



# IDRIS

## MPI

Dimitri Lecas - Rémi Lacroix - Serge Van Criekingen - Myriam Peyrounette

*CNRS — IDRIS*

v5.3 13 mars 2025



# Plan I

## Introduction

- A propos
- Introduction
- Concepts de l'échange de messages
- Mémoire distribuée
- Historique
- Bibliographie

## Environnement

### Communications point à point

- Notions générales
- Opérations d'envoi et de réception bloquantes
- Types de données de base
- Autres possibilités

### Communications collectives

- Notions générales
- Synchronisation globale : `MPI_Barrier()`
- Diffusion générale : `MPI_Bcast()`
- Diffusion sélective : `MPI_Scatter()`
- Collecte : `MPI_Gather()`
- Collecte générale : `MPI_Allgather()`
- Collecte : `MPI_Gatherv()`

## Plan II

- Collectes et diffusions sélectives : `MPI_Alltoall()`
- Réductions réparties
- Compléments

## Modèles de communication

- Modes d'envoi point à point
- Appels bloquants
  - Envois synchrones
  - Envois *bufferisés*
  - Envois standards
- Nombre d'éléments reçus
- Appels non bloquants
- Communications mémoire à mémoire (RMA)

## Types de données dérivés

- Introduction
- Types contigus
- Types avec un pas constant
- Validation des types de données dérivés
- Exemples
  - Type « colonne d'une matrice »
  - Type « ligne d'une matrice »
  - Type « bloc d'une matrice »
- Types homogènes à pas variable

## Plan III

- Taille des types de données
- Types hétérogènes
- Conclusion
- Memento

## Communicateurs

- Introduction
- Exemple
- Communicateur par défaut
- Groupes et communicateurs
- Partitionnement d'un communicateur
- Topologies
  - Topologies cartésiennes
  - Subdiviser une topologie cartésienne

## MPI-IO

- Introduction
- Ouverture et fermeture d'un fichier
- Lectures/écritures : généralités
- Lectures/écritures individuelles
  - Via des déplacements explicites
  - Via des déplacements implicites individuels
  - Via des déplacements implicites partagés
- Lectures/écritures collectives

## Plan IV

- Via des déplacements explicites
- Via des déplacements implicites individuels
- Via des déplacements implicites partagés

Positionnement explicite des pointeurs dans un fichier

Lectures/écritures non bloquantes

- Via des déplacements explicites
- Via des déplacements implicites individuels
- Lectures/écritures collectives et non bloquantes

## MPI 4.x

### MPI-IO Vues

- Définition des vues
- Construction de sous-tableaux
- Lecture d'un fichier par blocs de deux éléments
- Utilisation successive de plusieurs vues
- Gestion des trous dans les types de données
- Conseils

## Conclusion

# Introduction

# Introduction

## A propos

Ce document est mis à jour régulièrement. La version la plus récente est disponible sur le site Web de l'IDRIS : <http://www.idris.fr/formations/mpi/>

- IDRIS  
Institut du développement et des ressources en informatique scientifique  
Rue John Von Neumann  
Bâtiment 506  
BP 167  
91403 ORSAY CEDEX  
France  
<http://www.idris.fr>

# Introduction

## Parallélisme

L'intérêt de faire de la programmation parallèle est :

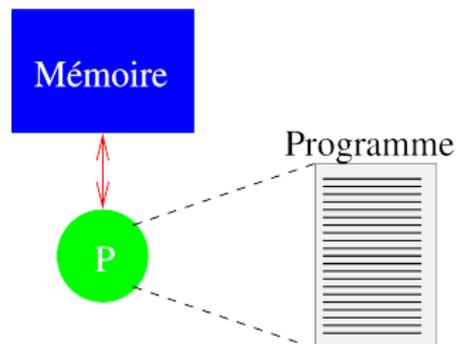
- De réduire le temps de restitution ;
- D'effectuer de plus gros calculs ;
- D'exploiter le parallélisme des processeurs modernes (multi-coeurs, multithreading).

Mais pour travailler à plusieurs, la coordination est nécessaire. [MPI](#) est une bibliothèque permettant de coordonner des processus en utilisant le paradigme de l'échange de messages.

# Introduction

## Modèle de programmation séquentiel

- le programme est exécuté par un et un seul processus ;
- toutes les variables et constantes du programme sont allouées dans la mémoire allouée au processus ;
- un processus s'exécute sur un processeur physique de la machine.

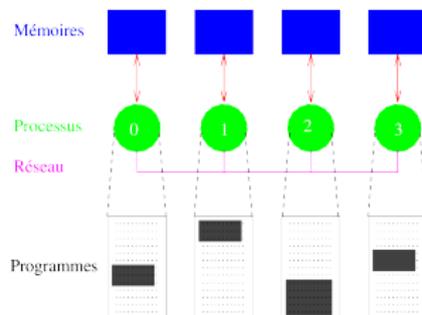


**Figure 1** – Modèle de programmation séquentiel

# Introduction

## Modèle de programmation par échange de messages

- le programme est écrit dans un langage classique ([Fortran](#), [C](#), [C++](#), etc.);
- toutes les variables du programme sont privées et résident dans la mémoire locale allouée à chaque processus ;
- chaque processus exécute éventuellement des parties différentes d'un programme ;
- une donnée est échangée entre deux ou plusieurs processus via un appel, dans le programme, à des sous-programmes particuliers.

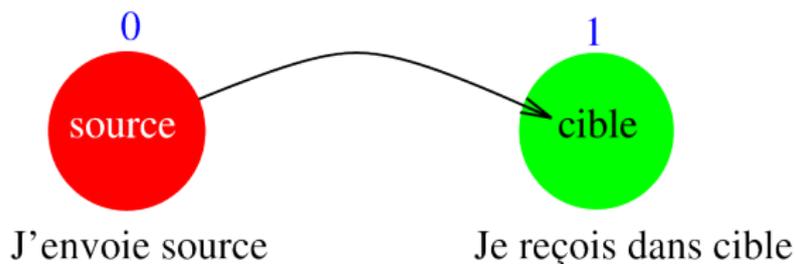


**Figure 2** – Modèle de programmation par échange de messages

# Introduction

## Concepts de l'échange de messages

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



**Figure 3** – échange d'un message

# Introduction

## Constitution d'un message

- 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.

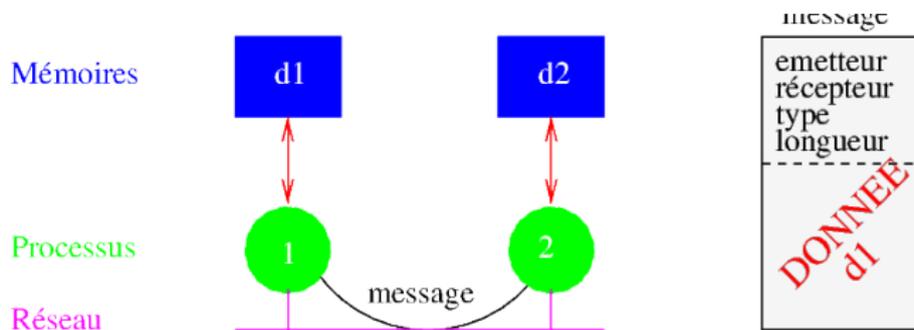


Figure 4 – Constitution d'un message

# Introduction

## Environnement

- Les messages échangés sont interprétés et gérés par un environnement qui peut être comparé à la téléphonie, au courrier postal, à la 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 (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

# Introduction

## Architecture des supercalculateurs

La plupart des supercalculateurs sont des machines à mémoire distribuée. Ils sont composés d'un ensemble de nœud, à l'intérieur d'un nœud la mémoire est partagée.

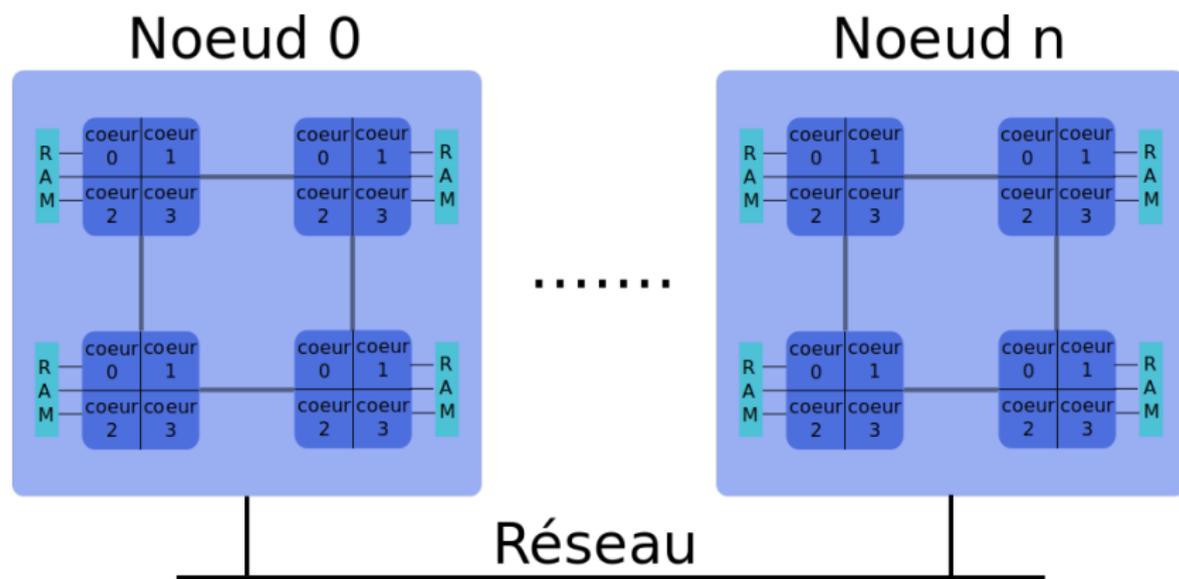


Figure 5 – Architecture des supercalculateurs

# Introduction

## Jean Zay

- 1 116 nœuds
- 2 processeurs Intel Cascade Lake (20 cœurs)  
à 2,5 Ghz par nœud
- 4 GPU Nvidia V100 par nœud (sur 391 nœud)
- 44 640 cœurs
- 214 To (192 Go par nœud)



# Introduction

## MPI vs OpenMP

OpenMP utilise un schéma à mémoire partagée, tandis que pour MPI la mémoire est distribuée.

Processus 0 ..... Processus n

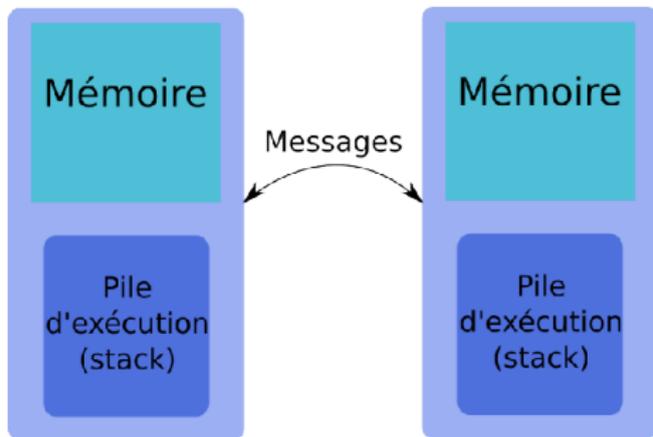


Figure 6 – Schéma MPI

Processus

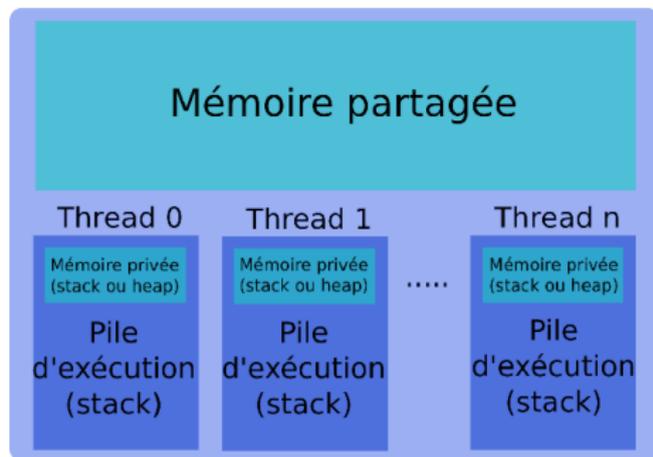
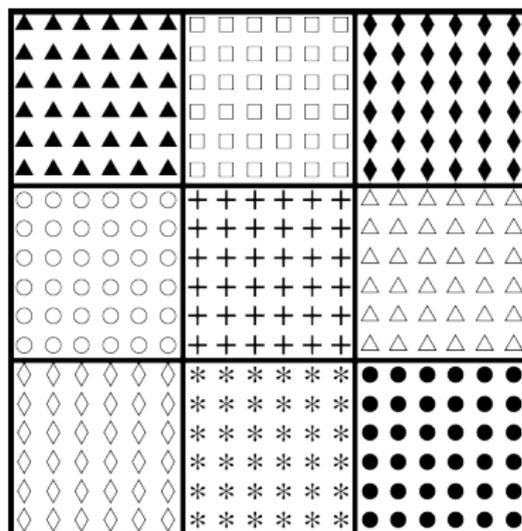


Figure 7 – Schéma OpenMP

# Introduction

## Décomposition de domaine

Un schéma que l'on rencontre très souvent avec **MPI** est la décomposition de domaine. Chaque processus possède une partie du domaine global, et effectue principalement des échanges avec ses processus voisins.



**Figure 8** – Découpage en sous-domaines

# Introduction

## Historique

- **Version 1.0** : en juin 1994, le forum MPI, avec la participation d'une quarantaine d'organisations, aboutit à la définition d'un ensemble de sous-programmes concernant la bibliothèque d'échanges de messages MPI
- **Version 1.1** : juin 1995, avec seulement des changements mineurs
- **Version 1.2** : en 1997, avec des changements mineurs pour une meilleure cohérence des dénominations de certains sous-programmes
- **Version 1.3** : septembre 2008, avec des clarifications dans MPI 1.2, en fonction des clarifications elles-mêmes apportées par MPI-2.1
- **Version 2.0** : apparue en juillet 1997, cette version apportait des compléments importants volontairement non intégrés dans MPI 1.0 (gestion dynamique de processus, copies mémoire à mémoire, entrées-sorties parallèles, etc.)
- **Version 2.1** : juin 2008, avec seulement des clarifications dans MPI 2.0 mais aucun changement
- **Version 2.2** : septembre 2009, avec seulement de petites additions
- **Version 3.0** : septembre 2012, cette version apportait les communications collectives non bloquantes, nouvelle interface Fortran, etc.
- **Version 3.1** : juin 2015, avec des corrections et des petites additions

# Introduction

## MPI 4.0

Version 4.0 : juin 2021

- Grand nombre
- Communication par morceaux
- MPI Session

Version 4.1 : novembre 2023

# Introduction

## Bibliographie

- Site du MPI Forum <http://www.mpi-forum.org>
- Normes disponible en PDF sur <http://www.mpi-forum.org/docs/>
- William Gropp, Ewing Lusk et Anthony Skjellum : *Using MPI, third edition Portable Parallel Programming with the Message-Passing Interface*, MIT Press, 2014.
- William Gropp, Torsten Hoefler, Rajeev Thakur et Erwing Lusk : *Using Advanced MPI Modern Features of the Message-Passing Interface*, MIT Press, 2014.
- Victor Eijkhout : The Art of HPC <http://theartofhpc.com>

# Introduction

## Implémentations MPI *open source*

Elles peuvent être installées sur un grand nombre d'architectures mais leurs performances sont en général en dessous de celles des implémentations constructeurs.

- **MPICH** : <http://www.mpich.org>
- **Open MPI** : <http://www.open-mpi.org>

# Introduction

## Outils

- Débogueurs
  - Totalview  
<https://totalview.io>
  - DDT  
<https://www.linaroforge.com/linaro-ddt>
- Outils de mesure de performances
  - FPMPI : *FPMPI*  
<http://www.mcs.anl.gov/research/projects/fpmapi/WWW/>
  - Scalasca : *Scalable Performance Analysis of Large-Scale Applications*  
<http://www.scalasca.org>
  - MUST : *MPI Runtime Correctness Analysis*  
<https://itc.rwth-aachen.de/must/>

# Introduction

## Bibliothèques scientifiques parallèles *open source*

- **ScaLAPACK** : résolution de problèmes d'algèbre linéaire par des méthodes directes.  
<http://www.netlib.org/scalapack/>
- **PETSc** : résolution de problèmes d'algèbre linéaire et non-linéaire par des méthodes itératives.  
<https://petsc.org/release/>
- **PaStiX** : résolution de grands systèmes linéaires creux.  
<https://solverstack.gitlabpages.inria.fr/pastix/>
- **FFTW** : transformées de Fourier rapides.  
<http://www.fftw.org>
- **HDF5** : Lecture et écriture sur fichiers.  
<https://www.hdfgroup.org/solutions/hdf5/>

# Environnement

# Environnement

## Description

- Toute unité de programme appelant des fonctions MPI doit inclure le fichier d'en-têtes `mpi.h`.
- Le sous-programme `MPI_Init()` permet d'initialiser l'environnement nécessaire :

```
int MPI_Init(int *argc, char ***argv)
```

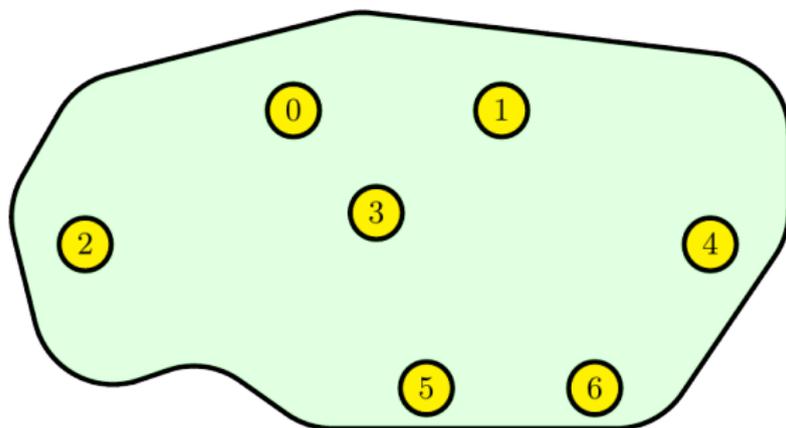
- Réciproquement, le sous-programme `MPI_Finalize()` désactive cet environnement :

```
int MPI_Finalize(void)
```

# Environnement

## Communicateurs

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



**Figure 9** – Communicateur `MPI_COMM_WORLD`

## Arrêt d'un programme

Parfois un programme se trouve dans une situation où il doit s'arrêter sans attendre la fin normale. C'est typiquement le cas si un des processus ne peut pas allouer la mémoire nécessaire à son calcul. Dans ce cas il faut utiliser le sous-programme `MPI_Abort()` et non l'instruction Fortran `stop` (Ou `exit` in C).

```
int MPI_Abort(MPI_Comm comm, int erreur)
```

- `comm` : tous les processus appartenant à ce communicateur seront stoppés, il est donc conseillé d'utiliser `MPI_COMM_WORLD` ;
- `erreur` : numéro d'erreur retourné à l'environnement UNIX.

## Code

Il n'est pas nécessaire de tester la valeur de `code` (valeur de retour en C) après des appels aux routines MPI. Par défaut, lorsque MPI rencontre un problème, le programme s'arrête comme lors d'un appel à `MPI_Abort()`.

## Rang et nombre de processus

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

```
int MPI_Comm_size(MPI_Comm comm, int *nb_procs)
```

- 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`) :

```
int MPI_Comm_rank(MPI_Comm comm, int *rang)
```

# Environnement

## Exemple

```
1 /* qui_je_suis */
2 #include <mpi.h>
3 #include <stdio.h>
4
5 int main(int argc, char *argv[]) {
6     int nb_procs, rang;
7
8     MPI_Init(&argc, &argv);
9
10    MPI_Comm_size(MPI_COMM_WORLD, &nb_procs);
11    MPI_Comm_rank(MPI_COMM_WORLD, &rang);
12
13    printf("Je suis le processus %d parmi %d\n", rang, nb_procs);
14
15    MPI_Finalize();
16 }
```

```
> mpiexec -n 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
```

## Compilation et exécution d'un code MPI

- Pour **compiler** un code MPI, on utilise un enrobeur (*wrapper*) de compilateur qui fait le lien avec la librairie MPI utilisée.
- Cet enrobeur diffère selon le langage de programmation, le compilateur et la librairie MPI utilisés. Par exemple : `mpif90`, `mpifort`, `mpicc`, ...

```
> mpicc <options> -c source.c  
> mpicc source.o -o mon_executable
```

- Pour **exécuter** un code MPI, on utilise un lanceur d'application MPI qui ordonne le lancement de l'exécution sur un nombre de processus choisi.
- Le lanceur défini par la norme MPI est `mpiexec`. Il existe également des lanceurs non standards, comme `mpirun`.

```
> mpiexec -n <nombre de processus> mon_executable
```

# Travaux pratiques MPI – Exercice 1 : Environnement MPI

- Implémenter un programme MPI dans lequel chaque processus affiche un message indiquant si son rang est **pair** ou **impair**. Par exemple :

```
> mpiexec -n 4 ./pair_impair
Moi, processus 0, je suis de rang pair
Moi, processus 2, je suis de rang pair
Moi, processus 3, je suis de rang impair
Moi, processus 1, je suis de rang impair
```

- Pour tester la parité, la fonction intrinsèque Fortran correspondant à l'opération *modulo* est **mod** :

```
mod(a, b)
```

(en C, utilisez le symbole **%** : `a%b`)

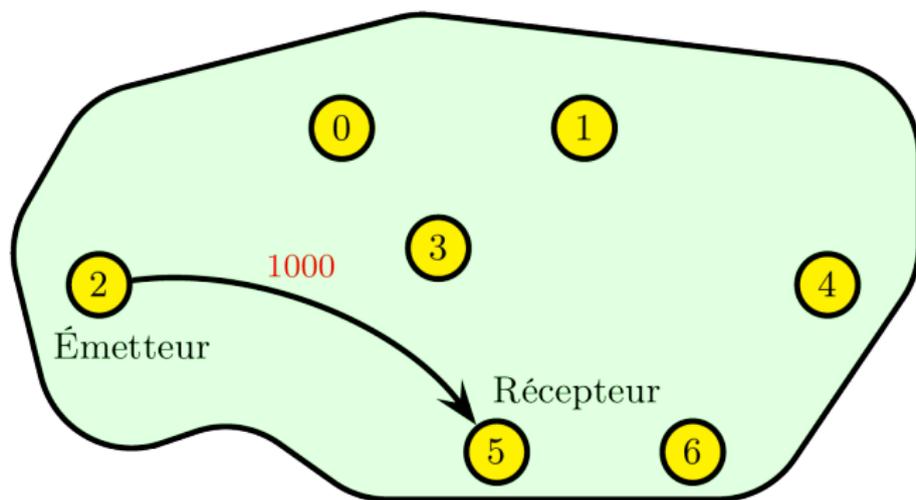
- Pour compiler votre programme, utilisez la commande **make**
- Pour exécuter votre programme, utilisez la commande **make exe**
- Pour être reconnu par le Makefile, le programme doit se nommer `pair_impair.f90` (ou `pair_impair.c`)

## Communications point à point

# Communications point à point

## Notions générales

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



**Figure 10** – Communication point à point

# Communications point à point

## Notions générales

- L'émetteur et le récepteur sont identifiés par leur **rang** dans le communicateur.
- L'entité transmise entre deux processus est appelée **message**.
- Un message est caractérisé par son **enveloppe**. Celle-ci est constituée :
  - du rang du processus émetteur ;
  - du rang du processus récepteur ;
  - de l'étiquette (*tag*) du message ;
  - du communicateur qui définit le groupe de processus et le contexte de communication.
- Les données échangées sont **typées** (entiers, réels, etc ou types dérivés personnels).
- Il existe dans chaque cas plusieurs **modes** de transfert, faisant appel à des protocoles différents.
- Si deux messages sont envoyés avec la même enveloppe, l'ordre de réception et d'envoi sont les mêmes.

# Communications point à point

## Opération d'envoi `MPI_Send`

```
int MPI_Send(const void *message, int longueur, MPI_Datatype type_message,  
            int rang_dest, int etiquette, MPI_Comm comm)
```

Envoi, à partir de l'adresse `message`, d'un message de taille `longueur`, de type `type_message`, étiqueté `etiquette`, au processus `rang_dest` dans le communicateur `comm`.

### Remarque :

Cette opération est bloquante : l'exécution reste bloquée jusqu'à ce que le contenu de `message` puisse être réécrit sans risque d'écraser la valeur qui devait être envoyée.

# Communications point à point

## Opération de réception `MPI_Recv`

```
int MPI_Recv(void *message, int longueur, MPI_Datatype type_message,  
            int rang_source, int etiquette, MPI_Comm comm, MPI_Status *status)
```

Réception, à partir de l'adresse `message`, d'un message de taille `longueur`, de type `type_message`, étiqueté `etiquette`, du processus `rang_source`.

### Remarques :

- `statut` stocke des informations sur la communication : `rang_source`, `etiquette`, `code`,...
- L'appel `MPI_Recv` ne pourra fonctionner avec une opération `MPI_Send` que si ces deux appels ont la même enveloppe (`rang_source`, `rang_dest`, `etiquette`, `comm`).
- Cette opération est bloquante : l'exécution reste bloquée jusqu'à ce que le contenu de `message` corresponde au message reçu.

# Communications point à point

## Exemple (voir Fig. 10)

```
1  /* point_a_point */
2  #include <mpi.h>
3  #include <stdio.h>
4
5  int main(int argc, char *argv[]) {
6      int rang, valeur;
7      int etiquette=100;
8      MPI_Status statut;
9
10     MPI_Init (&argc, &argv);
11
12     MPI_Comm_rank (MPI_COMM_WORLD, &rang);
13
14     if (rang == 2) {
15         valeur = 1000;
16         MPI_Send (&valeur, 1, MPI_INT, 5, etiquette, MPI_COMM_WORLD);
17     } else if (rang == 5) {
18         MPI_Recv (&valeur, 1, MPI_INT, 2, etiquette, MPI_COMM_WORLD, &statut);
19         printf("Moi, processus 5, ai recu %d du processus 2.\n", valeur);
20     }
21
22     MPI_Finalize();
23 }
```

```
> mpiexec -n 7 point_a_point
```

```
Moi, processus 5, ai recu 1000 du processus 2
```

# Communications point à point

## Types de données de base C

Type MPI	Type C
<code>MPI_CHAR</code>	signed char
<code>MPI_SHORT</code>	signed short
<code>MPI_INT</code>	signed int
<code>MPI_LONG</code>	signed long int
<code>MPI_UNSIGNED_CHAR</code>	unsigned char
<code>MPI_UNSIGNED_SHORT</code>	unsigned short
<code>MPI_UNSIGNED</code>	unsigned int
<code>MPI_UNSIGNED_LONG</code>	unsigned long int
<code>MPI_FLOAT</code>	float
<code>MPI_DOUBLE</code>	double
<code>MPI_LONG_DOUBLE</code>	long double
<code>MPI_BYTE</code>	

# Communications point à point

## Autres possibilités

- À la réception d'un message, le rang de l'émetteur et l'étiquette peuvent être des « *jokers* », respectivement `MPI_ANY_SOURCE` et `MPI_ANY_TAG`.
- Une communication impliquant le processus « fictif » de rang `MPI_PROC_NULL` n'a aucun effet.
- `MPI_STATUS_IGNORE` est une constante prédéfinie qui peut être utilisée à la place de la variable `statut`.
- On peut communiquer des structures de données plus complexes en créant ses propres types dérivés.
- Il existe d'autres opérations qui effectuent **simultanément** un envoi et une réception : `MPI_Sendrecv()` et `MPI_Sendrecv_replace()`.

# Communications point à point

## Opération d'envoi et de réception simultanés `MPI_Sendrecv`

```
int MPI_Sendrecv(const void *message_emis, int longueur_message_emis,
                 MPI_Datatype type_message_emis, int rang_dest, int etiq_message_emis,
                 void *message_recu, int longueur_message_recu,
                 MPI_Datatype type_message_recu, int rang_source, int etiq_message_recu,
                 MPI_Comm comm, MPI_Status *status)
```

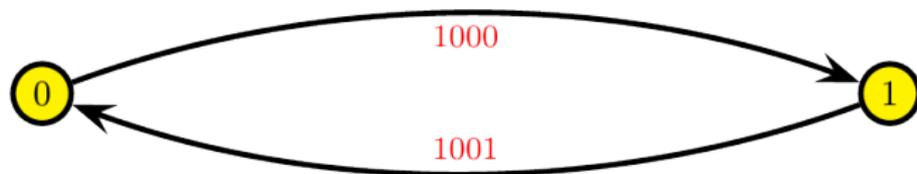
- Envoi, à partir de l'adresse `message_emis`, d'un message de taille `longueur_message_emis`, de type `type_message_emis`, étiqueté `etiq_message_emis`, au processus `rang_dest` dans le communicateur `comm` ;
- Réception, à partir de l'adresse `message_recu`, d'un message de taille `longueur_message_recu`, de type `type_message_recu`, étiqueté `etiq_message_recu`, du processus `rang_source` dans le communicateur `comm`.

### Remarque :

- La zone de réception `message_recu` doit différer de la zone d'envoi `message_emis`.

