

# Introduction au parallélisme par échange de message via MPI

Master DFE – année 2020/2021

Mathieu Lobet, Maison de la Simulation  
Mathieu.lobet@cea.fr

# Introduction au parallélisme par échange de message via MPI

## 1) Description de l'approche

# Cours et matériel supplémentaires sur internet

Selon moi, le cours le plus complet sur MPI en français et anglais:

<http://www.idris.fr/formations/mpi/>

Les implémentations fournissent en général une documentation en ligne complète :

<https://www.open-mpi.org/>

<https://www.mcs.anl.gov/research/projects/mpi/learning.html>

<http://mpi.deino.net/>

# Information concernant le cours

Il existe des implémentations de MPI pour Fortran95, Fortran08, C, C++ et python.

La syntaxe diffère pour chaque langage.

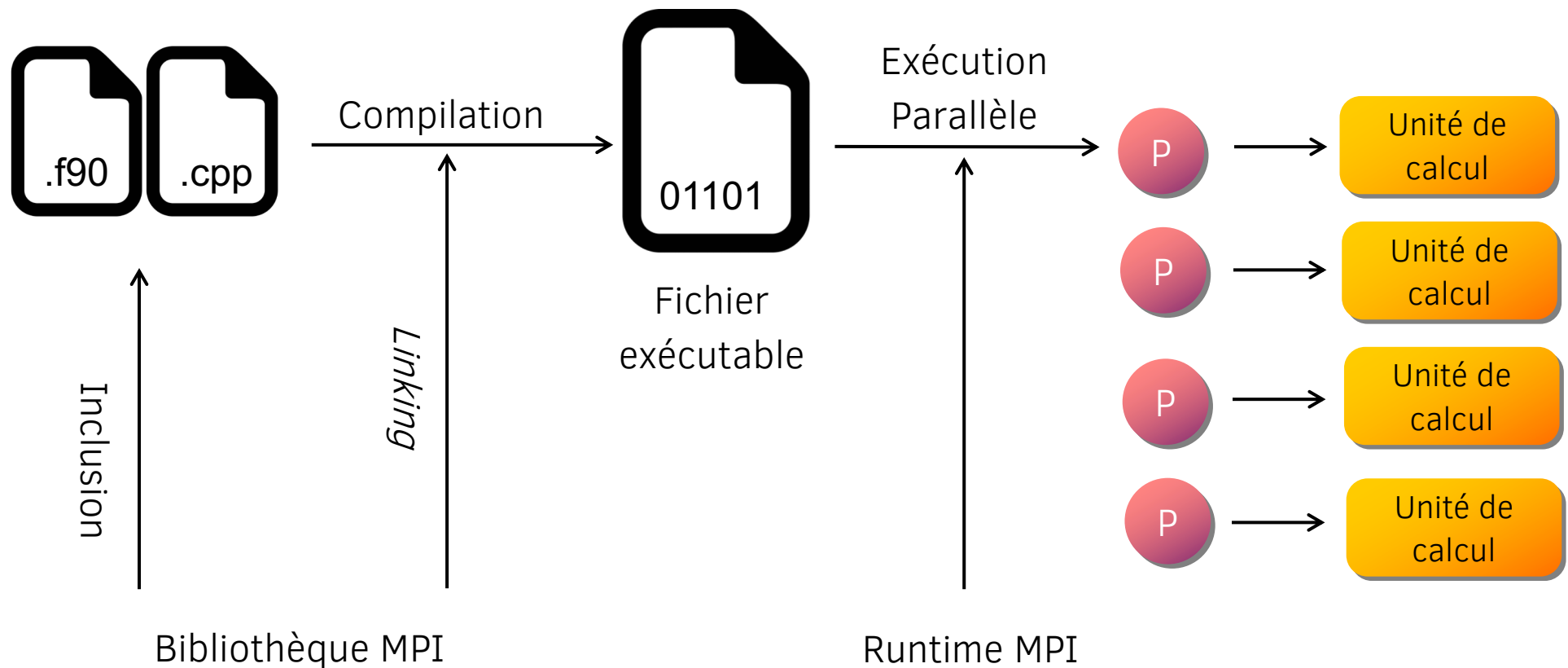
Dans ce cours, on programmera **en C++ en utilisant l'API C**.

Il existe également une version C++ avec classe pour les communicateurs ainsi qu'une version adaptée au Fortran moderne.

**La philosophie reste la même peu importe le langage choisi.**

# La chaîne de compilation et d'exécution d'un programme MPI

Programme parallèle avec appel aux fonctions MPI



# Implémentations MPI

Implémentations libres que vous pouvez vous procurer pour vos ordinateurs :

- OpenMPI
- MPICH
- Deino (Windows)

Implémentation propriétaire :

- IntelMPI

Disponible en Fortran, C, C++, Python



L'installation peut se faire par les sources ou via `apt-get install`

# Compilation d'un programme MPI (C++)

La compilation fait appel au **wrapper MPI mpic++**. Le wrapper utilise le compilateur C++ par défaut (par exemple g++) en y ajoutant des options de compilation supplémentaires et les chemins vers la bibliothèque MPI.

Pour connaître le contenu du wrapper :



```
> mpic++ -show
```

Pour compiler :



```
> mpic++ program.cpp -o executable
```

# Exécution d'un programme MPI

- L'exécution se fait par l'intermédiaire de la commande `mpirun`.
- `-np` représente le **nombre de processus MPI à lancer**.
- Si le nombre de processus MPI est inférieur au nombre d'unités de calcul disponibles, chaque processus est exécuté par une unité indépendante
- Il est possible de lancer plus de processus MPI que d'unités de calcul, dans ce cas, les ressources sont partagées.
- En OpenMPI, il faut spécifier « `--oversubscribe` » pour activer cette possibilité

Pour exécuter votre code en ligne de commande :



```
> mpirun -np 6 ./executable
```

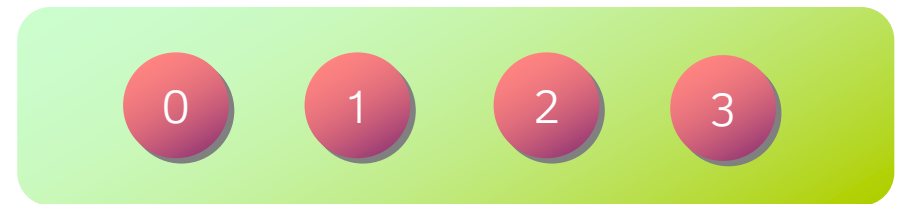
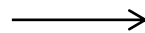


# Démarrage d'un programme MPI : notion de communicateur

## *Communicator*

- Un **communicateur** est un ensemble de processus MPI capables de communiquer entre eux.
- Au sein d'un communicateur, chaque processus MPI est représenté par **un rang (*rank*) unique sous forme d'un entier**.
- Le communicateur par défaut regroupe **l'ensemble des processus** et se nomme **MPI\_COMM\_WORLD**.

```
> mpirun -np 4 ./executable
```



MPI\_COMM\_WORLD : Communicateur composé de 4 rangs

# Démarrage d'un programme MPI (Fortran et C++) : inclure MPI

La première étape consiste à ne pas oublier d'**inclure le module MPI** (header en C/C++)



.f90

```
Program test
```

```
Use mpi
```

```
...
```

```
End program
```



.cpp

```
#include <mpi.h>
```

```
int main( int argc, char *argv[] )
```

```
{
```

```
}
```

# Démarrage d'un programme MPI : initialiser et finaliser MPI

- MPI est une bibliothèque qui fonctionne par **appel à des fonctions**
- La deuxième étape importante est l'**initialisation de MPI**
- Il ne faut pas oublier de **finaliser pour finir son programme proprement**



.cpp

```
#include <mpi.h>

int main( int argc, char *argv[] )
{
    int ierror ;

    ierror = MPI_Init() ;

    ...

    ierror = MPI_Finalize() ;

}
```

# Démarrage d'un programme MPI : initialiser MPI

- L'initialisation se fait avec la fonction **MPI\_INIT**.
- Toute fonction MPI renvoie en **dernier argument un code d'erreur noté ici ierror**.



.cpp

```
ierror = MPI_INIT();
```

- Le code d'erreur permet si besoin de vérifier qu'un appel s'est bien déroulé



[https://www.open-mpi.org/doc/v1.8/man3/MPI\\_Init.3.php](https://www.open-mpi.org/doc/v1.8/man3/MPI_Init.3.php)

# Démarrage d'un programme MPI (Fortran) : finaliser MPI

- La finalisation se fait avec la fonction **MPI\_FINALIZE**.
- Elle est appelée à la toute fin du programme



.cpp

```
lerror = MPI_Finalize();
```

# Démarrage d'un programme MPI : récupérer le nombre de rangs

Le **nombre de rangs** dans le communicateur MPI\_COMM\_WORLD se récupère via la fonction **MPI\_COMM\_SIZE** : c'est le **nombre total de processus** demandé.



.cpp

```
Int number_of_ranks ;
```

```
lerror = MPI_Comm_size(MPI_COMM_WORLD, &number_of_ranks) ;
```

- MPI\_COMM\_WORLD : communicateur (ici celui par défaut)
- Number\_of\_ranks : entier renvoyé contenant le nombre de rangs MPI



[https://www.open-mpi.org/doc/current/man3/MPI\\_Comm\\_size.3.php](https://www.open-mpi.org/doc/current/man3/MPI_Comm_size.3.php)

# Démarrage d'un programme MPI : récupérer le rang de chaque processus MPI

Chaque processus récupère son **rang** dans le communicateur `MPI_COMM_WORLD` via la fonction **`MPI_COMM_RANK`**.



.cpp

```
int rank ;  
  
lerror = MPI_Comm_rank(MPI_COMM_WORLD, &rank) ;
```

- `MPI_COMM_WORLD` : communicateur (ici celui par défaut)
- `rank` : entier renvoyé désignant le rang du processus qui appelle la fonction



[https://www.open-mpi.org/doc/v3.0/man3/MPI\\_Comm\\_rank.3.php](https://www.open-mpi.org/doc/v3.0/man3/MPI_Comm_rank.3.php)

# Mesure du temps : MPI\_Wtime

- **MPI\_Wtime** permet de récupérer le temps écoulé sur le processus courant en seconde



.cpp

```
double time ;  
Time = MPI_Wtime() ;
```

Par deux appels et une soustraction, cette fonction permet de déterminer le **temps passer dans une section du code**



.cpp

```
double time ;  
time = MPI_Wtime() ;  
  
! Des calculs...  
...  
  
! Ce temps est le temps passé entre les deux appels à MPI_Wtime  
time = MPI_Wtime() - time ;
```



[https://www.open-mpi.org/doc/v4.0/man3/MPI\\_Wtime.3.php](https://www.open-mpi.org/doc/v4.0/man3/MPI_Wtime.3.php)



# Exercice n°1 : votre premier programme MPI



- Rendez vous sur le GitHub des exercices : <https://github.com/Maison-de-la-Simulation/HPC-DFE-Paris-Saclay/>
- Télécharger les exercices sur votre session de travail
- Décompressez l'archive en ligne de commande



```
> tar xvf archivedossier.tar
```

•Rendez vous dans le dossier de l'exercice n°1 appelé 1\_initialization



```
> cd exercices/mpi/1_initialization
```

Ouvrez les instructions contenues dans le fichier README.md avec votre éditeur de fichier favori (vim, emacs, atom, gedit...) ou depuis la page GitHub

# Différencier du code pour des processus donnés



Le programme s'exécute simultanément autant de fois qu'il y a de processus en parallèle : chaque ligne de code est appelée par chaque processus.

Pour faire en sorte que certaines portions de code soient réservées à certains processus, on utilise des conditions `if` avec le numéro de rang comme condition.



.cpp

```
if (rank == 1) {  
    ! Cette portion de code ne sera exécutée que par le rang 1  
    ...  
}
```

# Votre premier programme MPI

A ce stade du cours, vous savez maintenant :

- Écrire un programme parallèle simple
- Compiler un programme MPI
- Exécuter un programme MPI
- Récupérer le nombre de rangs et le rang de chaque processus

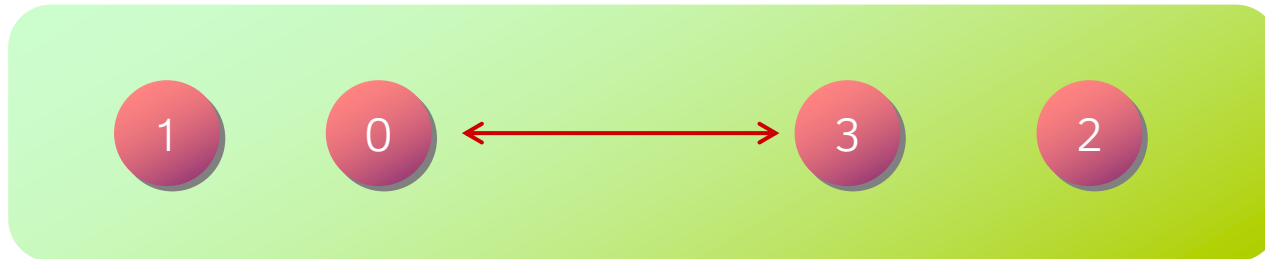
# Introduction au parallélisme par échange de message via MPI

## 2) Les communications point à point bloquantes

# Notion de communication point à point

## *Point to point communication*

- L'échange de message constitue la base du concept MPI
- L'échange de message se décompose toujours en **deux étapes** :
  - **Envoi** : Un processus envoie un message à un processus destinataire en spécifiant le rang
  - **Réception** : Un processus doit explicitement recevoir le message en connaissant le rang de l'expéditeur



← Le processus de rang 0 envoie un message au processus de rang 3

→ Le processus de rang 3 reçoit un message du processus de rang 0

# Notion de communication point à point : Envoi de données via MPI\_Send

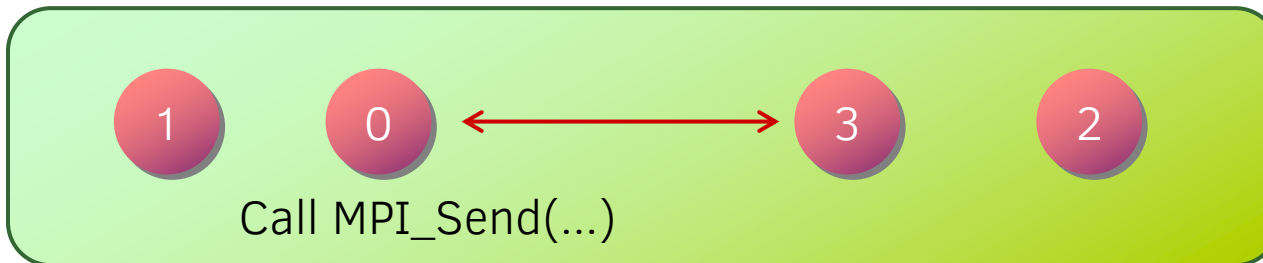
*Point to point communication : MPI\_Send*

• **MPI\_Send** est la fonction appelée par le processus expéditeur



.cpp

```
MPI_Send(message, size, data_type, destination_rank, tag, communicator)
```



[https://www.open-mpi.org/doc/v1.8/man3/MPI\\_Send.3.php](https://www.open-mpi.org/doc/v1.8/man3/MPI_Send.3.php)

# Notion de communication point à point : Envoi de données via MPI\_Send (C/C++)

*Point to point communication : MPI\_Send*

• **MPI\_Send** est la fonction appelée par le processus expéditeur



.cpp

```
MPI_Send(message, size, data_type, destination_rank, tag, communicator)
```

- **Message (const void \*)** : la variable contenant le message à envoyer (booléen, entier, double, caractère, chaîne, tableau, structure plus complexe...)
- **Size (int)** : nombre d'éléments constituant le message (> 1 uniquement pour une chaîne ou un tableau)
- **data\_type** : type de variable utilisée pour le message (MPI\_INT pour les integer, MPI\_DOUBLE pour les double, MPI\_FLOAT pour les float...)
- **Tag (int)** : numéro attribué à la communication si plusieurs coms vers le même processus
- **destination\_rank (int)** : rang du processus destinataire



[https://www.open-mpi.org/doc/v1.8/man3/MPI\\_Send.3.php](https://www.open-mpi.org/doc/v1.8/man3/MPI_Send.3.php)

# Tag MPI



La notion de tag permet de différencier des communications mais cet aspect ne sera pas exploité dans ce cours.

