquette,MPI_COMM_WORLD,requete1,code)
39       *** Ce calcul est indépendant du message
40       call sgetrf(nb, nb, b, nb, pivota, code)
41       call MPI_WAIT(requete1,statut,code)
42       *** Ce calcul dépend du message précédent
43       a(:, :) = transpose(c(1:na,1:na)) + a(:, :)
44     end do
45   end if

```

```
> mpiexec -np 2 Optimiser  
Temps : 235 secondes
```

L'utilisation d'un schéma persistant de communication permet de cacher la latence et de réduire les surcoûts induits par chaque appel aux sous-programmes de communication dans la boucle. Le gain peut être considérable lorsque ce mode de communication est réellement implémenté.

```

23 if (rang == 0) then
24   call MPI_SEND_INIT(c, m*m, MPI_REAL, 1, etiquette, MPI_COMM_WORLD, request0, code)
25   do k = 1, 1000
26     J'envoie un gros message
27     call MPI_START(request0, code)
28     call sgetrf(na, na, a, na, pivota, code)
29     call MPI_WAIT(request0, statut, code)
30     *** Ce calcul modifie le contenu du tableau C
31     c(1:nb, 1:nb) = matmul(a(1:nb, 1:nb), b)
32   end do
33   elseif (rang == 1) then
34     call MPI_RECV_INIT(c, m*m, MPI_REAL, 0, etiquette, MPI_COMM_WORLD, statut, request1, code)
35     do k = 1, 1000
36       call sgetrf(na, na, a, na, pivota, code)
37       *** Je reçois le gros message
38       call MPI_START(request1, code)
39       *** Ce calcul est indépendant du message
40       call sgetrf(nb, nb, b, nb, pivota, code)
41       call MPI_WAIT(request1, statut, code)
42       *** Ce calcul dépend du message précédent
43       a(:, :) = transpose(c(1:na, 1:na)) + a(:, :)
44     end do
45   end if

```

```
> mpiexec -np 2 Optimiser  
Temps : 235 secondes
```

L'infrastructure matérielle du T3E et du NEC SX-5 ne permet pas une implémentation efficace du mode persistant.

Remarques

→ Une communication activée par `MPI_START()` sur une requête créée par l'un des sous-programmes `MPI_XXXX_INIT()` est équivalente à une communication non bloquante `MPI_XXXX()`.

→ Pour redéfinir un nouveau schéma persistant avec la même requête, il faut auparavant libérer celle associée à l'ancien schéma en appelant le sous-programme `MPI_REQUEST_FREE()` (`requete`,code).

→ Ce sous-programme ne libérera la requête `requete` qu'une fois que la communication associée sera réellement terminée.

5.10 – Conseils 2

➡ Minimiser les surcoûts induits par des appels répétitifs aux sous-programmes de communication en utilisant une fois pour toute un schéma persistant de communication et activer celui-ci autant de fois qu'il est nécessaire à l'aide du sous-programme `MPI_START()`.

➡ Recouvrir les communications par des calculs tout en évitant la copie temporaire des messages car un schéma persistant (ex. `MPI_SEND_INIT()`) est forcément active d'une façon **non bloquante** à l'appel du sous-programme `MPI_START()`.

6 – Types de données dérivés

6.1 – Introduction

Dans les communications point à point, les données échangées sont types : `MPI_INTEGER`, `MPI_REAL`, `MPI_COMPLEX`, etc.

On peut créer des structures de données plus complexes à l'aide de sous-programmes tels que `MPI_TYPE_CONTIGUOUS()`, `MPI_TYPE_VECTOR()`, `MPI_TYPE_HVECTOR()`.

À chaque fois que l'on crée un type de données, il faut le valider à l'aide du sous-programme `MPI_TYPE_COMMIT()`.

Si on souhaite réutiliser le même type, on doit le libérer avec le sous-programme `MPI_TYPE_FREE()`.

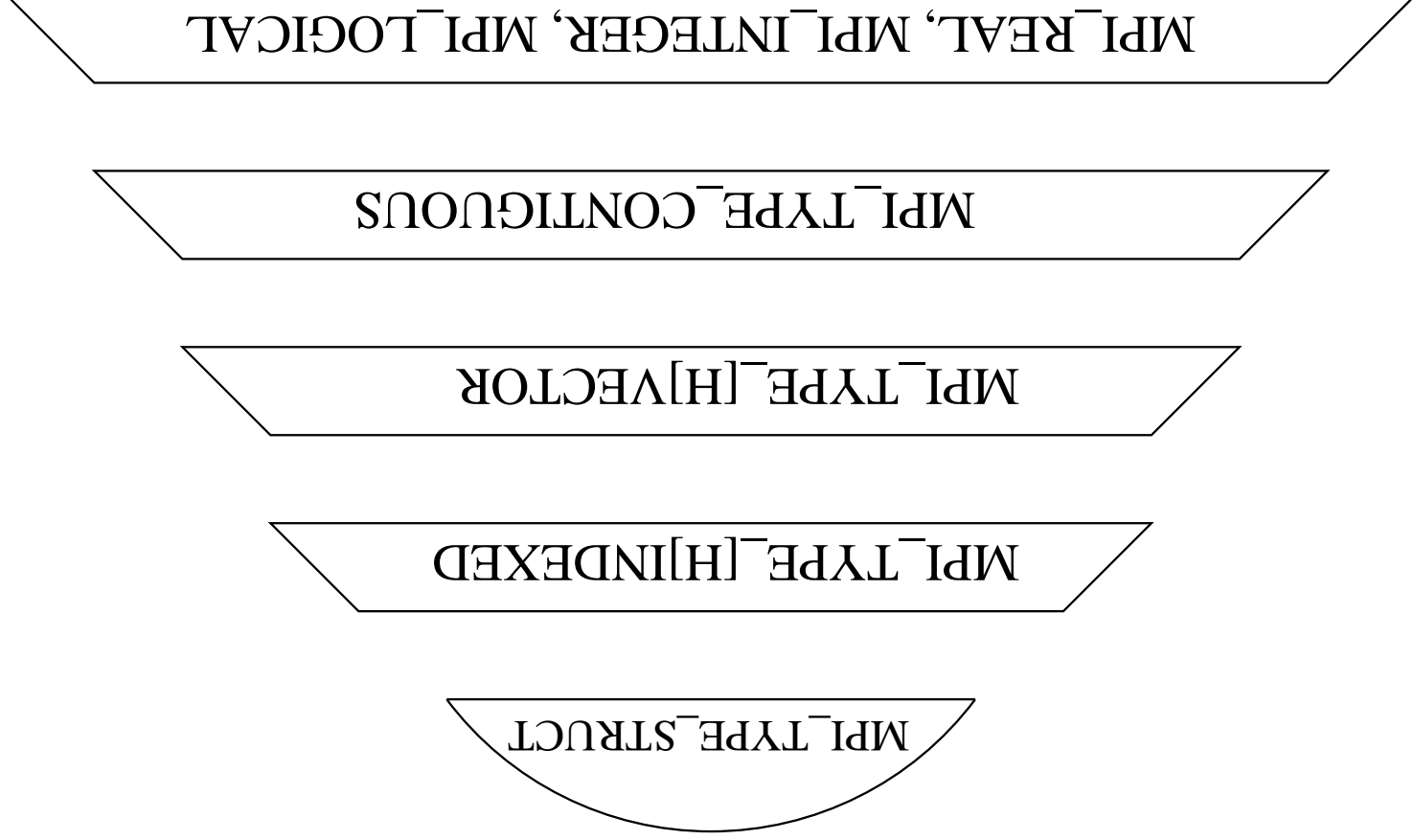


Fig. 26 – Hiérarchie des constructeurs de type MPI

6.2 – Types contigus

→ `MPI_TYPE_CONTIGUOUS()` crée une structure de données à partir d'un ensemble homogène de type prédéfini de données contigus en mémoire.

1	6	11	16	21	26
2	7	12	17	22	27
3	8	13	18	23	28
4	9	14	19	24	29
5	10	15	20	25	30

`call MPI_TYPE_CONTIGUOUS (5, MPI_REAL, nouveau_type, code)`

FIG. 27 – Sous-programme `MPI_TYPE_CONTIGUOUS()`

`integer, intent(in) :: nombre, ancien_type`
`integer, intent(out) :: nouveau_type, code`
`call MPI_TYPE_CONTIGUOUS (nombre, ancien_type, nouveau_type, code)`

6.3 – Types avec un pas constant

→ `MPI_TYPE_VECTOR()` crée une structure de données à partir d'un ensemble homogène de type prédéfini de données distantes d'un pas constant en mémoire.

Le pas est donné en nombre d'éléments.

```

integer, intent(in) :: nombre_bloc, longueur_bloc
integer, intent(in) :: pas ! donne en éléments
integer, intent(in) :: ancien_type
integer, intent(out) :: nouveau_type, code
call MPI_TYPE_VECTOR(nombre_bloc, longueur_bloc, pas, ancien_type, nouveau_type, code)
    
```

FIG. 28 – *Sous-programme MPI_TYPE_VECTOR()*

```

call MPI_TYPE_VECTOR(6, 1, 5, MPI_REAL, nouveau_type, code)
    
```

1	6	11	16	21	26
2	7	12	17	22	27
3	8	13	18	23	28
4	9	14	19	24	29
5	10	15	20	25	30

→ `MPI_TYPE_HVECTOR()` crée une structure de données à partir d'un ensemble

homogène de type prédéfini de données distantes d'un pas constant en

mémoire.

Le pas est donné en nombre d'octets.

→ Cette instruction est utile lorsque le type générique n'est plus un type de base (`MPI_INTEGER`, `MPI_REAL`, ...) mais un type plus complexe construit à l'aide des

sous-programmes *MPI* vus précédemment.

Le pas ne peut plus alors être exprimé en nombre d'éléments du type générique.

6.4 – Descriptif des sous-programmes

```

integer, intent(in) :: nombre_bloc, longueur_bloc
integer, intent(in) :: pas i donne en octets
integer, intent(in) :: ancien_type
integer, intent(out) :: nouveau_type, code
call MPI_TYPE_HVECTOR (nombre_bloc, longueur_bloc, pas, ancien_type, nouveau_type, code)
    
```

```

integer, intent(inout) :: nouveau_type
integer, intent(out) :: code
call MPI_TYPE_COMMIT (nouveau_type, code)
    
```

```

integer, intent(inout) :: nouveau_type
integer, intent(out) :: code
call MPI_TYPE_FREE (nouveau_type, code)
    
```

6.5 – Exemples

```

1 program colonne
2   implicit none
3   include 'mpif.h'
4
5   integer, parameter
6     :: nb_lignes=5, nb_colonnes=6
7     integer, parameter
8     :: dimension(nb_lignes, nb_colonnes)
9     integer, dimension
10    :: (MPI_STATUS_SIZE)
11    statut
12    :: rang, code, type_colonne
13
14    call MPI_INIT(code)
15    call MPI_COMM_RANK(MPI_COMM_WORLD, rang, code)
16
17    i Initialisation de la matrice sur chaque processus
18    a(:, :) = real(rang)
19
20    i Définition du type type_colonne
21    call MPI_TYPE_CONTIGUOUS(nb_lignes, MPI_REAL, type_colonne, code)
22
23    i Validation du type type_colonne
24    call MPI_TYPE_COMMIT(type_colonne, code)

```

```

22  i Envoi de la première colonne
23  if ( rang == 0 ) then
24  call MPI_SEND(a(1,1), 1, type_colonne, 1, etiquette, MPI_COMM_WORLD, code)
25
26  i Réception dans la dernière colonne
27  elseif ( rang == 1 ) then
28  call MPI_RECV(a(1,nb_colonnes), 1, type_colonne, 0, etiquette, &
29  MPI_COMM_WORLD, statut, code)
30  end if
31
32  i Libère le type
33  call MPI_TYPE_FREE(type_colonne, code)
34
35  call MPI_FINALIZE(code)
36
37  end program colonne

```



```

1  program ligne
2  implicit none
3  include 'mpi.h'
4
5  integer, parameter :: nb_lignes=5,nb_colonnes=6
6  integer, parameter :: etiquette=100
7  real, dimension(nb_lignes,nb_colonnes) :: a
8  integer, dimension(MPI_STATUS_SIZE) :: statut
9  integer :: rang,code,type_ligne
10
11  call MPI_INIT(code)
12  call MPI_COMM_RANK(MPI_COMM_WORLD,rang,code)
13
14  i Initialisation de la matrice sur chaque processus
15  a(:,:) = real(rang)
16
17  i Definition du type type_ligne
18  call MPI_TYPE_VECTOR(nb_colonnes,1,nb_lignes,MPI_REAL,type_ligne,code)
19
20  i Validation du type type_ligne
21  call MPI_TYPE_COMMIT(type_ligne,code)
    
```

Le type « ligne » d'une matrice »

```
22 i Envoi de la deuxième ligne
23 if ( rang == 0 ) then
24   call MPI_SEND(a(2,1), 1, type_ligne, 1, etiquette, MPI_COMM_WORLD, code)
25
26   i Réception dans l'avant-dernière ligne
27   elseif ( rang == 1 ) then
28     call MPI_RECV(a(nb_lignes-1,1), 1, type_ligne, 0, etiquette, &
29                 MPI_COMM_WORLD, statut, code)
30   end if
31
32   i Libère le type type_ligne
33   call MPI_TYPE_FREE(type_ligne, code)
34
35   call MPI_FINALIZE(code)
36
37 end program ligne
```

Le type « bloc d'une matrice »

```

1  program bloc
2  implicit none
3  include 'mpif.h'
4
5  integer, parameter
6  :: nb_lignes=5, nb_colonnes=6
7  integer, parameter
8  :: etiquette=100
9  integer, parameter
10 :: nb_lignes_bloc=2, nb_colonnes_bloc=3
11 real, dimension(nb_lignes, nb_colonnes) :: a
12 integer, dimension(MPI_STATUS_SIZE) :: statut
13 integer :: rang, code, type_bloc
14
15 call MPI_INIT(code)
16 call MPI_COMM_RANK(MPI_COMM_WORLD, rang, code)
17
18 i Initialisation de la matrice sur chaque processus
19 a(:, :) = real(rang)
20
21 i Creation du type type_bloc
22 call MPI_TYPE_VECTOR(nb_colonnes_bloc, nb_lignes_bloc, nb_lignes, &
23 MPI_REAL, type_bloc)
24 call MPI_TYPE_COMMIT(type_bloc, code)

```