# Communications point à point

Opération d'envoi et de réception simultanés `MPI_Sendrecv`



**Figure 11** – Communication `sendrecv` entre les processus 0 et 1

# Communications point à point

## Exemple (voir Fig. 11)

```
1  /* sendrecv */
2  #include <mpi.h>
3  #include <stdio.h>
4
5  int main(int argc, char *argv[]) {
6      int rang, valeur, num_proc, message;
7      int etiquette=110;
8
9      MPI_Init(&argc, &argv);
10     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
11
12     num_proc=(rang+1)%2;
13     message = rang+1000;
14     MPI_Sendrecv(&message, 1, MPI_INT, num_proc, etiquette, &valeur, 1, MPI_INT,
15                 num_proc, etiquette, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
16
17     printf("Moi, processus %d, ai recu %d du processus %d.\n", rang, valeur, num_proc);
18
19     MPI_Finalize();
20 }
```

```
> mpiexec -n 2 sendrecv
```

```
Moi, processus 1, ai recu 1000 du processus 0
Moi, processus 0, ai recu 1001 du processus 1
```

# Communications point à point

## Attention !

Dans le cas d'une implémentation **synchrone** de `MPI_Send()`, l'exemple précédent serait en situation de verrouillage si l'appel à `MPI_Sendrecv()` était remplacé par un `MPI_Send()` suivi d'un `MPI_Recv()`. En effet, chacun des deux processus attendrait un ordre de réception qui ne viendrait jamais, puisque les deux envois resteraient en suspens.

```
val = rang+1000;
MPI_Send(&val,1,MPI_INT,num_proc,etiquette,MPI_COMM_WORLD);
MPI_Recv(valeur,1,MPI_INT,num_proc,etiquette,MPI_COMM_WORLD,&statut);
```

# Communications point à point

## Opération d'envoi et de réception simultanés `MPI_Sendrecv_replace`

```
int MPI_Sendrecv_replace(void * message, int longueur, MPI_Datatype type_message,
                        int rang_dest, int etiq_message_emis,
                        int rang_source, int etiq_message_recu, MPI_Comm comm,
                        MPI_Status *statut)
```

- Envoi, à partir de l'adresse `message`, d'un message de taille `longueur`, de type `type_message`, étiqueté `etiq_message_emis`, au processus `rang_dest` dans le communicateur `comm` ;
- Réception d'un message à la même adresse, d'une taille et d'un type identique, étiqueté `etiq_message_recu`, du processus `rang_source` dans le communicateur `comm`.

### Remarque :

- Contrairement à l'usage imposée par `MPI_Sendrecv()`, la zone de réception coïncide ici avec la zone d'envoi `message`.

# Communications point à point

## Exemple

```
1  /* joker */
2  #include <mpi.h>
3  #include <stdio.h>
4
5  int main(int argc, char *argv[]) {
6      int nb_procs, rang;
7      int m=4, etiquette=11;
8      int A[m][m];
9      MPI_Status statut;
10
11     MPI_Init(&argc, &argv);
12     MPI_Comm_size(MPI_COMM_WORLD, &nb_procs);
13     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
14
15     if (rang == 0) {
16         A[0][0]=1; A[0][1]=2; A[0][2]=3; A[0][3]=4; A[1][0]=5; A[1][1]=6; A[1][2]=7;
17         A[1][3]=8; A[2][0]=9; A[2][1]=10; A[2][2]=11; A[2][3]=12; A[3][0]=13;
18         A[3][1]=14; A[3][2]=15; A[3][3]=16;
19         MPI_Send(A, 3, MPI_INT, 1, etiquette, MPI_COMM_WORLD);
20     } else {
21         MPI_Recv(&(A[0][1]), 3, MPI_INT, MPI_ANY_SOURCE, MPI_ANY_TAG,
22                 MPI_COMM_WORLD, &statut);
23         printf("Moi processus %d, ai reçu 3 elements du processus %d avec comme "
24               "etiquette %d les elements sont %d %d %d.\n",
25               rang, statut.MPI_SOURCE, statut.MPI_TAG, A[0][1], A[0][2], A[0][3]);
26     }
27     MPI_Finalize();
28 }
```

# Communications point à point

```
> mpiexec -n 2 joker  
Moi processus 1, ai reçu 3 elements du processus 0  
avec comme etiquette 11 les elements sont 1 2 3.
```

## Travaux pratiques MPI – Exercice 2 : Ping-pong

- Communications point à point : *Ping-Pong* entre deux processus
- L'exercice 2 est décomposé en 3 étapes :
  1. *Ping* : compléter le script `ping_pong_1.c` de manière à ce que le processus de rang 0 **envoie** un message contenant une série aléatoire de 1000 réels au rang 1.
  2. *Ping-Pong* : compléter le script `ping_pong_2.c` de manière à ce que le processus de rang 1 **renvoie** le message vers le processus de rang 0, et mesurer le temps pris par la communication à l'aide de la fonction `MPI_Wtime()`.
  3. *Match de Ping-Pong* : compléter le script `ping_pong_3.c` de manière à enchaîner 9 *Ping-Pong*, **en faisant varier la taille du message**, et mesurer les temps pris par chaque échange. Les débits correspondants seront affichés.

## Travaux pratiques MPI – Exercice 2 : Ping-pong

### Remarques :

- Pour compiler la première étape : `make ping_pong_1`
- Pour exécuter la première étape : `make exe1`
- Pour compiler la seconde étape : `make ping_pong_2`
- Pour exécuter la seconde étape : `make exe2`
- Pour compiler la dernière étape : `make ping_pong_3`
- Pour exécuter la dernière étape : `make exe3`
  
- La génération de nombres réels pseudo-aléatoires uniformément répartis dans l'intervalle  $[0,1[$  se fait en C par un appel au sous-programme `rand` :

```
rand() / (RAND_MAX+1.);
```

- Les mesures de temps peuvent s'effectuer de la façon suivante :

```
temps_debut=MPI_Wtime();  
.....  
temps_fin=MPI_Wtime();  
printf("... en %f secondes.\n",temps_fin-temps_debut);
```

## Communications collectives

# Communications collectives

## Notions générales

- Les communications **collectives** permettent de faire en une seule opération une série de communications point à point.
- Une communication collective concerne toujours **tous** 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).
- 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.

# Communications collectives

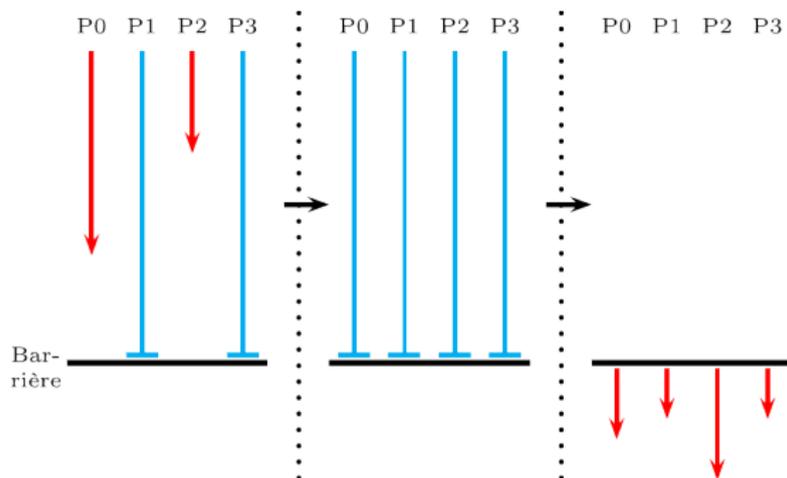
## Types de communications collectives

Il y a trois types de sous-programmes :

1. celui qui assure les synchronisations globales : `MPI_Barrier()`.
2. ceux 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()` ;
  - collecte et diffusion sélective, par tous les processus, de données réparties : `MPI_Alltoall()`.
3. ceux qui, en plus de la gestion des communications, effectuent des opérations sur les données transférées :
  - opérations de réduction (somme, produit, maximum, minimum, etc.), qu'elles soient d'un type prédéfini ou d'un type personnel : `MPI_Reduce()` ;
  - opérations de réduction avec diffusion du résultat (équivalent à un `MPI_Reduce()` suivi d'un `MPI_Bcast()`) : `MPI_Allreduce()`.

# Communications collectives

## Synchronisation globale : `MPI_Barrier()`



**Figure 12** – Synchronisation globale : `MPI_Barrier()`

```
int MPI_Barrier(MPI_Comm comm)
```

# Communications collectives

Diffusion générale : `MPI_Bcast()`

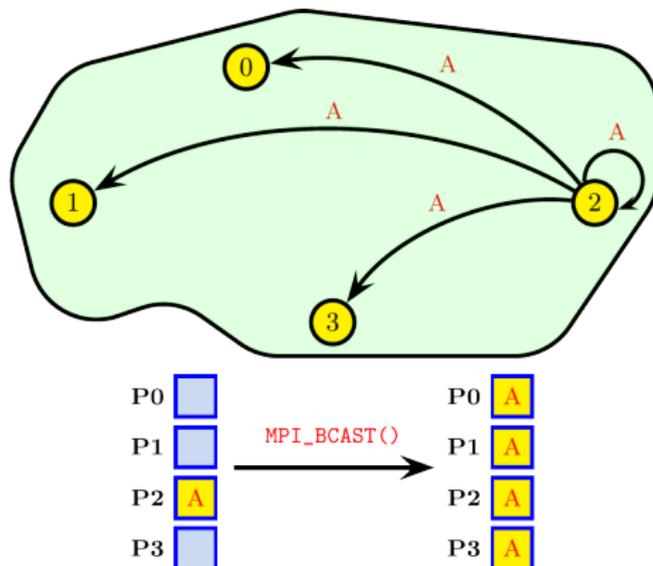


Figure 13 – Diffusion générale : `MPI_Bcast()`

## Diffusion générale : `MPI_Bcast()`

```
int MPI_Bcast(void *message, int longueur, MPI_Datatype type_message,  
             int rang_source, MPI_Comm comm)
```

1. Envoi, à partir de l'adresse `message`, d'un message constitué de `longueur` élément de type `type_message`, par le processus `rang_source`, à tous les autres processus du communicateur `comm`.
2. Réception de ce message à l'adresse `message` pour les processus autre que `rang_source`.

# Communications collectives

## Exemple de `MPI_Bcast()`

```
1  /* bcast */
2  #include <mpi.h>
3  #include <stdio.h>
4
5  int main(int argc, char *argv[]) {
6      int rang, valeur;
7
8      MPI_Init(&argc, &argv);
9      MPI_Comm_rank(MPI_COMM_WORLD, &rang);
10
11     if (rang == 2) valeur = rang+1000;
12
13     MPI_Bcast(&valeur, 1, MPI_INT, 2, MPI_COMM_WORLD);
14
15     printf("Moi, processus %d, ai recu %d du processus 2\n",
16           rang, valeur);
17
18     MPI_Finalize();
19 }
```

```
> mpiexec -n 4 bcast
```

```
Moi, processus 2, ai recu 1002 du processus 2
Moi, processus 0, ai recu 1002 du processus 2
Moi, processus 1, ai recu 1002 du processus 2
Moi, processus 3, ai recu 1002 du processus 2
```

# Communications collectives

Diffusion sélective : `MPI_Scatter()`

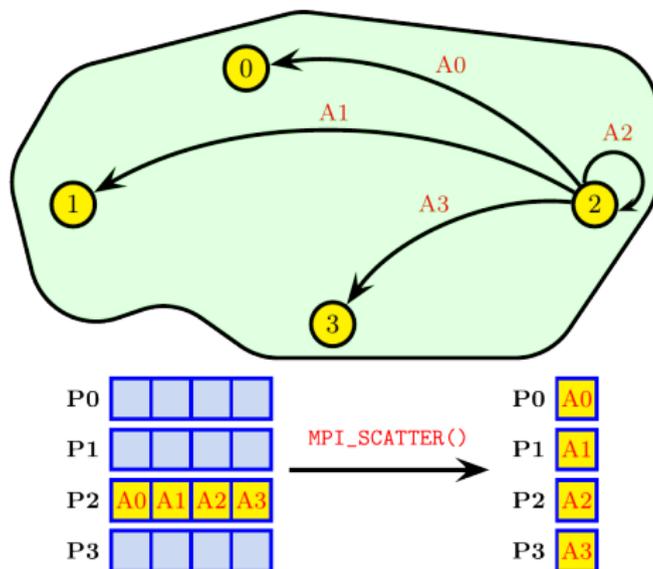


Figure 14 – Diffusion sélective : `MPI_Scatter()`

## Diffusion sélective : `MPI_Scatter()`

```
int MPI_Scatter(const void *message_a_repartir, int longueur_message_emis,  
               MPI_Datatype type_message_emis, void *message_recu,  
               int longueur_message_recu, MPI_Datatype type_message_recu,  
               int rang_source, MPI_Comm comm)
```

1. Distribution, par le processus `rang_source`, à partir de l'adresse `message_a_repartir`, d'un message de taille `longueur_message_emis`, de type `type_message_emis`, à tous les processus du communicateur `comm` ;
2. réception du message à l'adresse `message_recu`, de longueur `longueur_message_recu` et de type `type_message_recu` par tous les processus du communicateur `comm`.

### Remarques :

- Les couples (`longueur_message_emis`, `type_message_emis`) et (`longueur_message_recu`, `type_message_recu`) doivent être tels que les quantités de données envoyées et reçues soient égales.
- Les données sont distribuées en tranches égales, une tranche étant constituée de `longueur_message_emis` éléments du type `type_message_emis`.
- La *i*ème tranche est envoyée au *i*ème processus.

# Communications collectives

## Exemple de MPI\_Scatter()

```
1  /* scatter */
2  #include <mpi.h>
3  #include <stdio.h>
4  #include <stdlib.h>
5  int main(int argc, char *argv[]) {
6      int nb_valeurs=8, rang, nb_procs, longueur_tranche, i;
7      float *valeurs, *donnees;
8      MPI_Init(&argc, &argv);
9      MPI_Comm_size(MPI_COMM_WORLD, &nb_procs);
10     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
11     longueur_tranche = nb_valeurs/nb_procs;
12     donnees = (float *) malloc(longueur_tranche*sizeof(float));
13     if (rang == 2) {
14         valeurs = (float *) malloc(nb_valeurs*sizeof(float));
15         for (i=0; i<nb_valeurs;i++) valeurs[i]=1001.+i;
16         printf("Moi, processus %d envoie mon tableau valeurs : ", rang);
17         for (i=0; i<nb_valeurs;i++) {printf("%f ", valeurs[i]);} printf("\n"); }
18     MPI_Scatter(valeurs, longueur_tranche, MPI_FLOAT,
19               donnees, longueur_tranche, MPI_FLOAT, 2, MPI_COMM_WORLD);
20     printf("Moi, processus %d, ai recu ", rang);
21     for (i=0; i<longueur_tranche; i++) printf("%f ", donnees[i]);
22     printf("du processus 2\n");
23     MPI_Finalize(); }
```

```
> mpiexec -n 4 scatter
Moi, processus 2 envoie mon tableau valeurs :
1001. 1002. 1003. 1004. 1005. 1006. 1007. 1008.

Moi, processus 0, ai recu 1001. 1002. du processus 2
Moi, processus 1, ai recu 1003. 1004. du processus 2
Moi, processus 3, ai recu 1007. 1008. du processus 2
Moi, processus 2, ai recu 1005. 1006. du processus 2
```

# Communications collectives

Collecte : `MPI_Gather()`

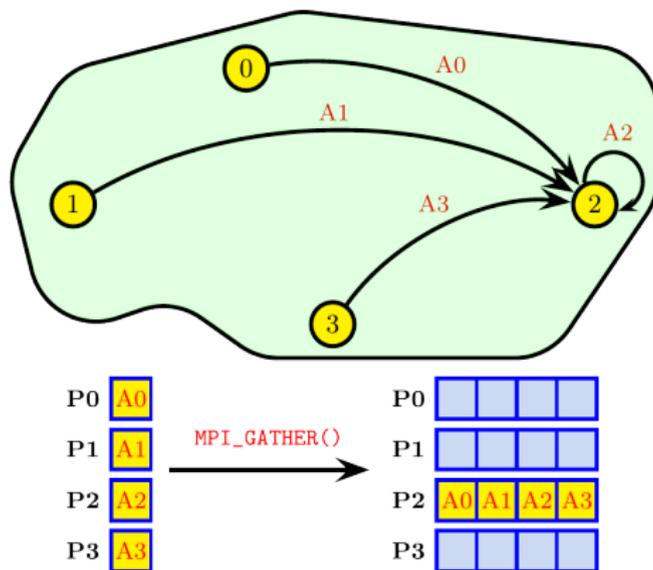


Figure 15 – Collecte : `MPI_Gather()`

## Collecte : `MPI_Gather()`

```
int MPI_Gather(const void*message_emis, int longueur_message_emis,  
             MPI_Datatype type_message_emis, void *message_recu,  
             int longueur_message_recu, MPI_Datatype type_message_recu,  
             int rang_dest, MPI_Comm comm)
```

1. Envoi de chacun des processus du communicateur `comm`, d'un message `message_emis`, de taille `longueur_message_emis` et de type `type_message_emis`.
2. Collecte de chacun de ces messages, par le processus `rang_dest`, à partir l'adresse `message_recu`, sur une longueur `longueur_message_recu` et avec le type `type_message_recu`.

### Remarques :

- Les couples (`longueur_message_emis`, `type_message_emis`) et (`longueur_message_recu`, `type_message_recu`) doivent être tels que les quantités de données envoyées et reçues soient égales.
- Les données sont collectées dans l'ordre des rangs des processus.

# Communications collectives

## Collecte : MPI\_Gather()

```
1  /* Gather */
2  #include <mpi.h>
3  #include <stdio.h>
4  #include <stdlib.h>
5  int main(int argc, char *argv[]) {
6      int nb_valeurs=8, rang, nb_procs, longueur_tranche, i;
7      float donnees[nb_valeurs], *valeurs;
8      MPI_Init(&argc, &argv);
9      MPI_Comm_size(MPI_COMM_WORLD, &nb_procs);
10     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
11     longueur_tranche = nb_valeurs/nb_procs;
12     valeurs = (float *) malloc(longueur_tranche*sizeof(float));
13     for (i=0; i<longueur_tranche; i++) valeurs[i]=1001.+rang*longueur_tranche+i;
14     printf("Moi, processus %d envoie mon tableau valeurs : ", rang);
15     for (i=0; i<longueur_tranche; i++) {printf("%f ", valeurs[i]);} printf("\n");
16     MPI_Gather(valeurs, longueur_tranche, MPI_FLOAT,
17              donnees, longueur_tranche, MPI_FLOAT, 2, MPI_COMM_WORLD);
18     if (rang==2) {
19         printf("Moi, processus %d, ai recu ", rang);
20         for (i=0; i<nb_valeurs; i++) { printf("%f ", donnees[i]); } printf("\n");
21     MPI_Finalize(); }
```

```
> mpiexec -n 4 gather
Moi, processus 1 envoie mon tableau valeurs :1003. 1004.
Moi, processus 0 envoie mon tableau valeurs :1001. 1002.
Moi, processus 2 envoie mon tableau valeurs :1005. 1006.
Moi, processus 3 envoie mon tableau valeurs :1007. 1008.

Moi, processus 2, ai recu 1001. 1002. 1003. 1004. 1005. 1006. 1007. 1008.
```

# Communications collectives

## Collecte générale : `MPI_Allgather()`

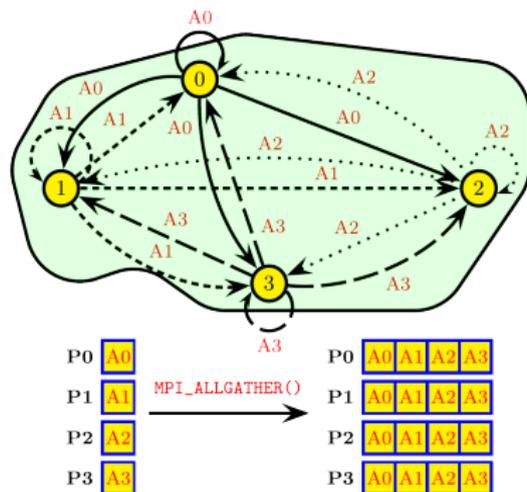


Figure 16 – Collecte générale : `MPI_Allgather()`

# Communications collectives

## Collecte générale : `MPI_Allgather()`

Correspond à un `MPI_Gather()` suivi d'un `MPI_Bcast()` :

```
int MPI_Allgather(const void *message_emis, int longueur_message_emis,
                 MPI_Datatype type_message_emis, void *message_recu,
                 int longueur_message_recu, MPI_Datatype type_message_recu,
                 MPI_Comm comm)
```

1. Envoi de chacun des processus du communicateur `comm`, d'un message `message_emis`, de taille `longueur_message_emis` et de type `type_message_emis`.
2. Collecte de chacun de ces messages, par tous les processus, à partir l'adresse `message_recu`, sur une longueur `longueur_message_recu` et avec le type `type_message_recu`.

### Remarques :

- Les couples (`longueur_message_emis`, `type_message_emis`) et (`longueur_message_recu`, `type_message_recu`) doivent être tels que les quantités de données envoyées et reçues soient égales.
- Les données sont collectées dans l'ordre des rangs des processus.

# Communications collectives

## Exemple de `MPI_Allgather()`

```
1  /* allgather */
2  #include <mpi.h>
3  #include <stdio.h>
4  #include <stdlib.h>
5  int main(int argc, char *argv[]) {
6      int nb_valeurs=8, rang, nb_procs, longueur_tranche, i;
7      float donnees[nb_valeurs], *valeurs;
8
9      MPI_Init(&argc, &argv);
10     MPI_Comm_size(MPI_COMM_WORLD, &nb_procs);
11     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
12     longueur_tranche = nb_valeurs/nb_procs;
13     valeurs = (float *) malloc(longueur_tranche*sizeof(float));
14     for (i=0; i<longueur_tranche; i++) valeurs[i]=1001.+rang*longueur_tranche+i;
15     MPI_Allgather(valeurs, longueur_tranche, MPI_FLOAT,
16                 donnees, longueur_tranche, MPI_FLOAT, MPI_COMM_WORLD);
17     printf("Moi, processus %d, ai recu ", rang);
18     for (i=0; i<nb_valeurs; i++) {printf("%f ", donnees[i]);} printf("\n");
19     MPI_Finalize(); }
```

```
> mpiexec -n 4 allgather
```

```
Moi, processus 1, ai recu 1001. 1002. 1003. 1004. 1005. 1006. 1007. 1008.
Moi, processus 3, ai recu 1001. 1002. 1003. 1004. 1005. 1006. 1007. 1008.
Moi, processus 2, ai recu 1001. 1002. 1003. 1004. 1005. 1006. 1007. 1008.
Moi, processus 0, ai recu 1001. 1002. 1003. 1004. 1005. 1006. 1007. 1008.
```

# Communications collectives

Collecte "variable" : `MPI_Gatherv()`

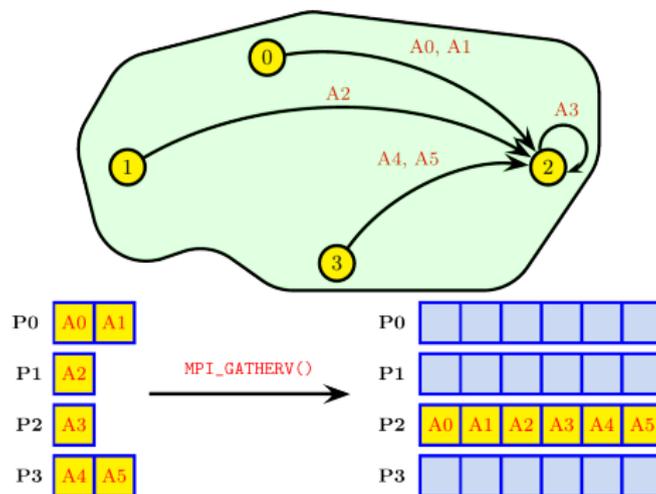


Figure 17 – Collecte : `MPI_Gatherv()`

## Collecte "variable" : `MPI_Gatherv()`

Correspond à un `MPI_Gather()` pour lequel la taille des messages varie :

```
int MPI_Gatherv(const void *message_emis, int longueur_message_emis,
               MPI_Datatype type_message_emis, void *message_recu,
               const int *nb_elts_recus, const int *deplts,
               MPI_Datatype type_message_recu, int rang_dest, MPI_Comm comm)
```

Le  $i$ ème processus du communicateur `comm` envoie au processus `rang_dest`, un message depuis l'adresse `message_emis`, de taille `longueur_message_emis`, de type `type_message_emis`, avec réception du message à l'adresse `message_recu`, de type `type_message_recu`, de taille `nb_elts_recus(i)` avec un déplacement de `deplts(i)`.

### Remarques :

- Les couples (`longueur_message_emis`, `type_message_emis`) du  $i$ ème processus et (`nb_elts_recus(i)`, `type_message_recu`) du processus `rang_dest` doivent être tels que les quantités de données envoyées et reçues soient égales.

# Communications collectives

## Exemple de `MPI_Gatherv()`

```
1  /* gatherv */
2  #include <mpi.h>
3  #include <stdio.h>
4  #include <stdlib.h>
5
6  int main(int argc, char *argv[]) {
7      int nb_valeurs=10, rang, nb_procs, longueur_tranche, reste, i;
8      float donnees[nb_valeurs];
9      float *valeurs;
10     int *nb_elements_recus, *deplacements;
11
12     MPI_Init (&argc, &argv);
13     MPI_Comm_size (MPI_COMM_WORLD, &nb_procs);
14     MPI_Comm_rank (MPI_COMM_WORLD, &rang);
15     longueur_tranche = nb_valeurs/nb_procs;
16     reste = nb_valeurs%nb_procs;
17     if (rang < reste) longueur_tranche = longueur_tranche+1;
18     valeurs = (float *) malloc(longueur_tranche*sizeof(float));
19     for (i=0; i<longueur_tranche; i++)
20         valeurs[i]=1001.+rang*(nb_valeurs/nb_procs)+(rang<reste?rang:reste)+i;
21     printf("Moi, processus %d envoie mon tableau valeurs : ", rang);
22     for (i=0; i<longueur_tranche; i++) {printf("%f ", valeurs[i]); }printf("\n");
23     if (rang == 2) {
24         nb_elements_recus = (int *) malloc(nb_procs*sizeof(int));
25         deplacements = (int *) malloc(nb_procs*sizeof(int));
26         nb_elements_recus[0] = nb_valeurs/nb_procs;
27         if (reste > 0) nb_elements_recus[0] = nb_elements_recus[0]+1;
28         deplacements[0] = 0;
29         for (i=1; i<nb_procs; i++) {
30             deplacements[i] = deplacements[i-1]+nb_elements_recus[i-1];
31             nb_elements_recus[i] = nb_valeurs/nb_procs;
32             if (i < reste) nb_elements_recus[i] = nb_elements_recus[i]+1;
33         } }
```

# Communications collectives

## Exemple de `MPI_Gatherv()` (suite)

```
MPI_Gatherv(valeurs, longueur_tranche, MPI_FLOAT,
            donnees, nb_elements_recus, deplacements, MPI_FLOAT, 2, MPI_COMM_WORLD);
if (rang==2) {
    printf("Moi, processus %d, ai recu ", rang);
    for (i=0; i<nb_valeurs; i++) printf("%f ", donnees[i]);
    printf("\n"); }
MPI_Finalize();
}
```

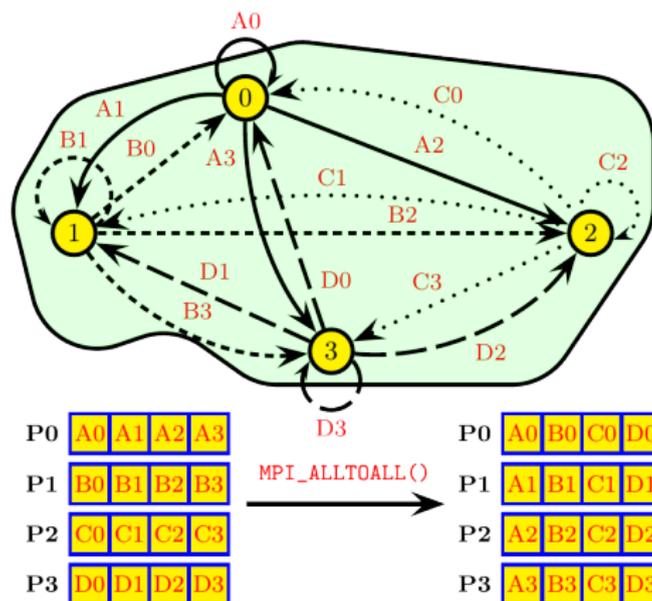
```
> mpiexec -n 4 gatherv
```

```
Moi, processus 0 envoie mon tableau valeurs : 1001. 1002. 1003.
Moi, processus 2 envoie mon tableau valeurs : 1007. 1008.
Moi, processus 3 envoie mon tableau valeurs : 1009. 1010.
Moi, processus 1 envoie mon tableau valeurs : 1004. 1005. 1006.

Moi, processus 2 ai recu 1001. 1002. 1003. 1004. 1005. 1006. 1007. 1008. 1009. 1010.
```

# Communications collectives

## Collectes et diffusions sélectives : `MPI_Alltoall()`



**Figure 18** – Collecte et diffusion sélectives : `MPI_Alltoall()`

## Collectes et diffusions sélectives : `MPI_Alltoall()`

```
int MPI_Alltoall(const void *message_emis, int longueur_message_emis,  
                MPI_Datatype type_message_emis, void *message_recu,  
                int longueur_message_recu, MPI_Datatype type_message_recu,  
                MPI_Comm comm)
```

Ici, le  $i$ ème processus envoie la  $j$ ème tranche au  $j$ ème processus qui le place à l'emplacement de la  $i$ ème tranche.