Une valeur de tag par défaut peut être donnée ou l'utilisation du paramètre **`MPI_ANY_TAG`** permet d'ignorer l'utilisation de ce dernier.



# Exemple d'utilisation de MPI\_Send

*Point to point communication : MPI\_Send*

- Envoi d'un message de type real(8) au processus de rang 3 par le processus de rang 2



.cpp

```
double message ;
int tag = 0 ;
ierror ;

message = 1245.76 ;

if (rank == 2) {
    ierror = MPI_Send(&message, 1, MPI_DOUBLE, 3, tag, &
        MPI_COMM_WORLD)
}
```

# Exemple d'utilisation de MPI\_Send

*Point to point communication : MPI\_Send*

Envoi d'un message de type tableau contenant 5 entiers au processus de rang 6 par le processus de rang 1



.cpp

```
int message[5] ;  
int tag, ierror ;  
  
message = { 12,45,37,43,59 } ;  
  
if (rank == 1) {  
    ierror = MPI_Send(message, 5, MPI_INT, 6, tag, &  
        MPI_COMM_WORLD) ;  
}
```

# Exemple d'utilisation de MPI\_SEND

*Point to point communication : MPI\_SEND*

Envoi d'un message au processus de rang 6 par le processus de rang 1 de type chaîne de caractère contenant 4 caractères et commençant au deuxième élément du message



.cpp

```
char message[10] = « abcdefghij » ;  
int tag, ierror ;  
  
if (rank == 1) {  
    ! Seulement « bcde » est envoyé  
    ierror = MPI_Send(&message[1], 4, MPI_CHAR, 6, tag, MPI_COMM_WORLD) ;  
}
```

# Notion de communication point à point : Réception de données via MPI\_

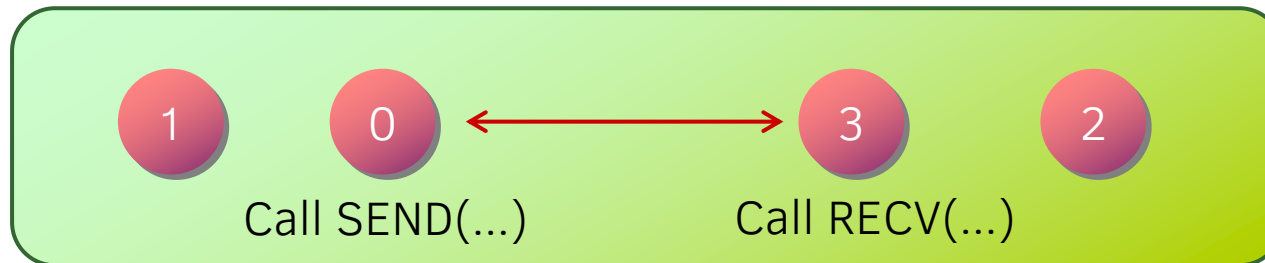
*Point to point communication : MPI\_RECV*

• **MPI\_RECV** est la fonction appelée par le processus destinataire



.f90

```
MPI_RECV(message, size, data_type, source_rank, tag, communicator, status, ierror)
```



[https://www.open-mpi.org/doc/v1.8/man3/MPI\\_Recv.3.php](https://www.open-mpi.org/doc/v1.8/man3/MPI_Recv.3.php)

# Notion de communication point à point : Réception de données via MPI\_

*Point to point communication : MPI\_RECV*

• **MPI\_RECV** est la fonction appelée par le processus destinataire



.f90

```
MPI_RECV(message, size, data_type, source_rank, tag, communicator, status, ierror)
```

- **Message** : la variable contenant le message à recevoir (booléen, entier, double, caractère, chaîne)
- **Size** : nombre d'éléments constituant le message (> 1 uniquement pour une chaîne ou un tableau)
- **data\_type** : type de variable utilisée pour le message (MPI\_INTEGER pour les integer, MPI\_DOUBLE pour les double)
- **source\_rank** : rang du processus expéditeur
- **Status** : état de la communication (en dehors de la portée de ce cours)

# Notion de communication point à point : Réception de données via MPI\_

*Point to point communication : MPI\_Recv*

• **RECV** est la fonction appelée par le processus destinataire



.cpp

```
MPI_Recv(message, size, data_type, source_rank, tag, communicator, status, ierror)
```

- **Message (void \*)** : la variable contenant le message à recevoir (booléen, entier, double, caractère, chaîne)
- **Size (int)** : nombre d'éléments constituant le message (> 1 uniquement pour une chaîne ou un tableau)
- **data\_type (MPI\_Datatype)** : type de variable utilisée pour le message (MPI\_INT pour les int, MPI\_DOUBLE pour les doubles)
- **source\_rank (int)** : rang du processus expéditeur
- **Status (MPI\_Status \*)** : état de la communication (en dehors de la portée de ce cours)



[https://www.open-mpi.org/doc/v1.8/man3/MPI\\_Recv.3.php](https://www.open-mpi.org/doc/v1.8/man3/MPI_Recv.3.php)

# Exemple d'utilisation de MPI\_SEND and MPI\_RECV (Fortran95)

*Point to point communication : MPI\_SEND and MPI\_RECV*

- Envoi d'un message de type real(8) au processus de rang 3 par le processus de rang 2
- Réception d'un message de type real(8) par le rang 2 venant du rang 3



.f90

```
Real(8) :: message
Integer :: tag
Integer :: ierror

If (rank == 2) then

    Message = 1245.76

    Call MPI_SEND(message, 1, MPI_DOUBLE_PRECISION, 3, tag, &
        MPI_COMM_WORLD, ierror)

End if

If (rank == 3) then

    Call MPI_RECV(message, 1, MPI_DOUBLE_PRECISION, 2, tag, &
        MPI_COMM_WORLD, MPI_STATUS_IGNORE, ierror)
```

# Exemple d'utilisation de MPI\_Send and MPI\_Recv (C/C++)

*Point to point communication : MPI\_Send and MPI\_Recv*

- Envoi d'un message de type **double** au processus de rang 3 par le processus de rang 2
- Réception d'un message de type **double** par le rang 2 venant du rang 3



.cpp

```
double message ;
int tag, ierror ;

if (rank == 2) {
    Message = 1245.76 ;
    MPI_Send(&message, 1, MPI_DOUBLE, 3, tag, MPI_COMM_WORLD) ;
}

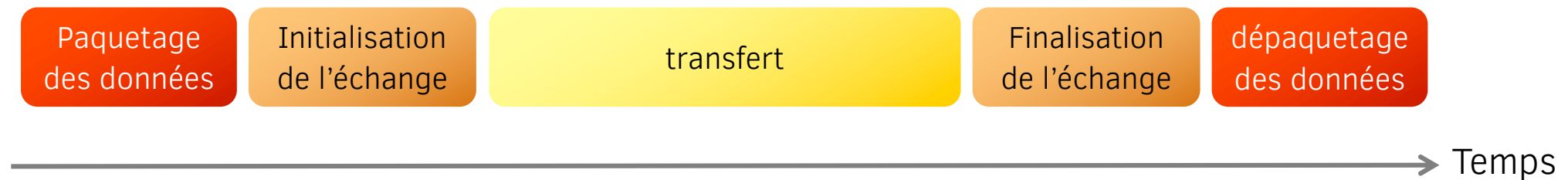
If (rank == 3) {
    MPI_Recv(&message, 1, MPI_DOUBLE, 2, tag,
            MPI_COMM_WORLD, MPI_STATUS_IGNORE) ;
}
```



# Mieux comprendre une communication

*Point to point communication : SEND and RECV*

• Une communication se compose d'un ensemble de sous-étapes :



MPI\_SEND

Réseau ou canal de communication

MPI\_RECV

## Exercice n°2 : Utilisation des communications point à point



•Rendez vous dans le dossier de l'exercice n°2 appelé 2\_blocking\_com



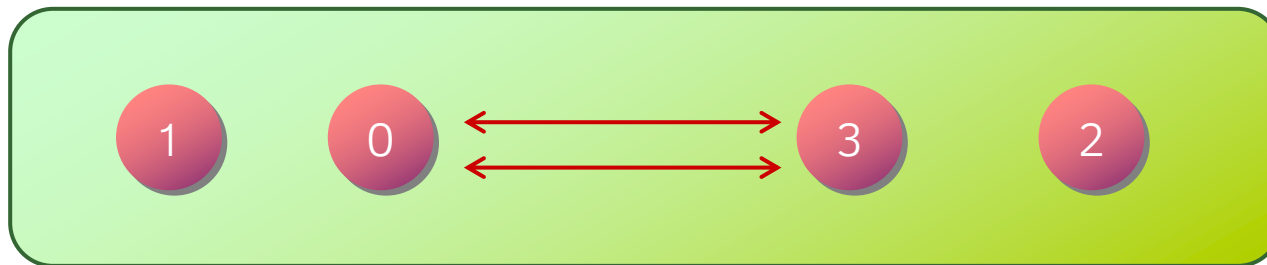
```
> cd exercises/mpi/2_blocking_com
```

•Ouvrez les instructions contenues dans le fichier **README.md** avec votre éditeur de fichier favori

# Communication point à point : MPI\_SENDRECV

*Point to point communication*

• Il est parfois nécessaire de faire un **échange mutuel** de données. Pour ce faire, il existe une fonction



Le processus de rang 0 envoie et reçoit un message au processus de rang 3

Le processus de rang 3 envoie et reçoit un message du processus de rang 0

# Communication point à point : MPI\_SENDRECV (Fortran95)

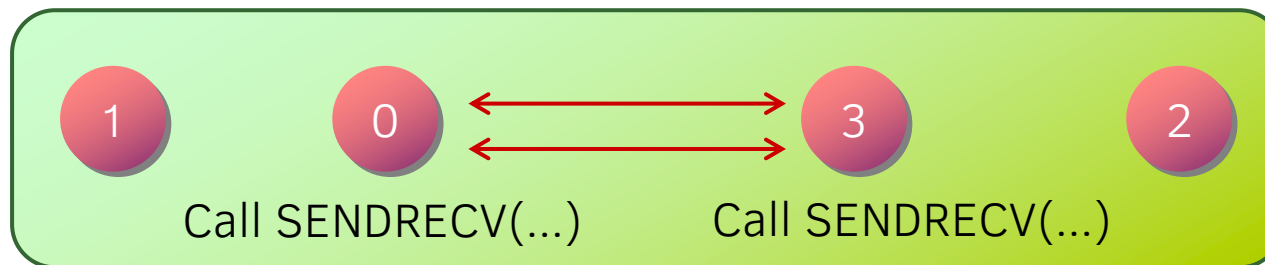
*Point to point communication : MPI\_SENDRECV*

• **SENDRECV** est appelée par les processus expéditeur et destinataire en même temps



.f90

```
MPI_SENDRECV(  
    &  
    send_message, send_size, send_type, destination, send_tag, & rcv_message, rcv_size, rcv_type
```



[https://www.open-mpi.org/doc/v1.8/man3/MPI\\_Sendrecv.3.php](https://www.open-mpi.org/doc/v1.8/man3/MPI_Sendrecv.3.php)

# Communication point à point : MPI\_Sendrecv (C/C++)

*Point to point communication : MPI\_Sendrecv*

- **MPI\_Sendrecv** est appelée par les processus expéditeur et destinataire en même temps



.cpp

```
ierror = MPI_Sendrecv(  
    send_message, send_size, send_type, destination, send_tag, rcv_message, rcv_size, rcv_type
```

- **send\_message (const void \*)** : données envoyées
- **rcv\_message (void \*)** : données reçues
- Les autres paramètres sont les mêmes que pour **MPI\_Send** et **MPI\_Recv**



[https://www.open-mpi.org/doc/v1.8/man3/MPI\\_Sendrecv.3.php](https://www.open-mpi.org/doc/v1.8/man3/MPI_Sendrecv.3.php)

# Exemple d'utilisation de MPI\_SENDRECV (Fortran 95)

*Point to point communication : MPI\_SENDRECV*

- Envoi et réception d'un message de type real(8) au processus de rang 3 par le processus de rang 2
- Envoi et réception d'un message de type real(8) par le rang 2 venant du rang 3



.f90

```
Real(8) :: send_message
Real(8) :: recv_message
Integer :: send_tag
Integer :: recv_tag
Integer :: ierror

If (rank == 2) then

    send_message = 1245.76

    Call MPI_SENDRECV(send_message, 1, MPI_DOUBLE_PRECISION, 3, send_tag, &
        recv_message, 1, MPI_DOUBLE_PRECISION, 3, recv_tag, &
        MPI_COMM_WORLD, MPI_STATUS_IGNORE, ierror)

End if

If (rank == 3) then

    send_message = 4567.32

    Call MPI_SENDRECV(send_message, 1, MPI_DOUBLE_PRECISION, 2, send_tag, &
        recv_message, 1, MPI_DOUBLE_PRECISION, 2, recv_tag, &
        MPI_COMM_WORLD, MPI_STATUS_IGNORE, ierror)
```

End if

# Exemple d'utilisation de MPI\_Sendrecv

*Point to point communication : MPI\_Sendrecv*

- Envoi et réception d'un message de type **double** au processus de rang 3 par le processus de rang 2
- Envoi et réception d'un message de type **double** par le rang 2 venant du rang 3