```

24 i Envoi d'un bloc
25 if ( rang == 0 ) then
26     call MPI_SEND(a(1,1), 1, type_bloc, 1, etiquette, MPI_COMM_WORLD, code)
27
28 i Réception du bloc
29 elseif ( rang == 1 ) then
30     call MPI_RECV(a(nb_lignes-1, nb_colonnes-2), 1, type_bloc, 0, etiquette, &
31                 MPI_COMM_WORLD, statut, code)
32 end if
33
34 i Libération du type type_bloc
35 call MPI_TYPE_FREE(type_bloc, code)
36
37 call MPI_FINALIZE(code)
38
39 end program bloc

```

6.6 – Types homogènes à pas variable

↳ `MPI_TYPE_INDEXED()` permet de créer une structure de données composée d'une

séquence de blocs contenant un nombre variable d'éléments et séparés par un pas variable en mémoire. Ce dernier est exprimé en **éléments**.

↳ `MPI_TYPE_INDEXED()` a la même fonctionnalité que `MPI_TYPE_INDEXED()` sauf que

le pas séparant deux blocs de données est exprimé en **octets**.

Cette instruction est utile lorsque le type générique n'est pas un type de base `MPI` (`MPI_INTEGER`, `MPI_REAL`, ...) mais un type plus complexe construit avec les

sous-programmes `MPI` vs précédemment. On ne peut exprimer alors le pas en nombre d'éléments du type générique d'où le recours à `MPI_TYPE_INDEXED()`.

↳ Attention à la portabilité avec `MPI_TYPE_INDEXED()` !

$nb=3$, $longueurs_blocs=(2,1,3)$, $déplacements=(0,3,7)$

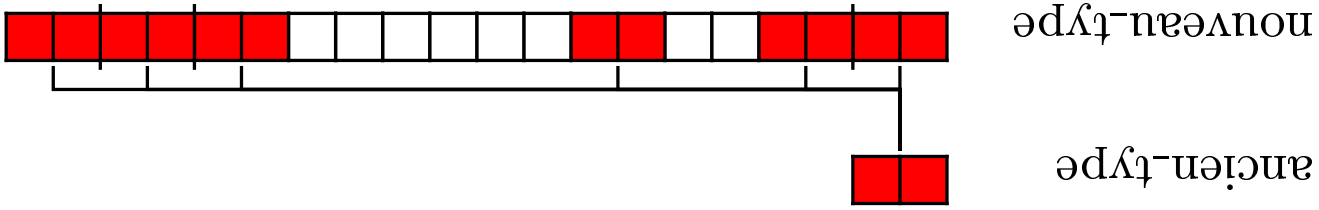


FIG. 29 – Le constructeur *MPI_TYPE_INDEXED*

```
integer, intent(in) :: nb
integer, intent(in), dimension(nb) :: longueurs_blocs
integer, intent(in), dimension(nb) :: déplacements
i Attention les déplacements sont donnés en éléments
integer, intent(in) :: ancien_type
integer, intent(out) :: nouveau_type, code
call MPI_TYPE_INDEXED(nb, longueurs_blocs, déplacements, ancien_type, nouveau_type, code)
```

$nb=4$, $longueurs_blocs=(2,1,2,1)$, $déplacements=(2,10,14,24)$

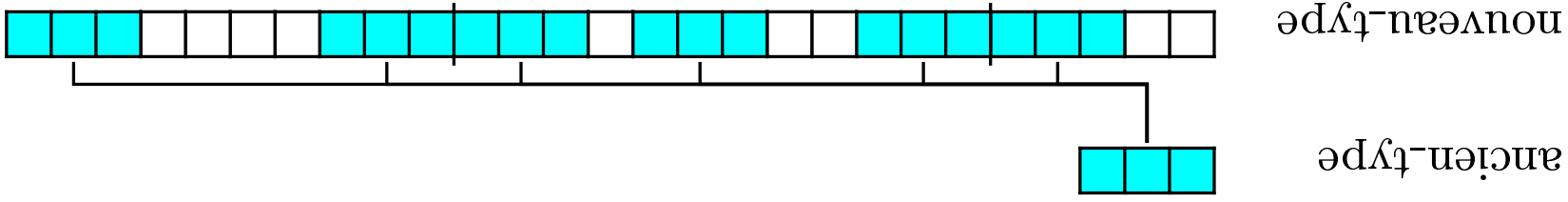


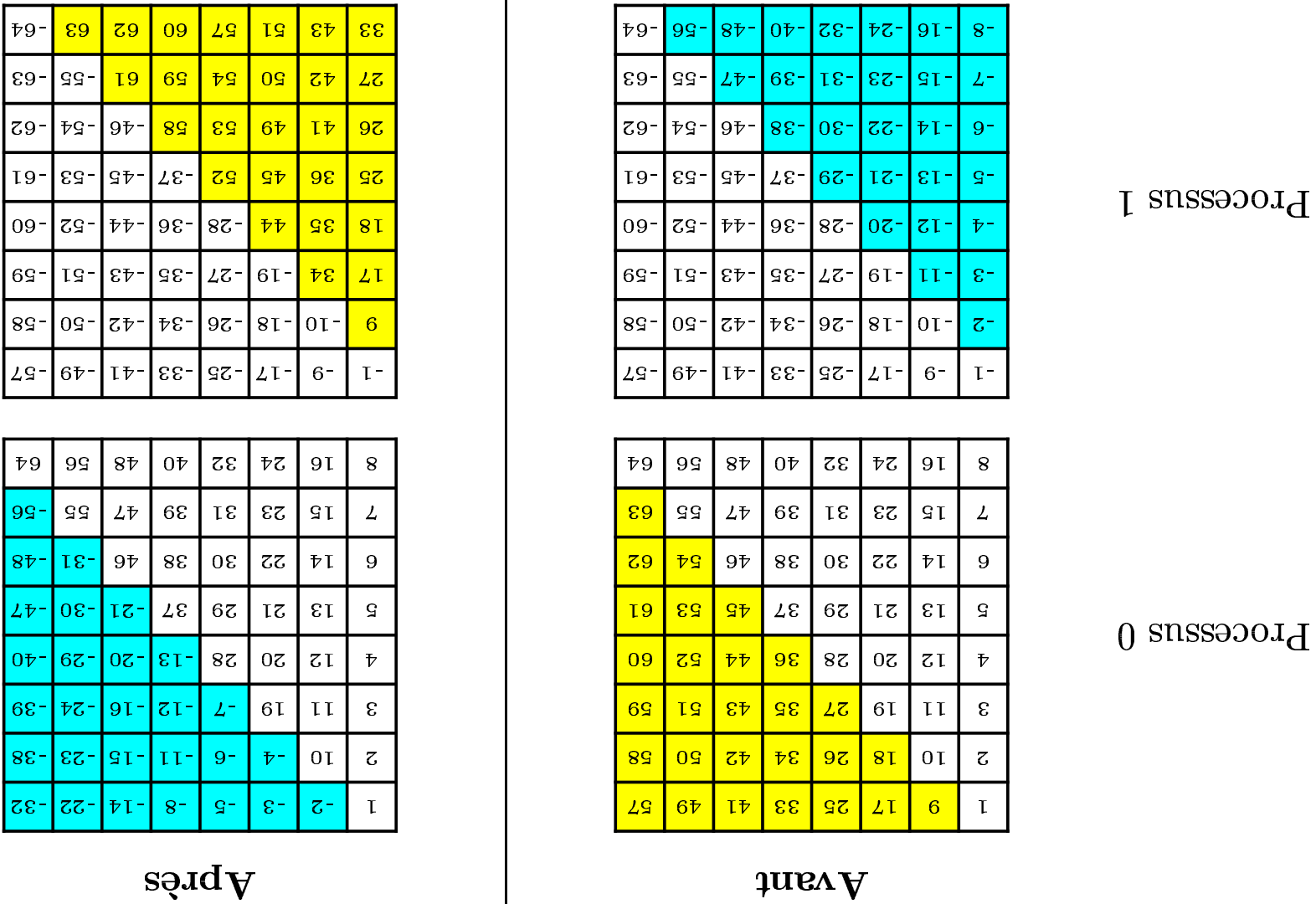
FIG. 30 – *Le constructeur MPI_TYPE_HINDEXED*

```
integer, intent(in) :: nb
integer, intent(in), dimension(nb) :: longueurs_blocs
integer, intent(in), dimension(nb) :: déplacements
i Attention les déplacements sont donnés en octets
integer, intent(in), dimension(nb) :: ancien_type
integer, intent(out) :: nouveau_type, code
call MPI_TYPE_HINDEXED(nb, longueurs_blocs, déplacements, ancien_type, nouveau_type, code)
```

Dans l'exemple suivant, chacun des deux processus :

- ① initialise sa matrice (nombres croissants positifs sur le processus 0 et négatifs décroissants sur le processus 1) ;
- ② construit son type de données (*datatype*) : matrice triangulaire (supérieure pour le processus 0 et inférieure pour le processus 1) ;
- ③ envoie sa matrice triangulaire à l'autre et reçoit une matrice triangulaire qu'il stocke à la place de celle qu'il a envoyée via l'instruction `MPI_SENDRECV_REPLACE()` ;
- ④ libère ses ressources et quitte MPI.

FIG. 31 – Échanges entre les 2 processus



```

1  program triangle
2  implicit none
3  include 'mpif.h'
4
5  integer, parameter :: n=8, etiquette=100
6  real, dimension(n,n) :: a
7  integer, dimension(MPI_STATUS_SIZE) :: statut
8  integer :: i, code
9  integer :: rang, type-triangle
10 integer, dimension(n) :: longueurs_blocs, deplacements
11
12 call MPI_INIT(code)
13 call MPI_COMM_RANK(MPI_COMM_WORLD, rang, code)
14
15 i Initialisation de la matrice sur chaque processus
16 a(:, :) = reshape( / (sign(i, -rang), i=1, n*n) / , (/n, n/))
17
18 i Creation du type matrice triangulaire sup pour le processus 0
19 i et du type matrice triangulaire inferieure pour le processus 1
20 if (rang == 0) then
21     longueurs_blocs(:) = (/ (i-1, i=1, n) /)
22     deplacements(:) = (/ (n*(i-1), i=1, n) /)
23 else
24     longueurs_blocs(:) = (/ (n-i, i=1, n) /)
25     deplacements(:) = (/ (n*(i-1)+i, i=1, n) /)
26 endif
    
```

```

27 call MPI_TYPE_INDEXED(n, longueurs_blocs, déplacements, MPI_REAL, type_triangle, code)
28 call MPI_TYPE_COMMIT(type_triangle, code)
29
30 i Permutation des matrices triangulaires supérieures et inférieures
31 call MPI_SENDRECV_REPLACE(a, 1, type_triangle, mod(rang+1, 2), etiquette, mod(rang+1, 2), &
32 etiquette, MPI_COMM_WORLD, statut, code)
33
34 i Libération du type triangle
35 call MPI_TYPE_FREE(type_triangle, code)
36
37 call MPI_FINALIZE(code)
38
39 end program triangle

```

6.7 – Types hétérogènes

⇒ Le sous-programme `MPI_TYPE_STRUCT()` est le constructeur de types le plus général. Il a les mêmes fonctionnalités que `MPI_TYPE_INDEXED()` mais permet en plus la réplique de blocs de données de types différents.

⇒ Les paramètres de `MPI_TYPE_STRUCT()` sont les mêmes que ceux de `MPI_TYPE_INDEXED()` avec en plus :

⇒ le champ *anciens-types* est maintenant un vecteur de types de données *MPI* ;
⇒ compte tenu de l'hétérogénéité des données et de leur alignement en mémoire, le calcul du déplacement entre deux éléments repose sur la différence de leurs adresses.

MPI, via `MPI_ADDRESS()`, fournit un sous-programme portable qui permet de retourner l'adresse d'une variable.

nb=5, longeurs_blocs=(3,1,5,1,1), déplacements=(0,7,11,21,26),

anciens_types=(type1,type2,type3,type1,type3)

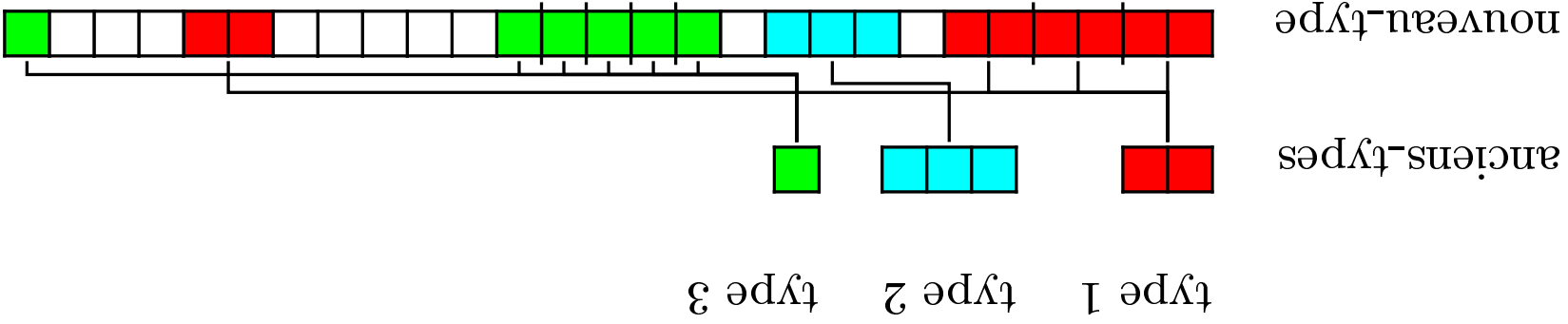


FIG. 32 – *Le constructeur MPI_TYPE_STRUCT*

```
integer, intent(in) :: nb
integer, intent(in), dimension(nb) :: longeurs_blocs
integer, intent(in), dimension(nb) :: déplacements
integer, intent(in), dimension(nb) :: anciens_types
integer, intent(out) :: nouveau_type, code

call MPI_TYPE_STRUCT(nb, longeurs_blocs, déplacements, anciens_types, nouveau_type, code)
```

```
<type>, intent(in) :: variable  
integer, intent(out) :: adresse_variable, code  
call MPI_ADDRESS(variable, adresse_variable, code)
```

6 — Types de données dérivés : hétérogènes

```

1  program Interaction_Particules
2  implicit none
3  include 'mpif.h'
4
5  integer, parameter :: n=100, etiquette=100
6  integer, dimension(MPI_STATUS_SIZE) :: statut
7  integer :: rang, code, type_particule, i
8  integer, dimension(4) :: types, longeurs_blocs, déplacements
9
10 type Particule
11   character(len=5) :: type
12   integer :: masse
13   real, dimension(3) :: coords
14   logical :: classe
15 end type Particule
16 type(Particule), dimension(n) :: p, temp_p
17
18 call MPI_INIT(code)
19 call MPI_COMM_RANK(MPI_COMM_WORLD, rang, code)
20
21 i Construction du type de données
22 types = (/MPI_CHARACTER, MPI_INTEGER, MPI_REAL, MPI_LOGICAL/)
23 longeurs_blocs = (/5,1,3,1/)

```

6 — Types de données dérivés : hétérogènes

```

24  call MPI_ADDRESS(p(1)%type,addresses(1),code)
25  call MPI_ADDRESS(p(1)%masse,addresses(2),code)
26  call MPI_ADDRESS(p(1)%coords,addresses(3),code)
27  call MPI_ADDRESS(p(1)%classe,addresses(4),code)
28
29  i Calcul des déplacements relatifs à l'adresse de départ
30  do i=1,4
31  déplacements(i)=addresses(i) - addresses(1)
32  end do
33  call MPI_TYPE_STRUCT(4,longueurs_blocs,déplacements,types,type-particule,code)
34  i Validation du type structure
35  call MPI_TYPE_COMMIT(type-particule,code)
36  i Initialisation des particules pour chaque processus
37  ...
38  i Envoi des particules de 0 vers 1
39  if (rang == 0) then
40  call MPI_SEND(p(1)%type,n,type-particule,1,etiquette,MPI_COMM_WORLD,code)
41  else
42  call MPI_RECV(temp_p(1)%type,n,type-particule,0,etiquette,MPI_COMM_WORLD,statut,&
43  code)
44  endif
45
46  i Libération du type
47  call MPI_TYPE_FREE(type-particule,code)
48  call MPI_FINALIZE(code)
49  end program Interaction_Particules
    
```


6.8 – Sous-programmes annexes

☞ La taille totale d'un type de données : `MPI_TYPE_SIZE()`.

✓ Attention à ne pas le confondre avec `MPI_TYPE_EXTENT()` qui, lui, renvoie la taille d'un type **aligné** en mémoire.

☞ Les bornes inférieures et supérieures d'un type de données relativement à l'origine : `MPI_TYPE_LB()` et `MPI_TYPE_UB()`

☞ $MPI_TYPE_SIZE() \leq MPI_TYPE_UB() - MPI_TYPE_LB()$

☞ `MPI_LB`, `MPI_UB` sont des *pseudo* types de données (de taille nulle) permettant un alignement correct en mémoire au sein d'une structure créée avec `MPI_TYPE_STRUCT()`