### Remarque :

- Les couples (`longueur_message_emis`, `type_message_emis`) et (`longueur_message_recu`, `type_message_recu`) doivent être tels que les quantités de données envoyées et reçues soient égales.

# Communications collectives

## Exemple de `MPI_Alltoall()`

```
1  /* alltoall */
2  #include <mpi.h>
3  #include <stdio.h>
4  #include <stdlib.h>
5
6  int main(int argc, char *argv[]) {
7      int nb_valeurs=8;
8      int rang, nb_procs, longueur_tranche, i;
9      float donnees[nb_valeurs], valeurs[nb_valeurs];
10
11
12     MPI_Init(&argc, &argv);
13     MPI_Comm_size(MPI_COMM_WORLD, &nb_procs);
14     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
15
16     for (i=0; i<nb_valeurs; i++) valeurs[i]=1001.+rang*nb_valeurs+i;
17     longueur_tranche = nb_valeurs/nb_procs;
18
19     printf("Moi, processus %d envoie mon tableau valeurs : ", rang);
20     for (i=0; i<nb_valeurs; i++) printf("%f ", valeurs[i]);
21     printf("\n");
22
23     MPI_Alltoall(valeurs, longueur_tranche, MPI_FLOAT,
24                donnees, longueur_tranche, MPI_FLOAT, MPI_COMM_WORLD);
25
26     printf("Moi, processus %d, ai recu ", rang);
27     for (i=0; i<nb_valeurs; i++) printf("%f ", donnees[i]);
28     printf("\n");
29     MPI_Finalize();
30 }
```

## Exemple de `MPI_Alltoall()` (suite)

```
> mpiexec -n 4 alltoall
Moi, processus 1 envoie mon tableau valeurs :
1009. 1010. 1011. 1012. 1013. 1014. 1015. 1016.
Moi, processus 0 envoie mon tableau valeurs :
1001. 1002. 1003. 1004. 1005. 1006. 1007. 1008.
Moi, processus 2 envoie mon tableau valeurs :
1017. 1018. 1019. 1020. 1021. 1022. 1023. 1024.
Moi, processus 3 envoie mon tableau valeurs :
1025. 1026. 1027. 1028. 1029. 1030. 1031. 1032.

Moi, processus 0, ai recu 1001. 1002. 1009. 1010. 1017. 1018. 1025. 1026.
Moi, processus 2, ai recu 1005. 1006. 1013. 1014. 1021. 1022. 1029. 1030.
Moi, processus 1, ai recu 1003. 1004. 1011. 1012. 1019. 1020. 1027. 1028.
Moi, processus 3, ai recu 1007. 1008. 1015. 1016. 1023. 1024. 1031. 1032.
```

## 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. Le résultat est obtenu sur un seul processus (`MPI_Reduce()`) ou bien sur tous (`MPI_Allreduce()`), qui est en fait équivalent à 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 (par exemple à tous les éléments d'un vecteur).

# Communications collectives

## Réductions réparties : `MPI_Reduce()`

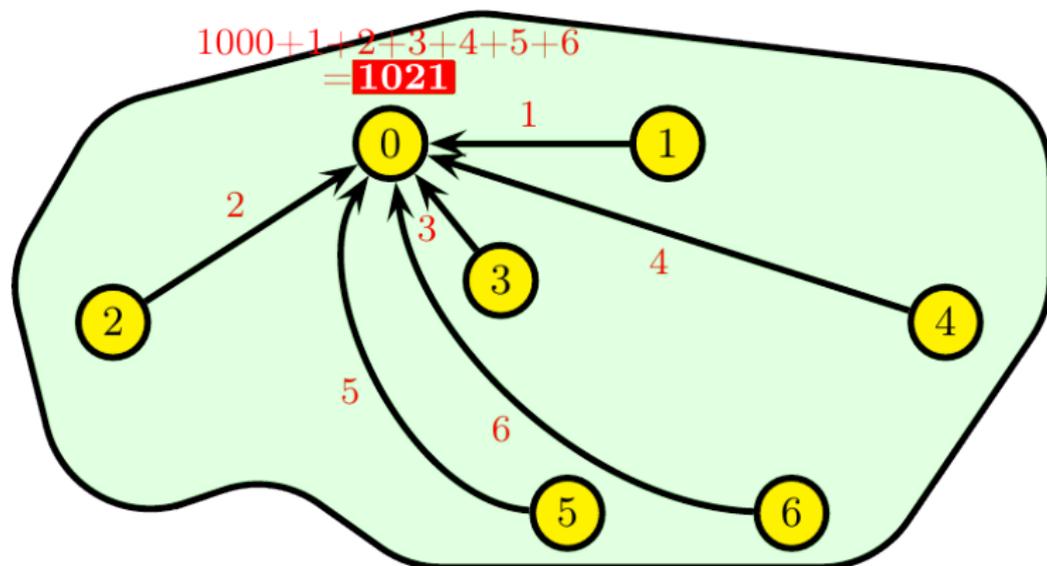


Figure 19 – Réduction répartie : `MPI_Reduce()` avec l'opérateur somme

## Opérations pour réductions réparties

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

## Réductions réparties : `MPI_Reduce()`

```
int MPI_Reduce(const void*message_emis,void *message_recu,int longueur,  
              MPI_Datatype type_message,MPI_Op operation,int rang_dest,  
              MPI_Comm comm)
```

1. Réduction répartie des éléments situés à partir de l'adresse `message_emis`, de taille `longueur`, de type `type_message`, pour les processus du communicateur `comm`,
2. Écrit le résultat à l'adresse `message_recu` pour le processus de rang `rang_dest`.

# Communications collectives

## Exemple de `MPI_Reduce()` (voir Fig.19)

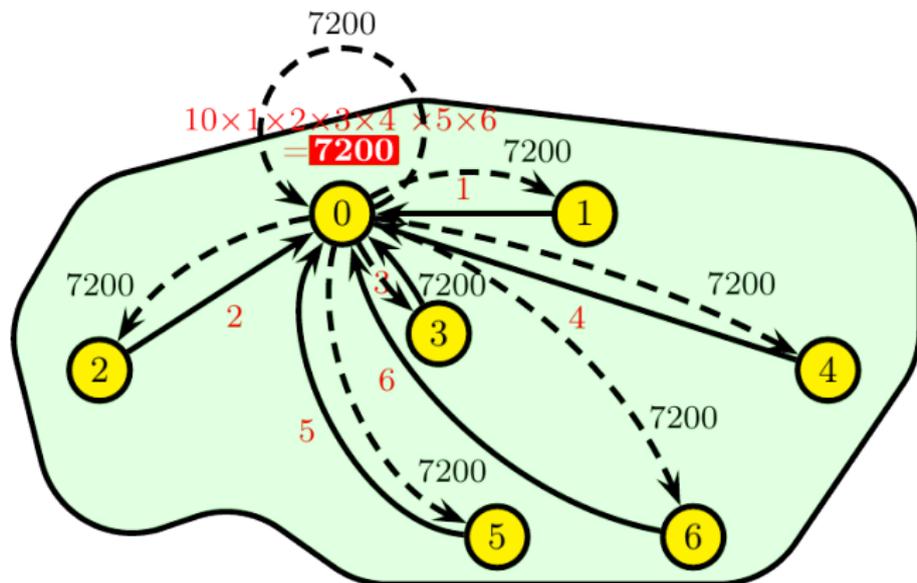
```
1  /* reduce */
2  #include <mpi.h>
3  #include <stdio.h>
4
5  int main(int argc, char *argv[]) {
6      int rang, nb_procs, valeur, somme, i;
7
8      MPI_Init(&argc, &argv);
9      MPI_Comm_size(MPI_COMM_WORLD, &nb_procs);
10     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
11
12     if (rang == 0)
13         valeur = 1000;
14     else
15         valeur = rang;
16
17     MPI_Reduce(&valeur, &somme, 1, MPI_INT, MPI_SUM, 0, MPI_COMM_WORLD);
18
19     if (rang == 0)
20         printf("Moi, processus 0, ma valeur de la somme globale est %d\n", somme);
21
22     MPI_Finalize();
23 }
```

```
> mpiexec -n 7 reduce
```

```
Moi, processus 0, ma valeur de la somme globale est 1021
```

# Communications collectives

## Réductions réparties avec diffusion du résultat : `MPI_Allreduce()`



**Figure 20** – Réduction répartie avec diffusion du résultat : `MPI_Allreduce` (utilisation de l'opérateur produit)

## Réductions réparties avec diffusion du résultat : `MPI_Allreduce()`

```
int MPI_Allreduce(const void *message_emis, void *message_recu, int longueur,  
                 MPI_Datatype type_message, MPI_Op operation, MPI_Comm comm)
```

1. Réduction répartie des éléments situés à partir de l'adresse `message_emis`, de taille `longueur`, de type `type_message`, pour les processus du communicateur `comm`,
2. Écrit le résultat à l'adresse `message_recu` pour tous les processus du communicateur `comm`.

## Exemple de `MPI_Allreduce()` (voir Fig.20)

```
1  /* allreduce */
2  #include <mpi.h>
3  #include <stdio.h>
4
5  int main(int argc, char *argv[]) {
6      int rang, nb_procs, valeur, produit, i;
7
8      MPI_Init(&argc, &argv);
9      MPI_Comm_size(MPI_COMM_WORLD, &nb_procs);
10     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
11
12     if (rang == 0)
13         valeur = 10;
14     else
15         valeur = rang;
16
17     MPI_Allreduce(&valeur, &produit, 1, MPI_INT, MPI_PROD, MPI_COMM_WORLD);
18
19     printf("Moi, processus %d, ai reçu la valeur du produit globale %d\n",
20           rang, produit);
21
22     MPI_Finalize();
23 }
```

## Exemple de `MPI_Allreduce()` (voir Fig.20) (suite)

```
> mpiexec -n 7 allreduce  
Moi, processus 6, ai recu la valeur du produit global 7200  
Moi, processus 2, ai recu la valeur du produit global 7200  
Moi, processus 0, ai recu la valeur du produit global 7200  
Moi, processus 4, ai recu la valeur du produit global 7200  
Moi, processus 5, ai recu la valeur du produit global 7200  
Moi, processus 3, ai recu la valeur du produit global 7200  
Moi, processus 1, ai recu la valeur du produit global 7200
```

## Compléments

- Le sous-programme `MPI_Scan()` permet d'effectuer des réductions partielles en considérant, pour chaque processus, les processus précédents du communicateur et lui-même. `MPI_Exscan()` est la version *exclusive* de `MPI_Scan()`, qui elle est inclusive.
- Les sous-programmes `MPI_Op_create()` et `MPI_Op_free()` permettent de définir des opérations de réduction personnelles.
- Pour toutes les opérations de réduction, le mot-clé `MPI_IN_PLACE` peut être utilisé pour que les données et résultats de l'opération soient stockés au même endroit (mais uniquement pour le ou les processus qui reçoivent les résultats).  
Exemple :

```
MPI_Allreduce(MPI_IN_PLACE, message_emis_et_recu, ...);
```

## Compléments

- De même que ce qui a été vu pour `MPI_Gatherv()` vis-à-vis de `MPI_Gather()`, les sous-programmes `MPI_Scatterv()`, `MPI_Allgatherv()` et `MPI_Alltoallv()` étendent `MPI_Scatter()`, `MPI_Allgather()` et `MPI_Alltoall()` au cas où le nombre d'éléments à diffuser ou collecter est différent suivant les processus.
- `MPI_Alltoallw()` est la version de `MPI_Alltoallv()` permettant de traiter des éléments hétérogènes (en exprimant les déplacements en octets et non en éléments).

## T.P. MPI – Exercice 3 : Communications collectives et réductions

- Il s'agit de calculer  $\pi$  par intégration numérique  $\pi = \int_0^1 \frac{4}{1+x^2} dx$ .
- On utilise la méthode des rectangles (point milieu).
- La fonction à intégrer est  $f(x) = \frac{4}{1+x^2}$ .
- *nbbloc* est le nombre de points.
- *largeur* =  $\frac{1}{nbbloc}$  est le pas de discrétisation et la largeur de chaque rectangle.
- La version séquentielle est disponible dans le fichier `pi.c`.
- Il vous faut écrire la version parallélisée avec MPI dans ce fichier.

## Modèles de communication

# Modèles de communication

## Modes d'envoi point à point

<i>Mode</i>	<i>Bloquant</i>	<i>Non bloquant</i>
Envoi standard	<code>MPI_Send()</code>	<code>MPI_Isend()</code>
Envoi synchrone	<code>MPI_Ssend()</code>	<code>MPI_Issend()</code>
Envoi <i>bufferisé</i>	<code>MPI_Bsend()</code>	<code>MPI_Ibsend()</code>
Réception	<code>MPI_Recv()</code>	<code>MPI_Irecv()</code>

## Appels bloquants

- Un appel est **bloquant** si l'espace mémoire servant à la communication peut être réutilisé immédiatement après la sortie de l'appel.
- Les données envoyées peuvent être modifiées après l'appel bloquant.
- Les données reçues peuvent être lues après l'appel bloquant.

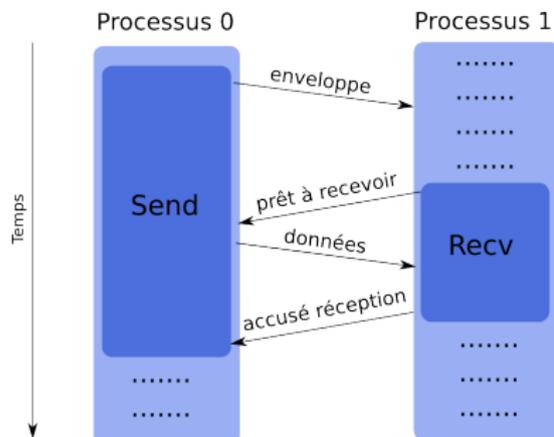
# Modèles de communication

## Envois synchrones

Un **envoi synchrone** implique une synchronisation entre les processus concernés. Un envoi ne pourra commencer que lorsque sa réception sera postée. Il ne peut y avoir de communication que si les deux processus sont prêts à communiquer.

## Protocole de *rendez-vous*

Le protocole de *rendez-vous* est généralement celui employé pour les envois en mode synchrone (dépend de l'implémentation). L'accusé de réception est optionnel.



# Modèles de communication

## Interface de `MPI_Ssend()`

```
int MPI_Ssend(const void* valeurs, int taille, MPI_Datatype type_message,  
             int dest, int etiquette, MPI_Comm comm)
```

## Avantages

- Consomment peu de ressources (pas de *buffer*)
- Rapides si le récepteur est prêt (pas de recopie dans un *buffer*)
- Connaissance de la réception grâce à la synchronisation

## Inconvénients

- Temps d'attente si le récepteur n'est pas là/pas prêt
- Risques d'inter-blocage

# Modèles de communication

## Exemple d'inter-blocage

Dans l'exemple suivant, on a un inter-blocage, car on est en mode synchrone, les deux processus sont bloqués sur le `MPI_Ssend()` car ils attendent le `MPI_Recv()` de l'autre processus. Or ce `MPI_Recv()` ne pourra se faire qu'après le déblocage du `MPI_Ssend()`.

```
1  /* ssendrecv */
2  #include <mpi.h>
3  #include <stdio.h>
4
5  int main(int argc, char *argv[]) {
6      int rang, num_proc, tmp, valeur;
7      int etiquette=110;
8
9      MPI_Init(&argc, &argv);
10     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
11
12     /* On suppose avoir exactement 2 processus */
13     num_proc = (rang+1)%2;
14
15     tmp = rang+1000;
16     MPI_Ssend(&tmp, 1, MPI_INT, num_proc, etiquette, MPI_COMM_WORLD);
17     MPI_Recv(&valeur, 1, MPI_INT, num_proc, etiquette, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
18
19     printf("Moi, processus %d ai reçu %d du processus %d\n", rang, valeur, num_proc);
20
21     MPI_Finalize();
22 }
```

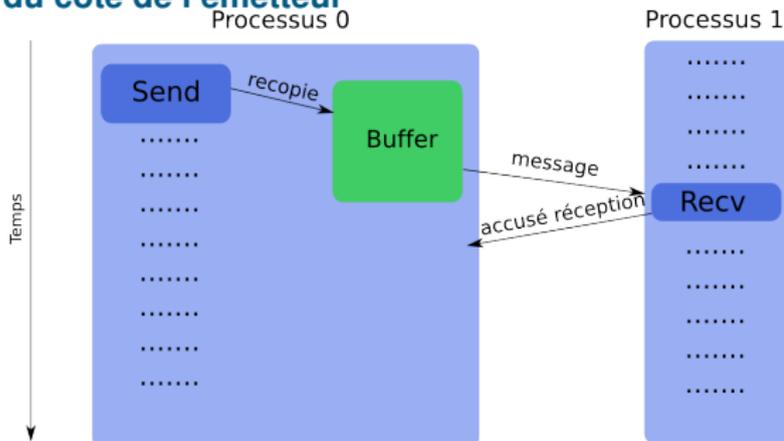
# Modèles de communication

## Envois *bufferisés*

Un *envoi bufferisé* implique la recopie des données dans un espace mémoire intermédiaire. Il n'y a alors pas de couplage entre les deux processus de la communication. La sortie de ce type d'envoi ne signifie donc pas que la réception a eu lieu.

## Protocole avec *buffer* utilisateur du côté de l'émetteur

Dans cette approche, le *buffer* se trouve du côté de l'émetteur et est géré explicitement par l'application. Un *buffer* géré par MPI peut exister du côté du récepteur. De nombreuses variantes sont possibles. L'accusé de réception est optionnel.



## Buffers

Les *buffers* doivent être gérés manuellement (avec appels à `MPI_Buffer_attach()` et `MPI_Buffer_detach()`). Ils doivent être alloués en tenant compte des surcoûts mémoire des messages (en ajoutant la constante `MPI_BSEND_OVERHEAD` pour chaque instance de message).

## Interfaces

```
int MPI_Buffer_attach(void *buf, int taille_buf)
int MPI_Buffer_detach(void *buf, int taille_buf)
int MPI_Bsend(const void *valeurs, int taille, MPI_Datatype type_message,
              int dest, int etiquette, MPI_Comm comm)
```

# Modèles de communication

## Avantages du mode bufferisé

- Pas besoin d'attendre le récepteur (recopie dans un *buffer*)
- Pas de risque de blocage (*deadlocks*)

## Inconvénients du mode bufferisé

- Consomment plus de ressources (occupation mémoire par les *buffers* avec risques de saturation)
- Les *buffers* d'envoi doivent être gérés manuellement (souvent délicat de choisir une taille adaptée)
- Un peu plus lent que les envois synchrones si le récepteur est prêt
- Pas de connaissance de la réception (découplage envoi-réception)
- Risque de gaspillage d'espace mémoire si les *buffers* sont trop sur-dimensionnés
- L'application plante si les *buffers* sont trop petits
- Il y a aussi souvent des *buffers* cachés gérés par l'implémentation MPI du côté de l'expéditeur et/ou du récepteur (et consommant des ressources mémoires)

## Absence d'inter-blocage

Dans l'exemple suivant, on a pas d'inter-blocage, car on est en mode bufferisé. Une fois la copie faite dans le *buffer*, l'appel `MPI_Bsend()` retourne et on passe à l'appel `MPI_Recv()`.

```
1  /* bsendrecv */
2  #include <mpi.h>
3  #include <stdio.h>
4  #include <stdlib.h>
5
6  int main(int argc, char *argv[]) {
7      int rang, num_proc, tmp, valeur, taille, surcout, taille_buf;
8      int etiquette=110, nb_elt=1, nb_msg=1;
9      int * buffer;
10
11     MPI_Init(&argc, &argv);
12     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
13     MPI_Type_size(MPI_INT, &taille);
14     /* Convertir taille MPI_BSEND_OVERHEAD (octets) en nombre d'integer */
15     surcout = (int) (1+(MPI_BSEND_OVERHEAD*1.)/taille);
16     buffer = (int *) malloc(nb_msg*(nb_elt+surcout)*sizeof(int));
17     taille_buf = taille*nb_msg*(nb_elt+surcout);
18     MPI_Buffer_attach(buffer, taille_buf);
19     /* On suppose avoir exactement 2 processus */
20     num_proc = (rang+1)%2;
21     tmp = rang+1000;
22     MPI_Bsend(&tmp, nb_elt, MPI_INT, num_proc, etiquette, MPI_COMM_WORLD);
23     MPI_Recv(&valeur, nb_elt, MPI_INT, num_proc, etiquette, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
24     printf("Moi, processus %d ai recu %d du processus %d\n", rang, valeur, num_proc);
25     MPI_Buffer_detach(&buffer, &taille_buf);
26     MPI_Finalize(); }
```

## Envois standards

Un envoi standard se fait en appelant le sous-programme `MPI_Send()`. Dans la plupart des implémentations, ce mode passe d'un mode *bufferisé* (*eager*) à un mode synchrone lorsque la taille des messages croît.

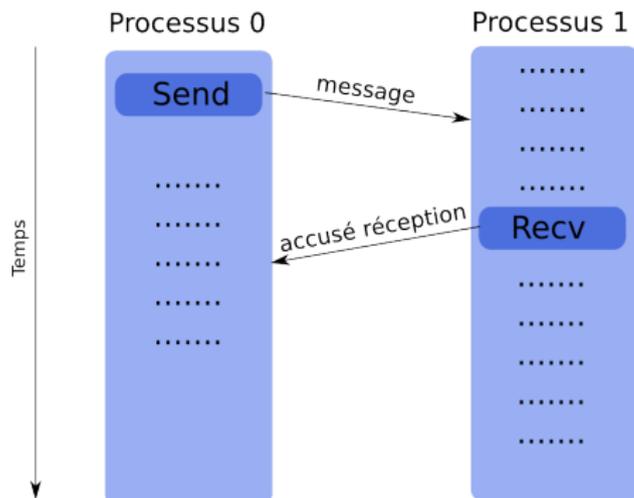
## Interfaces

```
int MPI_Send(const void *valeurs, int taille, MPI_Datatype type_message,  
            int dest, int etiquette, MPI_Comm comm)
```

# Modèles de communication

## Protocole *eager*

Le protocole *eager* est souvent employé pour les envois en mode standard (`MPI_Send()`) pour les messages de petites tailles. Il peut aussi être utilisé pour les envois avec `MPI_Bsend()` avec des petits messages (dépend de l'implémentation) et en court-circuitant le *buffer* utilisateur du côté de l'émetteur. Dans cette approche, le *buffer* se trouve du côté du récepteur. L'accusé de réception est optionnel.



# Modèles de communication

## Avantages du mode standard

- Souvent le plus performant (choix du mode le plus adapté par le constructeur)

## Inconvénients du mode standard

- Peu de contrôle sur le mode réellement utilisé (souvent accessible via des variables d'environnement)
- Risque de *deadlock* selon le mode réel
- Comportement pouvant varier selon l'architecture et la taille du problème

## Nombre d'éléments reçus

```
int MPI_Recv(void *message, int longueur, MPI_Datatype type_message,  
            int rang_source, int etiquette, MPI_Comm comm, MPI_Status *statut)
```

- Dans l'appel à `MPI_Recv()`, l'argument `longueur` correspond dans la norme au nombre d'éléments dans le buffer `message`.
- Ce nombre doit être supérieur au nombre d'éléments à recevoir.
- Quand c'est possible, pour des raisons de lisibilité, il est conseillé de mettre le nombre d'éléments à recevoir.
- On peut connaître le nombre d'éléments reçus avec `MPI_Get_count()` et à l'aide de l'argument `statut` retourné par l'appel à `MPI_Recv()`.

```
int MPI_Get_count(MPI_Status *statut, MPI_Datatype type_message, int *longueur)
```

## Nombre d'éléments reçus

`MPI_Probe` permet de vérifier les messages entrants sans les recevoir.

```
int MPI_Probe(int source, int tag, MPI_Comm comm, MPI_Status *status)
```

Une utilisation courante de `MPI_Probe` consiste à allouer de l'espace pour un message avant de le recevoir.

```
MPI_Probe(MPI_ANY_SOURCE, MPI_ANY_TAG, comm, &status);  
MPI_Get_count(&status, MPI_INT, &msgsize);  
buf = (int*) malloc(msgsize*sizeof(int));  
MPI_Recv(buf, msgsize, MPI_INT, status.MPI_SOURCE,  
         status.MPI_TAG, comm, MPI_STATUS_IGNORE);
```

# Modèles de communication

## Présentation

Le recouvrement des communications par des calculs est une méthode permettant de réaliser des opérations de communications en arrière-plan pendant que le programme continue de s'exécuter. Sur Jean Zay, la latence d'une communication inter-nœud est de  $1\mu\text{s}$  soit 2500 cycles processeur.

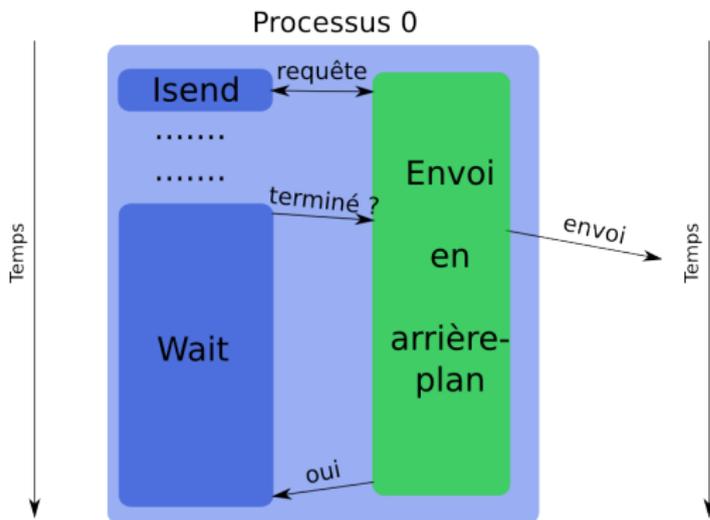
- Il est ainsi possible, si l'architecture matérielle et logicielle le permet, de masquer tout ou une partie des coûts de communications.
- Le recouvrement calculs-communications peut être vu comme un niveau supplémentaire de parallélisme.
- Cette approche s'utilise dans MPI par l'utilisation de sous-programmes non-bloquants (i.e. `MPI_Isend()`, `MPI_Irecv()` et `MPI_Wait()`).

## Appels non bloquants

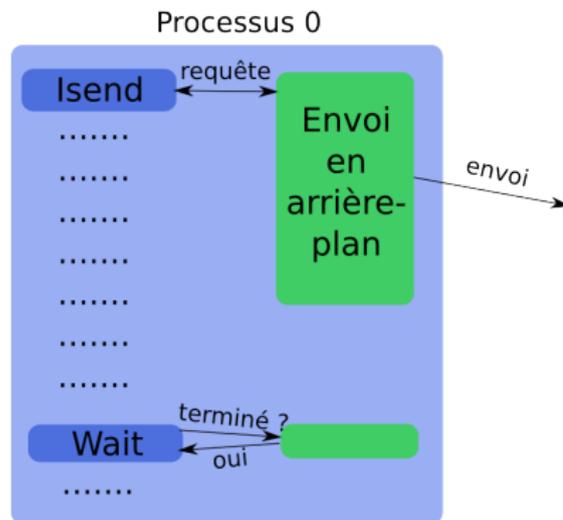
Un appel **non bloquant** rend la main très rapidement, mais n'autorise pas la réutilisation immédiate de l'espace mémoire utilisé dans la communication. Il est nécessaire de s'assurer que la communication est bien terminée (avec `MPI_Wait()` par exemple) avant de l'utiliser à nouveau.