.cpp

```
double send_message, rcv_message ;
int send_tag, rcv_tag, ierror ;

if (rank == 2) {

    send_message = 1245.76 ;

    ierror = MPI_Sendrecv(&send_message, 1, MPI_DOUBLE, 3, send_tag,
                        &rcv_message, 1, MPI_DOUBLE, 3, rcv_tag,
                        MPI_COMM_WORLD, MPI_STATUS_IGNORE) ;

}

If (rank == 3) {

    send_message = 4567.32 ;

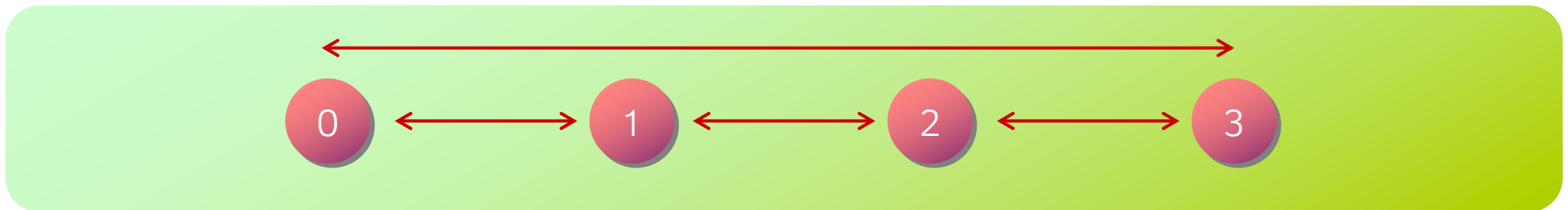
    ierror = MPI_Sendrecv(&send_message, 1, MPI_DOUBLE, 2, send_tag, &
                        &rcv_message, 1, MPI_DOUBLE, 2, rcv_tag, &
                        MPI_COMM_WORLD, MPI_STATUS_IGNORE) ;

}
```

# Communication point à point : MPI\_SENDRECV pour les communications

*Point to point communication*

- La fonction `MPI_SENDRECV` est également nécessaire pour effectuer des **communications chaînées**
- L'utilisation de `MPI_SEND` et `MPI_RECV` nécessiterait de gérer manuellement les synchronisations



Chaque processus reçoit un élément d'un processus A et envoie des données à un processus B



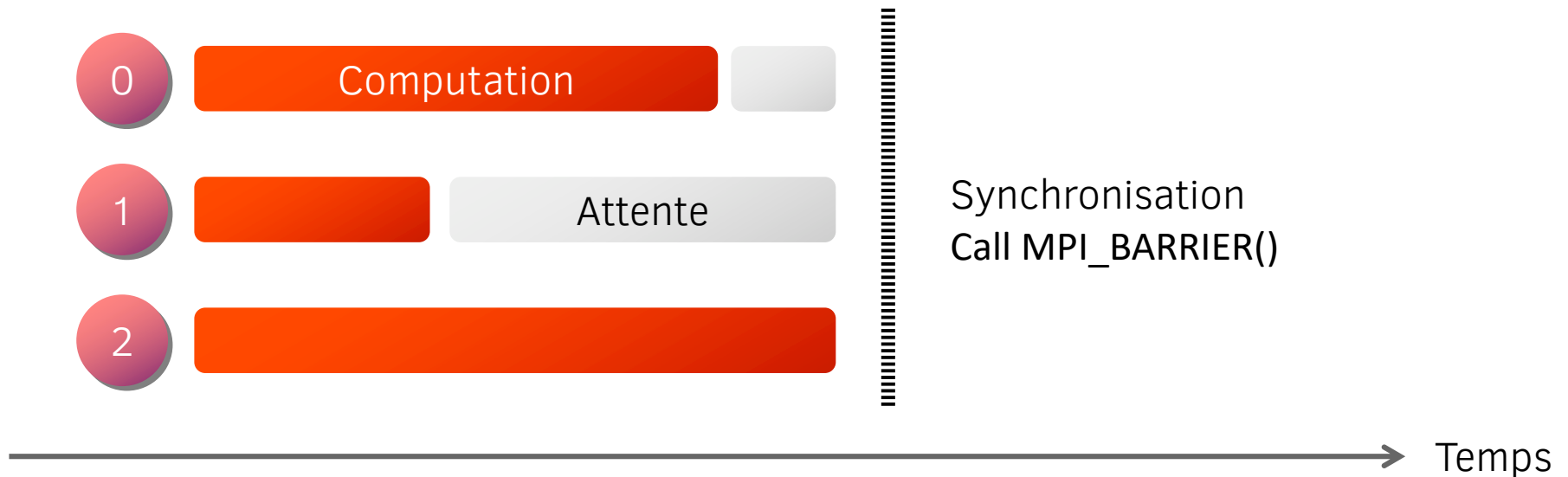
Lorsqu'une communication n'arrive pas à son terme, le programme attend et peut res



# Notion de barrière explicite

*Explicit barrier*

- Il est parfois nécessaire d'imposer une **étape de synchronisation ou barrière** qui ne sera pas
- La fonction **MPI\_BARRIER** est une **manière explicite d'exiger cette synchronisation** dans le code

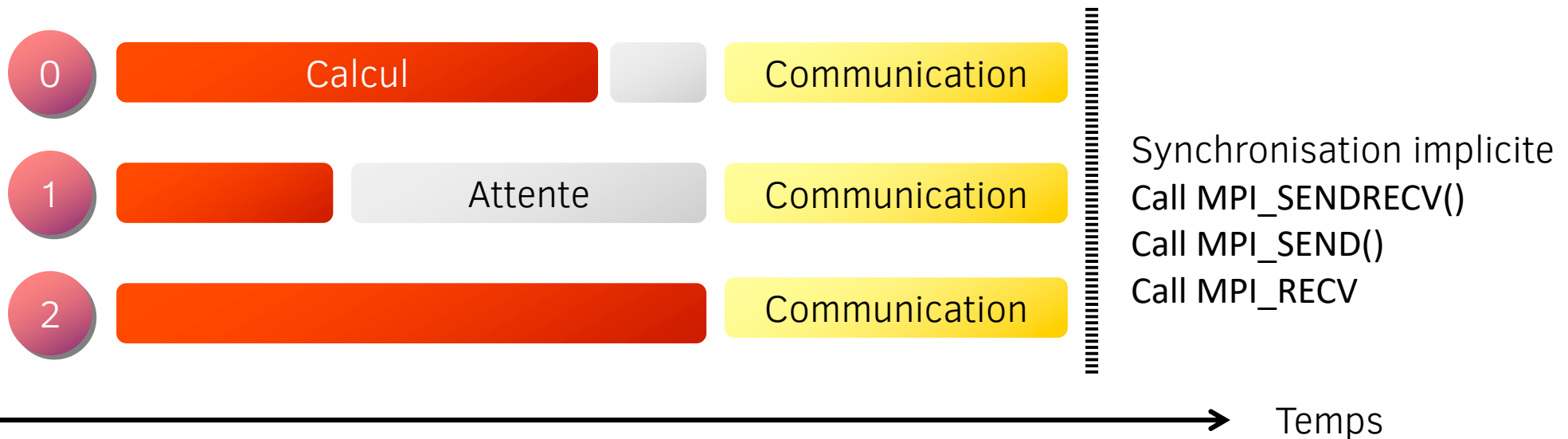


[https://www.open-mpi.org/doc/v1.5/man3/MPI\\_Barrier.3.php](https://www.open-mpi.org/doc/v1.5/man3/MPI_Barrier.3.php)

# Notion de barrière implicite

*Implicit barrier*

- Certaines fonctions d'échange induisent des barrières implicites au niveau des processus concernés.
- C'est le cas de `MPI_SEND`, `MPI_RECV`, `MPI_SENDRECV` d'où l'appellation de **communication bloquante**.



Si certains processus sont en avance, ils effectuent une attente active ou passive. Cet

## Exercice n°3 : Chaîne ou anneau de communication

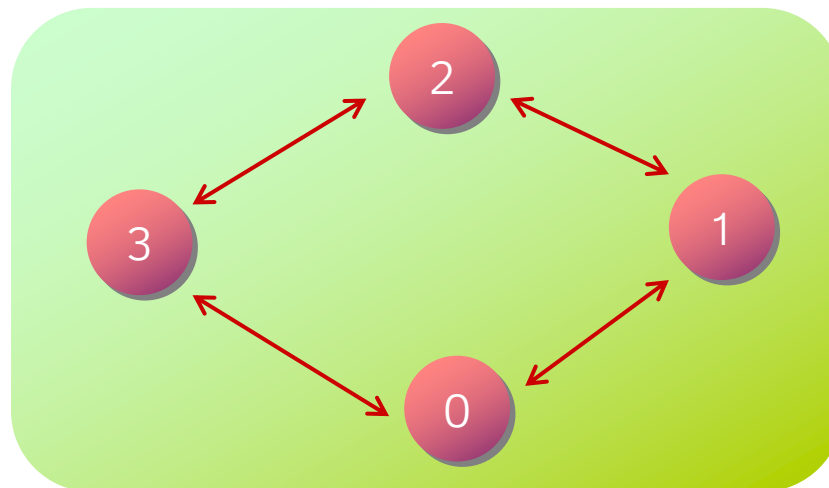


•Rendez vous dans le dossier de l'exercice n°3 appelé 3\_sendrecv



```
> cd exercices/mpi/3_sendrecv
```

•Ouvrez les instructions contenues dans le fichier **README.md** avec votre éditeur de fichier favori



Anneau de communication

# Premiers échanges MPI

A ce stade du cours, vous savez maintenant :

- Faire communiquer différents processus entre eux
- Gérer des chaînes de communication
- Demander une synchronisation explicite

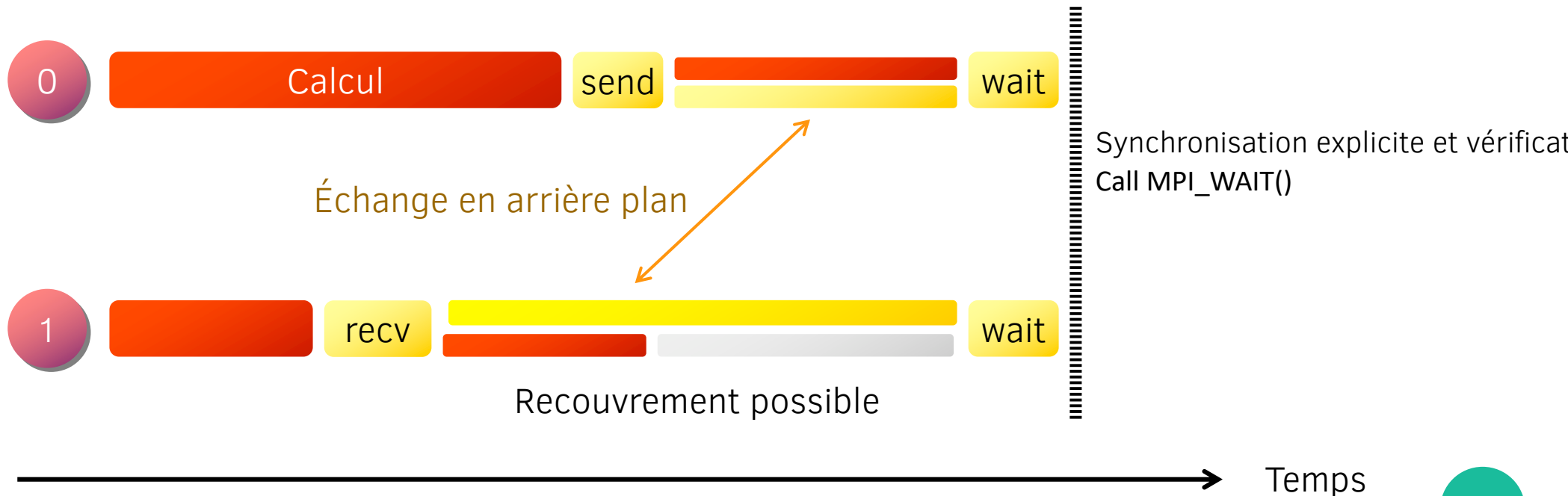
# Introduction au parallélisme par échange d

## 3) Les communications point à point non-bloquantes

# Communication point à point non-bloquante

*Non-blocking communication*

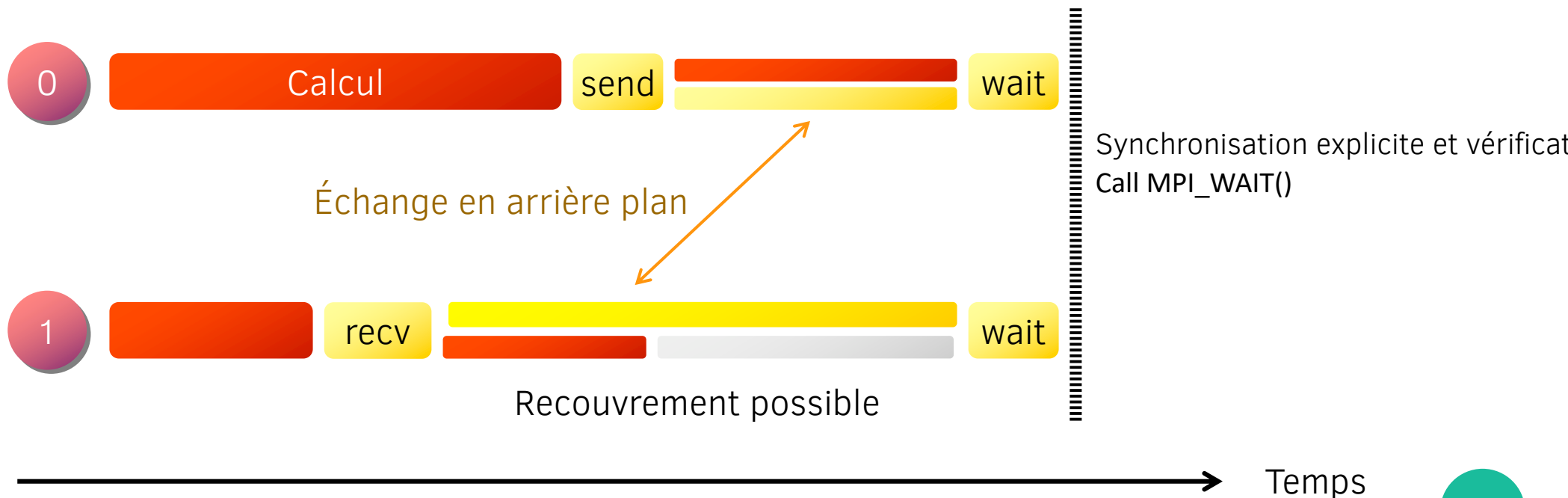
- Les **communications non-bloquantes** permettent d'éviter la **synchronisation implicite** des processeurs.
- Une **synchronisation explicite** est nécessaire pour s'assurer que les communications ont eu lieu.
- Ce type de communication permet de **recouvrir communication et calcul** : envoi et réception en



# Communication point à point non-bloquante

*Non-blocking communication*

- Les communications se font via les fonctions `MPI_ISEND` et `MPI_Irecv`.
- Les communications se voient attribuées **un identifiant unique**
- A un moment donné, il est nécessaire de **vérifier que ces communications ont eu lieu**, c'est le rôle de `MPI_WAIT`
- `MPI_WAIT` impose une barrière



# Les fonctions `MPI_Isend` et `MPI_Irecv` (Fortran95)

*MPI\_Isend and MPI\_Irecv*

- `MPI_ISEND` est la fonction appelée par le processus expéditeur



.f90

```
MPI_ISEND(message, size, data_type, destination_rank, tag, communicator, request, ierror)
```

- `MPI_Irecv` est la fonction appelée par le processus receveur



.f90

```
MPI_Irecv(message, size, data_type, source_rank, tag, communicator, request, ierror)
```

- Les paramètres sont les mêmes pour que pour les communications bloquantes.
- S'ajoute la variable `request` (`int`) utilisée par `MPI_WAIT` pour vérifier l'état de la communication.