```
integer, intent(in) :: type_donnee
integer, intent(out) :: adresse, code
call MPI_TYPE_LB(type_donnee, adresse, code)
call MPI_TYPE_UB(type_donnee, adresse, code)
```

```
integer, intent(in) :: type_donnee
integer, intent(out) :: taille, code
call MPI_TYPE_SIZE(type_donnee, taille, code)
```

6.9 – Conclusion

Les types dérivés *MPI* sont de puissants mécanismes portables de description de données.

➡ Ils permettent, lorsqu'ils sont associés à des instructions comme `MPI_SENDRECV()`, de simplifier l'écriture de sous-programmes d'échanges interprocessus.

➡ L'association des types dérivés et des topologies (décrites au chapitre suivant) fait de *MPI* l'outil idéal pour tous les problèmes de décomposition de domaines avec des maillages réguliers ou irréguliers.

- ➡ Dans la plupart des applications, plus particulièrement dans les méthodes de décomposition de domaine, il est intéressant de pouvoir disposer les processus suivant une topologie régulière.
- ➡ *MPI* permet de définir des topologies virtuelles du type cartésien ou graphe.

7.1 – Introduction

7 – Topologies


7.2 – Topologies de processus

- ↳ Topologies de type cartésien :
- ↳ chaque processus est défini dans une grille de processus ;
- ↳ la grille peut être périodique ou non ;
- ↳ les processus sont identifiés par leurs coordonnées dans la grille.
- ↳ Topologies de type graphe :
- ↳ généralisation à des topologies plus complexes.

7.3 – Topologies cartésiennes


→ Une topologie cartésienne est définie lorsqu'un ensemble de processus appartenant à un communicateur donne `comm_ancien` appellent le sous-programme `MPI_CART_CREATE()`.

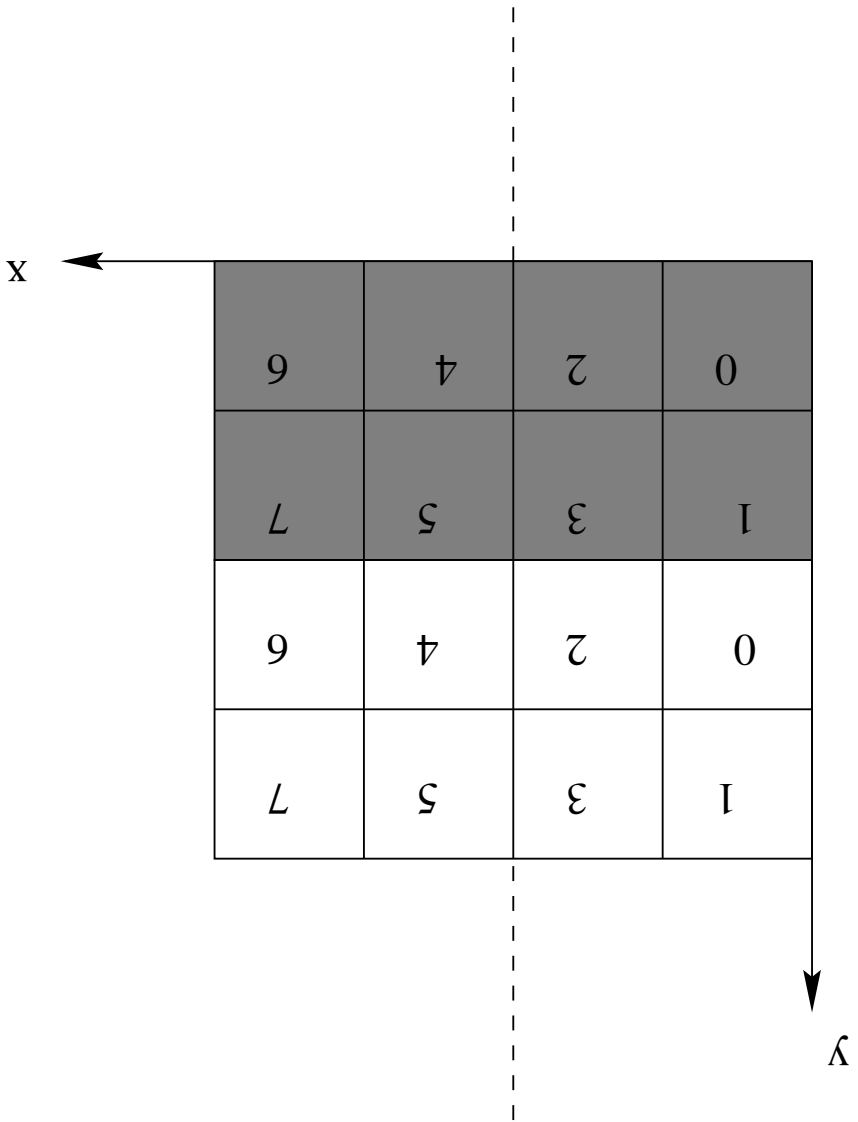
```
integer, intent(in) :: comm_ancien, ndims
integer, dimension(ndims), intent(in) :: dims
logical, dimension(ndims), intent(in) :: periods
logical, intent(in) :: reorganisation
integer, intent(out) :: comm_nouveau, code
call MPI_CART_CREATE(comm_ancien, ndims, dims, periods, reorganisation, comm_nouveau, code)
```

 Exemple sur une grille comportant 4 domaines suivant x et 2 suivant y, périodique en y.

```

include 'mpi.h'
integer, parameter :: ndims = 2
integer, dimension(ndims) :: dims
integer, dimension(ndims) :: logical
logical :: reorganisation
.....
dims(1) = 4
dims(2) = 2
periods(1) = .false.
periods(2) = .true.
reorganisation = .false.
call MPI_CART_CREATE(MPI_COMM_WORLD, ndims, dims, periods, reorganisation, code)
    
```

 Si `reorganisation = .false.` alors le rang des processus dans le nouveau communicateur (`comm_2D`) est le même que dans l'ancien communicateur (`MPI_COMM_WORLD`). Si `reorganisation = .true.`, MPI choisit l'ordre des processus.


 FIG. 33 – Topologie cartésienne 2D périodique en y

➔ Exemple sur une grille 3D comportant 4 domaines suivant x, 2 suivant y et 2 suivant z, non périodique.

```

include 'mpif.h'
integer :: comm_3D, code
integer, parameter :: ndims = 3
integer, dimension(ndims) :: dims
logical, dimension(ndims) :: periods
logical :: reorganisation

.....

dims(1) = 4
dims(2) = 2
dims(3) = 2
periods(:) = .false.
reorganisation = .false.

call MPI_CART_CREATE(MPI_COMM_WORLD, ndims, dims, periods, reorganisation, comm_3D, code)

```

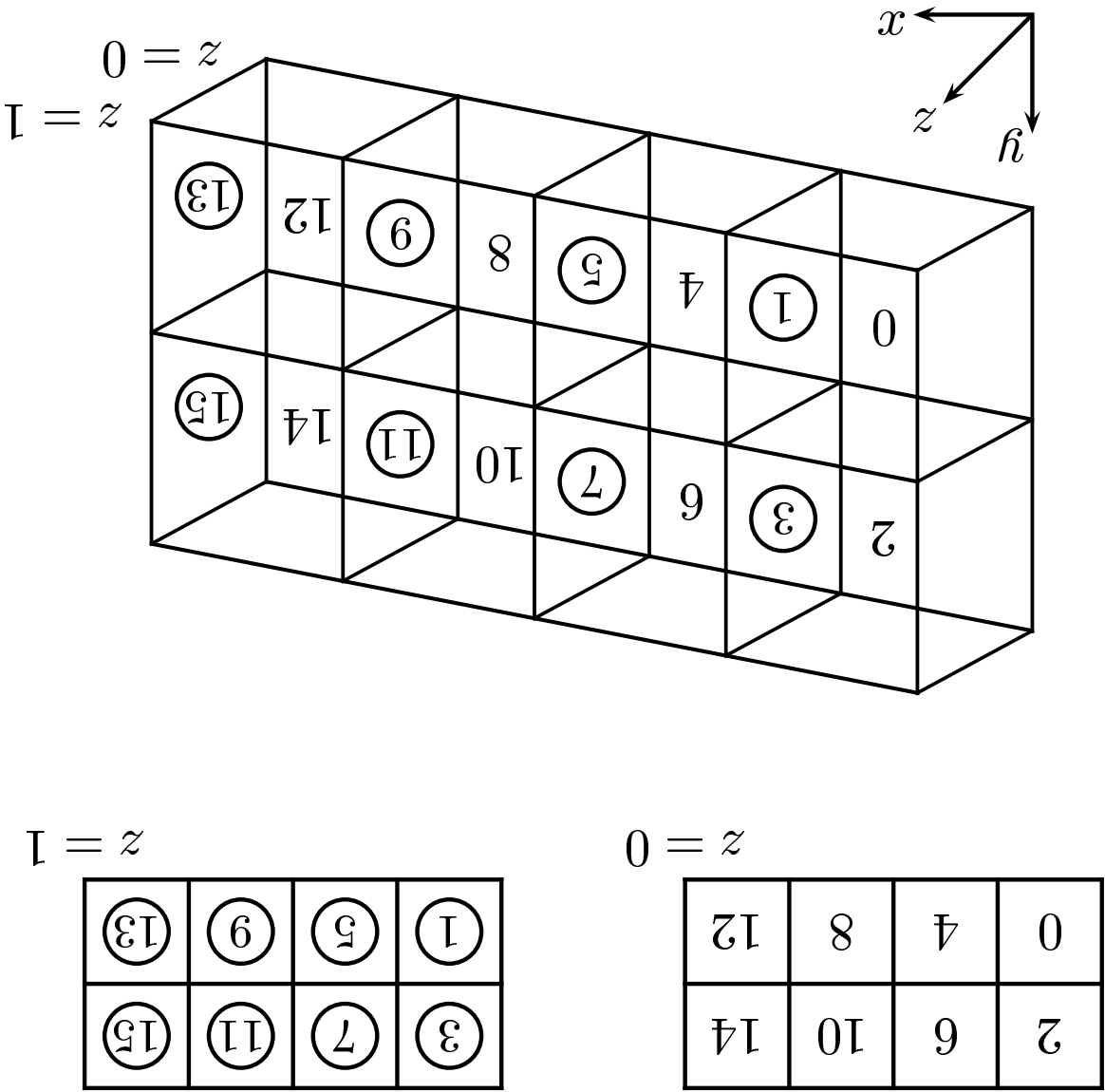


FIG. 34 — Topologie cartésienne 3D non périodique

➡ Dans une topologie cartésienne, le sous-programme `MPI_DIMS_CREATE()` retourne le nombre de processus dans chaque dimension de la grille en fonction du nombre total de processus.

```
integer, intent(in) :: nprocs, ndims
integer, dimension(ndims), intent(inout) :: dims
integer, intent(out) :: code
call MPI_DIMS_CREATE(nprocs, ndims, dims, code)
```

➡ Remarque : si les valeurs de **dims** en entrée valent toutes 0, cela signifie qu'on laisse à MPI le choix du nombre de processus dans chaque direction en fonction du nombre total de processus.

dims en entrée	call MPI_DIMS_CREATE	dims en sortie
(0,0)	(8,2,dims,code)	(4,2)
(0,0,0)	(16,3,dims,code)	(4,2,2)
(0,4,0)	(16,3,dims,code)	(2,4,2)
(0,3,0)	(16,3,dims,code)	error

➔ Dans une topologie cartésienne, le sous-programme `MPI_CART_RANK()` retourne le rang du processus associé aux coordonnées dans la grille.

```
integer, intent(in) :: comm_nouveau
integer, dimension(ndims), intent(in) :: coords
integer, intent(out) :: rang, code
call MPI_CART_RANK(comm_nouveau, coords, rang, code)
```

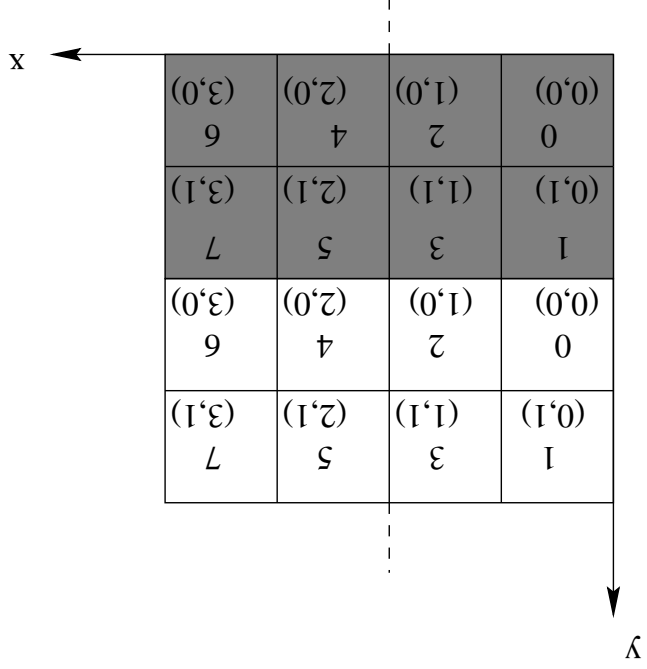


FIG. 35 – Topologie cartésienne 2D périodique en y

```

coords(1)=dims(1)-1
do i=0,dims(2)-1
  coords(2) = i
  call MPI_CART_RANK(comm_2D,coords,rang(i),code)
end do
.....
i=0, en entrée coords=(3,0), en sortie rang(0)=6.
i=1, en entrée coords=(3,1), en sortie rang(1)=7.
    
```

➔ Dans une topologie cartésienne, le sous-programme `MPI_CART_COORDS()` retourne les coordonnées d'un processus de rang donné dans la grille.

```
integer, intent(in) :: comm_nouveau, rang, ndims
integer, dimension(ndims), intent(out) :: coords
integer, intent(out) :: code
call MPI_CART_COORDS(comm_nouveau, rang, ndims, coords, code)
```

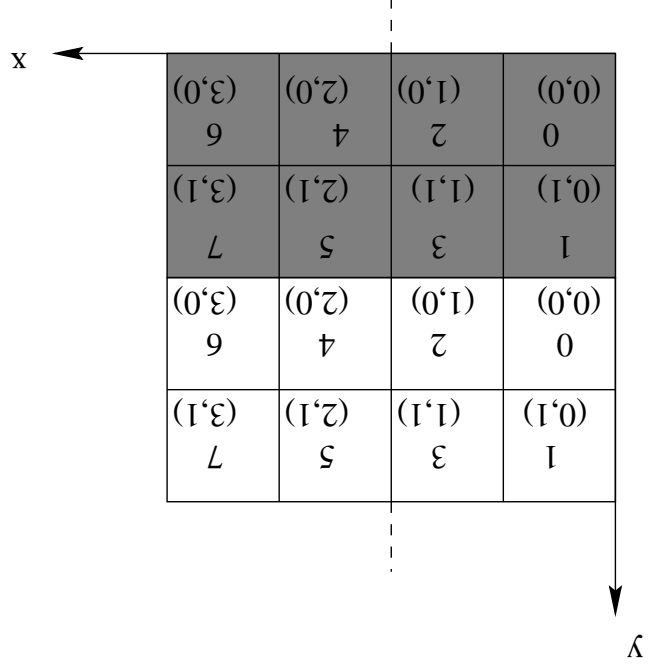


FIG. 36 – *Topologie cartésienne 2D périodique en y*

```

if (mod(rang,2) == 0) then
    call MPI_CART_COORDS(comm_2D,rang,2,coords,code)
end if
.....
En entrée, les valeurs de rang sont : 0,2,4,6.
En sortie, les valeurs de coords sont :
(0,0),(1,0),(2,0),(3,0).
    