# Modèles de communication

## Recouvrement partiel



## Recouvrement total



# Modèles de communication

## Avantages des appels non bloquants

- Possibilité de masquer tout ou une partie des coûts des communications (si l'architecture le permet)
- Pas de risques de *deadlock*

## Inconvénients des appels non bloquants

- Surcoûts plus importants (plusieurs appels pour un seul envoi ou réception, gestion des requêtes)
- Complexité plus élevée et maintenance plus compliquée
- Peu performant sur certaines machines (par exemple avec transfert commençant seulement à l'appel de `MPI_Wait()`)
- Risque de perte de performance sur les noyaux de calcul (par exemple gestion différenciée entre la zone proche de la frontière d'un domaine et la zone intérieure entraînant une moins bonne utilisation des caches mémoires)
- Limité aux communications point à point (a été étendu aux collectives dans MPI 3.0)

## Interfaces

`MPI_Isend()` `MPI_Issend()` et `MPI_Ibsend()` pour les envois non bloquants

```
int MPI_Isend(const void*valeurs, int taille, MPI_Datatype type_message,
             int dest, int etiquette, MPI_Comm comm, MPI_Request *req)
int MPI_Issend(const void*valeurs, int taille, MPI_Datatype type_message,
              int dest, int etiquette, MPI_Comm comm, MPI_Request *req)
int MPI_Ibsend(const void*valeurs, int taille, MPI_Datatype type_message,
              int dest, int etiquette, MPI_Comm comm, MPI_Request *req)
```

`MPI_Irecv()` pour les réceptions non bloquantes.

```
int MPI_Irecv(void *valeurs, int taille, MPI_Datatype type_message, int source,
             int etiquette, MPI_Comm comm, MPI_Request *req)
```

## Interfaces

`MPI_Wait()` attend la fin d'une communication. `MPI_Test()` est la version non bloquante.

```
int MPI_Wait(MPI_Request *req, MPI_Status *statut)
int MPI_Test(MPI_Request *req, int *flag, MPI_Status *statut)
```

`MPI_Waitall()` attend la fin de toutes les communications. `MPI_Testall()` est la version non bloquante.

```
int MPI_Waitall(int taille, MPI_Request reqs[], MPI_Status statuts[])
int MPI_Testall(int taille, MPI_Request reqs[], int *flag, MPI_Status statuts[])
```

## Interfaces

`MPI_Waitany()` attend la fin d'une communication parmi plusieurs. `MPI_Testany()` est la version non bloquante.

```
int MPI_Waitany(int taille, MPI_Request reqs[], int *indice, MPI_Status *statut)
int MPI_Testany(int taille, MPI_Request reqs[], int *indice, int *flag, MPI_Status *statut)
```

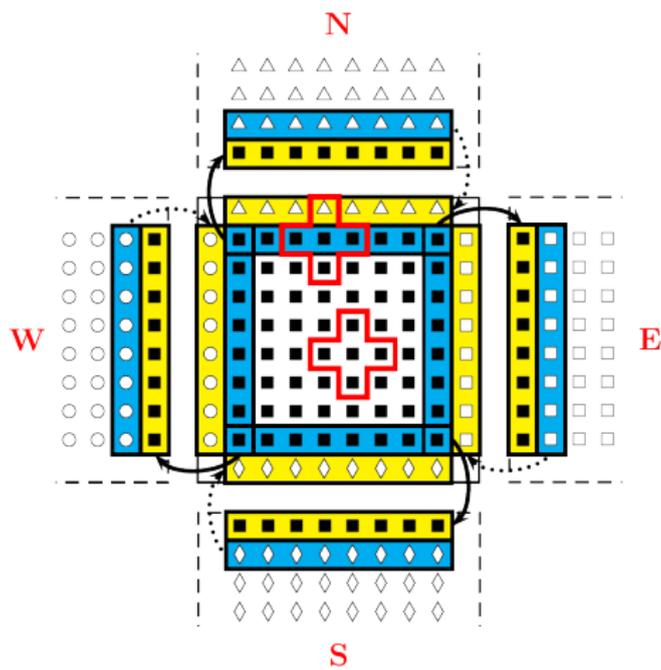
`MPI_Waitsome()` attend la fin d'une ou plusieurs communications.  
`MPI_Testsome()` est la version non bloquante.

```
int MPI_Waitsome(int taille, MPI_Request reqs[], int *nbfin, int *indices, MPI_Status *statuts)
int MPI_Testsome(int taille, MPI_Request reqs[], int *nbfin, int *indices, MPI_Status *statuts)
```

## Gestion des requêtes

- Après un appel aux fonctions bloquantes d'attente (`MPI_Wait()`, `MPI_Waitall()`, ...), la requête vaut `MPI_REQUEST_NULL`.
- De même après un appel aux fonctions non bloquantes d'attente lorsque le *flag* est à vrai.
- Une attente avec une requête qui vaut `MPI_REQUEST_NULL` ne fait rien.

# Modèles de communication



# Modèles de communication

```
1 void debut_communication(double *u) {
2     /* Envoi au voisin N et reception du voisin S */
3     MPI_Irecv(&(u[0]),1,type_ligne,voisin[S],etiquette,comm2d,&(requete[0]));
4     MPI_Isend(&(u[0]),1,type_ligne,voisin[N],etiquette,comm2d,&(requete[1]));
5
6     /* Envoi au voisin S et reception du voisin N */
7     MPI_Irecv(&(u[1]),1,type_ligne,voisin[N],etiquette,comm2d,&(requete[2]));
8     MPI_Isend(&(u[1]),1,type_ligne,voisin[S],etiquette,comm2d,&(requete[3]));
9
10    /* Envoi au voisin W et reception du voisin E */
11    MPI_Irecv(&(u[2]),1,type_colonne,voisin[E],etiquette,comm2d,&(requete[4]));
12    MPI_Isend(&(u[2]),1,type_colonne,voisin[W],etiquette,comm2d,&(requete[5]));
13
14    /* Envoi au voisin E et reception du voisin W */
15    MPI_Irecv(&(u[3]),1,type_colonne,voisin[W],etiquette,comm2d,&(requete[6]));
16    MPI_Isend(&(u[3]),1,type_colonne,voisin[E],etiquette,comm2d,&(requete[7]));
17 }
18 void fin_communication(double *u) {
19     MPI_Waitall(2*NB_VOISINS,requete,tab_statut);
20 }
```

# Modèles de communication

```
1 while(!(convergence) && (it < it_max) ) {
2   it = it+1;
3   /* Echanges des pointeurs */
4   temp = u; u = u_nouveau; u_nouveau = temp;
5
6   debut_communication(u);
7   calcul(u,u_nouveau, sx+1, ex-1, sy+1, ey-1);
8   fin_communication(u);
9
10  /* Nord */
11  calcul(u,u_nouveau,sx,sx,sy,ey);
12  /* Sud */
13  calcul(u,u_nouveau,ex,ex,sy,ey);
14  /* Ouest */
15  calcul(u,u_nouveau,sx,ex,sy,sy);
16  /* Est */
17  calcul(u,u_nouveau,sx,ex,ey,ey);
18
19  /* Calcul de l'erreur globale */
20  diffnorm = erreur_globale (u, u_nouveau);
21
22  /* Arrêt du programme si on a atteint la precision machine obtenue */
23  convergence = (diffnorm < eps);
24 }
```

# Modèles de communication

## Niveau de recouvrement sur différentes machines

<i>Machine</i>	<i>Niveau</i>
Zay(IntelMPI)	43%
Zay(IntelMPI) I_MPI_ASYNC_PROGRESS=yes	95%

Mesures faites en recouvrant un noyau de calcul et un noyau de communication de mêmes durées.

Un recouvrement de 0% signifie que la durée totale d'exécution vaut 2x la durée d'un noyau de calcul (ou communication).

Un recouvrement de 100% signifie que la durée totale vaut 1x la durée d'un noyau de calcul (ou communication).

## Communications collectives non bloquantes

- Version non bloquante des communications collectives
- Avec un I (*immediate*) devant : `MPI_Ireduce()`, `MPI_Ibcast()`, ...
- Attente avec les appels `MPI_Wait()`, `MPI_Test()` et leurs variantes
- Pas de correspondance bloquant et non bloquant
- Le *status* récupéré par `MPI_Wait()` contient une valeur non définie pour `MPI_SOURCE` et `MPI_TAG`
- Pour les processus d'un communicateur donné, l'ordre des appels doit être le même (comme en version bloquante)

```
int MPI_Ibarrier(MPI_Comm comm, MPI_Request *request)
```

# Modèles de communication

## Exemple d'utilisation du `MPI_Ibarrier`

Comment gérer les communications quand on ne sait pas à chaque itération si nos voisins vont envoyer un message.

```
int isAllFinish=0, isMySendFinish=0;

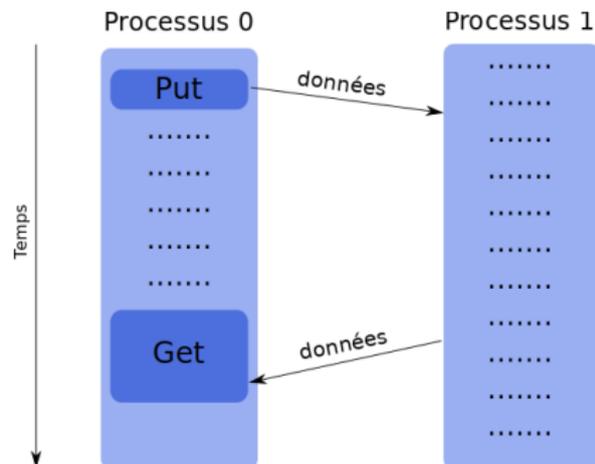
for (i=0; i<m; i++) {
    /* Envoi synchrone */
    MPI_Issend(sbuf[i], size[i], type, dst[i], tag, comm, &reqs[i]);
}

while (!isAllFinish) {
    /* Avons nous un message pret a etre recu */
    MPI_Iprobe(MPI_ANY_SOURCE, tag, comm, &flag, &astat);
    if (flag) {
        /* Recoit le message */
        MPI_Recv(rbuf, rsize, type, astat%MPI_SOURCE, tag, comm, &rstat);
    }
    if (!isMySendFinish) {
        /* Verifie si tous nos ssend sont termines */
        MPI_Testall(m, reqs, &flag, MPI_STATUSES_IGNORE);
        if (flag) {
            /* Si c'est le cas on demarre la ibarrier */
            MPI_Ibarrier(comm, &reqb);
            isMySendFinish=1;
        }
    } else {
        /* Test si tout le monde a fait la ibarrier */
        MPI_Test(&reqb, isAllFinish, MPI_STATUS_IGNORE);
    }
}
```

# Modèles de communication

## Communications mémoire à mémoire (RMA)

Les communications mémoire à mémoire (ou RMA pour *Remote Memory Access* ou *one sided communications*) consistent à accéder en écriture ou en lecture à la mémoire d'un processus distant sans que ce dernier doive gérer cet accès explicitement. Le processus cible n'intervient donc pas lors du transfert.



# Modèles de communication

## RMA - Approche générale

- Création d'une fenêtre mémoire avec `MPI_Win_create()` pour autoriser les transferts RMA dans cette zone.
- Accès distants en lecture ou écriture en appelant `MPI_Put()`, `MPI_Get()`, `MPI_Accumulate()`, `MPI_Fetch_and_op()`, `MPI_Get_accumulate()` et `MPI_Compare_and_swap()`.
- Libération de la fenêtre mémoire avec `MPI_Win_free()`.

## RMA - Méthodes de synchronisation

Pour s'assurer d'un fonctionnement correct, il est obligatoire de réaliser certaines synchronisations. 3 méthodes sont disponibles :

- Communication à cible active avec synchronisation globale (`MPI_Win_fence()`);
- Communication à cible active avec synchronisation par paire (`MPI_Win_start()` et `MPI_Win_complete()` pour le processus origine; `MPI_Win_post()` et `MPI_Win_wait()` pour le processus cible);
- Communication à cible passive sans intervention de la cible (`MPI_Win_lock()` et `MPI_Win_unlock()`).

# Modèles de communication

```
1  /* ex_fence */
2  #include <mpi.h>
3  #include <stdio.h>
4  #include <stdlib.h>
5
6  int main(int argc, char *argv[]) {
7      int rang, taille_reel, i;
8      int n=4, m=4, cible, nb_elements;
9      MPI_Aint dim_win, deplacement;
10     double *win_local, *tab, sum;
11     MPI_Win win;
12     int assert=0;
13
14     MPI_Init(&argc, &argv);
15     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
16     MPI_Type_size(MPI_DOUBLE, &taille_reel);
17
18     if (rang == 0) {
19         n = 0;
20         tab = (double *) malloc(m*sizeof(double)); }
21
22     win_local = (double *) malloc(n*sizeof(double));
23     dim_win = taille_reel*n;
24
25     MPI_Win_create(win_local, dim_win, taille_reel, MPI_INFO_NULL, MPI_COMM_WORLD, &win);
```

# Modèles de communication

```
26  if (rang == 0) {
27      for(i=0;i<m;i++) tab[i]=i+1;
28  } else {
29      for(i=0;i<n;i++) win_local[i]=0.0;
30  }
31
32  MPI_Win_fence(assert,win);
33  if (rang == 0) {
34      cible = 1; nb_elements = 2; deplacement =1;
35      MPI_Put (tab,nb_elements,MPI_DOUBLE,cible,
36              deplacement,nb_elements,MPI_DOUBLE,win);}
37
38  MPI_Win_fence(assert,win);
39  sum = 0.;
40  if (rang == 0) {
41      for (i=0;i<m-1;i++) sum=sum+tab[i];
42      tab[m-1] = sum;
43  } else {
44      for (i=0;i<n-1;i++) sum=sum+win_local[i];
45      win_local[n-1] = sum; }
46
47  MPI_Win_fence(assert,win);
48  if (rang == 0) {
49      nb_elements=1;deplacement=m-1;
50      MPI_Get (tab,nb_elements,MPI_DOUBLE,cible,
51              deplacement,nb_elements,MPI_DOUBLE,win);}

```

# Modèles de communication

## Avantages des RMA

- Permet de mettre en place plus efficacement certains algorithmes.
- Plus performant que les communications point à point sur certaines machines (utilisation de matériels spécialisés tels que moteur DMA, coprocesseur, mémoire spécialisée...).
- Possibilité pour l'implémentation de regrouper plusieurs opérations.

## Inconvénients des RMA

- La gestion des synchronisations est délicate.
- Complexité et risques d'erreurs élevés.
- Moins performant que les communications point à point sur certaines machines.

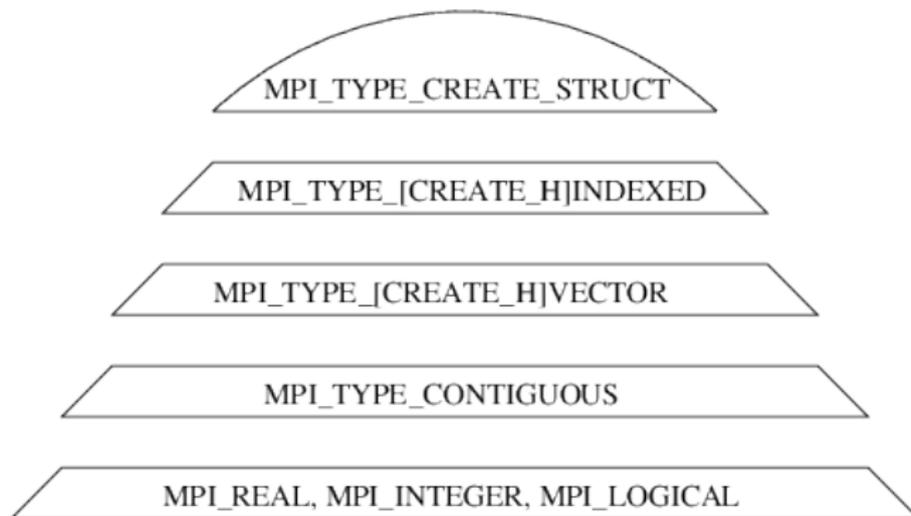
## Types de données dérivés

# Types de données dérivés

## Introduction

- Dans les communications, les données échangées sont typées : `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_indexed()` ou `MPI_Type_create_struct()` .
- Les types dérivés permettent notamment l'échange de données non contiguës ou non homogènes en mémoire et de limiter le nombre d'appels aux sous-programmes de communications.

## Types de données dérivés



**Figure 21** – Hiérarchie des constructeurs de type MPI

# Types de données dérivés

## Types contigus

- `MPI_Type_contiguous()` crée une structure de données à partir d'un ensemble homogène de type préexistant de données contiguës en mémoire.

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

```
MPI_Type_contiguous(5, MPI_FLOAT, &nouveau_type);
```

Figure 22 – Sous-programme `MPI_Type_contiguous`

```
int MPI_Type_contiguous(int nombre, MPI_Datatype ancien_type, MPI_Datatype *nouveau_type)
```

# Types de données dérivés

## Types avec un pas constant

- `MPI_Type_vector()` crée une structure de données à partir d'un ensemble homogène de type préexistant de données distantes d'un pas constant en mémoire. Le pas est donné en nombre d'éléments.

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

```
MPI_Type_vector(6,1,5,MPI_FLOAT,&nouveau_type);
```

Figure 23 – Sous-programme `MPI_Type_vector`

```
int MPI_Type_vector(int nombre,int nbr_elt_par_bloc,int pas,  
MPI_Datatype type_elt,MPI_Datatype *nouveau_type)
```

# Types de données dérivés

## Types avec un pas constant

- `MPI_Type_create_hvector()` crée une structure de données à partir d'un ensemble homogène de type préexistant 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_INT`, `MPI_FLOAT`, ...) mais un type plus complexe construit à l'aide des sous-programmes MPI, parce qu'alors le pas ne peut plus être exprimé en nombre d'éléments du type générique.

```
int MPI_Type_create_hvector(int nombre_bloc, int nbr_elt_par_bloc, MPI_Aint pas,
                           MPI_Datatype type_elt, MPI_Datatype *nouveau_type)
```

# Types de données dérivés

## Validation des types de données dérivés

- Les types dérivés doivent être validés avant d'être utilisés dans une communication. La validation s'effectue à l'aide du sous-programme `MPI_Type_commit()`.

```
int MPI_Type_commit(MPI_Datatype *nouveau_type)
```

- Si on souhaite réutiliser le même nom pour définir un autre type dérivé, on doit au préalable le libérer en utilisant le sous-programme `MPI_Type_free()`.

```
int MPI_Type_free(MPI_Datatype *nouveau_type)
```

# Types de données dérivés

```
1  /* ligne */
2  #include <mpi.h>
3  #include <stdio.h>
4  #include <stdlib.h>
5
6  int main(int argc, char *argv[]) {
7      int rang, i, j;
8      int nb_lignes=6, nb_colonnes=5, etiquette=100;
9      float a[nb_lignes][nb_colonnes];
10     MPI_Datatype type_ligne;
11     MPI_Status statut;
12
13     MPI_Init(&argc, &argv);
14     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
15
16     /* Initialisation de la matrice sur chaque processus */
17     for(i=0; i<nb_lignes; i++)
18         for(j=0; j<nb_colonnes; j++)
19             a[i][j]=rang;
20
21     /* Definition du type type_ligne */
22     MPI_Type_contiguous(nb_colonnes, MPI_FLOAT, &type_ligne);
23
24     /* Validation du type type_ligne */
25     MPI_Type_commit(&type_ligne);
```

# Types de données dérivés

```
26  /* Envoi de la premiere ligne */
27  if (rang == 0) {
28      MPI_Send(a,l,type_ligne,l,etiquette,MPI_COMM_WORLD);
29
30  /* Reception dans la derniere ligne */
31  } else {
32      MPI_Recv(&(a[nb_lignes-1][0]),nb_colonnes,MPI_FLOAT,0,etiquette,
33              MPI_COMM_WORLD,&statut); }
34
35  /* Libere le type */
36  MPI_Type_free(&type_ligne);
37
38  MPI_Finalize();
39  }
```

# Types de données dérivés

```
1  /* colonne */
2  #include <mpi.h>
3  #include <stdio.h>
4  #include <stdlib.h>
5
6  int main(int argc, char *argv[]) {
7      int rang, i, j;
8      int nb_lignes=6, nb_colonnes=5, etiquette=100;
9      float a[nb_lignes][nb_colonnes];
10     MPI_Datatype type_colonne;
11     MPI_Status statut;
12
13     MPI_Init(&argc, &argv);
14     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
15
16     /* Initialisation de la matrice sur chaque processus */
17     for(i=0; i<nb_lignes; i++)
18         for(j=0; j<nb_colonnes; j++)
19             a[i][j]=rang;
20
21     /* Definition du type type_colonne */
22     MPI_Type_vector(nb_lignes, 1, nb_colonnes, MPI_FLOAT, &type_colonne);
23
24     /* Validation du type type_colonne */
25     MPI_Type_commit(&type_colonne);
```

# Types de données dérivés

```
26  /* Envoi */
27  if (rang == 0) {
28      MPI_Send(&a[0][1], nb_lignes, MPI_FLOAT, 1, etiquette, MPI_COMM_WORLD);
29
30      /* Reception dans l'avant-derniere colonne */
31  } else {
32      MPI_Recv(&a[0][nb_colonnes-2], 1, type_colonne, 0, etiquette,
33              MPI_COMM_WORLD, &statut); }
34
35      /* Libere le type type_colonne */
36      MPI_Type_free(&type_colonne);
37
38      MPI_Finalize();
39  }
```

# Types de données dérivés

```
1  /* bloc */
2  #include <mpi.h>
3  #include <stdio.h>
4  #include <stdlib.h>
5
6  int main(int argc, char *argv[]) {
7      int rang, i, j;
8      int nb_lignes=6, nb_colonnes=5, etiquette=100;
9      int nb_lignes_bloc=3, nb_colonnes_bloc=2;
10     float a[nb_lignes][nb_colonnes];
11     MPI_Datatype type_bloc;
12     MPI_Status statut;
13
14     MPI_Init (&argc, &argv);
15     MPI_Comm_rank (MPI_COMM_WORLD, &rang);
16
17     /* Initialisation de la matrice sur chaque processus */
18     for (i=0; i<nb_lignes; i++)
19         for (j=0; j<nb_colonnes; j++)
20             a[i][j]=rang;
21
22     /* Definition du type type_bloc */
23     MPI_Type_vector (nb_lignes_bloc, nb_colonnes_bloc, nb_colonnes,
24                     MPI_FLOAT, &type_bloc);
25
26     /* Validation du type type_bloc */
27     MPI_Type_commit (&type_bloc);
```

# Types de données dérivés

```
28  /* Envoi d'un bloc */
29  if (rang == 0) {
30      MPI_Send(a,1,type_bloc,1,etiquette,MPI_COMM_WORLD);
31
32      /* Reception du bloc */
33  } else {
34      MPI_Recv (&(a[nb_lignes-3][nb_colonnes-2]),1,type_bloc,0,etiquette,
35              MPI_COMM_WORLD,&statut); }
36
37      /* Libere du type type_bloc */
38      MPI_Type_free (&type_bloc);
39
40      MPI_Finalize();
41  }
```

# Types de données dérivés

## 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_create_hindexed()` 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_INT`, `MPI_FLOAT`, ...) mais un type plus complexe construit avec les sous-programmes MPI vus précédemment. On ne peut exprimer alors le pas en nombre d'éléments du type générique d'où le recours à `MPI_Type_create_hindexed()`.
- Pour `MPI_Type_create_hindexed()`, comme pour `MPI_Type_create_hvector()`, utilisez `MPI_Type_size()` ou `MPI_Type_get_extent()` pour obtenir de façon portable la taille du pas en nombre d'octets.

## Types de données dérivés

nb=3, longueurs\_blocs=(2,1,3), déplacements=(0,3,7)

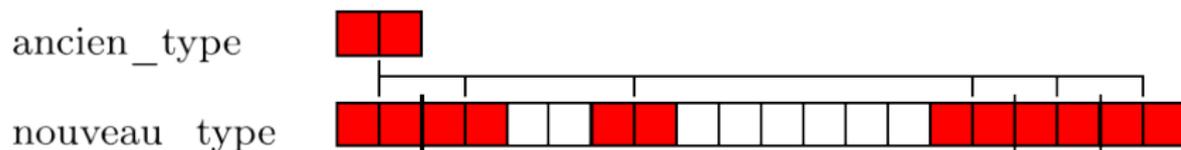
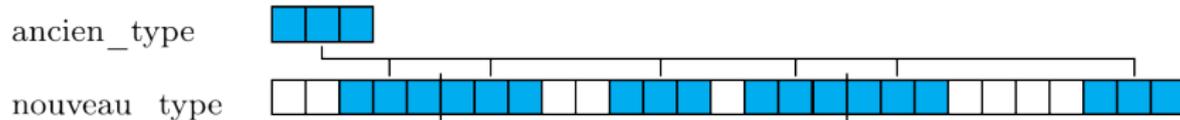


Figure 24 – Le constructeur `MPI_Type_indexed`

```
int MPI_Type_indexed(int nb, const int longueurs_blocs[], const int déplacements[],  
                    MPI_Datatype ancien_type, MPI_Datatype *nouveau_type)
```

## Types de données dérivés

nb=4, longueurs\_blocs=(2,1,2,1), déplacements=(2,10,14,24)



**Figure 25** – Le constructeur `MPI_Type_create_hindexed`

```
int MPI_Type_create_hindexed(int nb, const int longueurs_blocs[],
                             const MPI_Aint déplacements,
                             MPI_Datatype ancien_type,
                             MPI_Datatype *nouveau_type)
```

# Types de données dérivés

## Exemple : matrice triangulaire

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. Cela se fait avec l'instruction `MPI_Sendrecv_replace()` ;
4. libère ses ressources et quitte MPI.

# Types de données dérivés

AVANT

APRÈS

Processus 0

1	2	3	4	5	6	7	8
9	10	11	12	13	14	15	16
17	18	19	20	21	22	23	24
25	26	27	28	29	30	31	32
33	34	35	36	37	38	39	40
41	42	43	44	45	46	47	48
49	50	51	52	53	54	55	56
57	58	59	60	61	62	63	64

1	2	3	4	5	6	7	8
-2	10	11	12	13	14	15	16
-3	-4	19	20	21	22	23	24
-5	-6	-7	28	29	30	31	32
-8	-11	-12	-13	37	38	39	40
-14	-15	-16	-20	-21	46	47	48
-22	-23	-24	-29	-30	-31	55	56
-32	-38	-39	-40	-47	-48	-56	64

Processus 1

-1	-2	-3	-4	-5	-6	-7	-8
-9	-10	-11	-12	-13	-14	-15	-16
-17	-18	-19	-20	-21	-22	-23	-24
-25	-26	-27	-28	-29	-30	-31	-32
-33	-34	-35	-36	-37	-38	-39	-40
-41	-42	-43	-44	-45	-46	-47	-48
-49	-50	-51	-52	-53	-54	-55	-56
-57	-58	-59	-60	-61	-62	-63	-64

-1	9	17	18	25	26	27	33
-9	-10	34	35	36	41	42	43
-17	-18	-19	44	45	49	50	51
-25	-26	-27	-28	52	53	54	57
-33	-34	-35	-36	-37	58	59	60
-41	-42	-43	-44	-45	-46	61	62
-49	-50	-51	-52	-53	-54	-55	63
-57	-58	-59	-60	-61	-62	-63	-64

# Types de données dérivés

```
1  /* triangle */
2  #include <mpi.h>
3  #include <stdio.h>
4  #include <stdlib.h>
5
6  int main(int argc, char *argv[]) {
7      int rang, i, j;
8      int n=8, etiquette=100, sign=1;
9      float a[n][n];
10     MPI_Datatype type_triangle;
11     MPI_Status statut;
12     int longueurs_blocs[n], deplacements[n];
13
14     MPI_Init(&argc, &argv);
15     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
16
17     /* Initialisation de la matrice sur chaque processus */
18     if (rang == 1) sign=-1;
19     for (i=0; i<n; i++)
20         for (j=0; j<n; j++)
21             a[i][j]=sign*(1+i*n+j);
22
23     /* Creation du type matrice triangulaire inf pour le processus 0
24     et du type matrice triangulaire sup pour le processus 1 */
25     if (rang == 0) {
26         for (i=0; i<n; i++) longueurs_blocs[i] = i;
27         for (i=0; i<n; i++) deplacements[i] = n*i;
28     } else {
29         for (i=0; i<n; i++) longueurs_blocs[i] = n-i-1;
30         for (i=0; i<n; i++) deplacements[i] = (n+1)*i+1;
31     }
32
33     MPI_Type_indexed(n, longueurs_blocs, deplacements, MPI_FLOAT, &type_triangle);
34     MPI_Type_commit(&type_triangle);
35
36     /* Permutation des matrices triangulaires inferieure et superieure */
37     MPI_Sendrecv_replace(a, 1, type_triangle, (rang+1)%2, etiquette,
38                          (rang+1)%2, etiquette, MPI_COMM_WORLD, &statut);
39
40     /* Liberation du type triangle */
41     MPI_Type_free(&type_triangle);
42     MPI_Finalize();
43 }
```

# Types de données dérivés

## Taille des types de données

- `MPI_Type_size()` retourne le nombre d'octets nécessaire pour envoyer un type de données. Cette valeur ne tient pas compte des *trous* présents dans le type de données.

```
int MPI_Type_size(MPI_Datatype type_message, int *taille)
```

- L'étendue d'un type est l'espace mémoire occupé par le type (en octets). Cette valeur intervient directement pour calculer la position du prochain élément en mémoire (c'est-à-dire le **pas** entre des éléments successifs).

```
int MPI_Type_get_extent(MPI_Datatype type_message, MPI_Aint *borne_inf,  
MPI_Aint *etendue)
```

# Types de données dérivés

Exemple 1 : `MPI_Type_indexed(2, {2, 1}, {1, 4}, MPI_INT, &type)`

Type dérivé :



Deux éléments successifs :



`taille = 12` (3 entiers); `borne_inf = 4` (1 entier); `etendue = 16` (4 entiers)

Exemple 2 : `MPI_Type_vector(3, 1, nb_colonnes, MPI_INT, &type_demi_colonne)`

Vue 2D :

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

Vue 1D :

1	2	3	4	5	6	7	8	9	10	11
---	---	---	---	---	---	---	---	---	----	----

`taille = 12` (3 entiers); `borne_inf = 0`; `etendue = 44` (11 entiers)