[https://www.open-mpi.org/doc/v4.0/man3/MPI\\_Isend.3.php](https://www.open-mpi.org/doc/v4.0/man3/MPI_Isend.3.php)

[https://www.open-mpi.org/doc/v4.0/man3/MPI\\_Irecv.3.php](https://www.open-mpi.org/doc/v4.0/man3/MPI_Irecv.3.php)



# Les fonctions `MPI_Isend` et `MPI_Irecv` (C/C++)

*MPI\_Isend and MPI\_Irecv*

- `MPI_Isend` est la fonction appelée par le processus expéditeur



.cpp

```
MPI_Isend(&message, size, data_type, destination_rank, tag, communicator, request) ;
```

- `MPI_Irecv` est la fonction appelée par le processus receveur



.cpp

```
MPI_Irecv(&message, size, data_type, source_rank, tag, communicator, request) ;
```

- Les paramètres sont les mêmes pour que pour les communications bloquantes.
- S'ajoute la variable `request` (de type `MPI_Request *`) utilisé par `MPI_WAIT` pour vérifier l'état



[https://www.open-mpi.org/doc/v4.0/man3/MPI\\_Isend.3.php](https://www.open-mpi.org/doc/v4.0/man3/MPI_Isend.3.php)

[https://www.open-mpi.org/doc/v4.0/man3/MPI\\_Irecv.3.php](https://www.open-mpi.org/doc/v4.0/man3/MPI_Irecv.3.php)

# La fonction `MPI_Wait`

*MPI\_Wait*

**.MPI\_Wait** est la fonction appelée pour vérifier et attendre que la communication a bien été effectuée.



.f90

```
MPI_WAIT(request, status, ierror)
```



.cpp

```
MPI_Wait(request, status) ;
```



[https://www.open-mpi.org/doc/v4.0/man3/MPI\\_Wait.3.php](https://www.open-mpi.org/doc/v4.0/man3/MPI_Wait.3.php)

# Exemple d'utilisation des communications non-bloquantes (Fortran 95)

## *Non-blocking communication*

- Envoi d'un message de type real(8) au processus de rang 3 par le processus de rang 2
- Réception d'un message de type real(8) par le rang 2 venant du rang 3



.f90

```
Real(8) :: send_message, recv_message
Integer :: send_tag, recv_tag, ierror
Integer :: request

If (rank == 2) then

    send_message = 1245.76

    Call MPI_ISEND(send_message, 1, MPI_DOUBLE_PRECISION, 3, send_tag, &
        MPI_COMM_WORLD, request, ierror)

End if

If (rank == 3) then

    Call MPI_Irecv(recv_message, 1, MPI_DOUBLE_PRECISION, 2, recv_tag, &
        MPI_COMM_WORLD, request, ierror)

End if

... Calcul ...

Call MPI_WAIT(request, MPI_COMM_WORLD, ierror) ;
```

... Calcul suivant ...

# Exemple d'utilisation des communications non-bloquantes (C/C++)

## *Non-blocking communication*

- Envoi d'un message de type double au processus de rang 3 par le processus de rang 2
- Réception d'un message de type double par le rang 2 venant du rang 3



.cpp

```
double send_message, recv_message ;
int send_tag, recv_tag, ierror ;
MPI_Request request ;

If (rank == 2) {

    send_message = 1245.76

    Call MPI_Isend(send_message, 1, MPI_DOUBLE, 3, send_tag, &
        MPI_COMM_WORLD, request) ;

}

If (rank == 3) {

    Call MPI_Irecv(&recv_message, 1, MPI_DOUBLE, 2, recv_tag, &
        MPI_COMM_WORLD, request) ;

}

MPI_Wait(request, MPI_COMM_WORLD) ;
```

# La fonction `MPI_WAITALL`

*MPI\_Wait*

• `MPI_WAITALL` effectue l'action de `MPI_WAIT` sur un tableau de requêtes.



.f90

```
MPI_WAITALL(number_of_request, request, status, ierror)
```



.cpp

```
MPI_Waitall(number_of_request, request, status) ;
```

• `request` et `status` sont alors des tableaux (d'entiers en Fortran95, `MPI_Request` et `MPI_Status` en C++).



[https://www.open-mpi.org/doc/v4.0/man3/MPI\\_Waitall.3.php](https://www.open-mpi.org/doc/v4.0/man3/MPI_Waitall.3.php)

# Mélange des types de communication



Il est tout à fait possible de mélanger des appels bloquants à des appels non-bloquants.

Lorsque la requête est terminée, elle devient `MPI_REQUEST_NULL`. Il est également possi

## Exercice n°4 : Utilisation des communications non bloquantes



•Rendez vous dans le dossier de l'exercice n°4 appelé 4\_nonblocking\_com



```
> cd exercises/mpi/4_nonblocking_com
```

•Ouvrez les instructions contenues dans le fichier **README.md** avec votre éditeur de fichier favori

# Introduction au parallélisme par échange d

## 4) Les communications collectives

### 4.1) Communications collectives de base



# Communication collective : les différents types

## *Collective communication*

- Les **communications collectives** sont des communications qui font intervenir **plusieurs processus**
- 3 types de communication collective :
  - **Synchronisation** : c'est le `MPI_BARRIER`
  - **Transfert de données** (diffusion, collecte)
  - **Transfert et opérations sur les données** (opération de réduction)

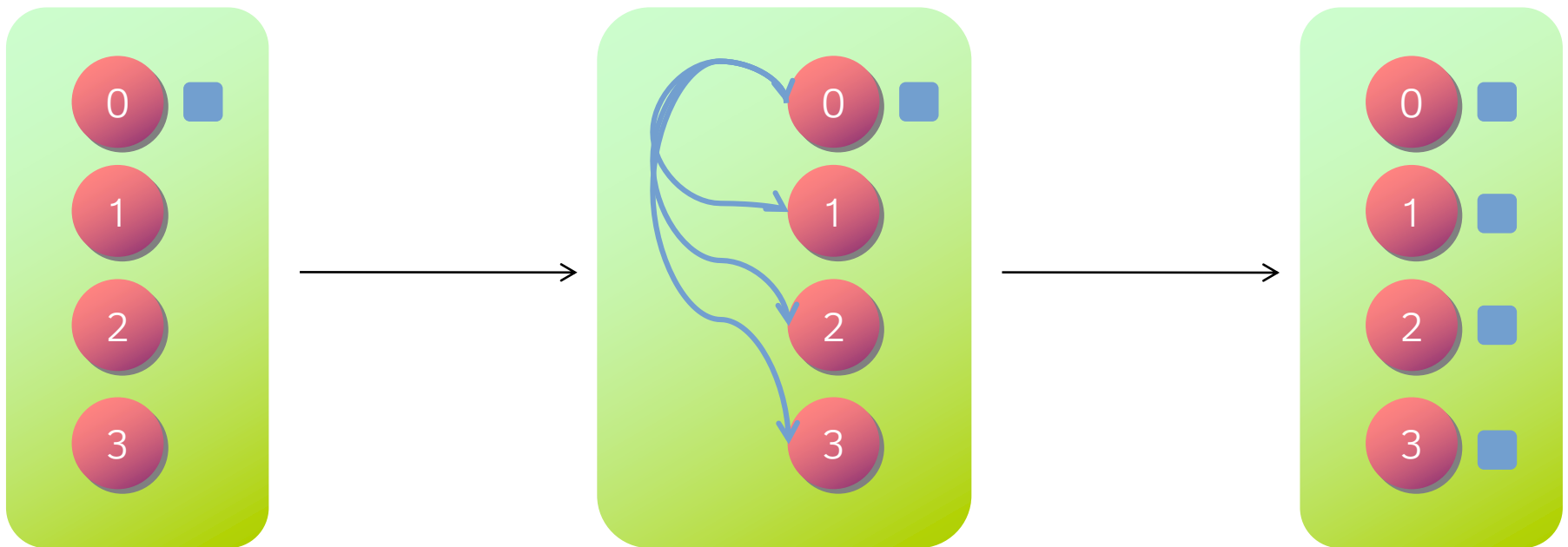


Les versions standards induisent des barrières implicites pour les processus concernés

# Communication collective : diffusion générale grâce à MPI\_Bcast

*Collective communication*

- Envoi d'une donnée depuis un processus vers tous les processus du communicateur



artition des données avant l'échange

MPI\_Bcast

Nouvelle répartition des données après

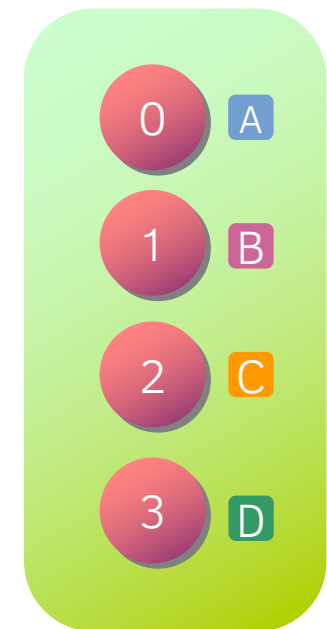
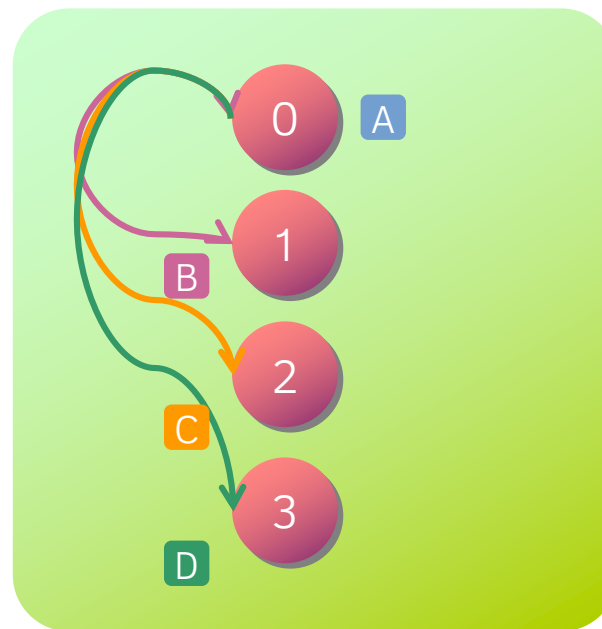
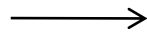
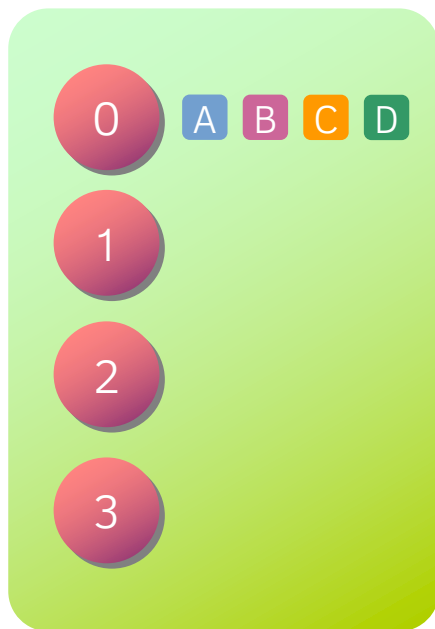


[https://www.open-mpi.org/doc/v4.0/man3/MPI\\_Bcast.3.php](https://www.open-mpi.org/doc/v4.0/man3/MPI_Bcast.3.php)

# Communication collective : diffusion sélective grâce à MPI\_Scatter

*Collective communication*

- Partage de données sélectif (selon les critères du développeur) depuis un processus vers tous les autres



Partition des données avant l'échange

MPI\_SCATTER

Nouvelle répartition des données

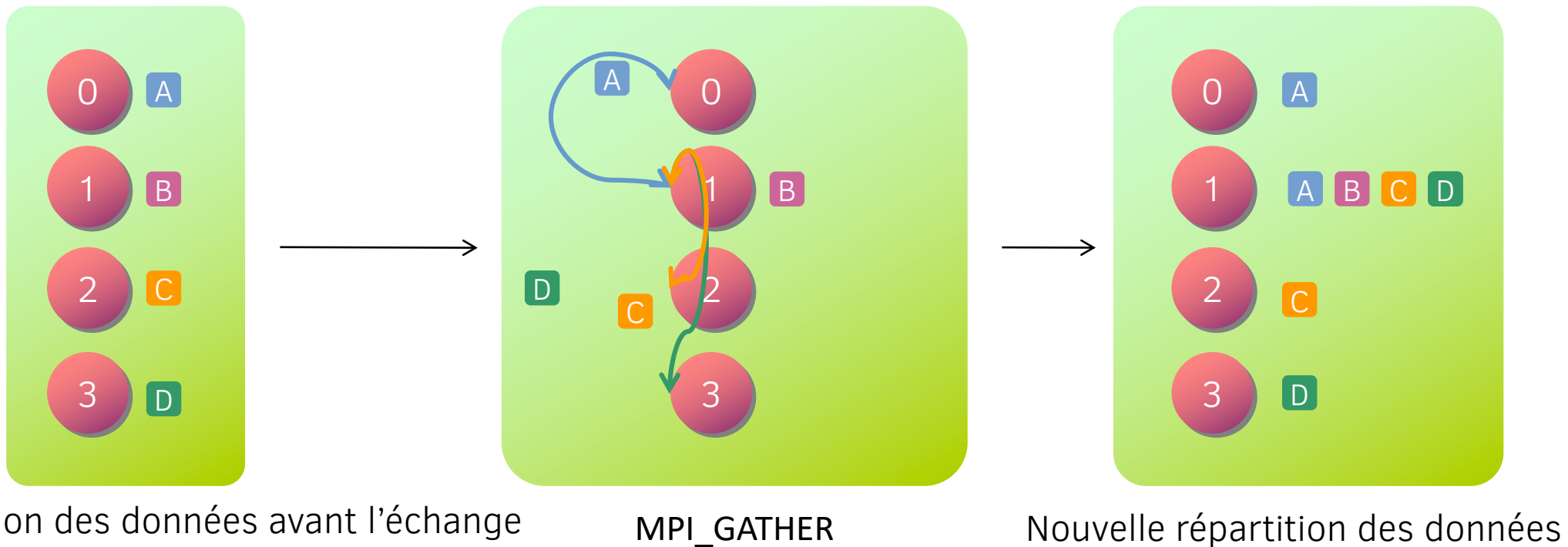


[https://www.open-mpi.org/doc/v4.0/man3/MPI\\_Scatter.3.php](https://www.open-mpi.org/doc/v4.0/man3/MPI_Scatter.3.php)

# Communication collective : collecte grâce à MPI\_Gather

*Collective communication*

- Envoi de données réparties sur plusieurs processus vers un processus unique



[https://www.open-mpi.org/doc/v4.0/man3/MPI\\_Gather.3.php](https://www.open-mpi.org/doc/v4.0/man3/MPI_Gather.3.php)

# Communication collective : collecte grâce à MPI\_Gather (Fortran95)

## *Collective communication*

• **MPI\_Gather** est appelée par les processus expéditeurs et destinataires en même temps



.f90

```
MPI_GATHER(send_buf, send_count, send_type, &
           rcv_buf, rcv_count, rcv_type, &
           destination, communicator, ierror)
```