```

➡ Dans une topologie cartésienne, un processus appelant le sous-programme

`MPI_CART_SHIFT()` se voit retourner le rang de ses processus voisins dans une direction donnée.

```
integer, intent(in) :: comm_nouveau, direction, pas
integer, intent(out) :: rang_precedent, rang_suisvant
integer, intent(out) :: code
call MPI_CART_SHIFT(comm_nouveau, direction, pas, rang_precedent, rang_suisvant, code)
```

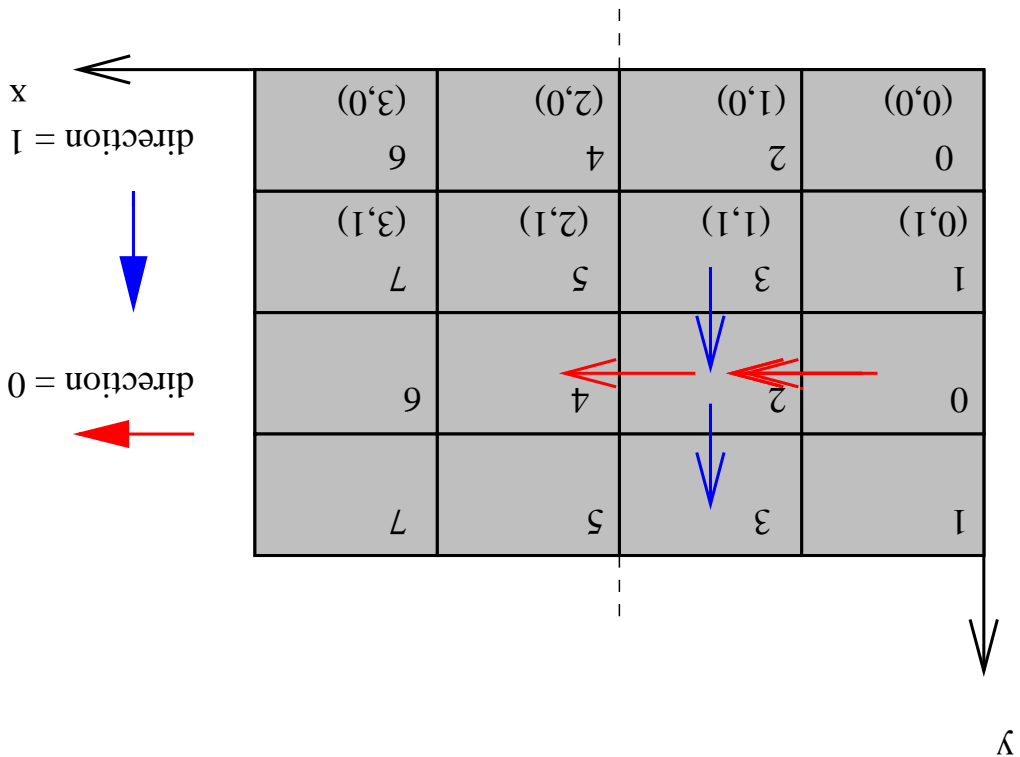
➡ Le paramètre **direction** correspond à l'axe du déplacement (xyz).
➡ Le paramètre **pas** correspond au pas du déplacement.


```

    Pour le processus 2, rang_gauche=0, rang_droit=4
    .....
    call MPI_CART_SHIFT(comm_2D,0,1,rang_gauche,rang_droit,code)

    Pour le processus 2, rang_bas=3, rang_haut=3
    .....
    call MPI_CART_SHIFT(comm_2D,1,1,rang_bas,rang_haut,code)
    
```

FIG. 37 – Appel du sous-programme *MPI_CART_SHIFT()*



```

call MPI_CART_SHIFT(comm_3D, 2, 1, rang_avant, rang_arriere, code)
.....
Pour le processus 0, rang_avant=-1, rang_arriere=1
    
```

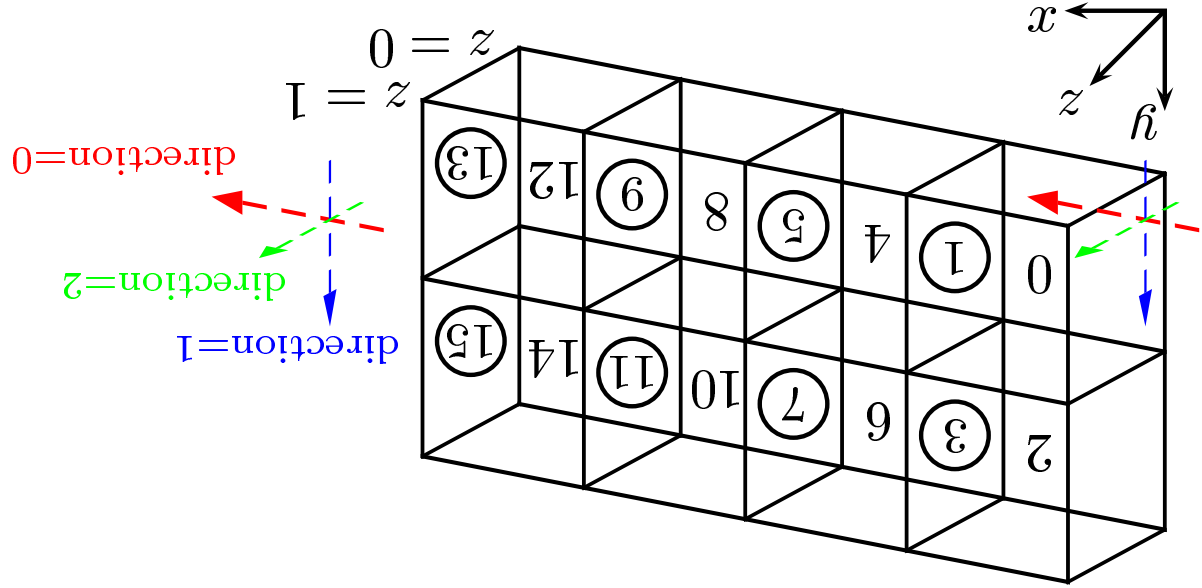
```

call MPI_CART_SHIFT(comm_3D, 1, 1, rang_bas, rang_haut, code)
.....
Pour le processus 0, rang_bas=-1, rang_haut=2
    
```

```

call MPI_CART_SHIFT(comm_3D, 0, 1, rang_gauche, rang_droit, code)
.....
Pour le processus 0, rang_gauche=-1, rang_droit=4
    
```

FIG. 38 – Appel du sous-programme *MPI_CART_SHIFT()*



👉 Exemple de programme :

```

1  program decomposition
2  implicit none
3  include 'mpif.h'
4
5  integer
6  integer, dimension(4)
7  integer, parameter (N=1,E=2,S=3,W=4)
8  integer, parameter (ndims = 2)
9  integer, dimension (ndims)
10 integer, dimension (ndims)
11 logical, dimension (ndims)
12 logical
13
14 call MPI_INIT(code)
15
16 call MPI_COMM_SIZE(MPI_COMM_WORLD, nb_procs, code)
17
18 i Connaitre le nombre de processus suivant x et y
19   dims(:) = 0
20
21 call MPI_DIMS_CREATE(nb_procs, ndims, dims, code)
    
```

```

22  i Creation grille 2D periodique en y
23  periods(1) = .false.
24  periods(2) = .true.
25  reorganisation = .false.
26
27  call MPI_CART_CREATE(MPI_COMM_WORLD, ndims, dims, periods, reorganisation, comm_2D, code)
28
29  i Initialisation du tableau voisin à la valeur MPI_PROC_NULL
30  voisin(:) = MPI_PROC_NULL
31
32  i Recherche de mes voisins Ouest et Est
33  call MPI_CART_SHIFT(comm_2D, 0, 1, voisin(W), voisin(E), code)
34
35  i Recherche de mes voisins Sud et Nord
36  call MPI_CART_SHIFT(comm_2D, 1, 1, voisin(S), voisin(N), code)
37
38  i Connaitre mes coordonnées dans la topologie
39  call MPI_COMM_RANK(comm_2D, rang_ds_topo, code)
40  call MPI_CART_COORDS(comm_2D, rang_ds_topo, ndims, coords, code)
41
42  call MPI_FINALIZE(code)
43
44  end program decomposition

```

```

1 integer, intent(in)      : comm_ancien, nb_procs
2 integer, dimension(nb_procs), intent(in) : index
3 logical, dimension(nb_voisins_max), intent(in) : liste_voisins
4 logical, intent(in)     : reorganisation
5 integer, intent(out)    : comm_nouveau, code
6
7
8 call MPI_GRAPH_CREATE(comm_ancien, nb_procs, index, liste_voisins, reorganisation, &
9                          comm_nouveau, code)
    
```

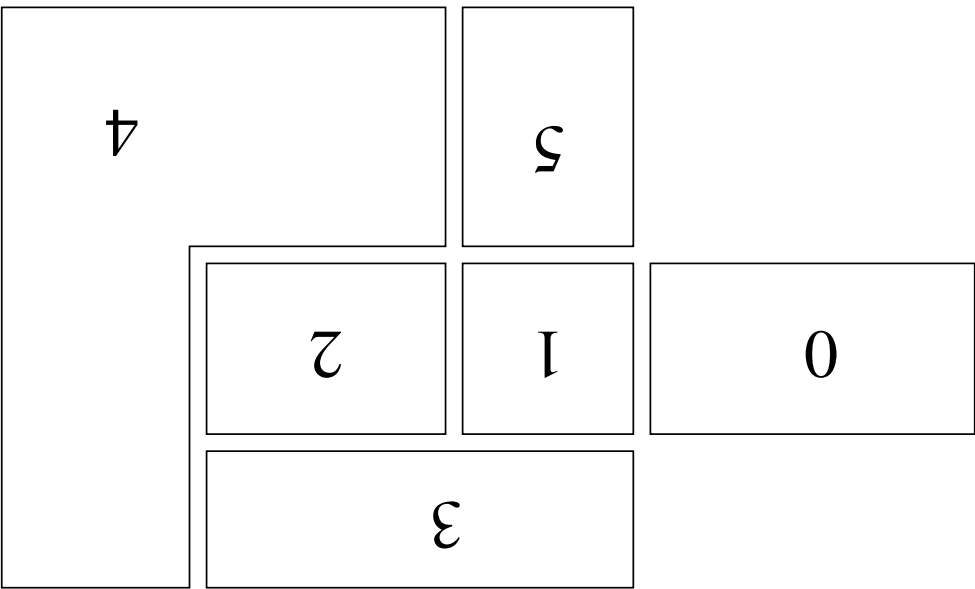
les voisins de chaque sous-domaine.

Il arrive cependant que dans certaines applications (géométries complexes), la décomposition de domaine ne soit plus une grille régulière mais un graphe dans lequel un sous-domaine peut avoir un ou plusieurs voisins quelconques. Le sous-programme `MPI_GRAPH_CREATE()` permet alors de définir une topologie de type graphe en indiquant

7.4 – Graphe de processus

```

index = (/ 1, 5, 8, 11, 14, 16 /)
voisins = (/ 1, 0,5,2,3, 1,3,4, 1,2,4, 3,2,5, 1,4 /)
    
```

 FIG. 39 – *Grappe de processus*


Les tableaux d'entiers index et liste-voisins permettent de définir la liste des voisins pour chacun des nœuds.

Numéro de processus	liste-voisins
0	1
1	0,5,2,3
2	1,3,4
3	1,2,4
4	3,2,5
5	1,4

7 — Topologies : graphe de processus

Deux autres fonctions sont utiles pour connaître :

👉 le nombre de voisins pour un processus donné :

```

integer, intent(in)
logical, intent(out)
integer, intent(in)
integer, intent(out)
:: comm_nouvel
rang ::
nb_voisins ::
code ::
call MPI_GRAPH_NEIGHBORS_COUNT(comm_nouvel, rang, nb_voisins, code)
    
```

👉 la liste des voisins pour un processus donné :

```

integer, intent(in)
integer, intent(in)
integer, intent(in)
integer, intent(in)
logical, intent(out)
integer, intent(out)
:: comm_nouvel
rang ::
nb_voisins ::
nb_voisins ::
code ::
call MPI_GRAPH_NEIGHBORS(comm_nouvel, rang, nb_voisins, voisins, code)
    
```

7 — Topologies : graphe de processus

```

1  program graphe
2
3  implicit none
4  include 'mpi.h'
5
6  integer :: rang,code,comm_graphe,nb_voisins,i,iteration=0
7  integer, parameter :: etiquette=100
8  integer, dimension(6) :: index
9  integer, dimension(16) :: liste_voisins
10 integer, allocatable, dimension(:) :: voisins
11 integer, dimension(MPI_STATUS_SIZE) :: statut
12
13 real :: propagation, & i Propagation du feu
14           & i depuis les voisins
15           feu=0.,
16           & i Valeur du feu
17           & i Rien n'a encore brûlé
18           bois=1.,
19           & i Tout a brûlé si arret <= 0.01
20
21 call MPI_INIT(code)
22
23 i On définit les voisins de chacune des parcelles
24 index = (/1, 5, 8, 11, 14, 16 /)
25 liste_voisins = (/1, 0,5,2,3, 1,3,4, 1,2,4, 3,2,5, 1,4 /)
26
27 call MPI_GRAPH_CREATE(MPI_COMM_WORLD,6,index,liste_voisins,.false.,comm_graphe,code)
28 call MPI_COMM_RANK(comm_graphe,rang,code)

```


7 — Topologies : graphe de processus

```
28 if (rang == 2) feu=1.
    i Le feu se declare arbitrairement sur la parcelle 2
30 call MPI_GRAPH_NEIGHBORS_COUNT(comm_graphe,rang,nb_voisins,code)
31 allocatable(voisins(nb_voisins)) i Allocated au tableau voisins
32 call MPI_GRAPH_NEIGHBORS(comm_graphe,rang,nb_voisins,voisins,code)
33
34 do while (arrêt > 0.01)
35     i On arrête des qu'il n'y a plus rien à brûler
36     do i=1,nb_voisins
37         i On propage le feu aux voisins
38         call MPI_SENDRECV(minval((/1.,feu/)),1,MPI_REAL,voisins(i),etiquette,&
39             propagation,1,MPI_REAL,voisins(i),etiquette,&
40             comm_graphe,statut,code)
41         i Le feu se développe en local sous l'influence des voisins
42         feu=1.2*feu + 0.2*propagation*bois
43         bois=bois/(1.+feu)
44         i On calcule ce qui reste de bois sur la parcelle
45     end do
46     call MPI_ALLREDUCE(bois,arrêt,1,MPI_REAL,MPI_SUM,comm_graphe,code)
47     iteration=iteration+1
48     print,'("Iteration ",i2," parcelle ",i2," bois=",f5.3)',iteration,rang,bois
49     call MPI_BARRIER(comm_graphe,code)
50     if (rang == 0) print,'("--")'
51 end do
52 deallocate(voisins)
53 call MPI_FINALIZE(code)
54 end program graphe
```

7 — Topologies : graphe de processus