# Types de données dérivés

## Changer l'étendue

- L'étendue est un paramètre du type de données. Par défaut, c'est généralement l'intervalle en mémoire entre le premier et le dernier composant du type (bornes incluses et en tenant compte de l'alignement mémoire). On peut modifier l'étendue d'un type pour créer un nouveau type adapté du précédent avec `MPI_Type_create_resized()`. Cela permet de choisir le pas entre des éléments successifs.

```
int MPI_Type_create_resized(MPI_Datatype ancien_type, MPI_Aint nouvelle_borne_inf,  
                             MPI_Aint nouvelle_etendue, MPI_Datatype *nouveau_type)
```

# Types de données dérivés

```
1  /* demi_colonne */
2  #include <mpi.h>
3  #include <stdio.h>
4  #include <stdlib.h>
5
6  int main(int argc, char *argv[]) {
7      int rang, i, j, taille_integer;
8      int nb_lignes=6, nb_colonnes=5, etiquette=100, sign=1;
9      int taille_demi_colonne=nb_lignes/2;
10     int a[nb_lignes][nb_colonnes];
11     MPI_Datatype type_demi_colonne1, type_demi_colonne2;
12     MPI_Status statut;
13     MPI_Aint borne_inf1, etendue1, borne_inf2, etendue2;
14
15     MPI_Init(&argc, &argv);
16     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
17
18     /* Initialisation de la matrice sur chaque processus */
19     if (rang == 1) sign=-1;
20     for (i=0; i<nb_lignes; i++)
21         for (j=0; j<nb_colonnes; j++)
22             a[i][j]=sign*(1+nb_colonnes*i+j);
23
24     /* Construction du type derive type_demi_colonne1 */
25     MPI_Type_vector(taille_demi_colonne, 1, nb_colonnes, MPI_INT, &type_demi_colonne1);
26
27     /* Connaitre la taille du type de base MPI_INT */
28     MPI_Type_size(MPI_INT, &taille_integer);
29
30     /* Informations sur le type type_demi_colonne1 */
31     MPI_Type_get_extent(type_demi_colonne1, &borne_inf1, &etendue1);
32     if (rang == 0) printf("Type_demi_colonne1: borne_inf=%d etendue %d\n",
33                          borne_inf1, etendue1);
34
35     /* Construction du type type_demi_colonne2 */
36     borne_inf2 = 0;
37     etendue2 = taille_integer;
38     MPI_Type_create_resized(type_demi_colonne1, borne_inf2, etendue2,
39                            &type_demi_colonne2);
```

# Types de données dérivés

```
40 /* Information sur le type type_demi_colonne2 */
41 MPI_Type_get_extent(type_demi_colonne2,&borne_inf2,&etendue2);
42 if (rang == 0) printf("Type_demi_colonne2: borne_inf=%d etendue %d\n",
43                       borne_inf2,etendue2);
44
45 /* Validation du type */
46 MPI_Type_commit(&type_demi_colonne2);
47
48 if (rang == 0) {
49     /* Envoi de la matrice a au processus 1 avec le type demi_colonne2 */
50     MPI_Send(a,2,type_demi_colonne2,1,etiquette,MPI_COMM_WORLD);
51 } else {
52     /* Reception pour le processus 1 dans la matrice a */
53     MPI_Recv(&a[nb_lignes-2][0],6,MPI_INT,0,etiquette,MPI_COMM_WORLD,&statut);
54     printf("Matrice A sur le processus 1\n");
55     for(i=0;i<nb_lignes;i++) {
56         for(j=0;j<nb_colonnes;j++)
57             printf("%d ",a[i][j]);
58         printf("\n"); } }
59
60 MPI_Finalize();
61 }
```

```
> mpiexec -n 2 demi_ligne
type_demi_colonne1: borne_inf=0, etendue=44
type_demi_colonne2: borne_inf=0, etendue=4
```

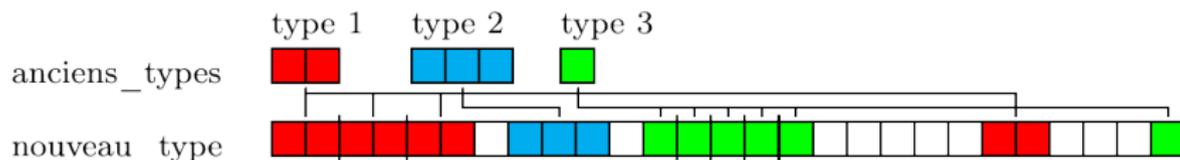
```
Matrice A sur le processus 1
-1 -2 -3 -4 -5
-6 -7 -8 -9 -10
-11 -12 -13 -14 -15
-16 -17 -18 -19 -20
 1  6 11  2  7
12 -27 -28 -29 -30
```

# Types de données dérivés

## Types hétérogènes

- Le sous-programme `MPI_Type_create_struct()` permet de créer une séquence de blocs de données en précisant le `type`, le `nombre d'éléments` et le `pas` de chaque bloc.
- Il s'agit du constructeur de types le plus complet. Il généralise `MPI_Type_indexed()` en permettant de définir un `type` différent pour chaque bloc.

`nb=5`, `longueurs_blocs=(3,1,5,1,1)`, `déplacements=(0,7,11,21,26)`,  
`anciens_types=(type1,type2,type3,type1,type3)`



```
int MPI_Type_create_struct(int nb, const int longueurs_blocs[], const MPI_Aint déplacements[],  
                          const MPI_Datatype anciens_types[], MPI_Datatype *nouveau_type)
```

# Types de données dérivés

## Calcul des déplacements

- `MPI_Type_create_struct()` est utile notamment pour créer des types MPI correspondant à des types dérivés Fortran ou à des structures C.
- L'alignement en mémoire des structures de données hétérogènes dépend de l'architecture et du compilateur.
- Attention, il faut corriger l'étendue du type MPI obtenu.
- Le sous-programme `MPI_Get_address()` permet de récupérer l'adresse d'une variable. C'est l'équivalent de l'opérateur de référencement (&) du C.
- Attention, même en C, il est préférable d'utiliser ce sous-programme MPI pour des raisons de portabilité.
- Il est conseillé d'utiliser `MPI_Aint_add()` et `MPI_Aint_diff()` pour faire des additions et soustractions sur des adresses.

```
int MPI_Get_address(const void *variable, MPI_Aint *adresse_variable)
MPI_Aint MPI_Aint_add(MPI_Aint base, MPI_Aint disp)
MPI_Aint MPI_Aint_diff(MPI_Aint addr1, MPI_Aint addr2)
```

# Types de données dérivés

```
1  /* Interaction_Particules */
2  #include <mpi.h>
3  #include <stdio.h>
4  #include <stdlib.h>
5  #include <stdbool.h>
6
7  struct Particule {
8      char categorie[5];
9      int masse;
10     float coords[3];
11     bool classe;
12 };
13
14 int main(int argc, char *argv[]) {
15     int rang, i;
16     int n=1000, etiquette=100;
17     int longueurs_blocs[4];
18     MPI_Datatype types[4], type_particule, temp;
19     MPI_Status statut;
20     MPI_Aint adresses[5], deplacements[5], lb, extent;
21     struct Particule p[n], temp_p[n];
22
23     MPI_Init(&argc, &argv);
24     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
25
26     /* Construction du type de donnees */
27     types[0] = MPI_CHARACTER; types[1] = MPI_INT;
28     types[2] = MPI_FLOAT; types[3] = MPI_LOGICAL;
29     longueurs_blocs[0]=5; longueurs_blocs[1]=1;
30     longueurs_blocs[2]=3; longueurs_blocs[3]=1;
```

# Types de données dérivés

```
31 MPI_Get_address(&(p[0]),&(adresses[0]));
32 MPI_Get_address(&(p[0].categorie),&(adresses[1]));
33 MPI_Get_address(&(p[0].masse),&(adresses[2]));
34 MPI_Get_address(&(p[0].coords),&(adresses[3]));
35 MPI_Get_address(&(p[0].classe),&(adresses[4]));
36 /* Calcul des déplacements relatifs a l'adresse de depart */
37 for (i=0;i<4;i++) deplacements[i] = MPI_Aint_diff(adresses[i+1],adresses[0]);
38 MPI_Type_create_struct(4,longueurs_blocs,deplacements,types,&temp);
39 MPI_Get_address(&(p[1]),&(adresses[1]));
40 lb=0;
41 extent = MPI_Aint_diff(adresses[1],adresses[0]);
42 MPI_Type_create_resized(temp,lb,extent,&type_particule);
43 /* Validation du type structure */
44 MPI_Type_commit(&type_particule);
45
46 /* Initialisation des particules pour chaque processus */
47
48 /* Envoi des particules de 0 vers 1 */
49 if (rang == 0) {
50     MPI_Send(&(p[0]),n,type_particule,1,etiquette,MPI_COMM_WORLD);
51 } else {
52     MPI_Recv(&(temp_p[0]),n,type_particule,0,etiquette,
53            MPI_COMM_WORLD,&statut);
54 }
55
56 /* Liberation du type */
57 MPI_Type_free(&type_particule);
58 MPI_Finalize();
59 }
```

# Types de données dérivés

## 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 dans l'un des prochains chapitres) fait de MPI l'outil idéal pour tous les problèmes de décomposition de domaine avec des maillages réguliers ou irréguliers.

# Types de données dérivés

## Memento

Sous-routines	longueurs_blocs	pas	types_anciens
<code>MPI_Type_Contiguous()</code>	constant*	constant*	constant
<code>MPI_Type_[Create_H]Vector()</code>	constant	constant	constant
<code>MPI_Type_[Create_H]Indexed()</code>	<i>variable</i>	<i>variable</i>	constant
<code>MPI_Type_Create_Struct()</code>	<i>variable</i>	<i>variable</i>	<i>variable</i>

(\*) paramètre caché, égal à 1

## Travaux pratiques MPI – Exercice 4 : Transposée d'une matrice

- Dans cet exercice, on se propose de se familiariser avec les types dérivés
- On se donne une matrice  $A$  de 4 lignes et 5 colonnes sur le processus 0
- Il s'agit pour le processus 0 d'envoyer au processus 1 cette matrice mais d'en faire automatiquement la transposition au cours de l'envoi

1.	2.	3.	4.	5.
6.	7.	8.	9.	10.
11.	12.	13.	14.	15.
16.	17.	18.	19.	20.

Processus 0



1.	6.	11.	16.
2.	7.	12.	17.
3.	8.	13.	18.
4.	9.	14.	19.
5.	10.	15.	20.

Processus 1

- Pour ce faire, on va devoir se construire deux types dérivés, un type `type_colonne` et un type `type_transpose`

## Travaux pratiques MPI – Exercice 5 : Produit réparti de matrices

- Communications collectives et réductions : produit de matrices  $C = A \times B$ 
  - On se limite au cas de matrices carrées dont l'ordre est un multiple du nombre de processus
  - Les matrices  $A$  et  $B$  sont sur le processus 0. Celui-ci distribue une tranche horizontale de la matrice  $A$  et une tranche verticale de la matrice  $B$  à chacun des processus. Chacun calcule alors un bloc diagonal de la matrice résultante  $C$ .
  - Pour calculer les blocs non diagonaux, chaque processus doit envoyer aux autres processus la tranche de  $A$  qu'il possède
  - Après quoi le processus 0 peut collecter les résultats et vérifier les résultats

# Travaux pratiques MPI – Exercice 5 : Produit réparti de matrices

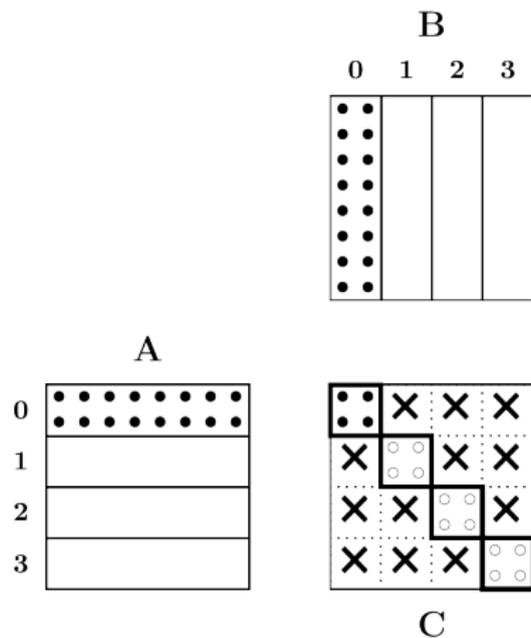


Figure 26 – Produit parallèle de matrices

## Travaux pratiques MPI – Exercice 5 : Produit réparti de matrices

- Toutefois, l'algorithme qui peut sembler le plus immédiat, et qui est le plus simple à programmer, consistant à faire envoyer par chaque processus sa tranche de la matrice A à chacun des autres, n'est pas performant parce que le schéma de communication n'est pas du tout équilibré. C'est très facile à voir en faisant des mesures de performances et en représentant graphiquement les traces collectées.



Figure 27 – Produit parallèle de matrices sur 16 processus, pour une taille de matrice de 1024 (premier algorithme)

## Travaux pratiques MPI – Exercice 5 : Produit réparti de matrices

- Mais en changeant l'algorithme pour faire *glisser* le contenu des tranches de processus à processus, on peut obtenir un équilibre parfait des calculs et des communications, et gagner ainsi un facteur 2.



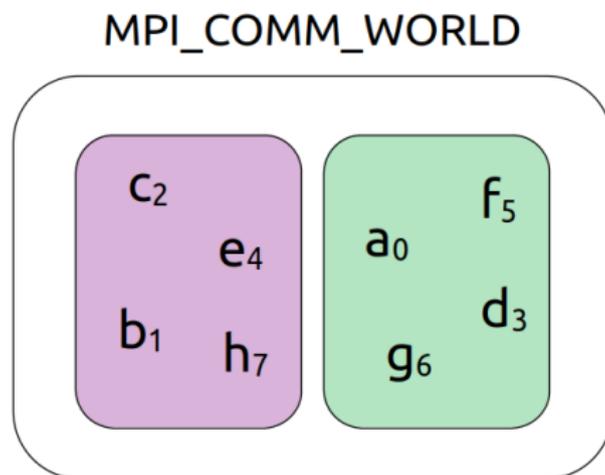
**Figure 28** – Produit parallèle de matrices sur 16 processus, pour une taille de matrice de 1024 (second algorithme)

# Communicateurs

# Communicateurs

## Introduction

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



**Figure 29** – Partitionnement d'un communicateur

## Exemple

Par exemple, on veut diffuser un message collectif aux processus de rang pair et un autre aux processus de rang impair.

- 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 émetteur doit envoyer le message est pair ou impair.
- Une solution est de créer un communicateur regroupant les processus pairs et un autre regroupant les processus impairs, puis d'initier les communications collectives à l'intérieur de ces groupes.

# Communicateurs

## Communicateur par défaut

- On ne peut créer un communicateur qu'à partir d'un autre communicateur. Le premier sera créé à partir de `MPI_COMM_WORLD`.
- En effet, suite à l'appel à `MPI_Init()`, un communicateur est créé pour toute la durée d'exécution du programme.
- Son identificateur `MPI_COMM_WORLD` est une variable définie dans les fichiers d'en-tête.
- Il ne peut être détruit que via l'appel à `MPI_Finalize()`
- Par défaut, il fixe donc la portée des communications point à point et collectives à tous les processus de l'application

## Groupes et communicateurs

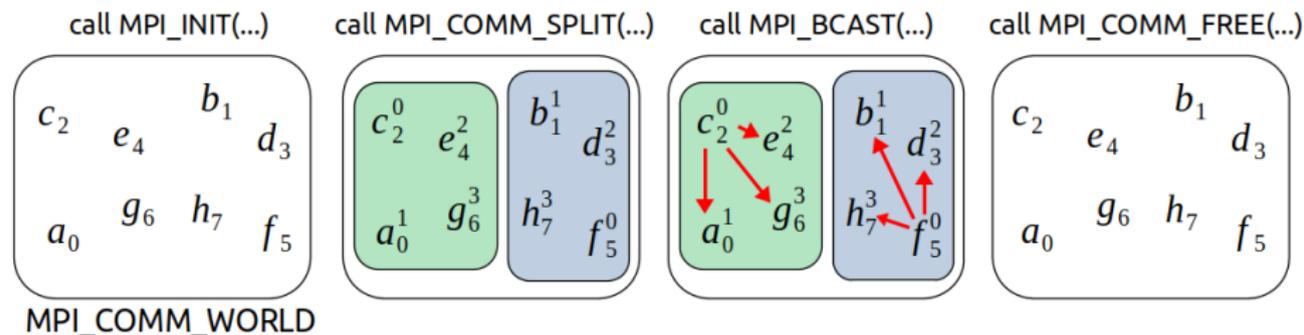
- Un communicateur est constitué :
  - d'un **groupe**, qui est un ensemble ordonné de processus ;
  - d'un **contexte** de communication mis en place à l'appel du sous-programme de construction du communicateur, qui permet de délimiter l'espace de communication.
- Les contextes de communication sont gérés par MPI (le programmeur n'a aucune action sur eux) : c'est un attribut « caché »
- Dans la bibliothèque MPI, divers sous-programmes existent pour construire des communicateurs : `MPI_Comm_create()`, `MPI_Comm_dup()`, `MPI_Comm_split()`
- Les **constructeurs de communicateurs** sont des **opérateurs collectifs** (qui engendrent des communications entre les processus)
- Les communicateurs que le programmeur crée peuvent être gérés dynamiquement et, de même qu'il est possible d'en créer, il est possible d'en détruire en utilisant le sous-programme `MPI_Comm_free()`

# Communicateurs

## Partitionnement d'un communicateur

Pour résoudre le problème de l'exemple, nous allons :

- partitionner le communicateur en processus de rang pair et d'autre part en processus de rang impair ;
- ne diffuser un message collectif qu'aux processus de rang pair et un autre qu'aux processus de rang impair.



**Figure 30** – Création/destruction d'un communicateur

# Communicateurs

## Partitionnement d'un communicateur avec `MPI_Comm_split()`

Le sous-programme `MPI_Comm_split()` permet de :

- partitionner un communicateur donné en autant de communicateurs que l'on veut
- donner le même nom à tous ces communicateurs : il aura la valeur du communicateur dans lequel se trouve le processus courant
- Méthode :
  1. définir une valeur couleur associant à chaque processus le numéro du communicateur auquel il appartiendra
  2. définir une valeur clef permettant de numéroter les processus dans chaque communicateur
  3. créer la partition où chaque communicateur s'appelle `nouveau_comm`

```
int MPI_Comm_split(MPI_Comm comm, int couleur, int clef, MPI_Comm *nouveau_comm)
```

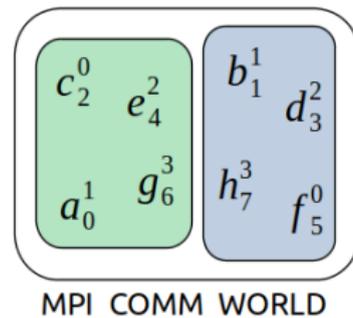
Un processus qui s'attribue une couleur égale à la valeur `MPI_UNDEFINED` aura pour `nouveau_com` le communicateur invalide `MPI_COMM_NULL`.

# Communicateurs

## Exemple

Voyons comment procéder pour construire le communicateur qui va subdiviser l'espace de communication entre processus de rangs pairs et impairs, via le constructeur `MPI_Comm_split()`.

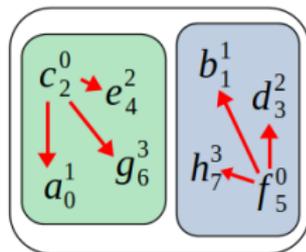
processus	a	b	c	d	e	f	g	h
rang_monde	0	1	2	3	4	5	6	7
couleur	0	1	0	1	0	1	0	1
clef	0	1	-1	3	4	-1	6	7
rang_pairs_imp	1	1	0	2	2	0	3	3



**Figure 31** – Construction du communicateur `CommPairsImpairs` avec `MPI_Comm_split()`

# Communicateurs

```
1  /* PairsImpairs */
2  #include <mpi.h>
3  #include <stdlib.h>
4
5  int main(int argc, char *argv[]) {
6      int rang, i, clef, rang_dans_monde;
7      int m=16;
8      float a[16];
9      MPI_Comm CommPairsImpairs;
10
11     MPI_Init (&argc, &argv);
12     MPI_Comm_rank (MPI_COMM_WORLD, &rang_dans_monde);
13
14     for (i=0; i<m; i++) a[i]=0.;
15     if (rang_dans_monde == 2) {for (i=0; i<m; i++) a[i]=2.;}
16     if (rang_dans_monde == 5) {for (i=0; i<m; i++) a[i]=5.;}
17
18     clef = rang_dans_monde;
19     if ((rang_dans_monde == 2) || (rang_dans_monde == 5)) {
20         clef = -1; }
21
22     /* Creation des communicateurs pair et impair en leur donnant une meme denomination */
23     MPI_Comm_split (MPI_COMM_WORLD, rang_dans_monde%2, clef, &CommPairsImpairs);
24
25     /* Diffusion du message par le processus 0 de chaque communicateur aux processus
26        de son groupe */
27     MPI_Bcast (a, m, MPI_FLOAT, 0, CommPairsImpairs);
28
29     /* Destruction des communicateurs */
30     MPI_Comm_free (&CommPairsImpairs);
31     MPI_Finalize ();
32 }
```

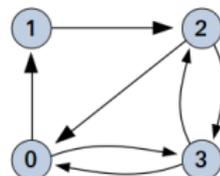


# Communicateurs

## Topologies

- Dans la plupart des applications, plus particulièrement dans les méthodes de décomposition de domaine où l'on fait correspondre le domaine de calcul à la grille de processus, il est intéressant de pouvoir disposer les processus suivant une topologie régulière
- MPI permet de définir des topologies virtuelles du type cartésien ou graphe
  - Topologies de type cartésien
    - ▶ chaque processus est défini dans une grille de processus ;
    - ▶ chaque processus a un voisin dans la grille ;
    - ▶ 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.

1	3	5	7
0	2	4	6



**Figure 32** – Topologie cartésienne 2D (gauche) et topologie de type graphe (droite)

# Communicateurs

## Topologies cartésiennes

- Une topologie cartésienne est définie à partir d'un communicateur donné `comm_ancien`, en appelant le sous-programme `MPI_Cart_create()`.
- On définit :
  - un entier `ndims` représentant le nombre de dimensions de la grille
  - un tableau d'entiers `dims` de dimension `ndims` indiquant le nombre de processus dans chaque dimension
  - un tableau de logiques de dimension `ndims` indiquant la périodicité dans chaque dimension
  - un logique `reorganisation` indiquant si la numérotation des processus peut être changé par MPI

```
int MPI_Cart_create(MPI_Comm comm_ancien, int ndims, const int *dims, const int *periods,  
int reorganisation, MPI_Comm *comm_nouveau)
```

## Exemple

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

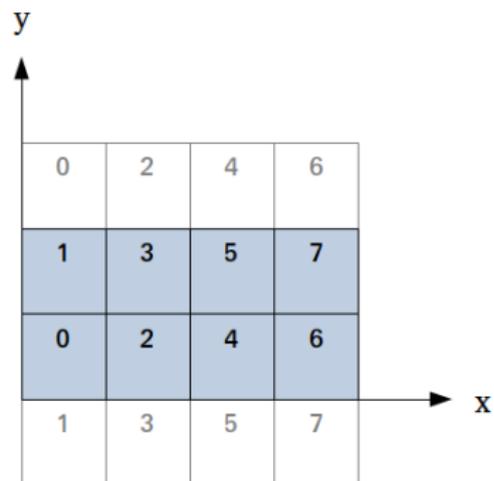
```
int ndims=2, reorganisation;
int dims[ndims], periods[ndims];
MPI_Comm comm_2D;

dims[0]=4; dims[1]=2;
periods[0]=false; periods[1]=true;
reorganisation=false;

MPI_Cart_create(MPI_COMM_WORLD, ndims, dims, periods, reorganisation, &comm_2D);
```

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`, l'implémentation MPI choisit l'ordre des processus.



**Figure 33** – Topologie cartésienne 2D périodique en y

## Exemple 3D

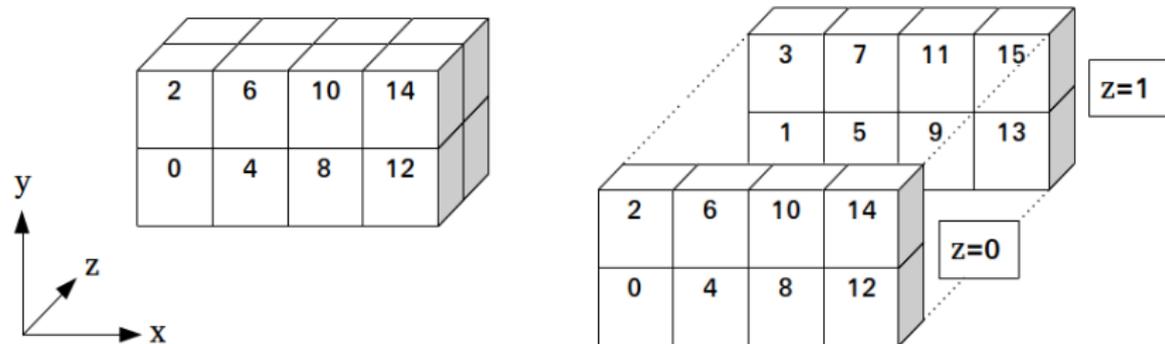
Exemple sur une grille 3D comportant 4 domaines suivant x, 2 suivant y et 2 suivant z, non périodique.

```
int ndims=3, reorganisation;
int dims[ndims], periods[ndims];
MPI_Comm comm_3D;

dims[0]=4; dims[1]=2; dims[2]=2;
periods[0]=false; periods[1]=false; periods[2]=false;
reorganisation=false;

MPI_Cart_create(MPI_COMM_WORLD, ndims, dims, periods, reorganisation, &comm_3D);
```

# Communicateurs



**Figure 34** – Topologie cartésienne 3D non périodique

## Distribution des processus

Le sous-programme `MPI_Dims_create()` retourne le nombre de processus dans chaque dimension de la grille en fonction du nombre total de processus.

```
int MPI_Dims_create(int nb_procs, int ndims, int *dims)
```

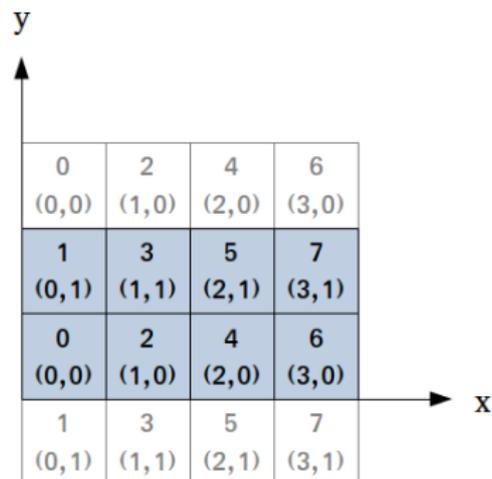
Remarque : si les valeurs de `dims` en entrée valent toutes 0, cela signifie qu'on laisse à MPI le choix du nombre de processus dans chaque direction en fonction du nombre total de processus.

dims en entrée	<code>MPI_Dims_create</code>	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

# Communicateurs

## Rang et coordonnées d'un processus

Dans une topologie cartésienne, le rang de chaque processus est associé à ses coordonnées dans la grille.

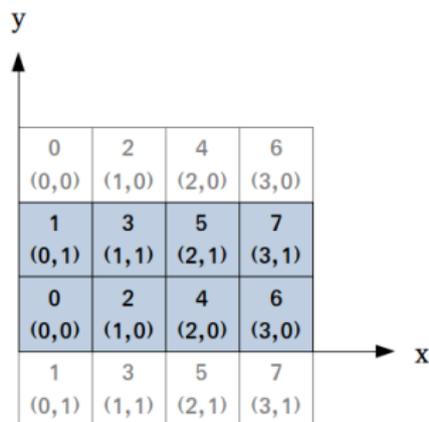


**Figure 35** – Topologie cartésienne 2D périodique en y

## Rang d'un processus connaissant ses coordonnées

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

```
int MPI_Cart_rank(MPI_Comm comm, const int coords[], int *rang)
```



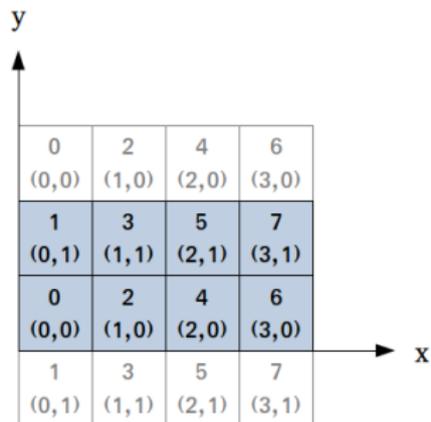
**Figure 36** – Topologie cartésienne 2D périodique en y

```
coords[0]=dims[0]-1;
for(i=0;i<dims[1];i++) {
    coords[1]=i;
    MPI_Cart_rank(comm_2D, coords, &(rang[i]));
}
.....
i=0, en entree coords=[3,0], en sortie rang[0]=6.
i=1, en entree coords=[3,1], en sortie rang[1]=7.
```

## Coordonnées d'un processus connaissant son rang

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

```
int MPI_Cart_coords(MPI_Comm comm, int rang, int ndims, int *coords)
```



**Figure 37** – Topologie cartésienne 2D périodique en y

```
if (rang%2 == 0)
    MPI_Cart_coords(comm_2D, rang, 2, &coords);
.....
En entree, les valeurs de rang sont : 0,2,4,6.
En sortie, les valeurs de coords sont :
(0,0), (1,0), (2,0), (3,0).
```

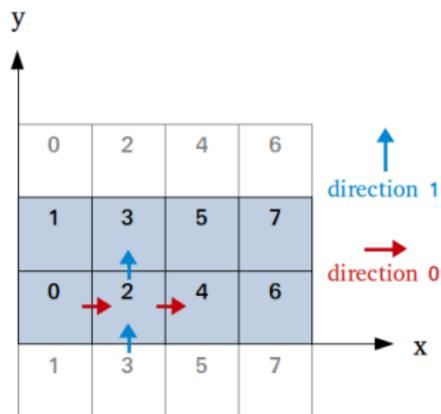
## Rang des voisins

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.

```
int MPI_Cart_shift(MPI_Comm comm, int direction, int pas, int *rang_precedent, int *rang_suivant)
```

- Le paramètre `direction` correspond à l'axe du déplacement (xyz).
- Le paramètre `pas` correspond au pas du déplacement.
- Si un rang n'a pas de voisin précédent (resp. suivant) dans la direction demandée, alors la valeur du rang précédent (resp. suivant) sera `MPI_PROC_NULL`.