- **send\_buf** : la valeur ou un ensemble de valeur
- **send\_count** : le nombre de valeur à envoyer
- **send\_type** : type MPI des valeurs envoyées
- **rcv\_buf** : le tableau réunissant les valeurs reçues
- **rcv\_count** : nombre d'éléments reçus
- **rcv\_type** : le type des données reçues
- **destination** : le processus qui reçoit les données



[https://www.open-mpi.org/doc/v4.0/man3/MPI\\_Gather.3.php](https://www.open-mpi.org/doc/v4.0/man3/MPI_Gather.3.php)

# Communication collective : collecte grâce à MPI\_Gather (C/C++)

## *Collective communication*

• **MPI\_Gather** est appelée par les processus expéditeurs et destinataires en même temps



.cpp

```
MPI_Gather(send_buf, send_count, send_type,  
          recv_buf, recv_count, recv_type,  
          destination, communicator) ;
```

- `send_buf (const void *)` : la valeur ou un ensemble de valeur
- `send_count (int)` : le nombre de valeur à envoyer
- `send_type (MPI_Datatype)` : type MPI des valeurs envoyées
- `recv_buf (void *)` : le tableau réunissant les valeurs reçues
- `recv_count (int)` : nombre d'éléments reçus
- `recv_type (MPI_Datatype)` : le type des données reçues
- `Destination (int)` : le processus qui reçoit les données



[https://www.open-mpi.org/doc/v4.0/man3/MPI\\_Gather.3.php](https://www.open-mpi.org/doc/v4.0/man3/MPI_Gather.3.php)

# Introduction au parallélisme par échange d

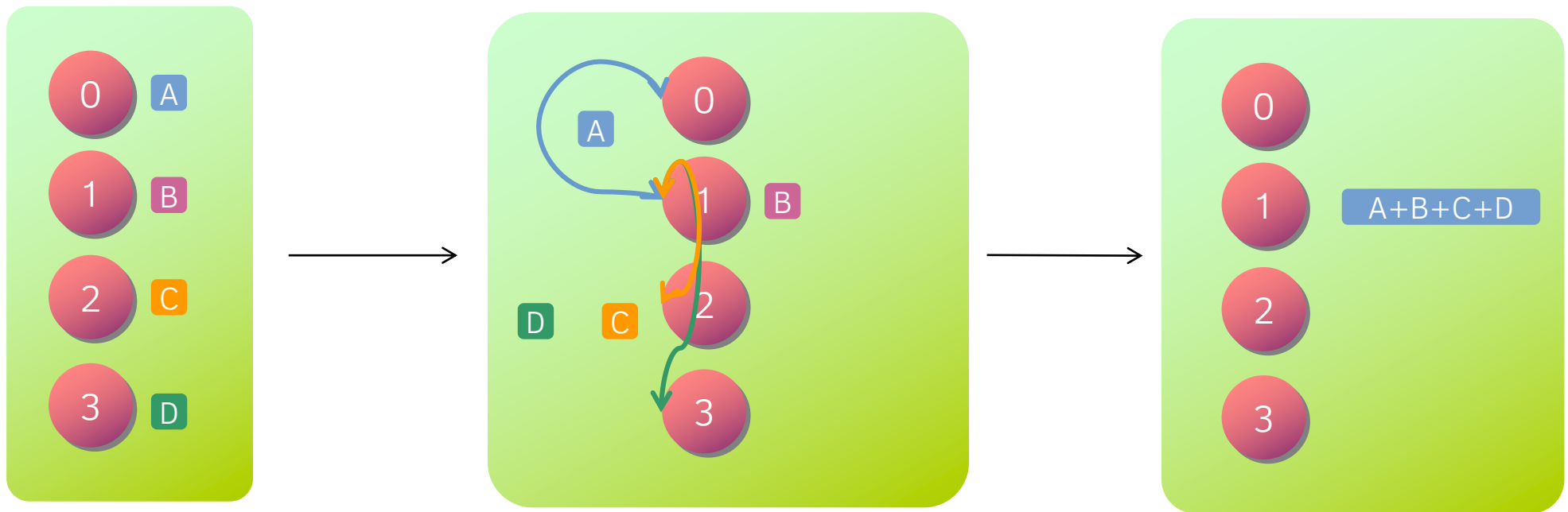
## 4) Les communications collectives

### 4.2) Les réductions

# Communication collective : réduction grâce à MPI\_Reduce

*Collective communication*

- Envoi de données réparties sur plusieurs processus vers un seul processus avec une opération



partition des données avant l'échange MPI\_REDUCE + MPI\_SUM Nouvelle répartition des données et ad



[https://www.open-mpi.org/doc/v4.0/man3/MPI\\_Reduce.3.php](https://www.open-mpi.org/doc/v4.0/man3/MPI_Reduce.3.php)



# Communication collective : réduction grâce à MPI\_Reduce (Fortran95)

## *Collective communication*

• **MPI\_Reduce** est appelée par les processus expéditeurs et destinataires en même temps



.f90

```
MPI_REDUCE(send_value, recv_value, size, MPI_data_type, MPI_reduction_operation, destination)
```

- **send\_value** : la valeur à envoyer par chaque processus
- **recv\_value** : la valeur reçue suite aux échanges et à la réduction par le destinataire seulement
- **Size** : nombre d'éléments (> 1 si tableau)
- **MPI\_data\_type** : le type de donnée (ex : **MPI\_INTEGER**)
- **MPI\_reduction\_operation** : type d'opération à effectuer pour la réduction (ex : **MPI\_SUM**)
- **Destination** : le processus qui va recevoir les données réduites



[https://www.open-mpi.org/doc/v4.0/man3/MPI\\_Reduce.3.php](https://www.open-mpi.org/doc/v4.0/man3/MPI_Reduce.3.php)

# Communication collective : réduction grâce à `MPI_Reduce` (C/C++)

*Collective communication*

• `MPI_Reduce` est appelée par les processus expéditeurs et destinataires en même temps



.cpp

```
MPI_Reduce(send_value, &recv_value, size, MPI_data_type, MPI_reduction_operation, destination)
```

- `send_value` (`const void *`) : la valeur à envoyer par chaque processus
- `recv_value` (`void *`) : la valeur reçue suite aux échanges et à la réduction par le destinataire
- `Size` : nombre d'éléments (>1 si tableau)
- `MPI_data_type` (`MPI_Datatype`) : le type de donnée (ex : `MPI_INT`)
- `MPI_reduction_operation` (`MPI_Op`) : type d'opération à effectuer pour la réduction (ex : `MPI_SUM`)
- `Destination` : le processus qui va recevoir les données réduites



[https://www.open-mpi.org/doc/v4.0/man3/MPI\\_Reduce.3.php](https://www.open-mpi.org/doc/v4.0/man3/MPI_Reduce.3.php)

# Communication collective : opération de réduction

## *Collective communication*

Il existe de multiple opérations de réduction disponibles (`MPI_reduction_operation`) :

- `MPI_SUM` : Somme l'ensemble des données
- `MPI_PROD` : multiplication des données
- `MPI_MAX` : maximum des valeurs
- `MPI_MIN` : minimum des valeurs
- ...



[https://www.open-mpi.org/doc/v4.0/man3/MPI\\_Reduce.3.php](https://www.open-mpi.org/doc/v4.0/man3/MPI_Reduce.3.php)

# Communication collective : exemple d'utilisation de MPI\_REDUCE

## *Collective communication*

- Réduction d'une simple variable réelle double précision



.f90

```
Real(8) :: rank_value
Real(8) :: reduction_value
Integer :: ierror
```

```
rank_value = rank
```

```
! Addition de l'ensemble des rank_value dans le processus 0
```

```
Call MPI_REDUCE(rank_value, reduction_value, 1 &
    MPI_DOUBLE_PRECISION, MPI_SUM, 0, &
    MPI_COMM_WORLD, ierror)
```

# Communication collective : exemple d'utilisation de MPI\_Reduce (C/C++)

*Collective communication*

- Réduction d'une simple variable réelle double précision



.cpp

```
double rank_value, reduction_value ;  
int ierror ;  
  
rank_value = rank ;  
  
// Addition de l'ensemble des rank_value dans le processus 0  
ierror = MPI_Reduce(rank_value, &reduction_value, 1  
                    MPI_DOUBLE, MPI_SUM, 0,  
                    MPI_COMM_WORLD) ;
```

# Communication collective : exemple d'utilisation de MPI\_REDUCE

## *Collective communication*

- Réduction d'un tableau d'entiers avec multiplication de toutes les valeurs



.f90

```
integer, dimension(4) :: rank_value  
integer, dimension(4) :: reduction_value  
Integer :: ierror
```

```
rank_value = (/ 18, 22, 43, 52/)
```

```
! Addition de l'ensemble des rank_value dans le processus 0
```

```
Call MPI_REDUCE(rank_value, reduction_value, 4, &  
               MPI_INTEGER, MPI_PROD, 0, &  
               MPI_COMM_WORLD, ierror)
```

# Communication collective : exemple d'utilisation de MPI\_Reduce (C/C++)

*Collective communication*

- Réduction d'un tableau d'entiers avec multiplication de toutes les valeurs



.cpp

```
int rank_value[4] ;
int reduction_value[4] ;
int ierror ;

rank_value = { 18, 22, 43, 52 } ;

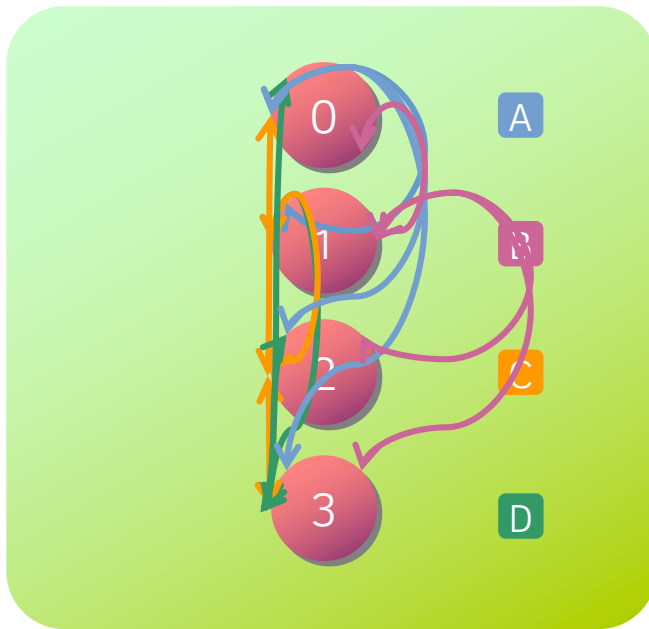
// Addition de l'ensemble des rank_value dans le processus 0

ierror = MPI_Reduce(rank_value, &reduction_value[0], 4
    MPI_INT, MPI_PROD, 0,
    MPI_COMM_WORLD) ;
```

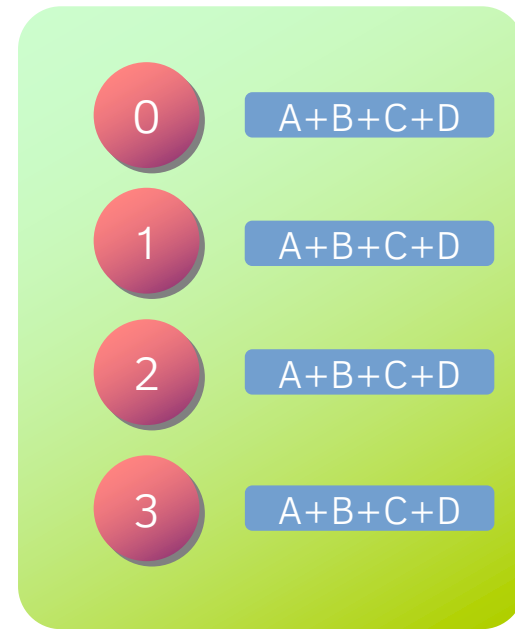
# Communication collective : réduction grâce à `MPI_Allreduce`

*Collective communication*

- Envoi de données réparties sur plusieurs processus **vers tous les processus** avec une **opération**



`MPI_ALLREDUCE + MPI_SUM`



Nouvelle répartition des données et addition



[https://www.open-mpi.org/doc/v4.0/man3/MPI\\_Allreduce.3.php](https://www.open-mpi.org/doc/v4.0/man3/MPI_Allreduce.3.php)



# Communication collective : réduction grâce à `MPI_ALLREDUCE` (Fortran90)

## *Collective communication*

- `MPI_ALLREDUCE` est appelée par les processus expéditeurs et destinataires en même temps



.f90

```
MPI_ALLREDUCE(send_value, recv_value, size, MPI_data_type, MPI_reduction_operation, comm)
```

- `send_value` : la valeur à envoyer par chaque processus
- `recv_value` : la valeur reçue suite aux échanges et à la réduction par le destinataire seulement
- `Size` : nombre d'éléments (> 1 si tableau)
- `MPI_data_type` : le type de donnée (ex : `MPI_INTEGER`)
- `MPI_reduction_operation` : type d'opération à effectuer pour la réduction (ex : `MPI_SUM`)



[https://www.open-mpi.org/doc/v4.0/man3/MPI\\_Allreduce.3.php](https://www.open-mpi.org/doc/v4.0/man3/MPI_Allreduce.3.php)

# Communication collective : réduction grâce à MPI\_Allreduce (C/C++)

## *Collective communication*

• **MPI\_Allreduce** est appelée par les processus expéditeurs et destinataires en même temps



.f90

```
MPI_Allreduce(send_value, &recv_value, size, MPI_data_type, MPI_reduction_operation, commu
```

- **send\_value** : la valeur à envoyer par chaque processus
- **recv\_value** : la valeur reçue suite aux échanges et à la réduction par le destinataire seulement
- **Size** : nombre d'éléments (> 1 si tableau)
- **MPI\_data\_type** : le type de donnée (ex : MPI\_INTEGER)
- **MPI\_reduction\_operation** : type d'opération à effectuer pour la réduction (ex : MPI\_SUM)



[https://www.open-mpi.org/doc/v4.0/man3/MPI\\_Allreduce.3.php](https://www.open-mpi.org/doc/v4.0/man3/MPI_Allreduce.3.php)

## Exercice n°5 : Utilisation de la communication collective MPI\_REDUCE



•Rendez vous dans le dossier de l'exercice n°5 appelé 5\_reduce\_com



```
> cd exercises/mpi/5_reduce_com
```

•Ouvrez les instructions contenues dans le fichier **README.md** avec votre éditeur de fichier favori

# Introduction au parallélisme par échange d

## 4) Les communications collectives

### 4.2) Les collectives avec des données de taille va

# Communication collective pour les données de taille variable



Certaines communications collectives (hors réduction) présentées précédemment imposent