```
< mpirun -np 6 graphe
Iteration 1 cellule 0 bois=1.000
Iteration 1 cellule 3 bois=0.602
Iteration 1 cellule 5 bois=0.953
Iteration 1 cellule 4 bois=0.589
Iteration 1 cellule 1 bois=0.672
Iteration 1 cellule 2 bois=0.068
--
Iteration 10 cellule 0 bois=0.008
Iteration 10 cellule 1 bois=0.000
Iteration 10 cellule 3 bois=0.000
Iteration 10 cellule 5 bois=0.000
Iteration 10 cellule 2 bois=0.000
Iteration 10 cellule 4 bois=0.000
--
.....
```

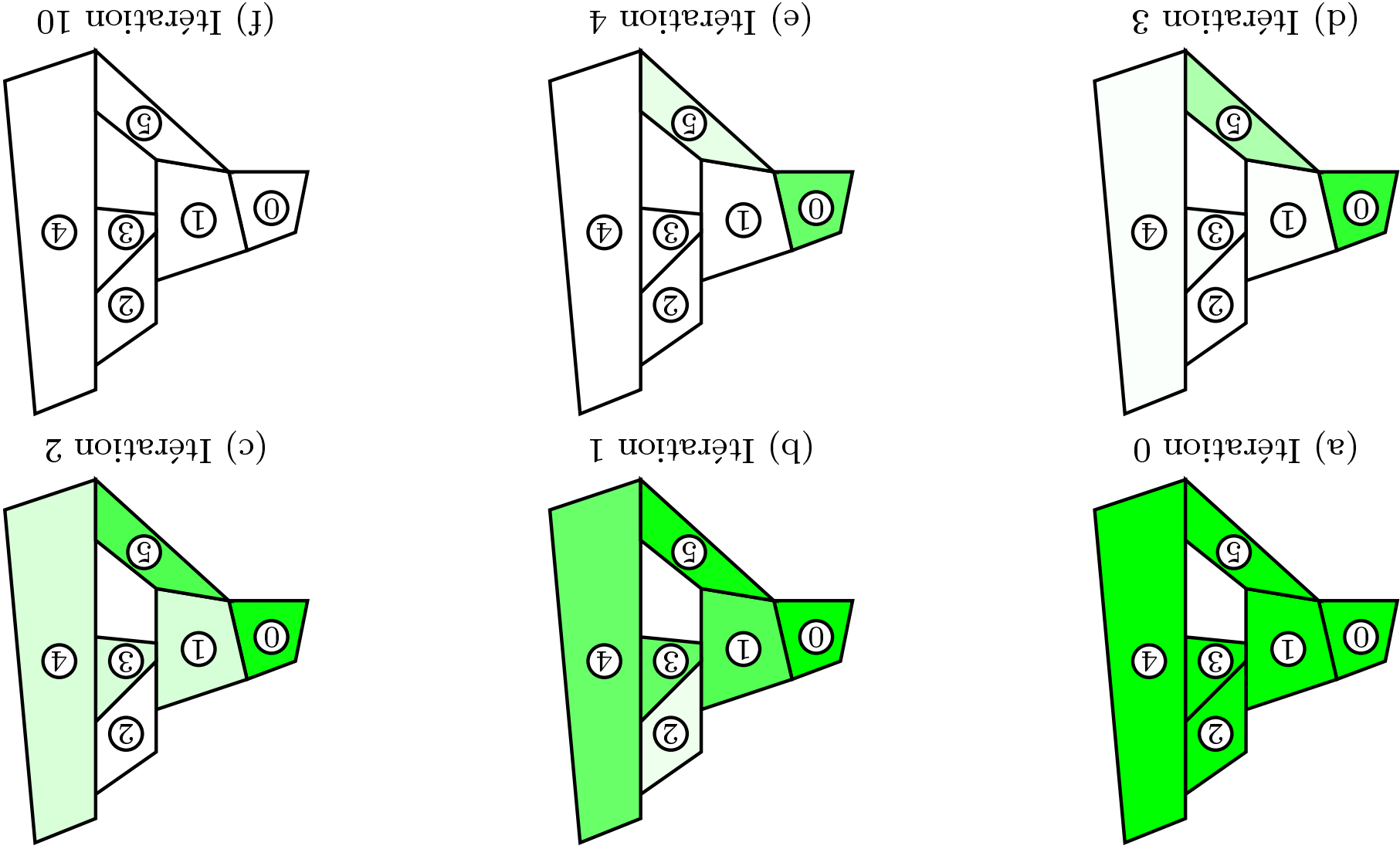


FIG. 40 – Définition d'une topologie quelconque via un graphe — Exemple de la propagation d'un feu de forêt

8 – Communicateurs

8.1 – Introduction

Il s'agit de partitionner un ensemble de processus afin de créer des sous-ensembles sur lesquels on puisse effectuer des opérations telles que des communications point à point, collectives, etc. Chaque sous-ensemble ainsi créé aura son propre espace de communication.

MPI_COMM_WORLD

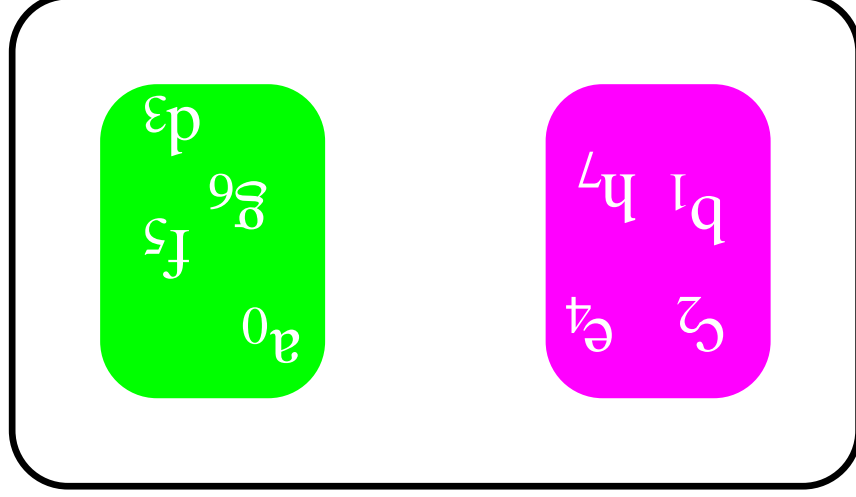


FIG. 41 – Partitionnement d'un communicateur

8.2 – Communicateur par défaut

C'est l'histoire de la poule et de l'œuf...

➡ On ne peut créer un communicateur qu'à partir d'un autre communicateur.

➡ Fort heureusement, cela a été résolu en postulant que la poule existait déjà. En

effet, un communicateur est fourni par défaut, dont l'identificateur `MPI_COMM_WORLD`

est un entier défini dans le fichier d'en-tête `mpi.h`.

➡ Ce communicateur initial `MPI_COMM_WORLD` est créé pour toute la durée d'exécution

du programme à l'appel du sous-programme `MPI_INIT()`.

➡ Ce communicateur ne peut être détruit.

➡ Par défaut, il fixe donc la portée des communications point à point et collectives à tous les processus de l'application.

Dans cet exemple, le processus 2 diffuse un message contenant un vecteur “a” à tous les processus du communicateur `MPI_COMM_WORLD` (donc de l'application) :

```

1  program monde
2  implicit none
3  include 'mpi.h'
4
5  integer, parameter :: m=16
6  integer :: rang,nb-procs,code
7  real, dimension(m) :: a
8
9  call MPI_INIT(code)
10 call MPI_COMM_SIZE(MPI_COMM_WORLD, nb-procs, code)
11 call MPI_COMM_RANK(MPI_COMM_WORLD, rang, code)
12 a(:)=0.
13 if(rang == 2) a(:) = 2.
14
15 call MPI_BCAST(a,m,MPI_REAL,2,MPI_COMM_WORLD,code)
16
17 call MPI_FINALIZE(code)
18
19 end program monde

```

8 — Communicateurs : par défaut

```
< mpirun -np 8 monde
```

MPI_COMM_WORLD

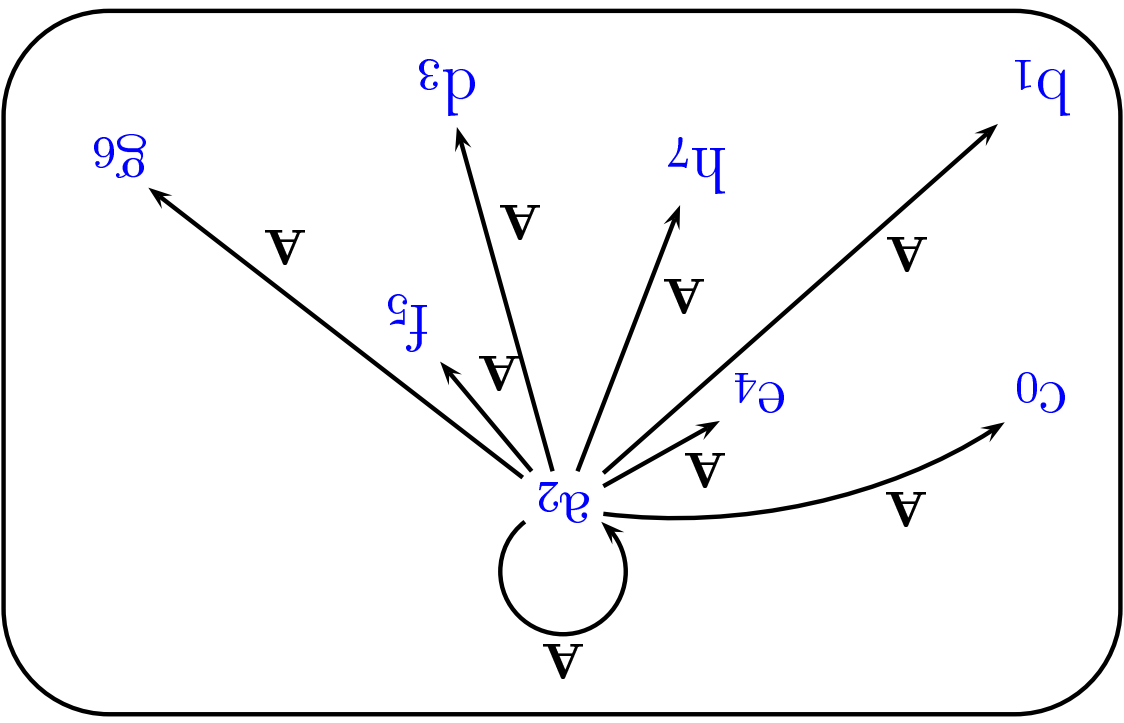


FIG. 42 – Le communicateur par défaut

Que faire pour que le processus 2 diffuse ce message au sous-ensemble de processus de rang pair, par exemple?

- ➡ Boucler sur des *send/recv* peut être très pénalisant surtout si le nombre de processus est élevé. De plus un test serait obligatoire dans la boucle pour savoir si le rang du processus auquel le processus 2 doit envoyer le message est pair ou impair.
- ➡ La solution est de créer un communicateur regroupant ces processus de sorte que le processus 2 diffuse le message à eux seuls.

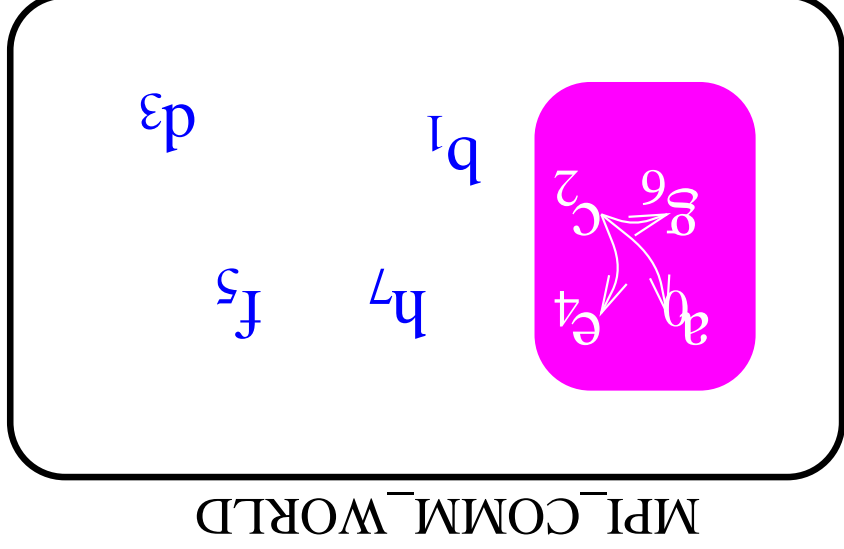


FIG. 43 – Un nouveau communicateur

8.3 – Groupes et communicateurs

Dans la bibliothèque *MPI*, il existe des sous-programmes pour :

① construire des groupes de processus ;

→ `MPI_GROUP_INCL()` ;

→ `MPI_GROUP_EXCL()` .

② construire des communicateurs :

→ `MPI_CART_CREATE()` ;

→ `MPI_CART_SUB()` ;

→ `MPI_COMM_CREATE()` ;

→ `MPI_COMM_DUP()` ;

→ `MPI_COMM_SPLIT()` .

- ➡ Les **constructeurs de groupes locaux** aux processus du groupe (qui n'engendrent pas de communications entre les processus).
- ➡ Les **constructeurs de communicateurs** sont des **opérateurs collectifs** (qui engendrent des communications entre les processus).
- ➡ Les groupes et les communicateurs que le programmeur crée peuvent être gérés dynamiquement. De même qu'il est possible de les créer, il est possible de les détruire : `MPI_GROUP_FREE()`, `MPI_COMM_FREE()`.

Un communicateur est constitué :

① d'un **groupe** de processus ;

② d'un **contexte** de communication mis en place à l'appel du sous-programme de construction du communicateur.

Sachant que :

👉 le **groupe** constitue un ensemble de processus ;

👉 le **contexte** de communication permet de délimiter l'espace de communication ;

👉 les contextes de communication sont gérés par *MPI* (le programmeur n'a aucune action sur eux : c'est un attribut « caché »).

En pratique, pour construire un communicateur, il existe deux façons de procéder :

① par l'intermédiaire d'un groupe de processus ;

② directement à partir d'un autre communicateur.

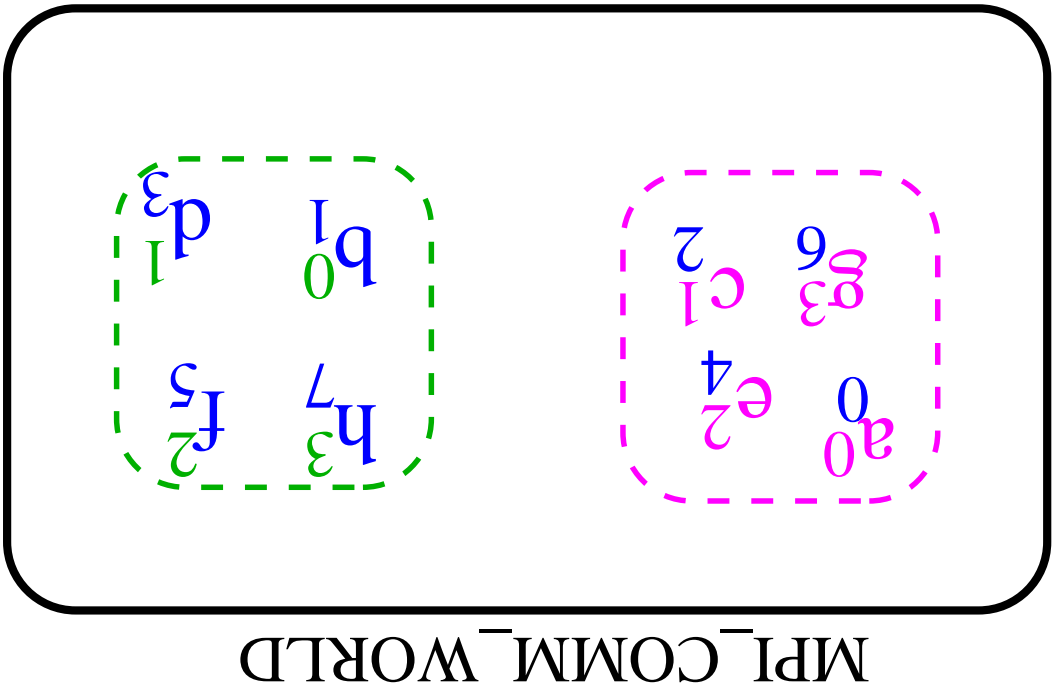
8.4 – Communicateur issu d'un groupe

- ➡ Un groupe est un ensemble ordonné de N processus.
- ➡ Chaque processus du groupe est identifié par un entier $0, 1, \dots, N - 1$ appelé **rang**.
- ➡ Initialement, tous les processus sont membres d'un groupe associé au communicateur par défaut `MPI_COMM_WORLD`.
- ➡ Des sous-programmes `MPI` (`MPI_GROUP_SIZE()`, `MPI_GROUP_RANK()`, etc.) permettent de connaître les attributs d'un groupe.
- ➡ Tout communicateur est associé à un groupe. Le sous-programme `MPI_COMM_GROUP()` permet de connaître le groupe auquel un communicateur est associé.

Dans l'exemple qui suit, nous allons :

- ☞ regrouper les processus de rang pair dans un communicateur (comm-pair) et les processus de rang impair dans un autre (comm-impair) ;
- ☞ ne diffuser un message qu'aux processus de chacun de ces deux groupes.

FIG. 44 – Deux groupes de processus dans un communicateur