# Communicateurs



**Figure 38** – Appel du sous-programme MPI\_Cart\_shift()

```
MPI_Cart_shift(comm_2D,0,1,&rang_gauche,&rang_droit);  
.....  
Pour le processus 2, rang_gauche=0, rang_droit=4
```

```
MPI_Cart_shift(comm_2D,1,1,&rang_bas,&rang_haut);  
.....  
Pour le processus 2, rang_bas=3, rang_haut=3
```

# Communicateurs

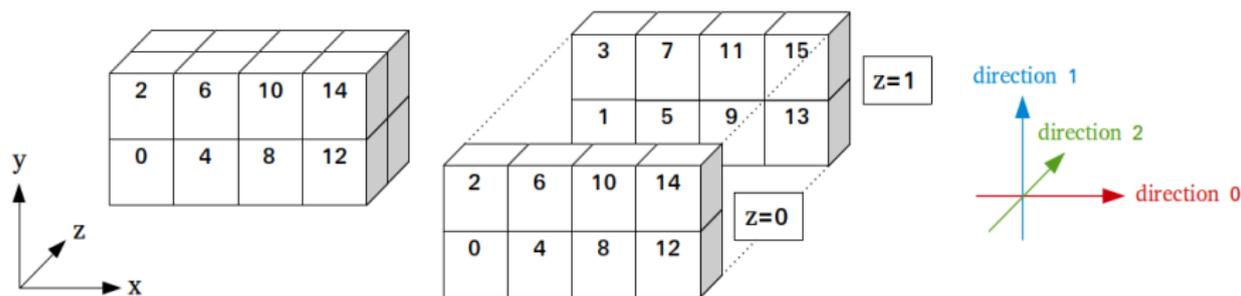


Figure 39 – Appel du sous-programme MPI\_Cart\_shift()

```
MPI_Cart_shift(comm_3D,0,1,&rang_gauche,&rang_droit)
.....
Pour le processus 0, rang_gauche=-1, rang_droit=4
```

```
MPI_Cart_shift(comm_3D,1,1,&rang_bas,&rang_haut)
.....
Pour le processus 0, rang_bas=-1, rang_haut=2
```

```
MPI_Cart_shift(comm_3D,2,1,&rang_avant,&rang_arriere)
.....
Pour le processus 0, rang_avant=-1, rang_arriere=1
```

## Exemple

- création d'une grille cartésienne 2D périodique en y
- récupération des coordonnées de chaque processus
- récupération des rangs voisins pour chaque processus

```
1  /* decomposition */
2  #include <mpi.h>
3  #include <stdlib.h>
4
5  int main(int argc, char *argv[]) {
6      int nb_procs, rang_ds_topo;
7      int ndims=2, N=1, E=2, S=3, W=4;
8      int dims[ndims], coords[ndims], voisin[4];
9      int periods[ndims], reorganisation;
10     MPI_Comm comm_2D;
11
12     MPI_Init(&argc, &argv);
13
14     MPI_Comm_size(MPI_COMM_WORLD, &nb_procs);
15
16     /* Connaitre le nombre de processus suivant x et y */
17     dims[0] = dims[1] = 0;
18
19     MPI_Dims_create(nb_procs, ndims, dims);
```

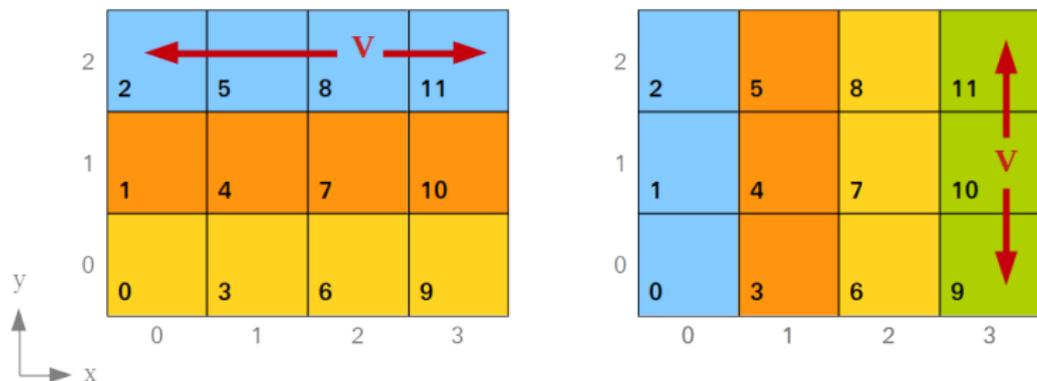
# Communicateurs

```
20  /* Creation grille 2D periodique en y */
21  periods[0] = 0;
22  periods[1] = 1;
23  reorganisation = 0;
24
25  MPI_Cart_create(MPI_COMM_WORLD, ndims, dims, periods, reorganisation, &comm_2D);
26
27  /* Connaitre mes coordonnees dans la topologie */
28  MPI_Comm_rank(comm_2D, &rang_ds_topo);
29  MPI_Cart_coords(comm_2D, rang_ds_topo, ndims, coords);
30
31  /* Recherche de mes voisins Ouest et Est */
32  MPI_Cart_shift(comm_2D, 0, 1, &(voisin[W]), &(voisin[E]));
33
34  /* Recherche des mes voisins Sud et Nord */
35  MPI_Cart_shift(comm_2D, 1, 1, &(voisin[S]), &(voisin[N]));
36
37  MPI_Finalize();
38 }
```

## Subdiviser une topologie cartésienne

- La question est de savoir comment dégénérer une topologie cartésienne de processus 2D ou 3D 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.

# Communicateurs



**Figure 40** – Deux exemples de distribution de données dans une topologie 2D dégénérée

# Communicateurs

## Subdiviser une topologie cartésienne

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.

```
int MPI_Cart_sub(MPI_Comm CommCart, const int conserve_dims[], MPI_Comm *CommCartD)
```

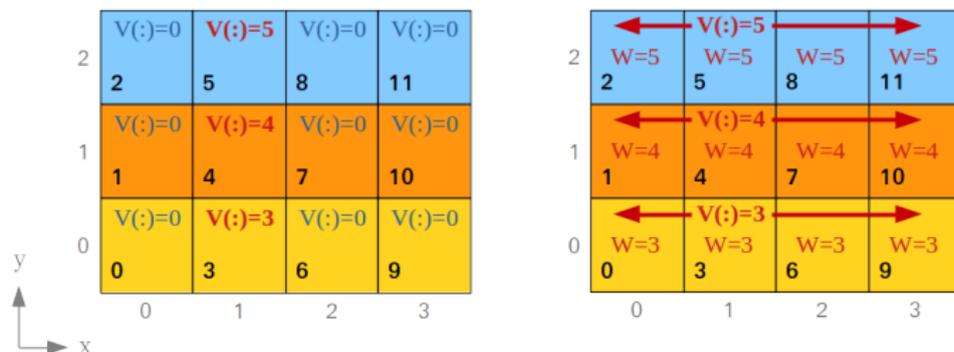


Figure 41 – Distribution d'un tableau  $V$  sur la grille 2D dégénérée

# Communicateurs

```
1  /* CommCartSub */
2  #include <mpi.h>
3  #include <stdlib.h>
4  #include <stdio.h>
5
6  int main(int argc, char *argv[]) {
7      int NDim2D=2, m=4;
8      int Dim2D[NDim2D], Periode[NDim2D], Coord2D[NDim2D], conserve_dims[NDim2D];
9      int ReOrdonne, rang, i;
10     MPI_Comm Comm2D, Comm1D;
11     float V[m], W;
12
13
14     MPI_Init(&argc, &argv);
15
16     /* Creation de la grille 2D initiale */
17     Dim2D[0] = 4;
18     Dim2D[1] = 3;
19     Periode[0] = 0; Periode[1] = 0;
20     ReOrdonne = 0;
21     MPI_Cart_create(MPI_COMM_WORLD, NDim2D, Dim2D, Periode, ReOrdonne, &Comm2D);
22     MPI_Comm_rank(Comm2D, &rang);
23     MPI_Cart_coords(Comm2D, rang, NDim2D, Coord2D);
```

# Communicateurs

```
24  /* Initialisation du vecteur V */
25  if (Coord2D[0] == 1) {for(i=0;i<m;i++) V[i]=rang; }
26
27  /* Chaque ligne de la grille doit etre une topologie cartesienne 1D */
28  conserve_dims[0] = 1;
29  conserve_dims[1] = 0;
30  /* Subdivision de la grille cartesienne 2D */
31  MPI_Cart_sub (Comm2D, conserve_dims, &Comm1D);
32
33  /* Les processus de la colonne 2 distribuent le vecteur V aux processus de leur ligne */
34  MPI_Scatter (V, 1, MPI_FLOAT, &W, 1, MPI_FLOAT, 1, Comm1D);
35
36  printf("Rang : %d ; Coordonnees : ( %d,%d); W = %f\n",
37         rang, Coord2D[0], Coord2D[1], W);
38
39  MPI_Finalize();
40 }
```

# Communicateurs

```
> mpiexec -n 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.
```

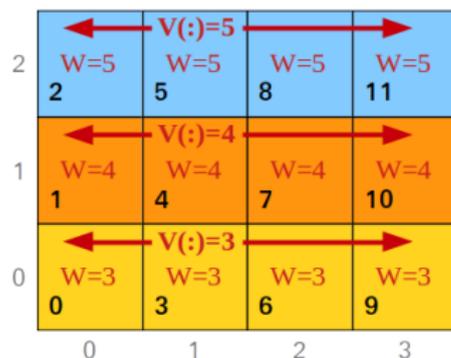
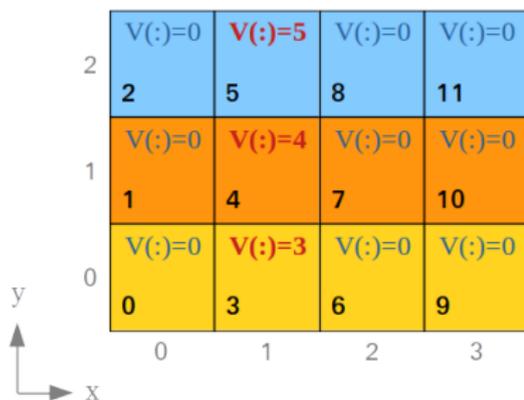
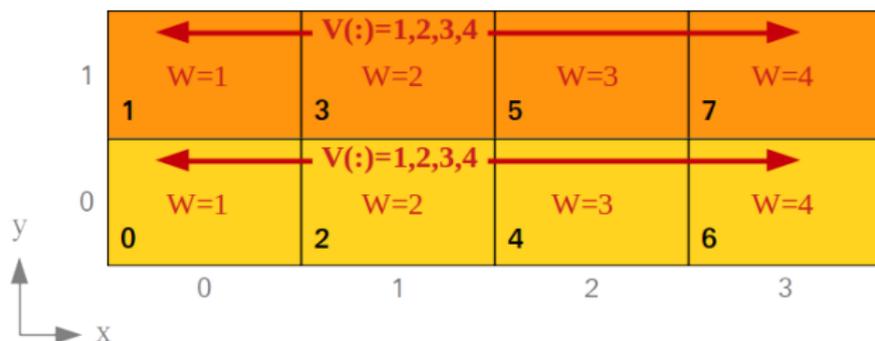


Figure 42 – Distribution d'un tableau  $V$  sur la grille 2D dégénérée

## Travaux pratiques MPI – Exercice 6 : Communicateurs

- En partant de la topologie cartésienne définie ci-dessous, subdiviser en 2 communicateurs suivant les lignes via `MPI_Comm_split()`. Ensuite, faire en sorte que les processus de la 2eme colonne diffusent sélectivement le vecteur  $V$  aux processus de leur ligne.



**Figure 43** – Subdivision d'une topologie 2D et communication suivant la topologie 1D obtenue

- Contrainte : définir les couleurs de chaque processus sans utiliser l'opération *modulo*.

# MPI-IO

## Optimisation des entrées-sorties

- Très logiquement, les applications qui font des calculs volumineux manipulent également des quantités importantes de données externes, et génèrent donc un nombre conséquent d'entrées-sorties.
- Le traitement efficace de celles-ci influe donc parfois très fortement sur les performances globales des applications.
- L'optimisation des entrées-sorties de codes parallèles se fait par la combinaison :
  - de leur **parallélisation**, pour éviter de créer un goulet d'étranglement en raison de leur sérialisation ;
  - de techniques mises en œuvre **explicitement** au niveau de la programmation (lectures / écritures non-bloquantes) ;
  - d'opérations spécifiques prises en charge par le **système d'exploitation** (regroupement des requêtes, gestion des tampons d'entrées-sorties, etc.).
- L'utilisation d'une bibliothèque facilite l'optimisation des entrées-sorties.

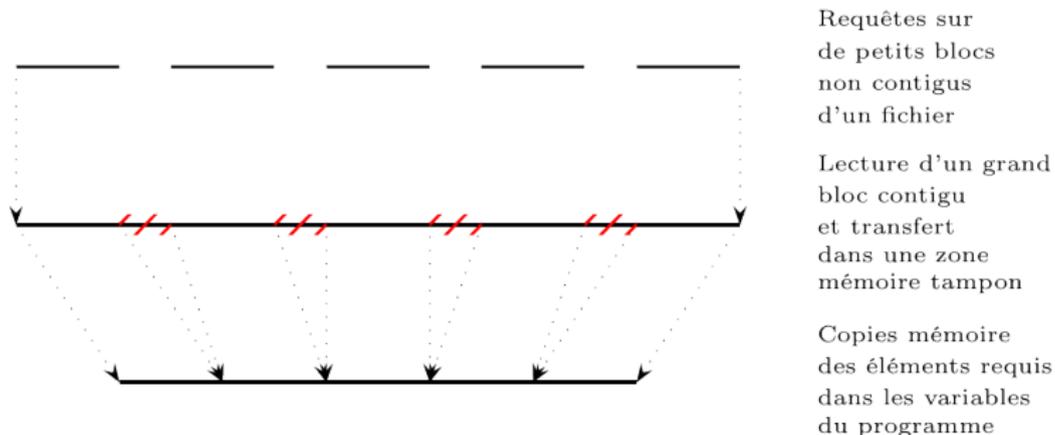
## L'interface MPI-IO

- La norme MPI-2 définit un ensemble de fonctions permettant de réaliser des entrées-sorties parallèles.
- L'interface est calquée sur celle utilisée pour l'échange de messages MPI. Par exemple, les **opérations collectives** et **non-bloquantes** sur les fichiers sont gérées de façon similaire à ce que propose **MPI** pour les messages entre processus. La définition des données accédées suivant les processus se fait par l'utilisation de **types de données** (de base ou bien dérivés).
- Bien sûr, de nombreux éléments (descripteurs de fichiers, attributs . . .) rappellent les interfaces d'entrées-sorties natives des langages de programmation.

# MPI-IO

## Exemple d'optimisation séquentielle implémentée par les bibliothèques

- Pour obtenir de bonnes performances, il est préférable de limiter le nombre de requêtes (latence) et de lire de larges blocs de données.
- Lorsqu'un seul processus accède à de nombreux petits blocs discontinus, il est possible de regrouper les requêtes pour plus de performances.
- Une bibliothèque MPI-IO peut implémenter cette optimisation de manière transparente.

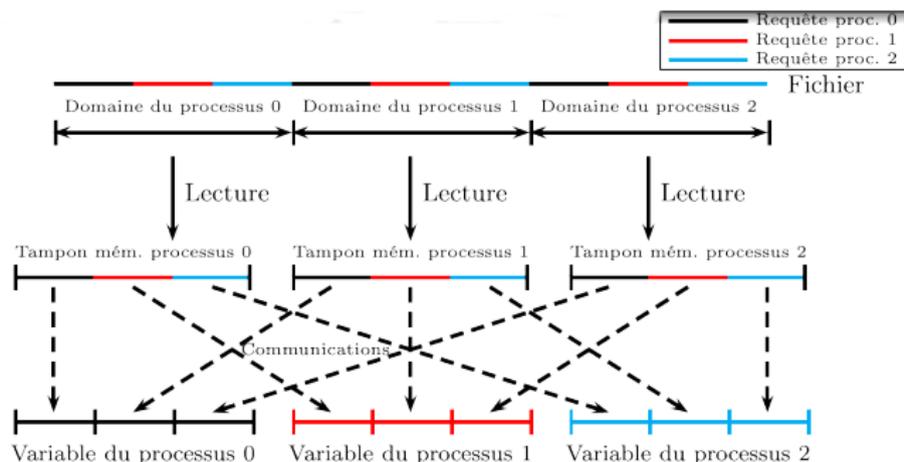


**Figure 44** – Mécanisme de *passoire* (*data sieving*) dans le cas d'accès nombreux, par un seul processus, à de petits blocs discontinus

# MPI-IO

## Exemple d'optimisation parallèle

Lorsqu'un ensemble de processus accède à des blocs discontinus (cas des tableaux distribués, par exemple), la bibliothèque d'I/O peut optimiser l'opération en maximisant l'accès aux données contiguës et en utilisant des communications collectives de redistribution.



**Figure 45** – Lecture en deux phases, par un ensemble de processus

## Ouverture d'un fichier

```
int MPI_File_open(MPI_Comm comm, const char *filename, int amode,  
MPI_Info info, MPI_File *descripteur)
```

- Ouvre le fichier `filename` avec les attributs `amode` ;
- `descripteur` est un objet opaque qui est ensuite utilisé comme référence dans toutes les opérations portant sur le fichier ;
- L'ouverture est une opération *collective* ;
- `filename` et `amode` doivent être identiques sur tous les rangs du communicateur `comm` ;
- Un objet de type `MPI_Info` est une base de donnée de type clé-valeur utile pour l'optimisation. `MPI_INFO_NULL` permet d'utiliser une valeur par défaut.

## Attributs

Attribut	Signification
<code>MPI_MODE_RDONLY</code>	seulement en lecture
<code>MPI_MODE_RDWR</code>	en lecture et écriture
<code>MPI_MODE_WRONLY</code>	seulement en écriture
<code>MPI_MODE_CREATE</code>	création du fichier s'il n'existe pas
<code>MPI_MODE_EXCL</code>	erreur si le fichier existe
<code>MPI_MODE_UNIQUE_OPEN</code>	le fichier n'est pas ouvert ailleurs
<code>MPI_MODE_SEQUENTIAL</code>	accès séquentiel
<code>MPI_MODE_APPEND</code>	pointeurs en fin de fichier (mode ajout)
<code>MPI_MODE_DELETE_ON_CLOSE</code>	destruction après la fermeture

Les attributs peuvent être combinés en utilisant l'opérateur `|`.

## Fermeture d'un fichier

```
int MPI_File_close(MPI_File *descripteur)
```

- Ferme le fichier ;
- La fermeture est une opération *collective*.

# MPI-IO

```
1  /* open01 */
2  #include <mpi.h>
3  #include <stdlib.h>
4  #include <stdio.h>
5
6  int main(int argc, char *argv[]) {
7      MPI_File descripteur;
8      int code, texte_longueur;
9      char texte_erreur[MPI_MAX_ERROR_STRING];
10
11     MPI_Init(&argc, &argv);
12
13     code = MPI_File_open(MPI_COMM_WORLD, "fichier.txt",
14                         MPI_MODE_RDWR | MPI_MODE_CREATE, MPI_INFO_NULL, &descripteur);
15     if (code != MPI_SUCCESS) {
16         MPI_Error_string(code, texte_erreur, &texte_longueur);
17         printf("%s\n", texte_erreur);
18         MPI_Abort(MPI_COMM_WORLD, 42);
19     }
20
21     code = MPI_File_close(&descripteur);
22     if (code != MPI_SUCCESS) {
23         printf("Erreur fermeture fichier\n");
24         MPI_Abort(MPI_COMM_WORLD, 2);
25     }
26
27     MPI_Finalize();
28 }
```

```
> ls -l fichier.txt
-rw-----  1 nom      grp      0 Feb 08 12:13 fichier.txt
```

## Gestion des erreurs

- Le comportement concernant l'argument code est différent pour la partie IO de MPI;
- Il est nécessaire de tester la valeur de cet argument ;
- Il est possible de changer ce comportement avec `MPI_File_set_errhandler()` ;
- Deux gestionnaires d'erreurs sont disponibles : `MPI_ERRORS_ARE_FATAL` et `MPI_ERRORS_RETURN` ;
- `MPI_Comm_set_errhandler()` permet de changer la gestion des erreurs pour les communications.

```
int MPI_File_set_errhandler(MPI_File descripteur, MPI_Errhandler gestionnaire)
```

Pour changer le comportement par défaut, il faut utiliser `MPI_FILE_NULL` comme descripteur.

## Généralités

- Les transferts de données entre fichiers et zones mémoire des processus se font via des appels explicites à des sous-programmes de lecture et d'écriture.
- On distingue trois propriétés des accès aux fichiers :
  - le **positionnement**, qui peut être explicite (en spécifiant un déplacement par rapport au début du fichier) ou implicite, via des pointeurs gérés par le système (ces pointeurs peuvent être de deux types : soit **individuels** à chaque processus, soit **partagés** par tous les processus) ;
  - la **synchronisation**, les accès pouvant être de type bloquants ou non bloquants ;
  - le **regroupement**, les accès pouvant être collectifs (c'est-à-dire effectués par tous les processus du communicateur au sein duquel le fichier a été ouvert) ou propres seulement à un ou plusieurs processus.
- Il est possible de mélanger les types d'accès effectués à un même fichier au sein d'une application.

Positionnement	Synchronisation	individuel	collectif
adresses explicites	bloquantes	MPI_File_read_at MPI_File_write_at	MPI_File_read_at_all MPI_File_write_at_all
	non bloquantes	MPI_File_iread_at MPI_File_iwrite_at	MPI_File_iread_at_all MPI_File_iwrite_at_all
pointeurs implicites individuels	bloquantes	MPI_File_read MPI_File_write	MPI_File_read_all MPI_File_write_all
	non bloquantes	MPI_File_iread MPI_File_iwrite	MPI_File_iread_all MPI_File_iwrite_all
pointeurs implicites partagés	bloquantes	MPI_File_read_shared MPI_File_write_shared	MPI_File_read_ordered MPI_File_write_ordered
	non bloquantes	MPI_File_iread_shared MPI_File_iwrite_shared	MPI_File_read_ordered_begin MPI_File_read_ordered_end MPI_File_write_ordered_begin MPI_File_write_ordered_end

## Le mécanisme de vue

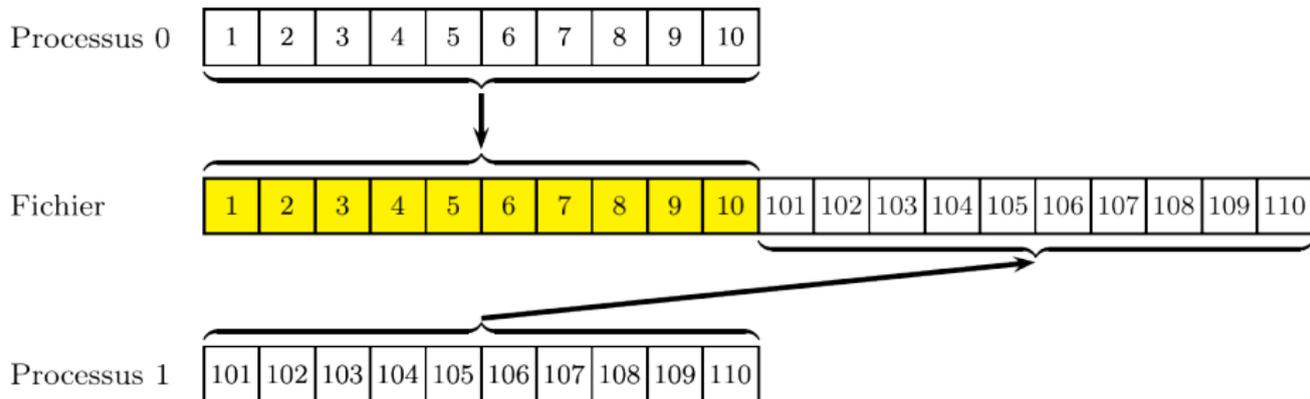
- Par défaut, les fichiers sont lus comme une simple suite d'octets mais MPI-IO dispose d'un mécanisme permettant une abstraction de plus haut niveau du contenu des fichiers : il est possible de décrire des structures de données complexes et de s'en servir comme gabarit lors de l'accès aux fichiers.
- Pour l'instant, il faut seulement savoir qu'un type élémentaire de données sert d'unité de base à ces constructions et que, par défaut, le type élémentaire est l'octet.
- Ce mécanisme de **vue** sera décrit en détail plus tard.