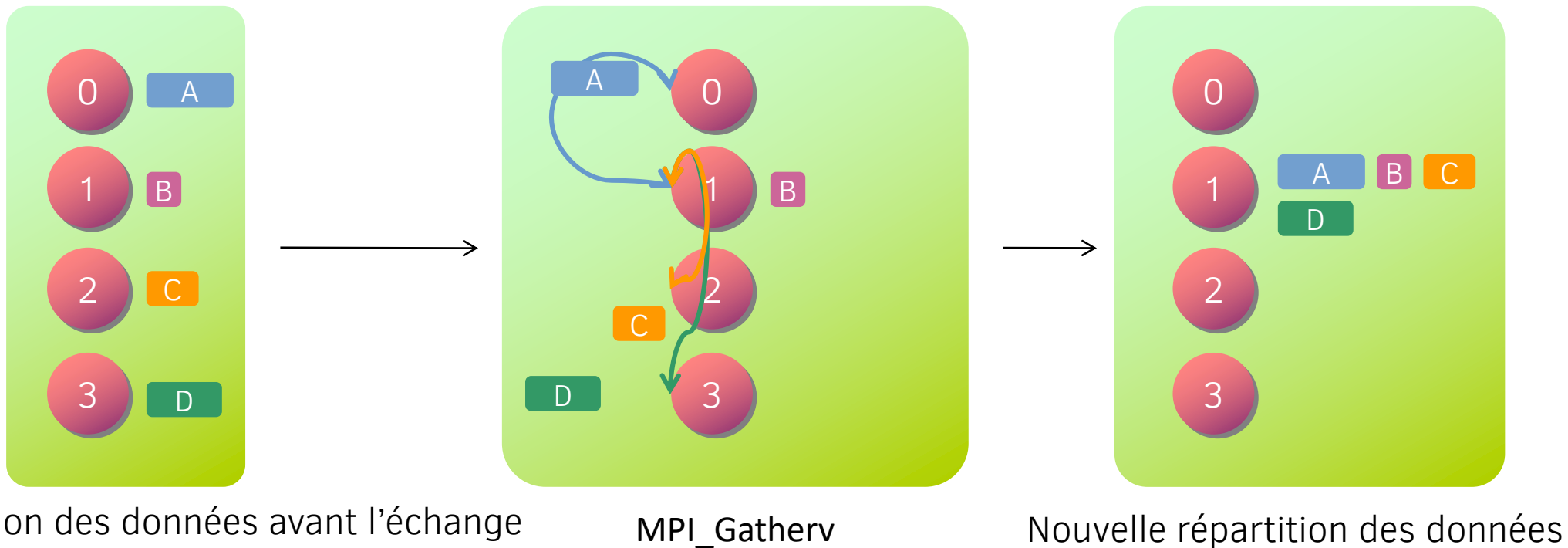
Elles possèdent le même nom suivi d'un v à la fin :

- MPI\_Gatherv
- MPI\_Allgatherv
- MPI\_Alltoallv
- MPI\_Scatterv

# Communication collective : collecte grâce à MPI\_Gatherv

*Collective communication*

- Envoi de données de taille différente réparties sur plusieurs processus vers un processus unique



[https://www.open-mpi.org/doc/v4.0/man3/MPI\\_Gatherv.3.php](https://www.open-mpi.org/doc/v4.0/man3/MPI_Gatherv.3.php)

# Communication collective pour les données de taille variable



`MPI_Gatherv` peut également être utilisée pour changer l'ordre des données une fois ra

# Communication collective : collecte grâce à MPI\_Gatherv (Fortran95)

## *Collective communication*

• **MPI\_Gatherv** est appelée par les processus expéditeurs et destinataires en même temps



.f90

```
MPI_GATHERV(send_array, send_count, send_type, &  
            recv_array, recv_count, displacement, recv_type, &  
            destination, communicator, ierror)
```

- **send\_array** : la valeur ou un ensemble de valeur
- **send\_count** : le nombre de valeur à envoyer, ce nombre peut être différent sur chaque processus
- **send\_type** : type MPI des valeurs envoyées
- **recv\_array** : le tableau réunissant les valeurs reçues
- **recv\_count** (tableau) : nombre d'éléments reçus de chaque rang
- **displacement** (tableau) : où placer chaque contribution dans **recv\_array**
- **recv\_type** : le type des données reçues
- **destination** : le processus qui reçoit les données



[https://www.open-mpi.org/doc/v4.0/man3/MPI\\_Gatherv.3.php](https://www.open-mpi.org/doc/v4.0/man3/MPI_Gatherv.3.php)



# Communication collective : collecte grâce à MPI\_Gatherv (C/C++)

## *Collective communication*

• **MPI\_Gatherv** est appelée par les processus expéditeurs et destinataires en même temps



.cpp

```
MPI_Gatherv(send_array, send_count, send_type,  
            recv_array, recv_count, displacement, recv_type,  
            destination, communicator) ;
```

- **send\_array** (const void \*) : la valeur ou un ensemble de valeur
- **send\_count** (int) : le nombre de valeur à envoyer
- **send\_type** (MPI\_Datatype) : type MPI des valeurs envoyées
- **recv\_array** (void \*) : le tableau réunissant les valeurs reçues
- **recv\_count** (const int \*) : nombre d'éléments reçus pour chaque rang
- **displacement** (const int \*) : où placer chaque contribution dans **recv\_array**
- **recv\_type** (MPI\_Datatype) : le type des données reçues
- **Destination** (int) : le processus qui reçoit les données

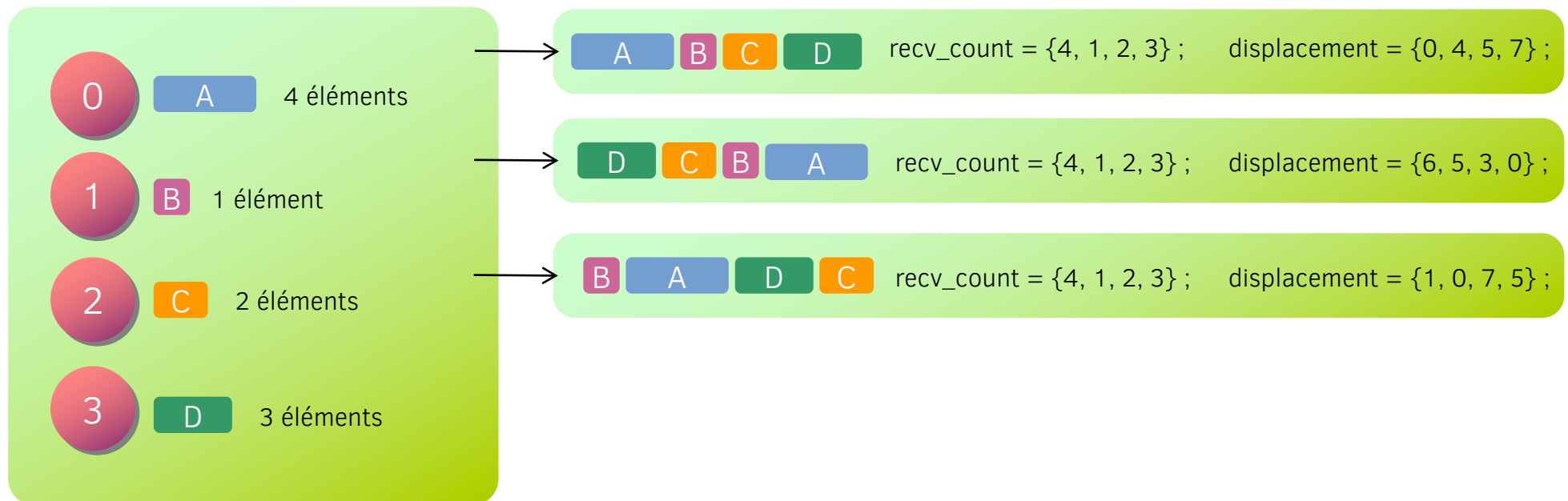


[https://www.open-mpi.org/doc/v4.0/man3/MPI\\_Gatherv.3.php](https://www.open-mpi.org/doc/v4.0/man3/MPI_Gatherv.3.php)

# Communication collective : notion de déplacement

*Collective communication*

• Le tableau **displacement** permet de **placer les contribution de chaque rang dans le tableau qui**



[https://www.open-mpi.org/doc/v4.0/man3/MPI\\_Gatherv.3.php](https://www.open-mpi.org/doc/v4.0/man3/MPI_Gatherv.3.php)

# Fonctions supplémentaires



Toutes les communications collectives présentées ont également un équivalent non-bloc

`.MPI_Igather`

`.MPI_Igatherv`

`.MPI_Iscatter`

`.MPI_Ibcast`

`.MPI_Ialltoall`

`.MPI_Ireduce`

`...`

D'autres variantes de communications collectives sont à découvrir dans le cours de l'IDF

## Exercice n°6 : Utilisation de la communication collective MPI\_GATHER



•Rendez vous dans le dossier de l'exercice n°6 appelé 6\_gather\_com



```
> cd exercises/mpi/6_gather_com
```

•Ouvrez les instructions contenues dans le fichier **README.md** avec votre éditeur de fichier favori

# Communications collectives MPI

A ce stade du cours, vous savez maintenant :

- Effectuer des communications collectives
- Effectuer des réductions

# Introduction au parallélisme par échange d

## 5) Topologie cartésienne

# Décomposition de domaine cartésienne

En calcul scientifique, il est courant de décomposer le domaine d'étude (grille, matrice) en sous-domaines.

Sur grille régulière et structurée, une approche simple et classique consiste à diviser le domaine en sous-domaines.

› Méthode de décomposition cartésienne

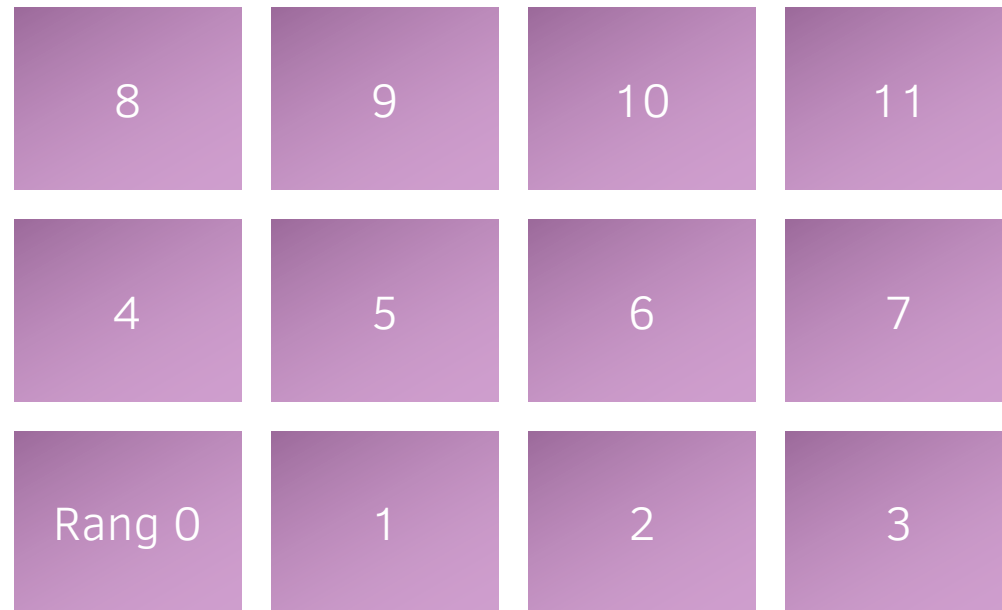


Il s'agit d'un exemple typique de parallélisme de donnée

# Décomposition de domaine cartésienne



Décomposition 2D

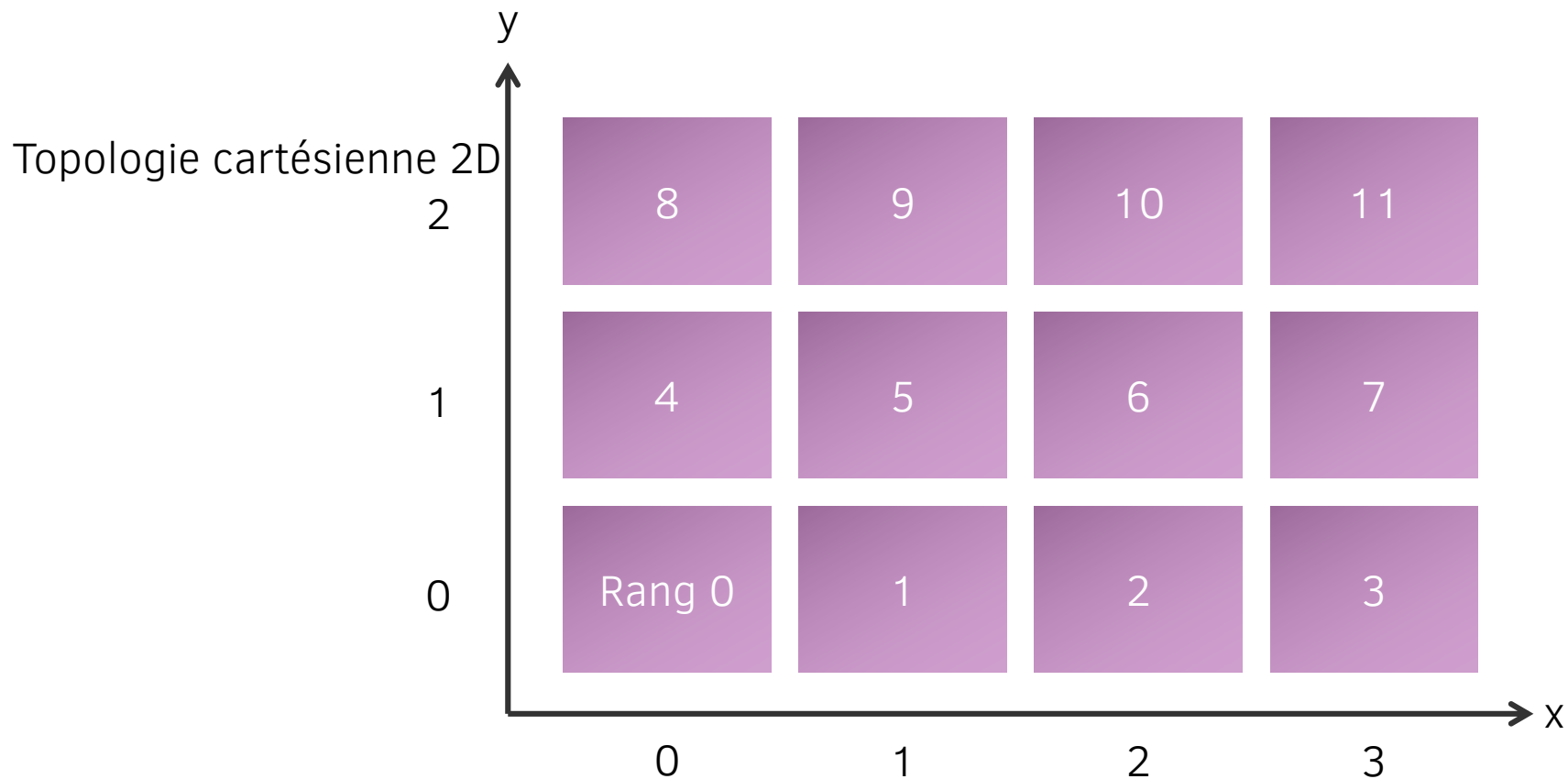




# Décomposition de domaine cartésienne : coordonnées

Une **topologie** cartésienne a besoin de :

- **Coordonnées** pour situer les processus (bloc) dans l'espace cartésien
- De **rangs** pour chaque processus en adéquation avec la topologie cartésienne
- exemple : le rang 5 a pour coordonnées (1,1)