```

1 program GroupePairImpair
2 include 'mpi.h'
3 integer, parameter :: m=16
4 integer :: i=0, dim_rangs_pair, rang, nb_procs, iproc, code, &
5 grp_monde, grp_pair, grp_impair
6 integer :: comm_pair, comm_impair
7 integer, dimension(:), allocatable :: rangs_pair
8 real, dimension(m) :: a
9 integer :: rang_ds_monde, rang_ds_pair, rang_ds_impair
10
11 call MPI_INIT(code)
12 call MPI_COMM_SIZE(MPI_COMM_WORLD, nb_procs, code)
13 call MPI_COMM_RANK(MPI_COMM_WORLD, rang, code)
14
15 *** Initialisation du vecteur A
16 a(:)=0.
17 if(rang == 2) a(:)=2.
18 if(rang == 3) a(:)=3.
19
20 *** Enregistrer le rang des processus pairs
21 dim_rangs_pair = int((nb_procs+1)/2)
22 allocate(rangs_pair(dim_rangs_pair))
23 do iproc = 0, nb_procs-1, 2
24   i = i + 1
25   rangs_pair(i) = iproc
26 end do

```

```

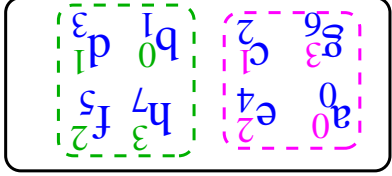
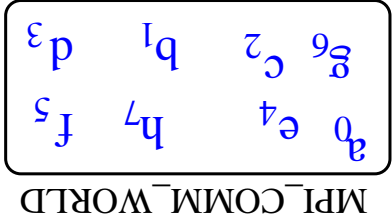
27 *** Connaître le groupe associé au communicateur MPI_COMM_WORLD
28 call MPI_COMM_GROUP(MPI_COMM_WORLD, grp_monde, code)
29
30 *** Créer le groupe des processus pairs
31 call MPI_GROUP_INCL(grp_monde, dim_rangs_pair, rangs_pair, grp_pair, code)
32
33 *** Créer le communicateur des processus pairs
34 call MPI_COMM_CREATE(MPI_COMM_WORLD, grp_pair, comm_pair, code)
35
36 *** Créer le groupe des processus impairs
37 call MPI_GROUP_EXCL(grp_monde, dim_rangs_pair, rangs_pair, grp_impair, code)
38
39 *** Créer le communicateur des processus impairs
40 call MPI_COMM_CREATE(MPI_COMM_WORLD, grp_impair, comm_impair, code)

```

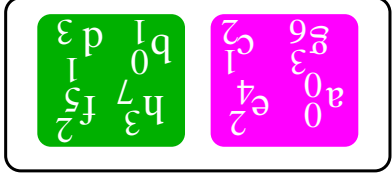
```

41 if (mod(rang,2) == 0) then
42     *** Trouver le rang du processus 2 dans le groupe pair
43     rang_ds_monde=2
44     call MPI_GROUP_TRANSLATE_RANKS (grp_monde, 1, rang_ds_monde, grp_pair, rang_ds_pair, code)
45     *** Diffuser le message seulement aux processus de rangs pairs
46     call MPI_BCAST (a,m,MPI_REAL, rang_ds_pair, comm_pair, code)
47     *** Destruction du communicateur comm_pair
48     call MPI_COMM_FREE (comm_pair, code)
49 else
50     *** Trouver le rang du processus 3 dans le groupe impair
51     rang_ds_monde=3
52     call MPI_GROUP_TRANSLATE_RANKS (grp_monde, 1, rang_ds_monde, grp_impair, rang_ds_impair, &
53     *** Diffuser le message seulement aux processus de rangs impairs
54     call MPI_BCAST (a,m,MPI_REAL, rang_ds_impair, comm_impair, code)
55     *** Destruction du communicateur comm_impair
56     call MPI_COMM_FREE (comm_impair, code)
57 end if
58 call MPI_FINALIZE (code)
59 end program GroupePairImpair

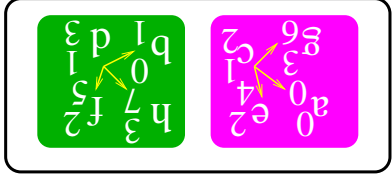
```

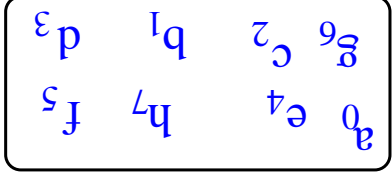
call MPI_GROUP_INCL (...)
call MPI_GROUP_EXCL (...)



call MPI_COMM_CREATE(...)



call MPI_BCAST(...)



call MPI_COMM_FREE(...)

```
$ mpirun -np 8 GroupePairImpair
```

FIG. 45 – *Création/destruction des groupes et des communicateurs pair et impair*

Une fois les groupes constitués, il est possible :

→ de les comparer :

`MPI_GROUP_COMPARE`(group1,group2,resultat).

→ d'appliquer des opérateurs ensemblistes sur deux groupes :

`MPI_GROUP_UNION`(group1,group2,nouveau-groupe),

`MPI_GROUP_INTERSECTION`(group1,group2,nouveau-groupe),

`MPI_GROUP_DIFFERENCE`(group1,group2,nouveau-groupe)

où nouveau-groupe peut être l'ensemble vide, auquel cas il prend la valeur

`MPI_GROUP_EMPTY`.

8.5 – Communicateur issu d'un autre

Le programme précédent présente cependant quelques inconvénients. Nous allons le réécrire pour :

→ ne pas nommer différemment ces deux communicateurs ;

→ ne pas passer par les groupes pour construire les communicateurs `comm_pair` et `comm_impair` ;

→ ne pas laisser le soin à *MPI* d'ordonner le rang des processus dans les communicateurs `comm_pair` et `comm_impair` ;

→ éviter les tests conditionnels, en particulier lors de l'appel au sous-programme `MPI_BCAST()`.

Le sous-programme `MPI_COMM_SPLIT()` permet de partitionner un communicateur donné en autant de communicateurs que l'on veut...

```
integer, intent(in) :: comm, couleur, c1ef
integer, intent(out) :: nouveau_comm, code
call MPI_COMM_SPLIT(comm, couleur, c1ef, nouveau_comm, code)
```

rang	0	1	2	3	4	5	6	7
processus	a	b	c	d	e	f	g	h
couleur	0	2	3	0	3	0	2	3
clef	2	15	0	0	1	3	11	1

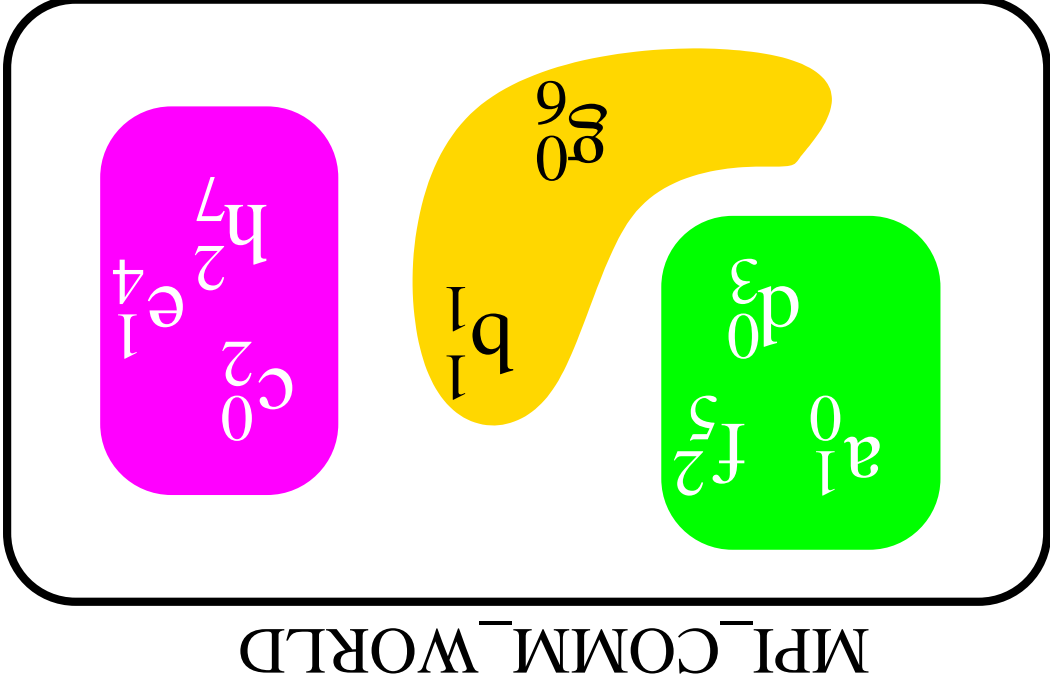


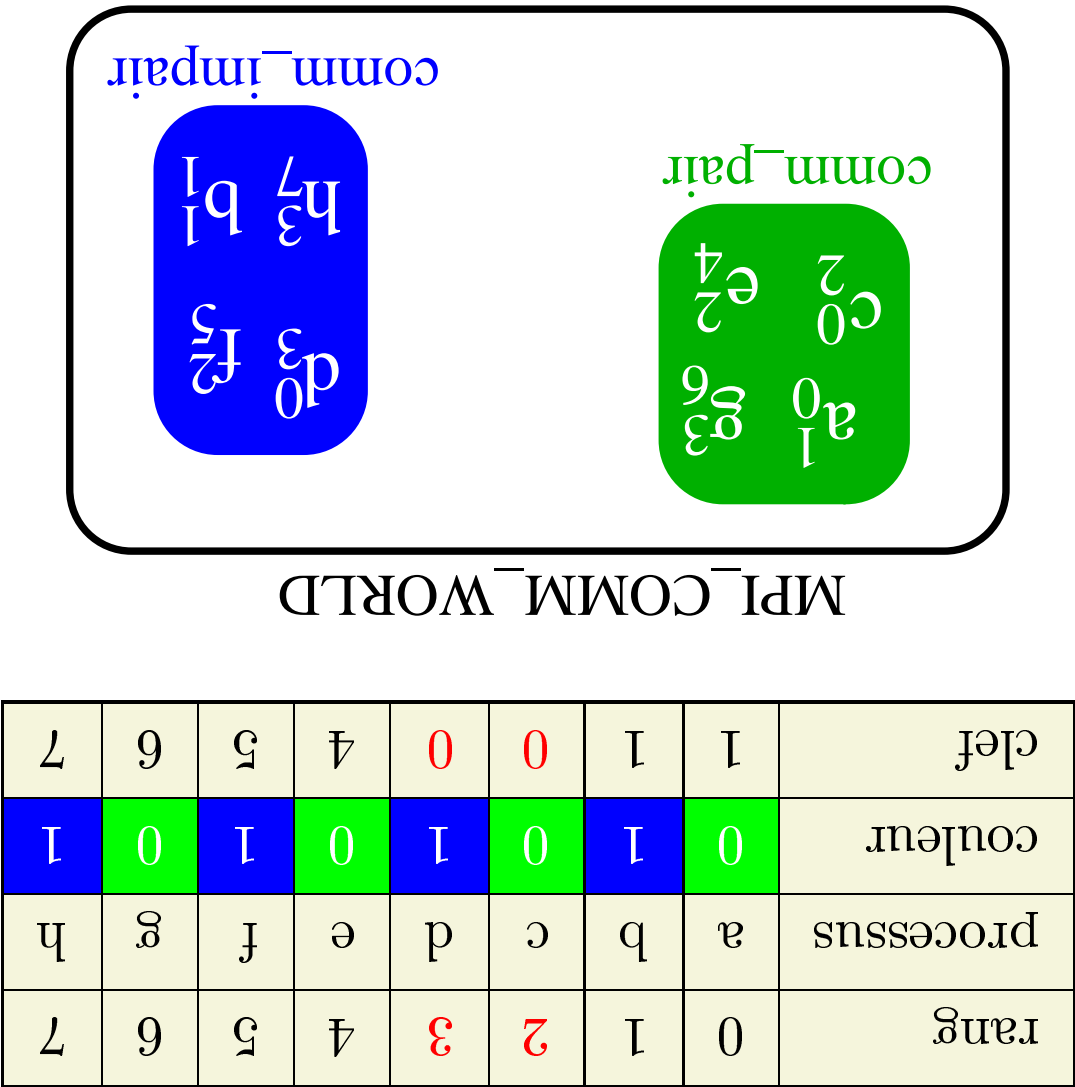
FIG. 46 – Construction de communicateurs avec `MPI_COMM_SPLIT()`

Un processus qui se voit attribuer une couleur égale à la valeur `MPI_UNDEFINED`, n'appartiendra qu'à son communicateur initial.

Voyons comment procéder pour construire nos deux communicateurs `pair` et `impair` avec le constructeur `MPI_COMM_SPLIT()`...

En pratique, ceci se met en place très simplement...

FIG. 47 – Construction des communicateurs pair et impair avec `MPI_COMM_SPLIT()`