## Déplacements explicites

```
int MPI_File_read_at(MPI_File descripteur, MPI_Offset offset,  
void *buf, int count, MPI_Datatype datatype, MPI_Status *statut)  
  
int MPI_File_write_at(MPI_File descripteur, MPI_Offset offset,  
const void *buf, int count, MPI_Datatype datatype, MPI_Status *statut)
```

- Écrit/lit à la position `offset` dans le fichier `descripteur`, `count` élément de type `datatype` depuis l'adresse `buf` ;
- La `position` dans le fichier s'exprime toujours comme un multiple du type élémentaire de la vue courante. Par défaut, la position s'exprime donc en octets ;
- La taille du `datatype` doit être un multiple du type élémentaire.

```
1  /* write_at */
2  #include <mpi.h>
3  #include <stdio.h>
4
5  int main(int argc, char *argv[]) {
6      int rang, i, nb_octets_entier, nb_valeurs=10, code;
7      int valeurs[nb_valeurs];
8      MPI_File descripteur;
9      MPI_Offset position_fichier;
10     MPI_Status statut;
11
12     MPI_Init(&argc, &argv);
13     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
14     for(i=0; i<nb_valeurs; i++) { valeurs[i] = i+rang*100; }
15     printf("Ecriture processus %d :", rang);
16     for(i=0; i<nb_valeurs; i++) { printf("%d ", valeurs[i]); }
17     printf("\n");
18
19     code = MPI_File_open(MPI_COMM_WORLD, "donnees.dat",
20                         MPI_MODE_WRONLY | MPI_MODE_CREATE, MPI_INFO_NULL, &descripteur);
21     if (code != MPI_SUCCESS) {
22         printf("Erreur ouverture fichier\n");
23         MPI_Abort(MPI_COMM_WORLD, 42); }
24     MPI_Type_size(MPI_INT, &nb_octets_entier);
25     position_fichier = rang*nb_valeurs*nb_octets_entier;
26
27     MPI_File_set_errhandler(descripteur, MPI_ERRORS_ARE_FATAL);
28     MPI_File_write_at(descripteur, position_fichier, valeurs, nb_valeurs, MPI_INT,
29                      &statut);
30     MPI_File_close(&descripteur);
31     MPI_Finalize();
32 }
```

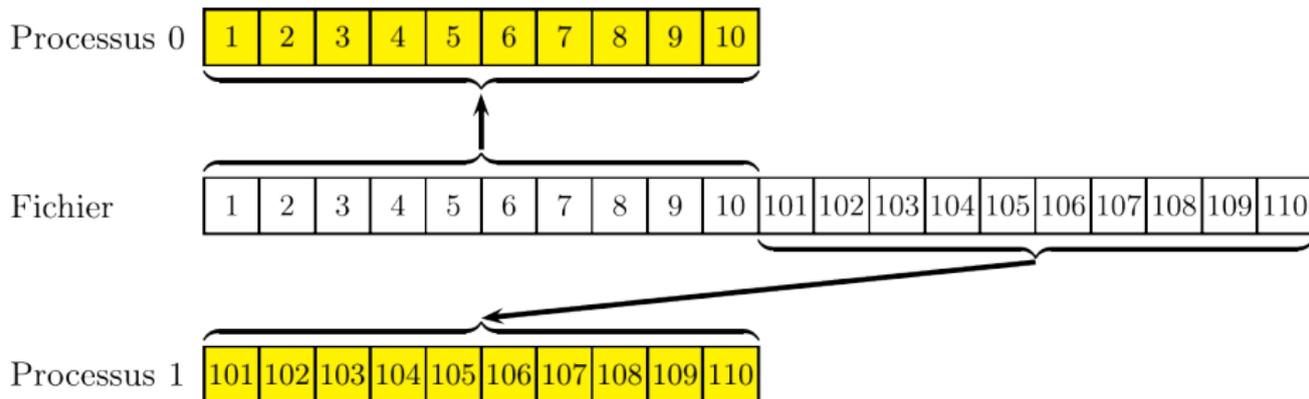


**Figure 46** – Exemple d'utilisation de `MPI_File_write_at()`

```
> mpiexec -n 2 write_at
```

```
Ecriture processus 0 : 1, 2, 3, 4, 5, 6, 7, 8, 9, 10  
Ecriture processus 1 : 101, 102, 103, 104, 105, 106, 107, 108, 109, 110
```

```
1  /* read_at */
2  #include <mpi.h>
3  #include <stdio.h>
4
5  int main(int argc, char *argv[]) {
6      int rang, i, nb_octets_entier, nb_valeurs=10, code;
7      int valeurs[nb_valeurs];
8      MPI_File descripteur;
9      MPI_Offset position_fichier;
10     MPI_Status statut;
11
12     MPI_Init(&argc, &argv);
13     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
14     MPI_File_set_errhandler(MPI_FILE_NULL, MPI_ERRORS_ARE_FATAL);
15     code = MPI_File_open(MPI_COMM_WORLD, "donnees.dat",
16                         MPI_MODE_RDONLY, MPI_INFO_NULL, &descripteur);
17     MPI_Type_size(MPI_INT, &nb_octets_entier);
18     position_fichier = rang*nb_valeurs*nb_octets_entier;
19     MPI_File_read_at(descripteur, position_fichier, valeurs, nb_valeurs, MPI_INT,
20                    &statut);
21     printf("Lecture processus %d : ", rang);
22     for (i=0; i<nb_valeurs; i++) {printf("%d ", valeurs[i]);} printf("\n");
23     MPI_File_close(&descripteur);
24     MPI_Finalize();
25 }
```



**Figure 47** – Exemple d'utilisation de `MPI_File_read_at()`

```
> mpiexec -n 2 read_at
```

```
Lecture processus 0 : 1, 2, 3, 4, 5, 6, 7, 8, 9, 10  
Lecture processus 1 : 101, 102, 103, 104, 105, 106, 107, 108, 109, 110
```

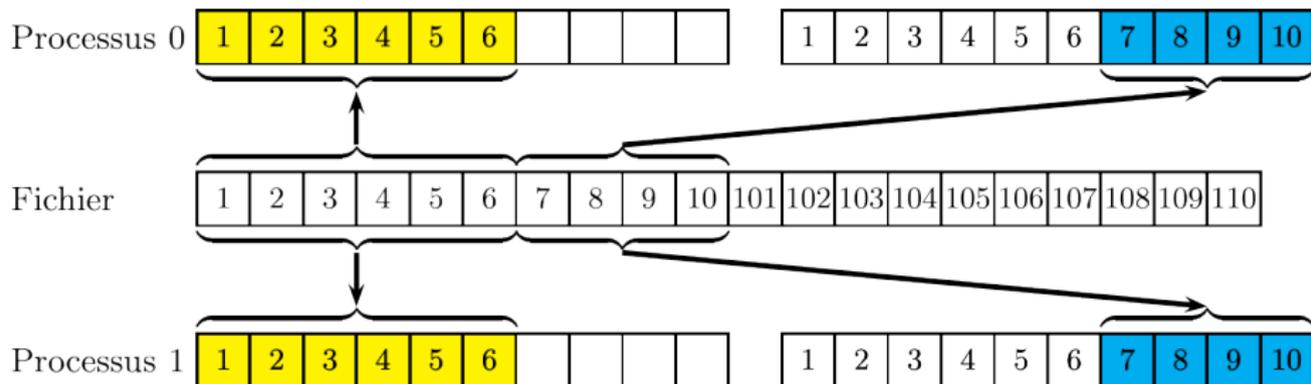
## Déplacements implicites individuels

```
int MPI_File_read(MPI_File descripteur, void *buf, int count,
MPI_Datatype datatype, MPI_Status *statut)

int MPI_File_write(MPI_File descripteur, const void *buf, int count,
MPI_Datatype datatype, MPI_Status *statut)
```

- Écrit/lit dans le fichier `descripteur`, `count` élément de type `datatype` depuis l'adresse `buf` ;
- Un pointeur individuel est géré par MPI, et ceci `par fichier` et `par processus`.
- Après chaque accès, le pointeur est positionné sur l'élément suivant.
- Dans tous ces sous-programmes, les pointeurs partagés ne sont jamais accédés ou modifiés explicitement.

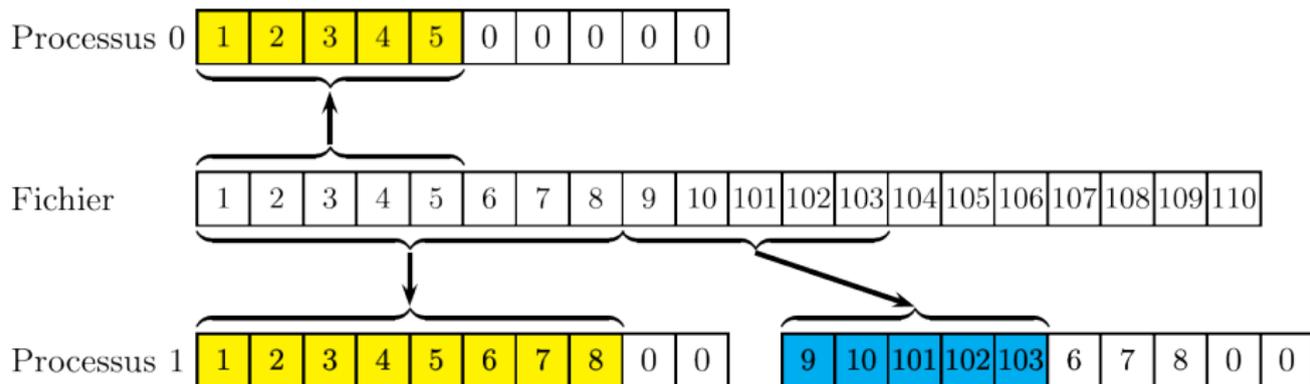
```
1  /* read01 */
2  #include <mpi.h>
3  #include <stdio.h>
4
5  int main(int argc, char *argv[]) {
6      int rang, i, nb_valeurs=10;
7      int valeurs[nb_valeurs];
8      MPI_File descripteur;
9      MPI_Status statut;
10
11     MPI_Init(&argc, &argv);
12     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
13     MPI_File_open(MPI_COMM_WORLD, "donnees.dat",
14                 MPI_MODE_RDONLY, MPI_INFO_NULL, &descripteur);
15     MPI_File_read(descripteur, valeurs, 6, MPI_INT, &statut);
16     MPI_File_read(descripteur, &(valeurs[6]), 4, MPI_INT, &statut);
17     printf("Lecture processus %d : ", rang);
18     for (i=0; i<nb_valeurs; i++) {printf("%d ", valeurs[i]);} printf("\n");
19     MPI_File_close(&descripteur);
20     MPI_Finalize();
21 }
```



**Figure 48** – Exemple 1 d'utilisation de `MPI_File_read()`

```
> mpiexec -n 2 read01
Lecture processus 1 : 1, 2, 3, 4, 5, 6, 7, 8, 9, 10
Lecture processus 0 : 1, 2, 3, 4, 5, 6, 7, 8, 9, 10
```

```
1  /* read02 */
2  #include <mpi.h>
3  #include <stdio.h>
4
5  int main(int argc, char *argv[]) {
6      int rang, i, nb_valeurs=10;
7      int valeurs[nb_valeurs];
8      MPI_File descripteur;
9      MPI_Status statut;
10
11     MPI_Init(&argc, &argv);
12     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
13     MPI_File_open(MPI_COMM_WORLD, "donnees.dat",
14                  MPI_MODE_RDONLY, MPI_INFO_NULL, &descripteur);
15     if (rang == 0) {
16         MPI_File_read(descripteur, valeurs, 5, MPI_INT, &statut);
17     } else {
18         MPI_File_read(descripteur, valeurs, 8, MPI_INT, &statut);
19         MPI_File_read(descripteur, valeurs, 5, MPI_INT, &statut); }
20     printf("Lecture processus %d : ", rang);
21     for (i=0; i<8; i++) {printf("%d ", valeurs[i]);} printf("\n");
22     MPI_File_close(&descripteur);
23     MPI_Finalize();
24 }
```



**Figure 49** – Exemple 2 d'utilisation de `MPI_File_read()`

```
> mpiexec -n 2 read02
```

```
Lecture processus 0 : 1, 2, 3, 4, 5, 0, 0, 0  
Lecture processus 1 : 9, 10, 101, 102, 103, 6, 7, 8
```

## Déplacements implicites partagés

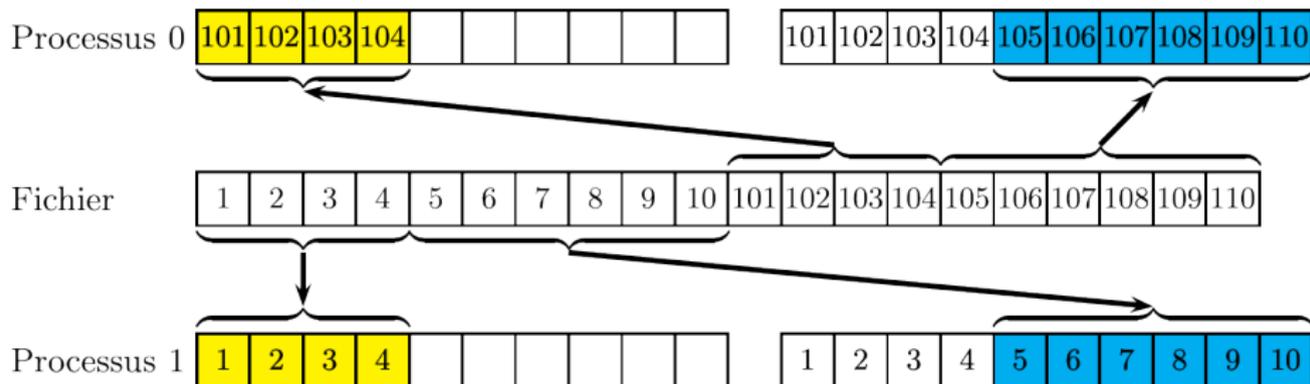
```
int MPI_File_read_shared(MPI_File descripteur, void *buf, int count,
    MPI_Datatype datatype, MPI_Status *statut)

int MPI_File_write_shared(MPI_File descripteur, const void *buf, int count,
    MPI_Datatype datatype, MPI_Status *statut)
```

- Écrit/lit dans le fichier **descripteur**, **count** élément de type **datatype** depuis l'adresse **buf** ;
- Il existe **un et un seul** pointeur partagé par fichier, commun à tous les processus du communicateur dans lequel le fichier a été ouvert.
- Tous les processus qui font une opération d'entrée-sortie utilisant le pointeur partagé doivent employer **la même vue** du fichier.
- Si on utilise les variantes non collectives des sous-programmes, l'ordre **n'est pas déterministe**. Si le traitement doit être déterministe, il faut explicitement gérer l'ordonnancement des processus ou utiliser les variantes collectives.
- Après chaque accès, le pointeur est positionné sur l'élément suivant.
- Dans tous ces sous-programmes, les pointeurs individuels ne sont jamais accédés ou modifiés.

```
1  /* read_shared01 */
2  #include <mpi.h>
3  #include <stdio.h>
4
5  int main(int argc, char *argv[]) {
6      int rang, i, nb_valeurs=10;
7      int valeurs[nb_valeurs];
8      MPI_File descripteur;
9      MPI_Status statut;
10
11     MPI_Init(&argc, &argv);
12     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
13     MPI_File_open(MPI_COMM_WORLD, "donnees.dat",
14                 MPI_MODE_RDONLY, MPI_INFO_NULL, &descripteur);
15     MPI_File_read_shared(descripteur, valeurs, 4, MPI_INT, &statut);
16     MPI_File_read_shared(descripteur, &(valeurs[4]), 6, MPI_INT, &statut);
17     printf("Lecture processus %d : ", rang);
18     for (i=0; i<nb_valeurs; i++) {printf("%d ", valeurs[i]);} printf("\n");
19     MPI_File_close(&descripteur);
20     MPI_Finalize();
21 }
```

# MPI-IO



**Figure 50** – Exemple 2 d'utilisation de `MPI_File_read_shared()`

```
> mpiexec -n 2 read_shared01
```

```
Lecture processus 1 : 1, 2, 3, 4, 5, 6, 7, 8, 9, 10  
Lecture processus 0 : 101, 102, 103, 104, 105, 106, 107, 108, 109, 110
```

## Lectures/écritures collectives

- Tous les processus du **communicateur** au sein duquel un fichier est ouvert participent aux opérations collectives d'accès aux données.
- Les opérations collectives sont généralement **plus performantes** que les opérations individuelles, parce qu'elles autorisent davantage de techniques d'optimisation mises en œuvre automatiquement ;
- Les accès sont effectués **dans l'ordre** des rangs des processus : le traitement est donc ici **déterministe**.

## Interfaces

```
int MPI_File_read_at_all(MPI_File descripteur, MPI_Offset offset,
void *buf, int count, MPI_Datatype datatype, MPI_Status *statut)

int MPI_File_write_at_all(MPI_File descripteur, MPI_Offset offset,
const void *buf, int count, MPI_Datatype datatype, MPI_Status *statut)

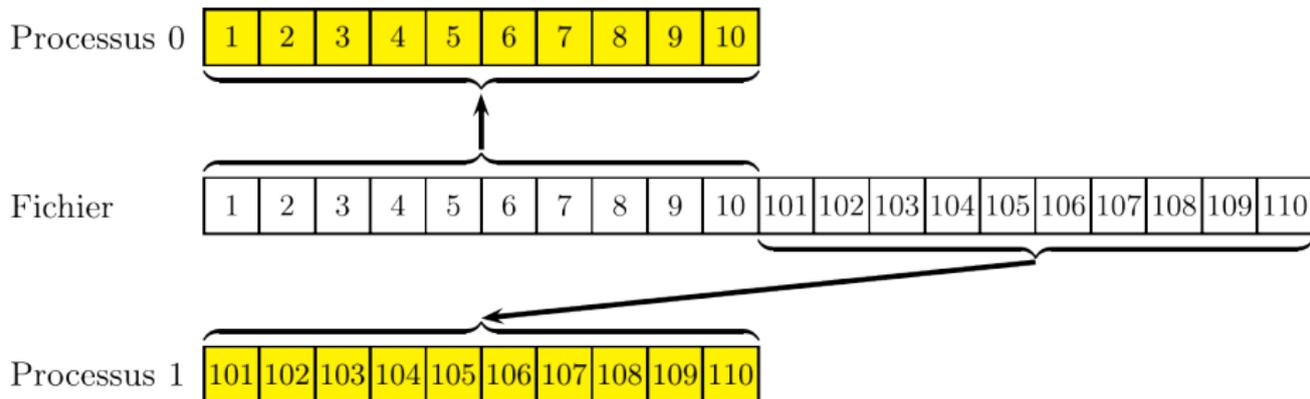
int MPI_File_read_all(MPI_File descripteur, void *buf, int count,
MPI_Datatype datatype, MPI_Status *statut)

int MPI_File_write_all(MPI_File descripteur, const void *buf, int count,
MPI_Datatype datatype, MPI_Status *statut)

int MPI_File_read_ordered(MPI_File descripteur, void *buf, int count,
MPI_Datatype datatype, MPI_Status *statut)

int MPI_File_write_ordered(MPI_File descripteur, const void *buf, int count,
MPI_Datatype datatype, MPI_Status *statut)
```

```
1  /* read_at_all */
2  #include <mpi.h>
3  #include <stdio.h>
4
5  int main(int argc, char *argv[]) {
6      int rang, i, nb_octets_entier, nb_valeurs=10, code;
7      int valeurs[nb_valeurs];
8      MPI_File descripteur;
9      MPI_Offset position_fichier;
10     MPI_Status statut;
11
12     MPI_Init(&argc, &argv);
13     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
14     code = MPI_File_open(MPI_COMM_WORLD, "donnees.dat",
15                         MPI_MODE_RDONLY, MPI_INFO_NULL, &descripteur);
16     MPI_Type_size(MPI_INT, &nb_octets_entier);
17     position_fichier = rang*nb_valeurs*nb_octets_entier;
18     MPI_File_read_at_all(descripteur, position_fichier,
19                          valeurs, nb_valeurs, MPI_INT, &statut);
20     printf("Lecture processus %d : ", rang);
21     for (i=0; i<nb_valeurs; i++) {printf("%d ", valeurs[i]);} printf("\n");
22     MPI_File_close(&descripteur);
23     MPI_Finalize();
24 }
```

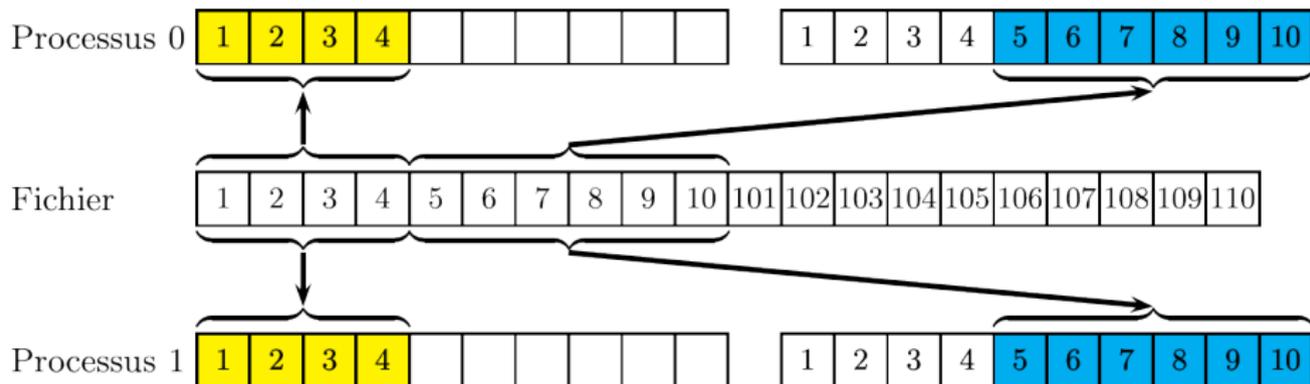


**Figure 51** – Exemple d'utilisation de `MPI_File_read_at_all()`

```
> mpiexec -n 2 read_at_all
```

```
Lecture processus 0 : 1, 2, 3, 4, 5, 6, 7, 8, 9, 10  
Lecture processus 1 : 101, 102, 103, 104, 105, 106, 107, 108, 109, 110
```

```
1  /* read_all01 */
2  #include <mpi.h>
3  #include <stdio.h>
4
5  int main(int argc, char *argv[]) {
6      int rang, i, nb_valeurs=10;
7      int valeurs[nb_valeurs];
8      MPI_File descripteur;
9      MPI_Status statut;
10
11     MPI_Init(&argc, &argv);
12     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
13     MPI_File_open(MPI_COMM_WORLD, "donnees.dat",
14                 MPI_MODE_RDONLY, MPI_INFO_NULL, &descripteur);
15     MPI_File_read_all(descripteur, valeurs, 4, MPI_INT, &statut);
16     MPI_File_read_all(descripteur, &(valeurs[4]), 6, MPI_INT, &statut);
17     printf("Lecture processus %d : ", rang);
18     for (i=0; i<nb_valeurs; i++) {printf("%d ", valeurs[i]);} printf("\n");
19     MPI_File_close(&descripteur);
20     MPI_Finalize();
21 }
```



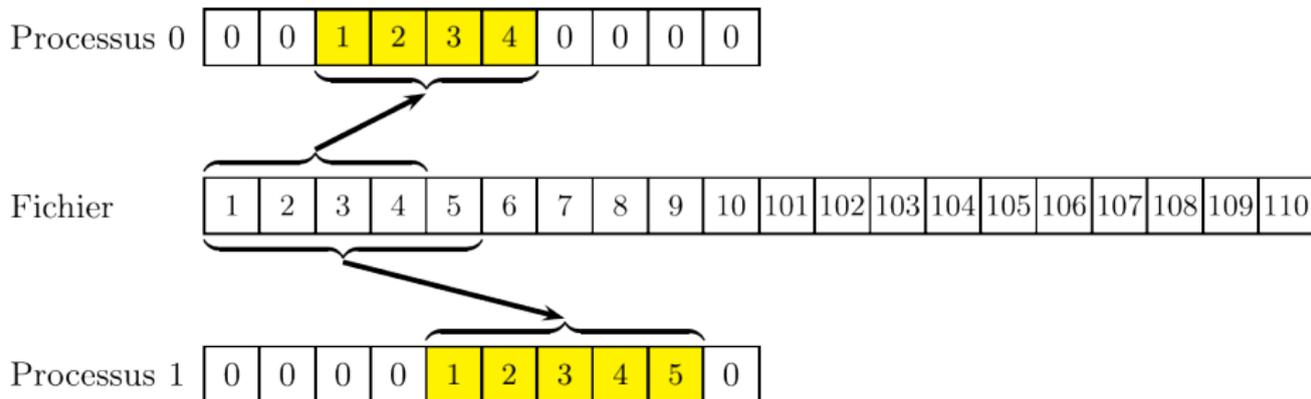
**Figure 52** – Exemple 1 d'utilisation de `MPI_File_read_all()`

```
> mpiexec -n 2 read_all01
```

```
Lecture processus 0 : 1, 2, 3, 4, 5, 6, 7, 8, 9, 10
```

```
Lecture processus 1 : 1, 2, 3, 4, 5, 6, 7, 8, 9, 10
```

```
1  /* read_all102 */
2  #include <mpi.h>
3  #include <stdio.h>
4
5  int main(int argc, char *argv[]) {
6      int rang, i, nb_valeurs=10, indicel, indice2;
7      int valeurs[nb_valeurs];
8      MPI_File descripteur;
9      MPI_Status statut;
10
11     MPI_Init(&argc, &argv);
12     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
13     MPI_File_open(MPI_COMM_WORLD, "donnees.dat",
14                 MPI_MODE_RDONLY, MPI_INFO_NULL, &descripteur);
15
16     if (rang == 0) {
17         indicel=2;
18         indice2=5;
19     } else {
20         indicel=4;
21         indice2=8;
22     }
23     MPI_File_read_all(descripteur, &(valeurs[indicel]),
24                     indice2-indicel+1, MPI_INT, &statut);
25     printf("Lecture processus %d : ", rang);
26     for (i=0; i<nb_valeurs; i++) {printf("%d ", valeurs[i]);} printf("\n");
27     MPI_File_close(&descripteur);
28     MPI_Finalize();
29 }
```



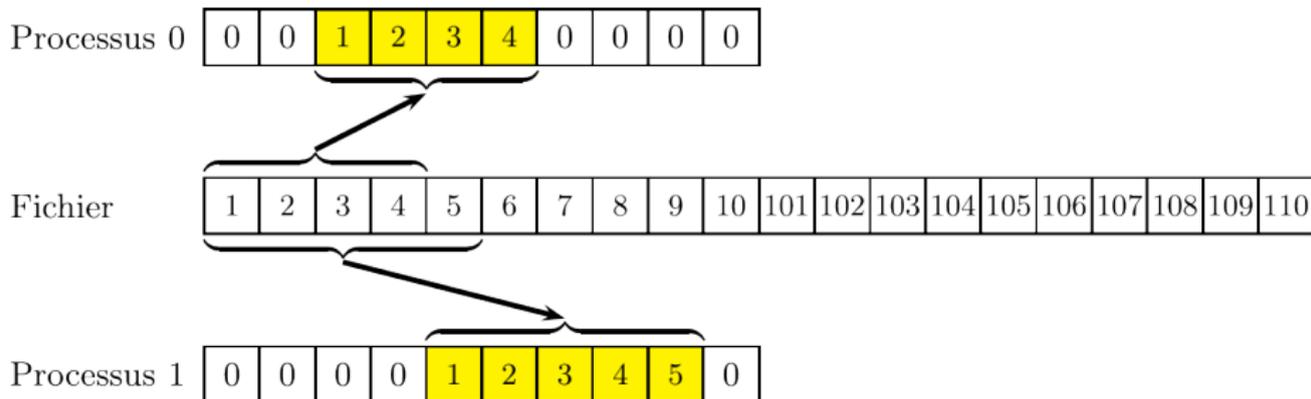
**Figure 53** – Exemple 2 d'utilisation de `MPI_File_read_all()`

```
> mpiexec -n 2 read_all02
```

```
Lecture processus 1 : 0, 0, 0, 0, 1, 2, 3, 4, 5, 0
```

```
Lecture processus 0 : 0, 0, 1, 2, 3, 4, 0, 0, 0, 0
```

```
1  /* read_all03 */
2  #include <mpi.h>
3  #include <stdio.h>
4
5  int main(int argc, char *argv[]) {
6      int rang, i, nb_valeurs=10, indicel, indice2;
7      int valeurs[nb_valeurs];
8      MPI_File descripteur;
9      MPI_Status statut;
10
11     MPI_Init(&argc, &argv);
12     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
13     MPI_File_open(MPI_COMM_WORLD, "donnees.dat",
14                  MPI_MODE_RDONLY, MPI_INFO_NULL, &descripteur);
15     if (rang == 0) {
16         MPI_File_read_all(descripteur, &(valeurs[2]), 4, MPI_INT, &statut);
17     } else {
18         MPI_File_read_all(descripteur, &(valeurs[4]), 5, MPI_INT, &statut);
19     }
20     printf("Lecture processus %d : ", rang);
21     for (i=0; i<nb_valeurs; i++) {printf("%d ", valeurs[i]);} printf("\n");
22     MPI_File_close(&descripteur);
23     MPI_Finalize();
24 }
```



**Figure 54** – Exemple 3 d'utilisation de `MPI_File_read_all()`

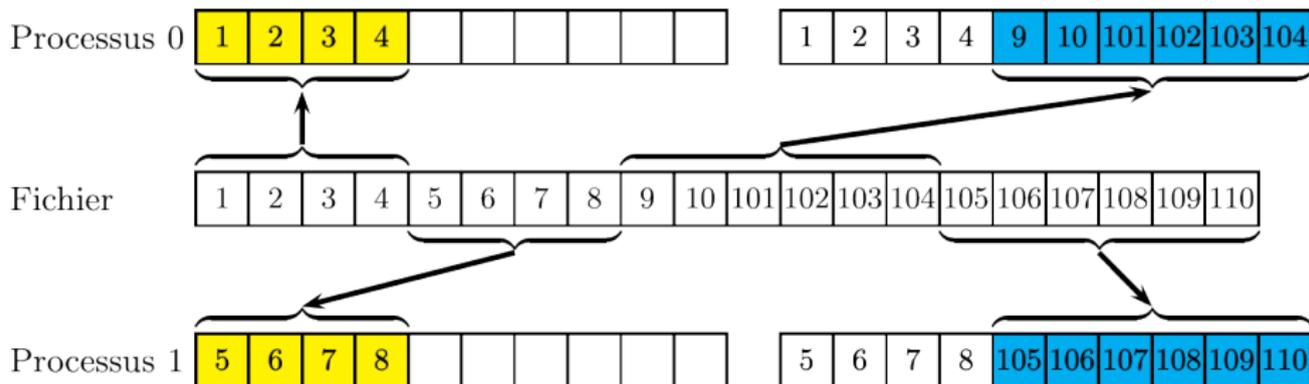
```
> mpiexec -n 2 read_all03
```

```
Lecture processus 1 : 0, 0, 0, 0, 1, 2, 3, 4, 5, 0
```

```
Lecture processus 0 : 0, 0, 1, 2, 3, 4, 0, 0, 0, 0
```

```
1  /* read_ordered */
2  #include <mpi.h>
3  #include <stdio.h>
4
5  int main(int argc, char *argv[]) {
6      int rang, i, nb_valeurs=10;
7      int valeurs[nb_valeurs];
8      MPI_File descripteur;
9      MPI_Status statut;
10
11     MPI_Init(&argc, &argv);
12     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
13     MPI_File_open(MPI_COMM_WORLD, "donnees.dat",
14                 MPI_MODE_RDONLY, MPI_INFO_NULL, &descripteur);
15     MPI_File_read_ordered(descripteur, valeurs, 4, MPI_INT, &statut);
16     MPI_File_read_ordered(descripteur, &(valeurs[4]), 6, MPI_INT, &statut);
17     printf("Lecture processus %d : ", rang);
18     for (i=0; i<nb_valeurs; i++) {printf("%d ", valeurs[i]);} printf("\n");
19     MPI_File_close(&descripteur);
20     MPI_Finalize();
21 }
```

# MPI-IO



**Figure 55** – Exemple d'utilisation de `MPI_File_ordered()`

```
> mpiexec -n 2 read_ordered
```

```
Lecture processus 1 : 5, 6, 7, 8, 105, 106, 107, 108, 109, 110
```

```
Lecture processus 0 : 1, 2, 3, 4, 9, 10, 101, 102, 103, 104
```

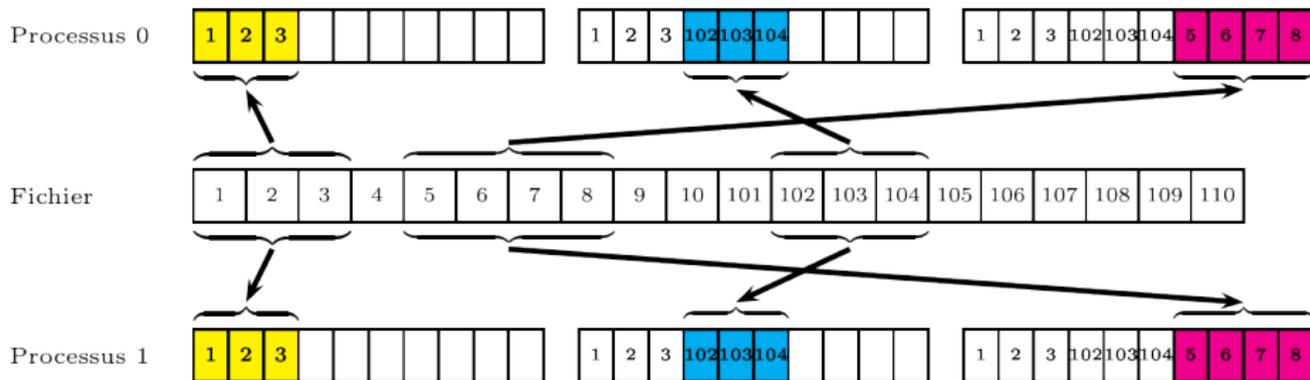
## Positionnement explicite des pointeurs dans un fichier

```
int MPI_File_seek(MPI_File descripteur, MPI_Offset offset, int mode)
int MPI_File_seek_shared(MPI_File descripteur, MPI_Offset offset, int mode)
```

- Il est possible de **positionner explicitement** les pointeurs individuels à l'aide du sous-programme `MPI_File_seek()`, et de même le pointeur partagé avec le sous-programme `MPI_File_seek_shared()`.
- Il y a **trois modes** possibles pour modifier la valeur d'un pointeur :
  - `MPI_SEEK_SET` permet de définir un déplacement absolu ;
  - `MPI_SEEK_CUR` permet un déplacement relativement à la position courante ;
  - `MPI_SEEK_END` positionne le pointeur à la fin du fichier, à laquelle un déplacement éventuel est ajouté.
- Avec `MPI_SEEK_CUR` et `MPI_SEEK_END`, on peut spécifier une valeur négative, ce qui permet de revenir **en arrière** dans le fichier.

```
1  /* seek */
2  #include <mpi.h>
3  #include <stdio.h>
4
5  int main(int argc, char *argv[]) {
6      int rang, i, nb_octets_entier, nb_valeurs=10;
7      int valeurs[nb_valeurs];
8      MPI_File descripteur;
9      MPI_Status statut;
10     MPI_Offset position_fichier;
11
12     MPI_Init(&argc, &argv);
13     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
14     MPI_File_open(MPI_COMM_WORLD, "donnees.dat",
15                 MPI_MODE_RDONLY, MPI_INFO_NULL, &descripteur);
16     MPI_File_read(descripteur, valeurs, 3, MPI_INT, &statut);
17     MPI_Type_size(MPI_INT, &nb_octets_entier);
18     position_fichier = 8*nb_octets_entier;
19     MPI_File_seek(descripteur, position_fichier, MPI_SEEK_CUR);
20     MPI_File_read(descripteur, &(valeurs[3]), 3, MPI_INT, &statut);
21     position_fichier = 4*nb_octets_entier;
22     MPI_File_seek(descripteur, position_fichier, MPI_SEEK_SET);
23     MPI_File_read(descripteur, &(valeurs[6]), 4, MPI_INT, &statut);
24     printf("Lecture processus %d : ", rang);
25     for (i=0; i<nb_valeurs; i++) {printf("%d ", valeurs[i]);} printf("\n");
26     MPI_File_close(&descripteur);
27     MPI_Finalize();
28 }
```

# MPI-IO



**Figure 56** – Exemple d'utilisation de `MPI_File_seek()`

```
> mpiexec -n 2 seek  
Lecture processus 1 : 1, 2, 3, 102, 103, 104, 5, 6, 7, 8  
Lecture processus 0 : 1, 2, 3, 102, 103, 104, 5, 6, 7, 8
```

## Accès à la fin du fichier

- Écrire à la fin du fichier augmente la taille du fichier.
- Lire à partir de la fin du fichier ne récupère aucune donnée. Lors d'une lecture, l'utilisation de `MPI_Get_count()` permet de connaître le nombre d'éléments réellement lus.