# Décomposition de domaine cartésienne : création



Deux solutions pour mettre en place une topologie cartésienne :

- Le faire à la main
- Faire appel aux fonctions MPI conçues pour ça

# Décomposition de domaine cartésienne : création via MPI\_Cart\_create (F90)

**MPI\_CART\_CREATE** permet de définir une topologie cartésienne à partir d'un ancien communicateur



.f90

```
MPI_CART_CREATE(old_communicator, dimension,  
                ranks_per_direction, periodicity,  
                reorganisation, cartesian_communicator, ierror)
```

- **old\_communicator** : ancien communicateur (MPI\_COMM\_WORLD par exemple)
- **dimension** (entier) : dimension de la topologie (2 pour 2D par exemple)
- **ranks\_per\_direction** (tableau d'entier) : le nombre de rangs dans chaque dimension
- **Periodicity** (tableau de booléens) : permet de définir les directions périodiques (true)
- **Reorganisation** (booléen) : réorganisation des rangs pour optimiser les échanges (true)
- **cartesian\_communicator** (entier) : nouveau communicateur renvoyé par la fonction qui vient remplir



[https://www.open-mpi.org/doc/v4.0/man3/MPI\\_Cart\\_create.3.php](https://www.open-mpi.org/doc/v4.0/man3/MPI_Cart_create.3.php)

# Décomposition de domaine cartésienne : création via MPI\_Cart\_create (C)

**MPI\_Cart\_create** permet de définir une topologie cartésienne à partir d'un ancien communicateur



.cpp

```
MPI_Cart_create(old_communicator, dimension,  
               ranks_per_adirection, periodicity,  
               reorganisation, cartesian_communicator) ;
```

- **old\_communicator** (MPI\_Comm) : ancien communicateur (MPI\_COMM\_WORLD par exemple)
- **dimension** (int) : dimension de la topologie (2 pour 2D par exemple)
- **ranks\_per\_direction** (int \*) : le nombre de rangs dans chaque dimension
- **Periodicity** (int \*) : permet de définir les directions périodiques
- **Reorganisation** (int) : réorganisation des rangs pour optimiser les échanges
- **cartesian\_communicator** (MPI\_Comm \*) : nouveau communicateur renvoyé par la fonction qui vient d'être appelée



[https://www.open-mpi.org/doc/v4.0/man3/MPI\\_Cart\\_create.3.php](https://www.open-mpi.org/doc/v4.0/man3/MPI_Cart_create.3.php)

# Décomposition de domaine cartésienne : création



Comme pour n'importe quel communicateur, on peut récupérer les rangs dans le con

# Décomposition de domaine cartésienne : récupérer les coordonnées d'un

• **MPI\_CART\_COORDS** permet de récupérer les coordonnées d'un rang donné dans la topologie cartésienne



.f90

```
CALL MPI_CART_COORDS(cartesian_communicator, rank, dimension, &  
rank_coordinates, ierror)
```

- Dimension (entier) : dimension de la topologie (2 pour 2D par exemple)
- rank\_coordinates (tableau d'entier) : les coordonnées du rang rank dans cartesian\_communicator



[https://www.open-mpi.org/doc/v2.0/man3/MPI\\_Cart\\_coords.3.php](https://www.open-mpi.org/doc/v2.0/man3/MPI_Cart_coords.3.php)

# Décomposition de domaine cartésienne : récupérer les coordonnées d'un

• **MPI\_Cart\_coords** permet de récupérer les coordonnées d'un rang donné dans la topologie cartésienne



.cpp

```
MPI_Cart_coords(cartesian_communicator, rank, dimension,  
                rank_coordinates) ;
```

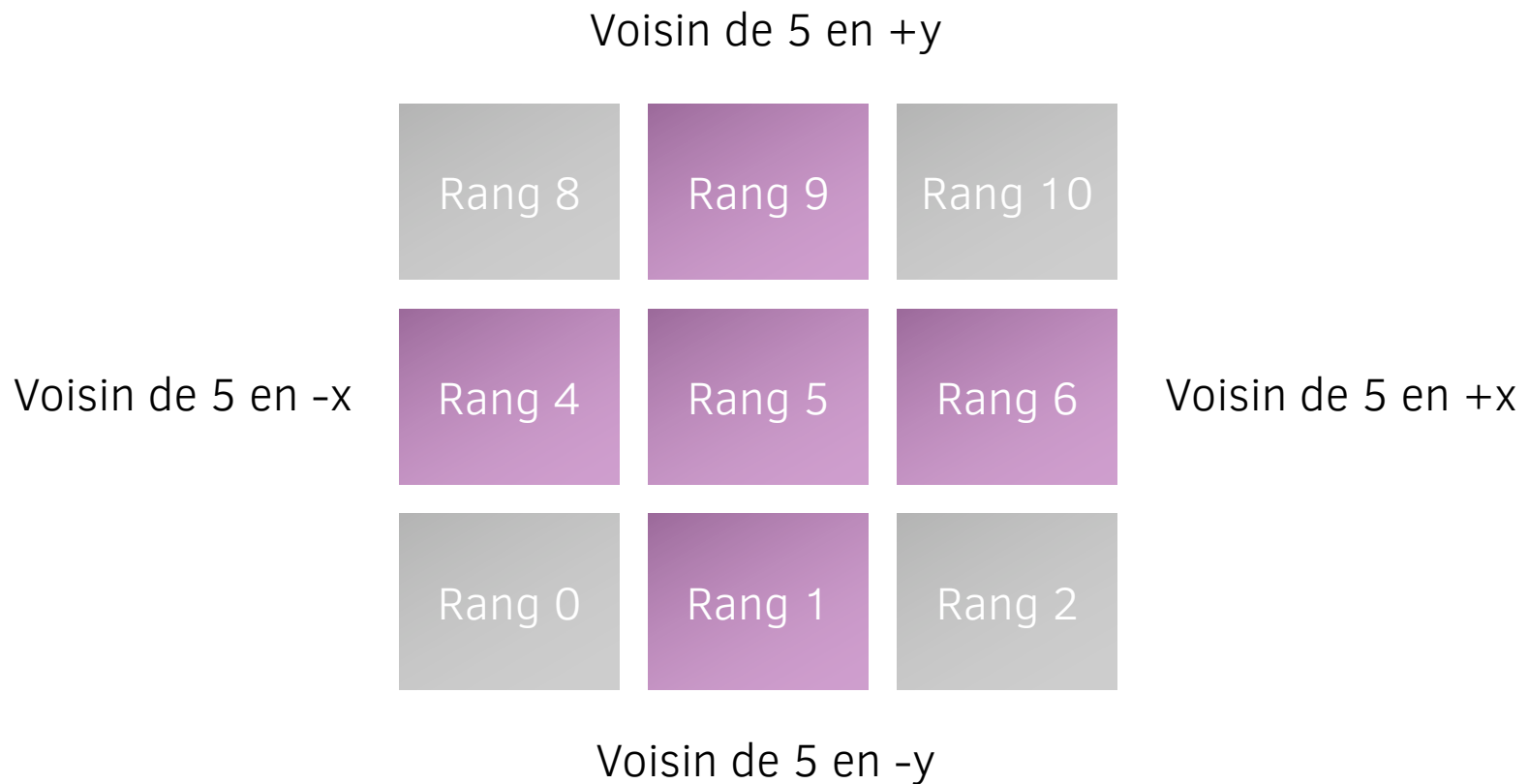
- Rank (int) : rang du processus MPI
- Dimension (int) : dimension de la topologie (2 pour 2D par exemple)
- ranks\_coordinates (int \*) : les coordonnées du rang rank dans cartesian\_communicator



[https://www.open-mpi.org/doc/v2.0/man3/MPI\\_Cart\\_coords.3.php](https://www.open-mpi.org/doc/v2.0/man3/MPI_Cart_coords.3.php)

# Décomposition de domaine cartésienne : les voisins

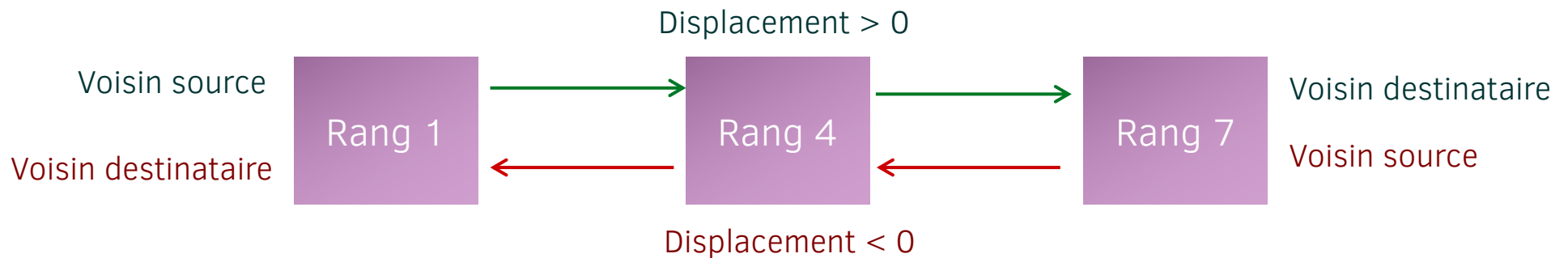
- Chaque processus doit être en mesure de récupérer **le rang de ses voisins** dans la topologie cartésienne





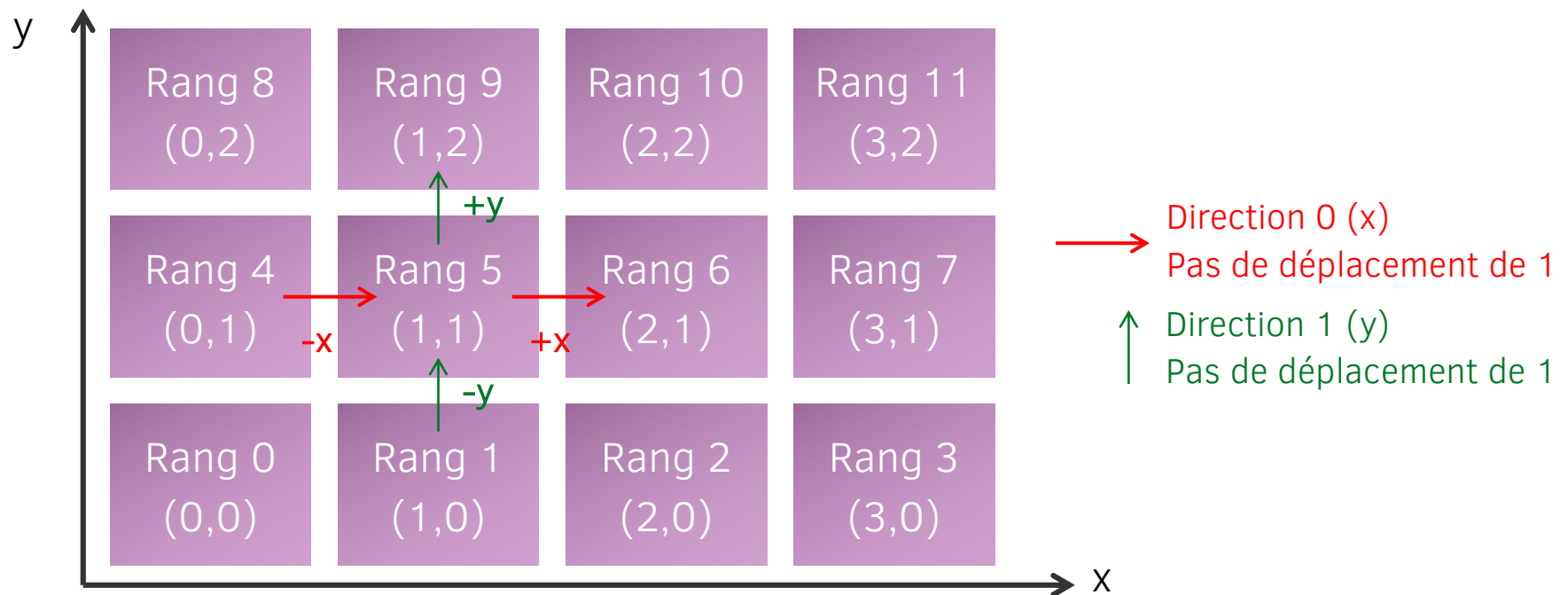
# Décomposition de domaine cartésienne : récupérer les rangs des voisins

**MPI\_CART\_SHIFT** permet de récupérer les rangs voisins d'un rang donné en spécifiant une direction

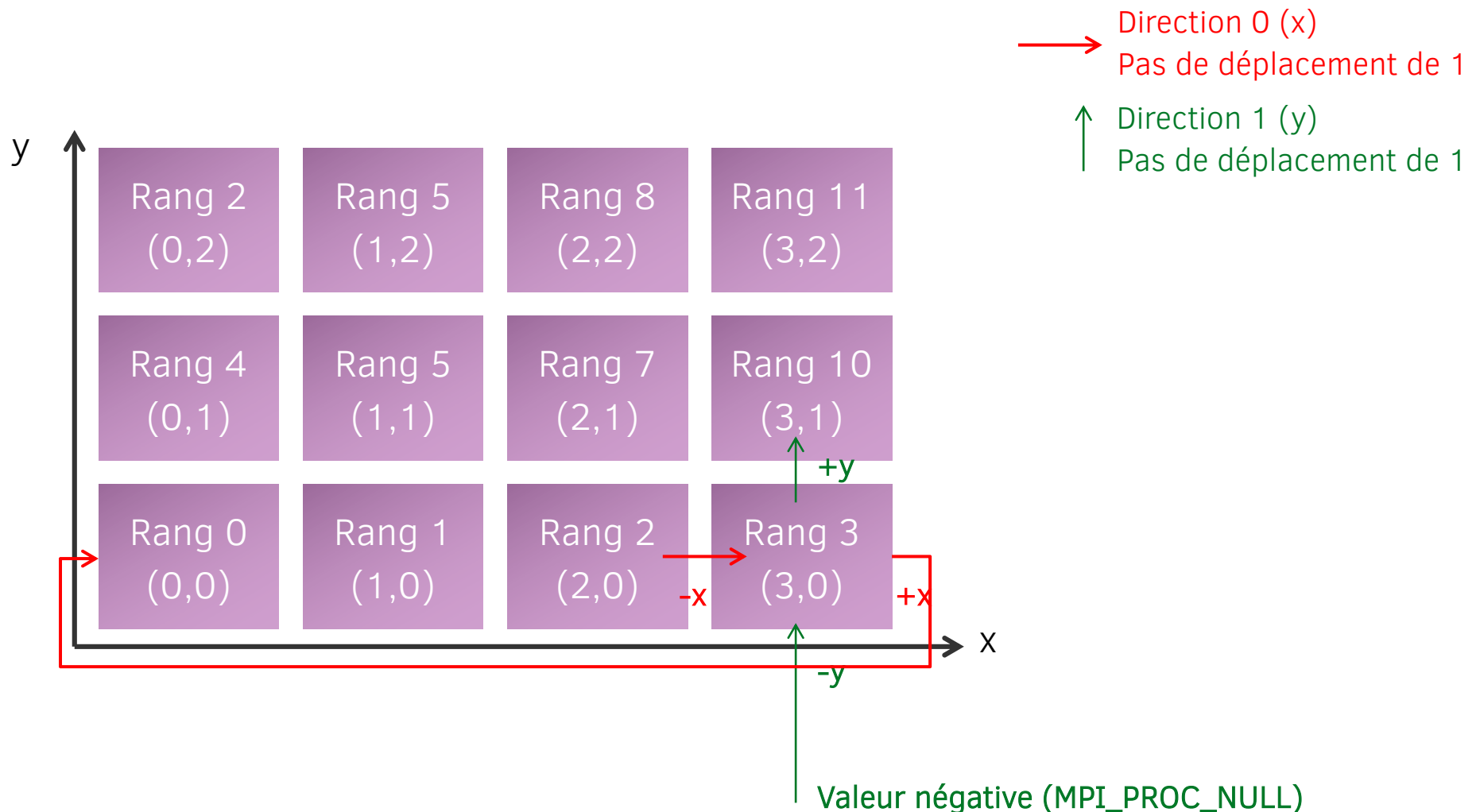


[https://www.open-mpi.org/doc/v2.0/man3/MPI\\_Cart\\_shift.3.php](https://www.open-mpi.org/doc/v2.0/man3/MPI_Cart_shift.3.php)