```

1  program ComPairImpair
2  implicit none
3  include 'mpif.h'
4
5  integer, parameter :: m=16
6  integer :: cléf, couleur, comm
7  integer :: nb_procs, rang, code
8  real, dimension(m) :: a
9
10 call MPI_INIT(code)
11 call MPI_COMM_SIZE(MPI_COMM_WORLD, nb_procs, code)
12 call MPI_COMM_RANK(MPI_COMM_WORLD, rang, code)
13
14 *** Initialisation du vecteur A
15 a(:)=0.
16 if(rang == 2) a(:)=2.
17 if(rang == 3) a(:)=3.

```

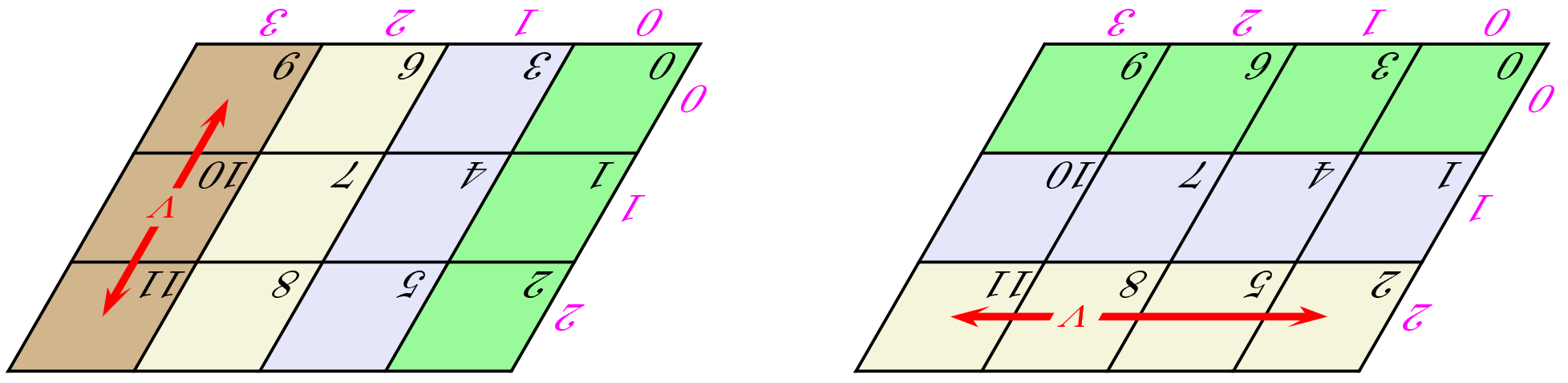
```

18  if ( mod(rang,2) == 0 ) then
19      ** Couleur et clefs des processus pairs
20      couleur = 0
21      if (rang == 2 ) then
22          clef = 0
23      else if ( rang == 0 ) then
24          clef = 1
25      else
26          clef = rang
27      end if
28  else
29      ** Couleur et clefs des processus impairs
30      couleur = 1
31      if (rang == 3) then
32          clef = 0
33      else
34          clef = rang
35      end if
36  end if
37
38  ** Créer les communicateurs pair et impair en leur donnant une même dénomination
39  call MPI_COMM_SPLIT(MPI_COMM_WORLD,couleur,clef,comm,code)
40  ** Le processus 0 de chaque communicateur diffuse son message aux processus
41  ** de son groupe
42  call MPI_BCAST(a,m,MPI_REAL,0,comm,code)
43  ** Destruction des communicateurs
44  call MPI_COMM_FREE(comm,code)
45  call MPI_FINALIZE(code)
46  end program ComPairImpair
    
```


8.6 – Subdiviser une topologie cartésienne

- ↳ La question est de savoir comment dégénérer une topologie cartésienne 2D ou 3D de processus en une topologie cartésienne respectivement 1D ou 2D.
- ↳ Pour *MPI*, dégénérer une topologie cartésienne 2D (ou 3D) revient à créer autant de communicateurs qu'il y a de lignes ou de colonnes (resp. de plans) dans la grille cartésienne initiale.
- ↳ L'intérêt majeur est de pouvoir effectuer des opérations collectives restreintes à un sous-ensemble de processus appartenant à :
 - ↳ une même ligne (ou colonne), si la topologie initiale est 2D ;
 - ↳ un même plan, si la topologie initiale est 3D.

FIG. 48 – Deux exemples de distribution de données dans une topologie 2D dé générée



Il existe deux façons de faire pour dégénérer une topologie :

→ en utilisant le sous-programme général `MPI_COMM_SPLIT()` ;

→ en utilisant le sous-programme `MPI_CART_SUB()` prévu à cet effet.

```

Logical, intent(in), dimension(NDim) :: Subdivision
Integer, intent(in) :: CommCart
Integer, intent(out) :: CommCartD, code
call MPI_CART_SUB(commCart, Subdivision, CommCartD, code)
    
```

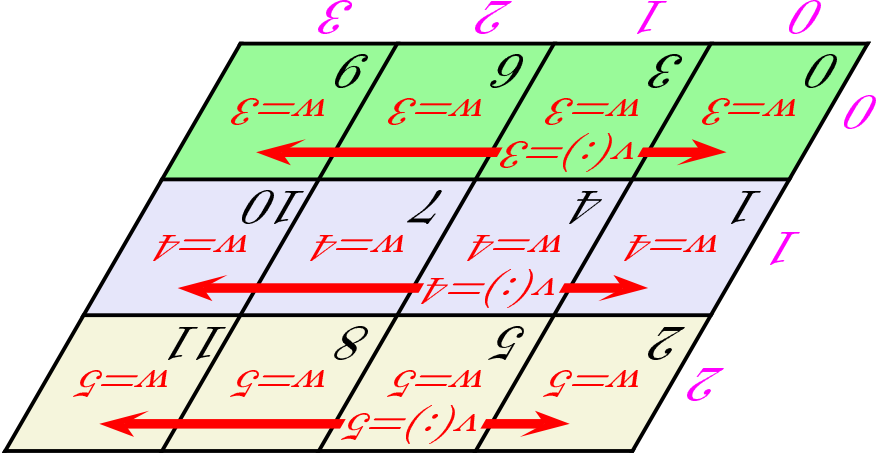
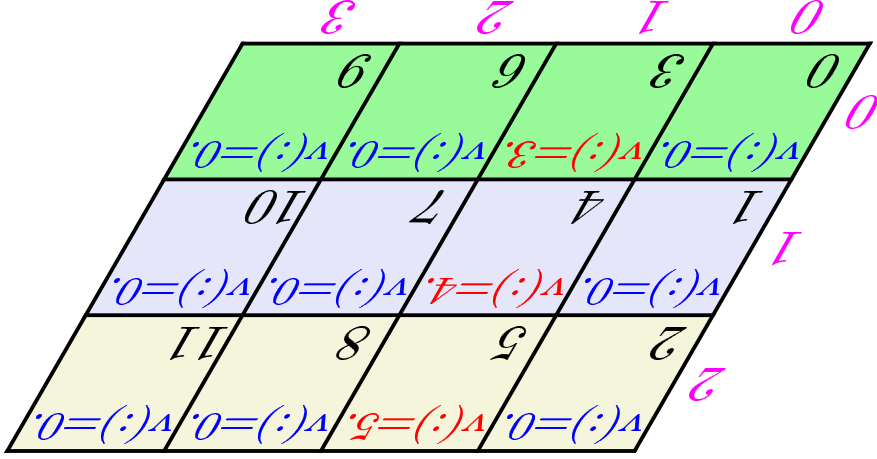


FIG. 49 – Représentation initiale d'un tableau V dans la grille 2D et représentation finale après la distribution de celui-ci sur la grille 2D dégénérée

```
1 program ComcartSub
2 implicit none
3 include 'mpi.h'
4
5 integer
6   :: Comm2D, Comm1D, rang, code
7   :: NDIm2D=2
8   :: integer, dimension(NDIm2D)
9   :: integer, dimension(NDIm2D)
10  :: Period, Subdivision
11  :: logical
12  :: Reordonne
13  :: integer, parameter
14  :: m=4
15  :: integer, dimension(m)
16  :: V(:)=0.
17  :: real
18  :: W=0.
```

```

13 call MPI_INIT(code)
14
15 *** Creation de la grille 2D initiale
16 Dim2D(1) = 4
17 Dim2D(2) = 3
18 Periode(:) = .false.
19 Redonne = .false.
20 call MPI_CART_CREATE(MPI_COMM_WORLD, NDim2D, Dim2D, Periode, Redonne, Comm2D, code)
21 call MPI_COMM_RANK(Comm2D, rang, code)
22 call MPI_CART_COORDS(Comm2D, rang, NDim2D, Coord2D, code)
23
24 *** Initialisation du vecteur V
25 if (Coord2D(1) == 1) V(:)=real(rang)
26 *** Nous voulons que chaque ligne de la grille soit une topologie cartésienne 1D
27 Subdivision(1) = .true.
28 Subdivision(2) = .false.
29
30 *** Subdivision de la grille cartésienne 2D.
31 call MPI_CART_SUB(Comm2D, Subdivision, Comm1D, code)
32 *** Les processus de la colonne 2 distribuent le vecteur V aux processus
33 *** de leur ligne
34 call MPI_SCATTER(V, 1, MPI_REAL, W, 1, MPI_REAL, 1, Comm1D, code)
35
36 print , ("Rang : ", I2, " ; Coordonnées : (", I1, ",", I1, ")", " ; W = ", F2.0), &
37 rang, Coord2D(1), Coord2D(2), W
38
39 call MPI_FINALIZE(code)
40 end program CommCartSub

```

```

< mpirun -np 12 commcartsub
Rang : 0 ; Coordonnees : (0,0) ; W = 3.
Rang : 1 ; Coordonnees : (0,1) ; W = 4.
Rang : 3 ; Coordonnees : (1,0) ; W = 3.
Rang : 8 ; Coordonnees : (2,2) ; W = 5.
Rang : 4 ; Coordonnees : (1,1) ; W = 4.
Rang : 5 ; Coordonnees : (1,2) ; W = 5.
Rang : 6 ; Coordonnees : (2,0) ; W = 3.
Rang : 10 ; Coordonnees : (3,1) ; W = 4.
Rang : 11 ; Coordonnees : (3,2) ; W = 5.
Rang : 9 ; Coordonnees : (3,0) ; W = 3.
Rang : 2 ; Coordonnees : (0,2) ; W = 5.
Rang : 7 ; Coordonnees : (2,1) ; W = 4.
    
```

8.7 – Intra et intercommunicateurs

- ↳ Les communicateurs que nous avons construits jusqu'à présent sont des **intracommunicateurs** (ex. `comm-pair` et `comm-impair`) car ils ne permettent pas que des processus appartenant à des communicateurs distincts puissent communiquer entre eux.
- ↳ Des processus appartenant à des intracommunicateurs distincts ne peuvent communiquer que s'il existe un lien de communication entre ces intracommunicateurs.
- ↳ Un **intercommunicateur** est un communicateur qui permet l'établissement de ce lien de communication.
- ↳ Une fois ce lien établi, seules sont possibles les communications point à point, les communications collectives au sein d'un **intercommunicateur** n'étant pas permises dans *MPI-1* (cette limitation a disparu dans *MPI-2*).
- ↳ Le sous-programme *MPI* `MPI_INTERCOMM_CREATE()` permet de construire des intercommunicateurs.
- ↳ Le couplage des modèles océan/atmosphère illustre bien l'utilité des intra et intercommunicateurs...

8.8 – Exemple récapitulatif

MPI_COMM_WORLD

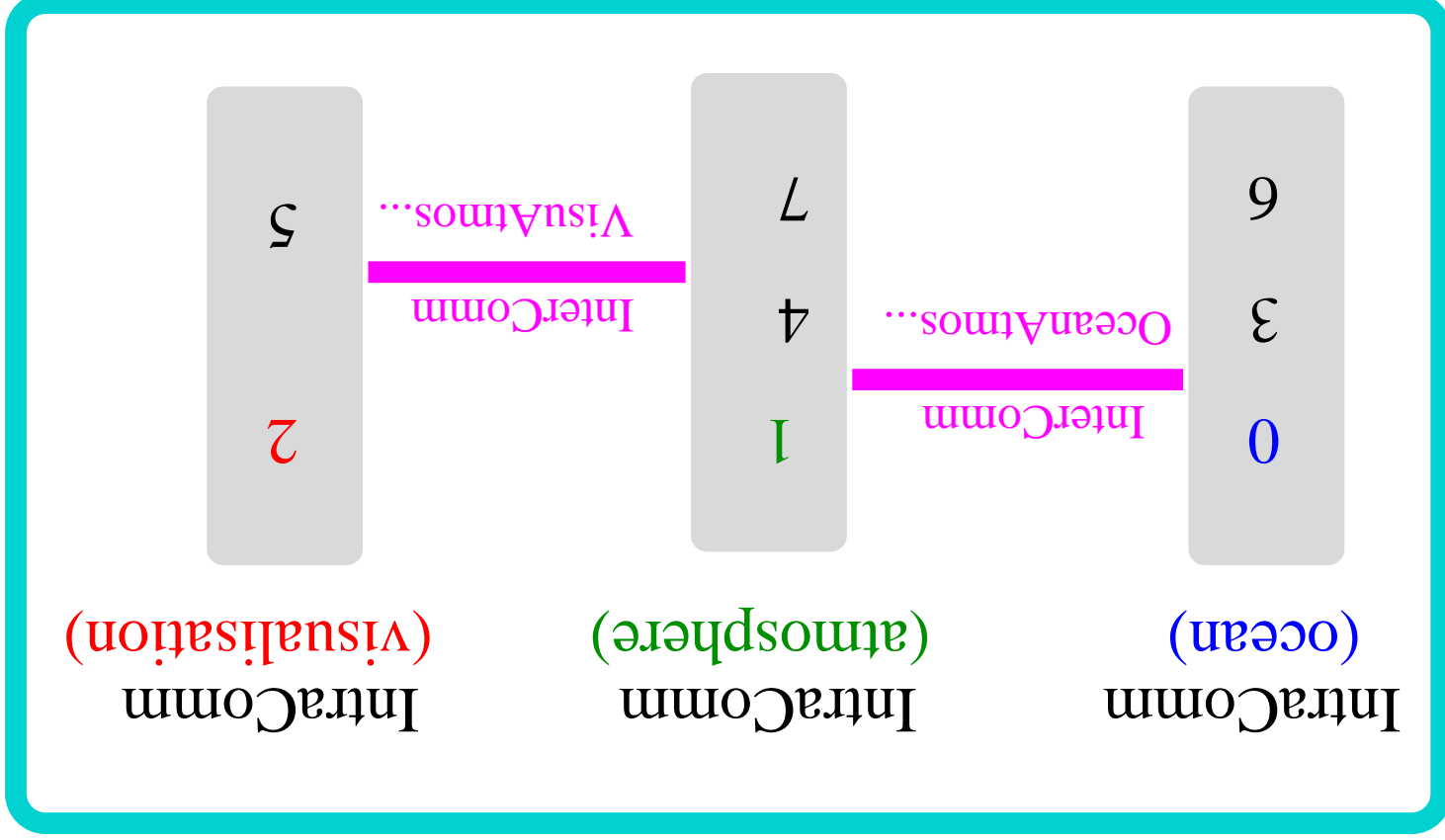


FIG. 50 – Couplage océan/atmosphère

8 — Communicateurs : exemple récapitulatif

```
1 program OceanAtmosphere
2   implicit none
3   include 'mpi.h'
4
5   integer, parameter :: tag1=111, tag2=222
6   integer :: NombreDeProcessus, RangMonde, NombreIntraComm, couleur, &
7   IntraComm, CommOceanAtmosphere, CommVisuAtmosphere
8
9   call MPI_INIT(code)
10  call MPI_COMM_SIZE(MPI_COMM_WORLD, NombreDeProcessus, code)
11  call MPI_COMM_RANK(MPI_COMM_WORLD, RangMonde, code)
12
13  *** Construction des 3 IntraCommunicateurs
14  NombreIntraComm = 3
15  couleur=mod(RangMonde, NombreIntraComm) i = 0,1,2
16  call MPI_COMM_SPLIT(MPI_COMM_WORLD, couleur, RangMonde, IntraComm, code)
```

8 — Communicateurs : exemple récursif