# Exemple de topologie cartésienne 2D



# Exemple de topologie cartésienne 2D : notion de périodicité



# Décomposition de domaine cartésienne : récupérer les rangs des voisins

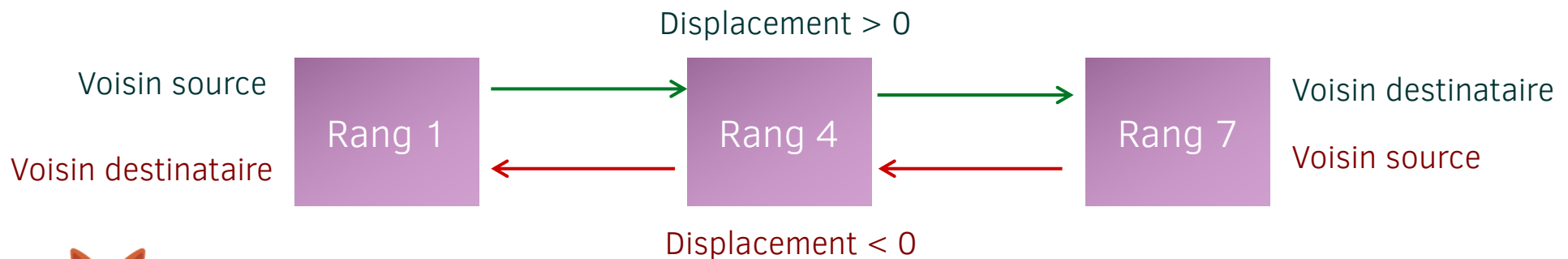
**MPI\_CART\_SHIFT** permet de récupérer les rangs voisins d'un rang donné



.f90

```
MPI_CART_SHIFT(cartesian_communicator, direction,  
displacement,  
rank_neighbors_src, rank_neighbors_dest, ierror)
```

- **direction** (integer) : 1 pour la première coordonnée, 2 pour la deuxième coordonnée
- **Displacement** (integer) : pas de déplacement dans la direction souhaitée, si  $> 0$  déplacement vers le voisin source, si  $< 0$  déplacement vers le voisin destinataire
- **rank\_neighbord\_source** : si  $\text{displacement} > 0$ , correspond au voisin source
- **rank\_neighbord\_dest** : si  $\text{displacement} < 0$ , correspond au voisin destinataire



[https://www.open-mpi.org/doc/v2.0/man3/MPI\\_Cart\\_shift.3.php](https://www.open-mpi.org/doc/v2.0/man3/MPI_Cart_shift.3.php)

# Décomposition de domaine cartésienne : récupérer les rangs des voisins

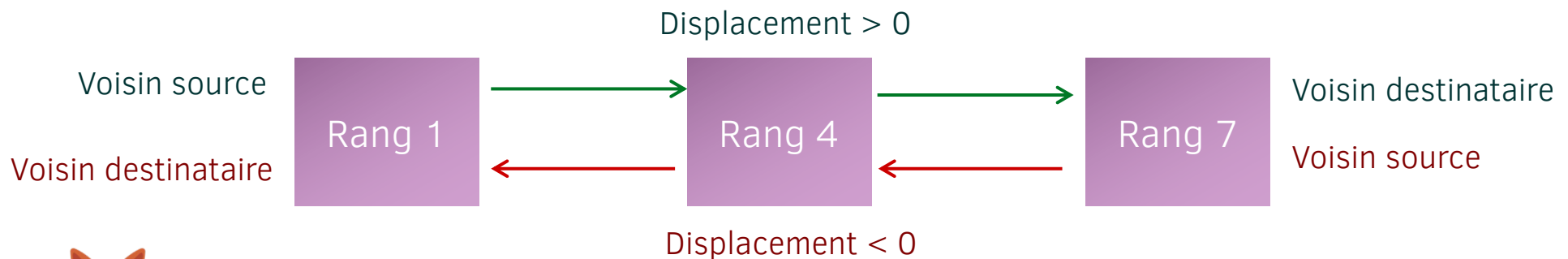
**MPI\_Cart\_shift** permet de récupérer les rangs voisins d'un rang donné



.cpp

```
MPI_Cart_shift(cartesian_communicator, direction,  
displacement,  
rank_neighbors_src, rank_neighbors_dest) ;
```

- **direction** (int) : 1 pour la première coordonnée, 2 pour la deuxième coordonnée
- **Displacement** (int) : pas de déplacement dans la direction souhaitée, si  $> 0$  déplacement vers les
- **rank\_neighbord\_source** : si  $\text{displacement} > 0$ , correspond au voisin source
- **rank\_neighbord\_dest** : si  $\text{displacement} > 0$ , correspond au voisin destinataire



[https://www.open-mpi.org/doc/v2.0/man3/MPI\\_Cart\\_shift.3.php](https://www.open-mpi.org/doc/v2.0/man3/MPI_Cart_shift.3.php)

# Décomposition de domaine cartésienne : récupérer les rangs des voisins



Lorsqu'un rang n'a pas de voisin (par exemple en non-périodique), `MPI_CART_SHIFT` re

Lorsqu'une communication a pour destinataire ou expéditeur `MPI_PROC_NULL`, elle pe



[https://www.open-mpi.org/doc/v2.0/man3/MPI\\_Cart\\_shift.3.php](https://www.open-mpi.org/doc/v2.0/man3/MPI_Cart_shift.3.php)

# Exemple de topologie cartésienne 2D (Fortran95)



.f90

```
Integer, dimension(2) :: ranks_per_direction = (/4, 3/)
Logical, dimension(2) :: periodicity = (/ .true., .true. /)
Logical                :: reorganisation = .true.
Integer                :: cartesian_communicator

Integer, dimension(2) :: rank_coordinates
Integer                :: rank_neighbors_mx, rank_neighbors_px
Integer                :: rank_neighbors_my, rank_neighbors_py

Call MPI_INIT(ierr)

Call MPI_CART_CREATE(MPI_COMM_WORLD, 2, ranks_per_direction, periodicity,
                    reorganisation, cartesian_communicator, ierr)

Call MPI_COMM_RANK(cartesian_communicator, rank, ierr)

Call MPI_CART_COORDS(cartesian_communicator, rank, 2, rank_coordinates, ierr)

CALL MPI_CART_SHIFT(cartesian_communicator, 1, 1, &
                    rank_neighbors_my, rank_neighbors_py, ierr)

CALL MPI_CART_SHIFT(cartesian_communicator, 0, 1, &
                    rank_neighbors_mx, rank_neighbors_px, ierr)
```

# Exemple de topologie cartésienne 2D (C/C++)



.cpp

```
int ranks_per_direction[2] = {4, 3};
int periodicity[2] = {1, 1};
int reorganisation = 1;
MPI_Comm cartesian_communicator;
int rank_coordinates[2];
int rank_neighbors_mx, rank_neighbors_px;
int rank_neighbors_my, rank_neighbors_py;

ierror = MPI_Init(&ierror)

ierror = MPI_Cart_create(MPI_COMM_WORLD, 2, ranks_per_direction, periodicity,
                        reorganisation, &cartesian_communicator);

ierror = MPI_Comm_rank(cartesian_communicator, &rank);

ierror = MPI_Cart_coords(cartesian_communicator, rank, 2, &rank_coordinates);

ierror = MPI_Cart_shift(cartesian_communicator, 1, 1,
                        &rank_neighbors_my, &rank_neighbors_py);

ierror = MPI_Cart_shift(cartesian_communicator, 0, 1,
                        &rank_neighbors_mx, &rank_neighbors_px);
```



## Exercice n°7 : Création d'une topologie cartésienne



•Rendez vous dans le dossier de l'exercice n°7 appelé 7\_cartesian\_com



```
> cd exercises/mpi/7_cartesian_com
```

•Ouvrez les instructions contenues dans le fichier **README.md** avec votre éditeur de fichier favori

# Décomposition de domaine cartésienne

A ce stade du cours, vous savez maintenant :

- Créer un communicateur cartésien pour décomposer vos données
- Utiliser les fonctions liées à la décomposition cartésienne

# Introduction au parallélisme par échange d

## 6) Types dérivés

# Types dérivés

Les types dérivés permettent de décrire des données plus complexes que les types classiques

- **MPI\_Type\_contiguous** : permet de sélectionner une portion contiguë d'un tableau
- **MPI\_Type\_vector** : permet de créer un sous-tableau à partir de sous-ensemble d'éléments contigus
- **MPI\_Type\_indexed** : permet de créer un sous-tableau à partir de sous-ensemble d'éléments non contigus
- **MPI\_Type\_create\_struct** : permet de créer l'équivalent d'une structure C en mélangeant les types



Les types dérivés sont ensuite utilisés dans les communications à la place des types classiques

# Types dérivés : MPI\_Type\_contiguous

**MPI\_Type\_contiguous** : permet de sélectionner une portion contiguë d'un tableau existant

Tableau existant



Sous-tableau contigu

# Types dérivés : MPI\_Type\_contiguous

La fonction `MPI_Type_contiguous` contient les arguments suivants :



.cpp

```
int MPI_Type_contiguous(int count, MPI_Datatype oldtype,  
                        MPI_Datatype *newtype)
```

- `count` (int) : nombre d'éléments qui compose le nouveau type
- `oldtype` (MPI\_Datatype) : le type de donnée qui compose le tableau initial
- `newtype` (MPI\_Datatype) : notre nouveau type dérivé

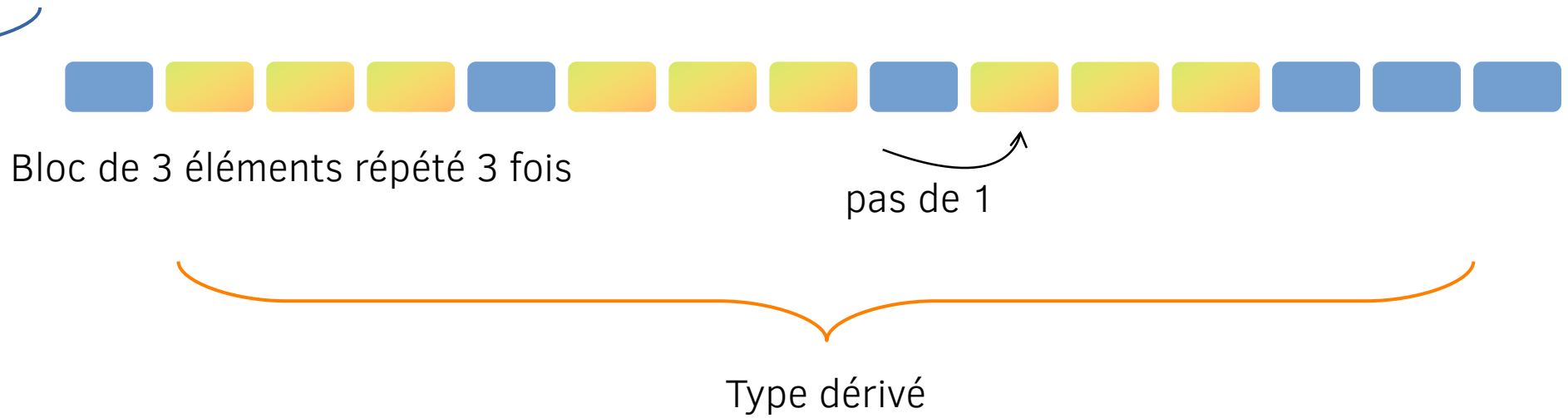


[https://www.open-mpi.org/doc/current/man3/MPI\\_Type\\_contiguous.3.php](https://www.open-mpi.org/doc/current/man3/MPI_Type_contiguous.3.php)

# Types dérivés : MPI\_Type\_vector

**MPI\_Type\_vector** : permet de créer un sous-tableau à partir de sous-ensemble d'éléments co

Tableau existant



# Types dérivés : MPI\_Type\_vector

La fonction `MPI_Type_vector` contient les arguments suivants :



.cpp

```
int MPI_Type_vector(int count, int blocklength, int stride,  
MPI_Datatype oldtype, MPI_Datatype *newtype)
```

- `count (int)` : nombre de blocs
- `Blocklength (int)` : la taille en nombre d'éléments de chaque bloc
- `Stride (int)` : distance (pas) entre chaque bloc
- `oldtype (MPI_Datatype)` : le type de donnée qui compose le tableau initial
- `newtype (MPI_Datatype)` : notre nouveau type dérivé



[https://www.rookiehpc.com/mpi/docs/mpi\\_type\\_vector.php](https://www.rookiehpc.com/mpi/docs/mpi_type_vector.php)



# Création d'un type dérivé

La fonction `MPI_Type_commit` permet officialiser la création du type :



Bloc de 3 éléments répété 3 fois

pas de 1



.cpp

```
MPI_Datatype column_type;  
  
MPI_Type_vector(3, 3, 1, MPI_INT, &column_type);  
  
MPI_Type_commit(&column_type);
```



[https://www.rookiehpc.com/mpi/docs/mpi\\_type\\_commit.php](https://www.rookiehpc.com/mpi/docs/mpi_type_commit.php)

# Exemple complet d'utilisation d'un type dérivé

Dans cette exemple, le rang 0 envoie de l'information au rang 1 à partir du tableau buffer et d'un



.cpp

```
If (rank == 0) {  
    MPI_Datatype vector_type;  
    MPI_Type_vector(3, 3, 1, MPI_INT, &vector_type);  
    MPI_Type_commit(&vector_type);  
  
    int buffer[12];  
  
    MPI_Send(&buffer[1], 1, vector_type, 1, 0, MPI_COMM_WORLD);  
}  
  
If (rank == 1) {  
    int received[9];  
    MPI_Recv(&received, 9, MPI_INT, 0, 0, MPI_COMM_WORLD, MPI_STATUS_IGNORE);  
}
```

## Exercice n°8 : Utilisation du type dérivé vector



•Rendez vous dans le dossier de l'exercice n°8 appelé 8\_type\_vector



```
> cd exercises/mpi/8_type_vector
```

•Ouvrez les instructions contenues dans le fichier **README.md** avec votre éditeur de fichier favori