```
17 *** Construction des deux InterCommunicateurs et calcul
18 select case(couleur)
19 case(0)
20 *** InterCommunicateur OceanAtmosphere pour que le groupe 0 communique
21 *** avec le groupe 1
22 call MPI_INTERCOMM_CREATE (IntraComm, 0, MPI_COMM_WORLD, 1, tag1, CommOceanAtmosphere, &
23 code)
24 call ocean(IntraComm, CommOceanAtmosphere)
25
26 case(1)
27 *** InterCommunicateur OceanAtmosphere pour que le groupe 1 communique
28 *** avec le groupe 0
29 call MPI_INTERCOMM_CREATE (IntraComm, 0, MPI_COMM_WORLD, 0, tag1, CommOceanAtmosphere, &
30 code)
31 *** InterCommunicateur CommVisuAtmosphere pour que le groupe 1 communique
32 *** avec le groupe 2
33 call MPI_INTERCOMM_CREATE (IntraComm, 0, MPI_COMM_WORLD, 2, tag2, CommVisuAtmosphere, code)
34 call atmosphere(IntraComm, CommOceanAtmosphere)
35
36 case(2)
37 *** InterCommunicateur CommVisuAtmosphere pour que le groupe 2 communique
38 *** avec le groupe 1
39 call MPI_INTERCOMM_CREATE (IntraComm, 0, MPI_COMM_WORLD, 1, tag2, CommVisuAtmosphere, code)
40 call visualisation(IntraComm, CommVisuAtmosphere)
41 end select
```

8 — Communicateurs : exemple récursif

```

42 subroutine ocean(IntraComm, CommOceanAtmosphere)
43   implicit none
44   include 'mpi.h'
45   integer, parameter :: n=1024, tag1=3333
46   real, dimension(n) :: a, b, c
47   integer :: rang, code, germe(1), IntraComm, CommOceanAtmosphere
48   integer, dimension(MPI_STATUS_SIZE) :: statut
49   integer, intrinsic :: irtc
50
51   *** Les processus 0, 3, 6 dédiés au modèle océanographique effectuent un calcul
52   germe(1)=irtc()
53   call random_seed(put=germe)
54   call random_number(a)
55   call random_number(b)
56   call random_number(c)
57   a(:) = b(:) * c(:)
58
59   *** Les processus impliqués dans le modèle océan effectuent une opération collective
60   call MPI_ALLREDUCE(a, c, n, MPI_REAL, MPI_SUM, IntraComm, code)
61
62   *** Rang du processus dans IntraComm
63   call MPI_COMM_RANK(IntraComm, rang, code)
64
65   *** Echange de messages avec les processus associés au modèle atmosphérique
66   call MPI_SENDRECV(c, n, MPI_REAL, rang, tag1, c, n, MPI_REAL, rang, tag1, &
67   CommOceanAtmosphere, statut, code)
68
69   *** Le modèle océanographique tient compte des valeurs atmosphériques
70   a = b(:) * c(:)
71 end subroutine ocean

```

```

72 subroutine atmosphere(IntraComm,CommVisuAtmosphere,CommVisuAtmosphere)
73   implicit none
74   include 'mpi.h'
75
76   integer,parameter
77     :: n=1024,tag1=3333,tag2=4444
78     :: a,b,c
79     :: rang,code,germe(1),IntraComm, &
80     CommOceanAtmosphere,CommVisuAtmosphere
81   integer,dimension(MPI_STATUS_SIZE) :: statut
82   integer,intrinsic :: irtc
83   *** Les processus 1, 4, 7 dédiés au modèle atmosphérique effectuent un calcul
84     germe(1)=irtc()
85   call random_seed(put=germe)
86
87   call random_number(a)
88   call random_number(b)
89   call random_number(c)
90
91   a(:) = b(:) + c(:)

```

```

92  *** Les processus dédiés au modèle atmosphère effectuent une opération collective
93  call MPI_ALLREDUCE(a,c,n,MPI_REAL,MPI_MAX,IntraComm,code)
94
95  *** Rang du processus dans IntraComm
96  call MPI_COMM_RANK(IntraComm,rang,code)
97
98  *** Échange de messages avec les processus dédiés au modèle océanographique
99  call MPI_SENDRECV(c,n,MPI_REAL,rang,tag1,c,n,MPI_REAL,rang,tag1,&
100  CommOceanAtmosphere,statut,code)
101
102  *** Le modèle atmosphère tient compte des valeurs océanographiques
103  a(:) = b(:) * c(:)
104
105  *** Envoi des résultats aux processus dédiés à la visualisation
106  if (rang == 0 .or. rang == 1) then
107    call MPI_SEND(a,n,MPI_REAL,rang,tag2,CommVisuAtmosphere,code)
108  end if
109
110 end subroutine atmosphere

```

```

111 subroutine visualisation(IntraComm,CommVisuAtmosphere)
112     implicit none
113     include 'mpi.h'
114     integer,parameter :: n=1024,tag2=4444
115     real,dimension(n) :: a,b,c
116     integer :: rang,code,IntraComm,CommVisuAtmosphere
117     integer(MPI_STATUS_SIZE) :: statut
118     *** Les processus 2 et 5 sont chargés de la visualisation
119     call MPI_COMM_RANK(IntraComm,rang,code)
120     *** Réception des valeurs du champ à tracer
121     call MPI_RECV(a,n,MPI_REAL,rang,tag2,CommVisuAtmosphere,statut,code)
122     print*,'Moi, processus ',rang,' je trace mon champ A : ',a(:)
123     end subroutine visualisation
124
125
126
127
128
    
```

8.9 – Conclusion

➡ **Groupes et contextes** définissent un objet appelé **communicateur**.

➡ Les communicateurs définissent la portée des communications.

➡ Ils sont utilisés pour dissocier les espaces de communication.

➡ Un communicateur doit être spécifié à l'appel de toute fonction d'échange de

messages.

➡ Ils permettent d'éviter les confusions lors de la sélection des messages, par exemple

au moment de l'appel à un sous-programme d'une bibliothèque scientifique qui

elle-même effectue des échanges de messages.

➡ Les communicateurs offrent une programmation modulaire du point de vue de

l'espace de communication. Dans le cadre de projets importants, chaque équipe

développe son module sans se soucier du choix du communicateur des autres

équipes (exemple : modèle couple océan/atmosphère).

9 – Évolution de MPI : MPI-2

- ➡ Historique :
- ➡ début des travaux en mars 1995 ;
- ➡ brouillon présenté pour *SuperComputing 96* ;
- ➡ version « officielle » disponible en juillet 1997 ;
- ➡ voir <http://www.erc.msstate.edu/mpi/mpi2.html>
- ➡ Principaux domaines nouveaux :
- ➡ gestion dynamique des processus :
- ➡ possibilité de développer des codes MPMD ;
- ➡ support multi plates-formes ;
- ➡ démarrage et arrêt dynamique de sous-tâches ;
- ➡ gestion de signaux système.
- ➡ communications de mémoire à mémoire ;
- ➡ entrées/sorties parallèles.

➡ Autres domaines où apparaissent des améliorations :

➡ extensions concernant les intracommunicateurs ;

➡ extensions concernant les intercommunicateurs ;

➡ divers autres apports :

➡ inter-opérabilité entre C et Fortran ;

➡ interfaçage avec C++ et Fortran 90 (avec des limitations dans ce dernier cas).

➡ Extensions proposées en dehors de MPI-2 : MPI (*Interoperable MPI*), MPI-RT (*real time extensions*)

Tab. 4 – *Changements de dénominations dans MPI-2*

Ancien nom	Nouveau nom
MPI_TYPE_HVECTOR()	MPI_TYPE_CREATE_HVECTOR()
MPI_TYPE_HINDEXED()	MPI_TYPE_CREATE_HINDEXED()
MPI_TYPE_STRUCT()	MPI_TYPE_CREATE_STRUCT()
MPI_ADDRESS()	MPI_GET_ADDRESS()
MPI_TYPE_EXTENT()	MPI_TYPE_EXTENT()
MPI_TYPE_LB()	MPI_TYPE_GET_EXTENT()
MPI_TYPE_UB()	MPI_TYPE_GET_EXTENT()
MPI_LB	MPI_TYPE_CREATE_RESIZED()
MPI_UB	MPI_TYPE_CREATE_RESIZED()

Il en va de même pour les versions C de ces sous-programmes.

Les anciens noms sont toujours acceptés pour des raisons de compatibilité.

barrière 43

bibliothèque 153, 183

bloquantes (communications) 24, 30, 74, 75, 80, 83, 84, 95

non-bloquantes (communications) 94

collectives (communications) 16

communicateur 16, 20, 21, 24, 41, 126, 148-155, 157, 162, 163, 165, 166, 169, 183

intercommunicateur 176

intracommunicateur 176

communication 16, 67, 70, 71, 74, 75, 78, 80, 83-85, 87-89, 94, 95, 148, 149, 155, 176, 183

contexte 24, 155, 183

envoi 74, 76-78, 83, 88, 89

étiquette 24, 41

groupe 57, 153-155, 157, 162, 163, 183

intercommunicateur 176

intracommunicateur 176

message 8, 12, 14, 15, 72, 75, 76, 78, 84, 88, 150, 152, 183

MPMD 10, 11

optimisation 67

performances 163

persistantes (communications) 88, 89, 94, 95

portabilité 18, 109, 123

processeur 7

processus 8, 10, 12, 14, 15, 20, 21, 23, 24, 28, 33, 33, 35, 41, 57, 66, 123-126, 131, 132, 134, 136, 148, 150, 152-155, 157, 163, 165, 176

rang 21, 24, 28, 132, 134, 136, 157, 163, 165

réception 74, 76, 77, 83, 88, 89

requête 89, 94

SMD 10, 11

surcoût 78, 95

synchrones (communications) 76, 78

topologie 16, 123-126, 131, 132, 134, 136, 169, 171

types dérivés 123

MPI_ANY_SOURCE	28
MPI_ANY_TAG	28
MPI_CHARACTER	119
MPI_COMM_WORLD	20–22, 25, 29, 32, 34, 36, 38, 39, 43, 45, 47, 50, 52, 55, 60, 62, 64, 68, 69, 79, 81, 82, 90, 92, 103–108, 114, 115, 119, 120, 127, 129, 139, 140, 144, 149, 150, 156, 158, 159, 167, 168, 173, 177, 178
MPI_COMPLEX	64, 96
MPI_DOUBLE_PRECISION	39, 69
MPI_GROUP_EMPTY	162
MPI_INTEGER	25, 29, 32, 34, 36, 39, 45, 60, 62, 96, 101, 109, 119
MPI_LIB	121
MPI_LOGICAL	119
MPI_MAX	69, 181
MPI_PACKED	39
MPI_PROC_NULL	28, 140

MPI_PROD	62
MPI_REAL 47, 50, 52, 55, 69, 79, 82, 90, 92, 96, 98, 100, 101, 103, 105, 107, 109, 115, 119, 145, 150, 160, 168, 173, 179, 181, 182	
MPI_STATUS_SIZE	25, 29, 32, 34, 38, 68, 81, 103, 105, 107, 114, 119, 144, 179, 180, 182
MPI_SUM	60, 145, 179
MPI_UB	121
MPI_UNDEFINED	165
mpif.h 22, 25, 29, 32, 34, 38, 45, 47, 50, 52, 55, 60, 62, 64, 68, 81, 103, 105, 107, 114, 119, 127, 129, 139, 144, 149, 150, 158, 167, 172, 177, 179, 180, 182	

Index des sous-programmes MPI

MPI_ADDRESS	116, 118, 120
MPI_ALLGATHER	42, 52, 66
MPI_ALLGATHERV	66
MPI_ALLREDUCE	42, 57, 62, 145, 179, 181
MPI_ALLTOALL	42, 55, 66
MPI_ALLTOALLV	66
MPI_BARRIER	42, 43, 145
MPI_BCAST	42, 45, 57, 150, 160, 163, 168
MPI_CART_COORDS	134, 135, 140, 173
MPI_CART_CREATE	126, 127, 129, 140, 153, 173
MPI_CART_RANK	132, 133
MPI_CART_SHIFT	136–138, 140
MPI_CART_SUB	153, 171, 173
MPI_COMM_CREATE	153, 159
MPI_COMM_DUP	153
MPI_COMM_FREE	154, 160, 168
MPI_COMM_GROUP	156, 159
MPI_COMM_RANK	21, 22, 25, 29, 32, 34, 38, 45, 47, 50, 52, 55, 60, 62, 64, 68, 81, 103, 105, 107, 114, 119, 140, 144, 150, 158, 167, 173, 177, 179, 181, 182

Index des sous-programmes MPI

MPI_COMM_SIZE	21, 22, 32, 34, 47, 50, 52, 55, 60, 62, 68, 81, 139, 150, 158, 167, 177
MPI_COMM_SPLIT	153, 163, 165, 168, 171, 177
MPI_DIMS_CREATE	131, 139
MPI_FINALIZE	19, 22, 25, 29, 32, 34, 39, 45, 47, 50, 52, 55, 60, 62, 64, 69, 104, 106, 108, 115, 120, 140, 145, 150, 160, 168, 173
MPI_GATHER	42, 50, 66
MPI_GATHERV	66
MPI_GRAPH_CREATE	141, 144
MPI_GRAPH_NEIGHBORS	143, 145
MPI_GRAPH_NEIGHBORS_COUNT	143, 145
MPI_GROUP_COMPARE	162
MPI_GROUP_DIFFERENCE	162
MPI_GROUP_EXCL	153, 159
MPI_GROUP_FREE	154
MPI_GROUP_INCL	153, 159
MPI_GROUP_INTERSECTION	162
MPI_GROUP_RANK	156
MPI_GROUP_SIZE	156
MPI_GROUP_TRANSLATE_RANKS	160

Index des sous-programmes MPI

MPI_GROUP_UNION 162
MPI_INIT . . . 19, 22, 25, 29, 32, 34, 38, 45, 47, 50, 52, 55, 60, 62, 64, 68, 81, 103, 105, 107, 114,
119, 139, 144, 149, 150, 158, 167, 173, 177

MPI_INTERCOMM_CREATE 175, 178

MPI_IPROBE 83

MPI_RECV 80, 82–84, 90

MPI_SEND 80, 82, 84, 90

MPI_XSEND 83

MPI_OP_CREATE 57, 64

MPI_OP_FREE 57

MPI_PACK 38, 39

MPI_PROBE 83

MPI_RECV 25, 32, 36, 39, 69, 72, 79, 104, 106, 108, 120, 182

MPI_RECV_INIT 92

MPI_REDUCE 42, 57, 60, 64, 69

MPI_REDUCE_FREE 94

MPI_SCAN 57

MPI_SCATTER 42, 47, 66, 173

MPI_SCATTERV 66

Index des sous-programmes MPI

MPI_SEND	25, 30, 32, 36, 39, 69, 72, 78, 104, 106, 108, 120
MPI_SENDRECV	29, 123, 145, 179, 181
MPI_SENDRECV_REPLACE	34, 112, 115
MPI_SEND	78, 79, 84, 181
MPI_SEND_INIT	92, 95
MPI_START	92, 94, 95
MPI_TEST	83
MPI_TYPE_COMMIT	96, 102, 103, 105, 107, 115, 120
MPI_TYPE_CONTIGUOUS	28, 96, 98, 103
MPI_TYPE_EXTENT	121
MPI_TYPE_FREE	96, 102, 104, 106, 108, 115, 120
MPI_TYPE_HINDEXED	109, 111
MPI_TYPE_HVECTOR	96, 101, 102
MPI_TYPE_INDEXED	28, 109, 110, 115, 116
MPI_TYPE_LB	121, 122
MPI_TYPE_SIZE	121, 122
MPI_TYPE_STRUCT	28, 116, 117, 120, 121
MPI_TYPEUB	121, 122
MPI_TYPE_VECTOR	28, 96, 99, 100, 105, 107

MPI_UNPACK	38, 39
MPI_WAIT	82, 83, 90, 92
MPI_WTIME	69, 71, 79, 82