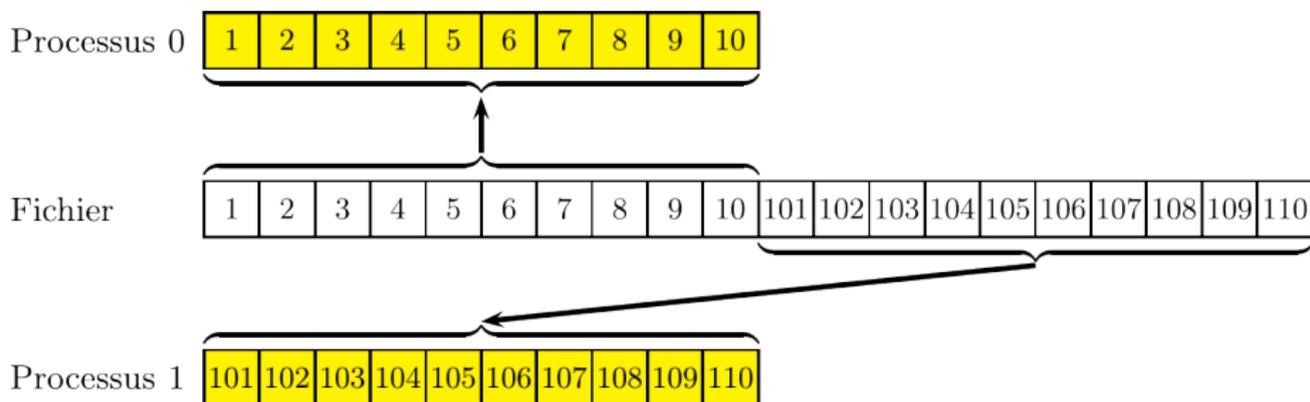
## Lectures/écritures non bloquantes

- L'intérêt est de faire un recouvrement entre les calculs et les entrées-sorties.
- Les entrées-sorties non bloquantes sont implémentées suivant le modèle utilisé pour les communications non bloquantes entre processus.
- Un accès non-bloquant doit donner lieu ultérieurement à un test explicite de complétude ou à une mise en attente (via `MPI_Test()`, `MPI_Wait()`, etc.)

```
1  /* iread_at */
2  #include <mpi.h>
3  #include <stdio.h>
4
5  int main(int argc, char *argv[]) {
6      int rang, i, nb_octets_entier, termine, nb_valeurs=10, nb_iterations=0;
7      int valeurs[nb_valeurs];
8      MPI_File descripteur;
9      MPI_Status statut;
10     MPI_Offset position_fichier;
11     MPI_Request requete;
12
13     MPI_Init (&argc, &argv);
14     MPI_Comm_rank (MPI_COMM_WORLD, &rang);
```

```
15 MPI_File_open(MPI_COMM_WORLD,"donnees.dat",
16               MPI_MODE_RDONLY,MPI_INFO_NULL,&descripteur);
17 MPI_Type_size(MPI_INT,&nb_octets_entier);
18 position_fichier = rang*nb_valeurs*nb_octets_entier;
19 MPI_File_iread_at(descripteur,position_fichier,
20                  valeurs,nb_valeurs,MPI_INT,&requete);
21 while( nb_iterations < 5000) {
22     nb_iterations = nb_iterations+1;
23     /* Calculs recouvrant le temps demande par l'operation de lecture */
24     /* */
25     MPI_Test(&requete,&termine,&statut);
26     if (termine) break;
27 }
28 if ( !termine) MPI_Wait (&requete,&statut);
29 printf("Apres %d iterations, lecture processus %d : ",nb_iterations, rang);
30 for (i=0;i<nb_valeurs;i++) {printf("%d ",valeurs[i]);} printf("\n");
31 MPI_File_close(&descripteur);
32 MPI_Finalize();
33 }
```

# MPI-IO



**Figure 57** – Exemple d'utilisation de `MPI_File_iread_at()`

```
> mpiexec -n 2 iread_at
```

```
Après 1 iterations, lecture processus 0 : 1, 2, 3, 4, 5, 6, 7, 8, 9, 10
```

```
Après 1 iterations, lecture processus 1 : 101, 102, 103, 104, 105, 106, 107, 108, 109, 110
```

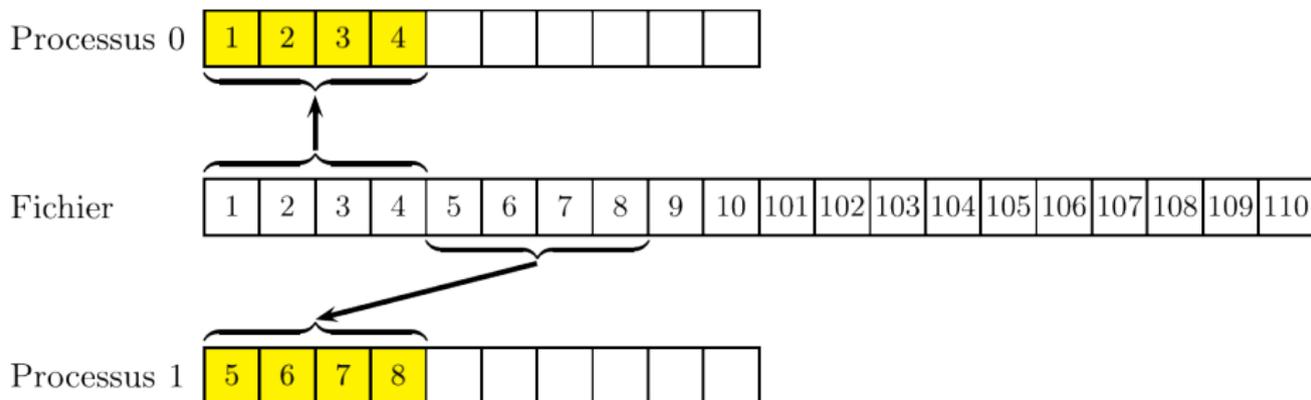
```
1  /* iwrite */
2  #include <mpi.h>
3  #include <stdio.h>
4
5  int main(int argc, char *argv[]) {
6      int rang, i, termine, nb_valeurs=10, nb_iterations=0;
7      int valeurs[nb_valeurs], temp[nb_valeurs];
8      MPI_File descripteur;
9      MPI_Request requete;
10
11     MPI_Init (&argc, &argv);
12     MPI_File_open (MPI_COMM_WORLD, "donnees.dat",
13                   MPI_MODE_WRONLY | MPI_MODE_CREATE, MPI_INFO_NULL, &descripteur);
14     for (i=0; i<nb_valeurs; i++) temp[i]=valeurs[i];
15     MPI_File_seek (descripteur, deplacement, MPI_SEEK_SET);
16     MPI_File_iwrite (descripteur, temp, nb_valeurs, MPI_INT, &requete);
17     while ( nb_iterations < 5000) {
18         nb_iterations = nb_iterations+1;
19         /* Calculs recouvrant le temps demande par l'operation de lecture */
20         /* */
21         MPI_Test (&requete, &termine, MPI_STATUS_IGNORE);
22         if (termine) {
23             for (i=0; i<nb_valeurs; i++) temp[i]=valeurs[i];
24             MPI_File_seek (descripteur, deplacement, MPI_SEEK_SET);
25             MPI_File_iwrite (descripteur, temp, nb_valeurs, MPI_INT, &requete); }
26     }
27     MPI_Wait (&requete, MPI_STATUS_IGNORE);
28     MPI_File_close (&descripteur);
29     MPI_Finalize();
30 }
```

## Lectures/écritures collectives et non bloquantes

- Il est possible d'effectuer des opérations qui soient à la fois **collectives** et **non bloquantes**.
- Il ne peut y avoir qu'une seule opération collective non bloquante en cours à la fois par processus.
- Entre les deux phases de l'opération collective non-bloquante, il est possible de faire des opérations non collectives sur le fichier, mais la zone mémoire concernée par l'opération collective ne peut être modifiée.

```
1 #include <mpi.h>
2 #include <stdio.h>
3
4 int main(int argc, char *argv[]) {
5     int rang, nb_valeurs=10, nb_iterations=0;
6     int valeurs[nb_valeurs], temp[nb_valeurs];
7     MPI_File descripteur;
8     MPI_Status statut;
9     MPI_Request req;
10
11     MPI_Init(&argc, &argv);
12     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
13     MPI_File_open(MPI_COMM_WORLD, "donnees.dat",
14                 MPI_MODE_RDONLY, MPI_INFO_NULL, &descripteur);
15     MPI_File_iread_all(descripteur, valeurs, 4, MPI_INT, &req);
16     printf("Processus numero : %d\n", rang);
17     MPI_Wait(&req, &statut);
18     printf("Lecture processus %d : %d %d %d %d\n",
19           rang, valeurs[0], valeurs[1], valeurs[2], valeurs[3]);
20     MPI_File_close(&descripteur);
21     MPI_Finalize();
22 }
```

# MPI-IO



**Figure 58** – Exemple d'utilisation de `MPI_File_iread_all()`

```
> mpiexec -n 2 iread_all
```

```
Processus numero : 0  
Lecture processus 0 : 1, 2, 3, 4  
Processus numero : 1  
Lecture processus 1 : 5, 6, 7, 8
```

## T.P. MPI – Exercice 7 : Lecture d'un fichier en mode parallèle

- On dispose du fichier binaire `donnees.dat`, constitué d'une suite de 484 valeurs entières
- En considérant un programme parallèle mettant en œuvre 4 processus, il s'agit de lire les 121 premières valeurs sur le processus 0, les 121 suivantes sur le processus 1, etc. et d'écrire celles-ci dans quatre fois quatre fichiers appelés `fichier_XXX0.dat` . . . `fichier_XXX3.dat`
- On emploiera pour ce faire 4 méthodes différentes, parmi celles présentées :
  - lecture via des déplacements explicites, en mode individuel ;
  - lecture via les pointeurs partagés, en mode collectif ;
  - lecture via les pointeurs individuels, en mode individuel ;
  - lecture via les pointeurs partagés, en mode individuel.
- Pour compiler utilisez la commande `make`, pour exécuter le code utilisez la commande `make exe` et pour vérifier les résultats utilisez la commande `make verification` qui génère des fichiers images correspondant aux quatre cas à traiter.

# MPI 4.x

## Ajout

- Grand nombre
- Communication par morceaux
- MPI Session
- Autres

# Grand nombre

- Les nombres d'éléments étaient en `integer` ou `int`.
- MPI 4.0 ajoute des fonctions avec `MPI_Count` à la place.
- En C ces nouvelles fonctions contiennent en plus `_c` à la fin de la fonction.

```
int MPI_Send(const void * buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm);
int MPI_Send_c(const void * buf, MPI_Count count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm);
```

- En *Fortran* les nombres en `integer` peuvent être remplacés par

```
integer(kind=MPI_COUNT_KIND).
```

- Uniquement disponible avec le module `mpi_f08`.
- Pas de changement de nom grâce au polymorphisme.

```
MPI_Send(buf, count, datatype, dest, tag, comm, ierror)
TYPE(*), DIMENSION(..), INTENT(IN) :: buf
INTEGER, INTENT(IN) :: count, dest, tag
TYPE(MPI_Datatype), INTENT(IN) :: datatype
TYPE(MPI_Comm), INTENT(IN) :: comm
INTEGER, OPTIONAL, INTENT(OUT) :: ierror
```

```
MPI_Send(buf, count, datatype, dest, tag, comm, ierror)
TYPE(*), DIMENSION(..), INTENT(IN) :: buf
INTEGER(KIND=MPI_COUNT_KIND), INTENT(IN) :: count
TYPE(MPI_Datatype), INTENT(IN) :: datatype
INTEGER, INTENT(IN) :: dest, tag
TYPE(MPI_Comm), INTENT(IN) :: comm
INTEGER, OPTIONAL, INTENT(OUT) :: ierror
```

## Communication par morceaux

- Contribution multiple à une communication.
- Utile pour l'hybride.
- Initialisation avec `MPI_Psend_init()` ou `MPI_Precv_init()` en fournissant le nombre d'éléments par partition et le nombre de partitions.
- `MPI_Start()` pour démarrer la communication.
- `MPI_Pready()` pour signaler qu'une partition est prête.
- Il n'est pas possible de faire un `MPI_Recv()` d'un `MPI_Psend_init()`
- `MPI_Wait()` pour terminer la communication
- `MPI_Parrived()` permet de savoir si une partition a été reçue

# Sessions

- Permettre de faire plusieurs `MPI_Init()`/`MPI_Finalize()`.
- `MPI_Session_init()` pour lancer une session.
- `MPI_Session_finalize()` pour terminer la session.
- Plus de notion de `MPI_COMM_WORLD`.
- Notion de *Process Sets* : `mpi://WORLD` et `mpi://SELF`.
- `MPI_Group_from_session_pset()` pour faire un groupe à partir d'un *pset*.
- `MPI_Comm_create_from_group()` pour faire un communicateur à partir d'un groupe.
- `MPI_Session_get_num_psets()` permet d'obtenir le nombre de *pset* disponible.
- `MPI_Session_get_nth_pset()` permet d'avoir le nom d'un *pset* disponible.

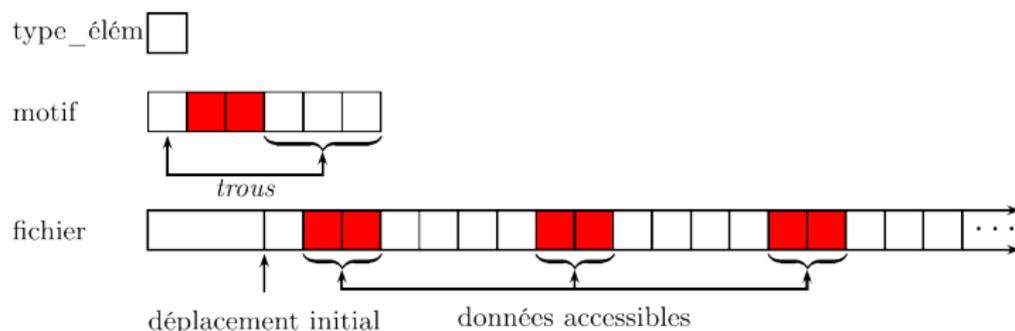
## Autres

- Ajout de `MPI_Isendrecv` et `MPI_Isendrecv_replace`.
- Ajout des communications collectives persistentes.
- Ajout de l'option `mpi_initial_errhandler` pour `mpiexec` afin de spécifier le gestionnaire d'erreur par défaut.

## MPI-IO Vues

## Définition des vues

- C'est un mécanisme permettant de décrire un schéma d'accès aux fichiers.
- Une vue est définie par trois variables : un **déplacement initial**, un **type élémentaire de données** et un **motif**.
- L'accès aux fichiers s'effectue par répétition du motif, une fois le positionnement initial effectué.



**Figure 59** – Type élémentaire de donnée et motif

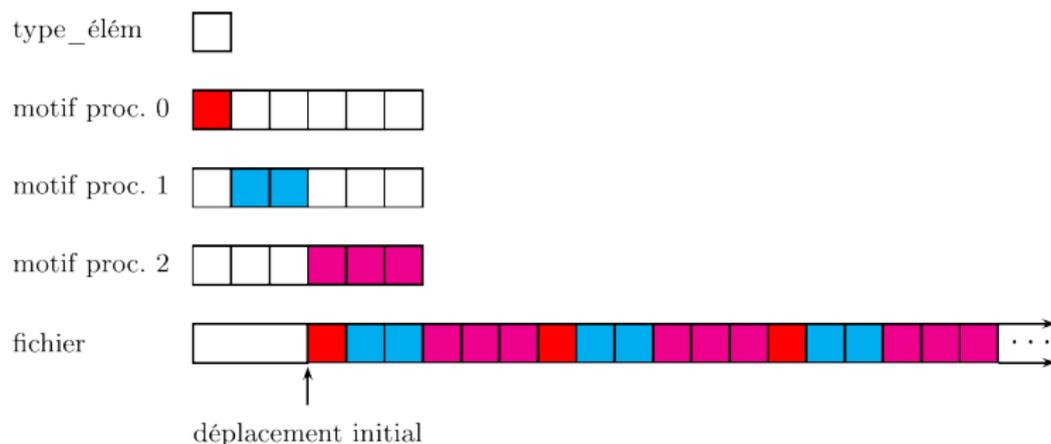
## Définition des vues

- Les vues sont construites à l'aide de **types dérivés** MPI.
- Il est possible de définir des **trous** dans une vue, de façon à ne pas tenir compte de certaines parties des données.
- La vue par défaut consiste en une simple suite d'octets (déplacement initial nul, **type\_élé**m et motif égaux à `MPI_BYTE`).

## Vues multiples

- Un processus donné peut définir et utiliser successivement **plusieurs vues** d'un même fichier.
- Les processus peuvent avoir des **vues différentes** du fichier, de façon à accéder à des parties complémentaires de celui-ci.

# MPI-IO Vues



**Figure 60** – Exemple de définition de motifs différents selon les processus

## Remarques :

- Un pointeur partagé n'est utilisable avec une vue que si tous les processus ont la même vue.
- Si le fichier est ouvert en écriture, les zones décrites par les différentes vues ne peuvent se recouvrir, même partiellement.

## Changement de la vue sur un fichier : `MPI_File_set_view()`

```
int MPI_File_set_view(MPI_File descripteur, MPI_Offset deplacement_initial, MPI_Datatype type_elem,
                     MPI_Datatype motif, char *mode, MPI_Info info)
```

- C'est une **opération collective** à l'ensemble des processus impliqués dans l'accès au fichier. Chaque processus peut définir **un déplacement initial et un motif différent**. L'étendue du type élémentaire doit être identique.
- Les pointeurs individuels et le pointeur partagé **sont réinitialisés au début de la vue**, en tenant compte du déplacement initial.

### Notes :

- Les types dérivés utilisés dans la vue doivent avoir été validés au préalable à l'aide du sous-programme `MPI_Type_commit()`.
- Il y a trois représentations possibles des données (mode) : "native", "internal" ou "external32".

### Construction de sous-tableaux

Un type dérivé utile pour créer un motif est le type “subarray”, qu’on introduit ici. Ce type permet de créer un sous-tableau à partir d’un tableau et se définit via le sous-programme `MPI_Type_create_subarray()`.

Le **profil** d’un tableau est un vecteur dont chaque élément est le nombre d’éléments dans chaque dimension. Soit par exemple le tableau `T(10, 0:5, -10:10)` (ou `T[10][6][21]`), son profil est le vecteur **(10,6,21)**.

# MPI-IO Vues / Types de données dérivés

```
int MPI_Type_create_subarray(int nb_dims, const int profil_tab[], const int profil_sous_tab[],
                             const int coord_debut[], int ordre, MPI_Datatype ancien_type,
                             MPI_Datatype *nouveau_type)
```

## Description des arguments

- `nb_dims` : nombre de dimension du tableau
- `profil_tab` : profil du tableau à partir duquel on va extraire un sous-tableau
- `profil_sous_tab` : profil du sous-tableau
- `coord_debut` : coordonnées de départ du sous-tableau dans le tableau, les indices du tableau commençant à 0. Par exemple, si on veut que les coordonnées de départ du sous-tableau soient `tab(2, 3)`, il faut que `coord_debut(:) = (/ 1, 2 /)`
- `ordre` : ordre de stockage des éléments
  - `MPI_ORDER_FORTRAN` spécifie le mode de stockage en Fortran, c.-à-d. suivant les colonnes
  - `MPI_ORDER_C` spécifie le mode de stockage en C, c.-à-d. suivant les lignes

# MPI-IO Vues / Types de données dérivés

## Échanges entre 2 processus avec subarray

AVANT

1	2	3	4
5	6	7	8
9	10	11	12

Processus 0

-1	-2	-3	-4
-5	-6	-7	-8
-9	-10	-11	-12

Processus 1

APRÈS

1	-7	-8	4
5	-11	-12	8
9	10	11	12

Processus 0

-1	-2	-3	-4
-5	-6	2	3
-9	-10	6	7

Processus 1

## Échanges entre 2 processus avec subarray : code

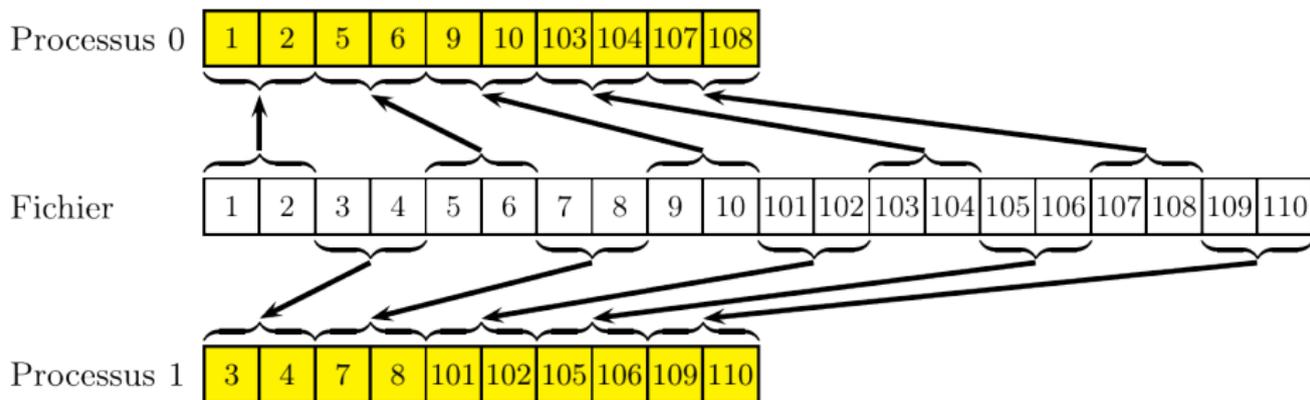
```
1  /* subarray */
2  #include <mpi.h>
3  #include <stdio.h>
4  #include <stdlib.h>
5
6  int main(int argc, char *argv[]) {
7      int rang, i, j;
8      int nb_lignes=3, nb_colonnes=4, sign=1, nb_dims=2, etiquette=1000;
9      int tab[nb_lignes][nb_colonnes], profil_tab[nb_dims];
10     int profil_sous_tab[nb_dims], coord_debut[nb_dims];
11     MPI_Datatype type_sous_tab;
12     MPI_Status statut;
13
14     MPI_Init (&argc, &argv);
15     MPI_Comm_rank (MPI_COMM_WORLD, &rang);
16
17     /* Initialisation du tableau tab sur chaque processus */
18     if (rang == 1) sign=-1;
19     for (i=0; i<nb_lignes; i++) {
20         for (j=0; j<nb_colonnes; j++) {
21             tab[i][j] = sign*(1+i*nb_colonnes+j); } }
```

## Échanges entre 2 processus avec subarray : code (suite)

```
22  /* Profil du tableau tab a partir duquel on va extraire un sous tableau */
23  profil_tab[0] = nb_lignes; profil_tab[1] = nb_colonnes;
24  /* Profil du sous tableau */
25  profil_sous_tab[0] = 2; profil_sous_tab[1] = 2;
26  /* Coordonnees de depart du sous tableau */
27  coord_debut[0]= rang; coord_debut[1]= rang+1;
28  /* Creation du type type_sous_tab */
29  MPI_Type_create_subarray(nb_dims,profil_tab,profil_sous_tab,coord_debut,
30                          MPI_ORDER_C,MPI_INT,&type_sous_tab);
31  MPI_Type_commit(&type_sous_tab);
32
33  /* Permutation du sous-tableau */
34  MPI_Sendrecv_replace(tab,1,type_sous_tab,(rang+1)%2,etiquette,
35                      (rang+1)%2,etiquette,MPI_COMM_WORLD,&statut);
36
37  MPI_Type_free(&type_sous_tab);
38  MPI_Finalize();
39 }
```

# MPI-IO Vues

## Exemple 1 : lecture d'un fichier par blocs de deux éléments



**Figure 61** – Exemple 1 : lecture d'un fichier par blocs de deux éléments

```
> mpiexec -n 2 read_view01
```

```
Lecture processus 1 : 3, 4, 7, 8, 101, 102, 105, 106, 109, 110  
Lecture processus 0 : 1, 2, 5, 6, 9, 10, 103, 104, 107, 108
```

## Exemple 1 (suite)

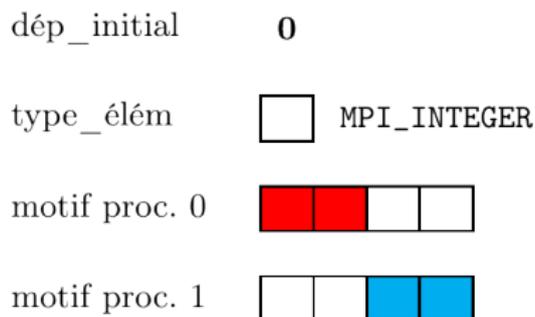


Figure 62 – Exemple 1 (suite) : création des vues sur le fichier

```
1 profil_tab[0] = 4;profil_sous_tab[0] = 2;  
2 if (rang == 0) coord[0]=0;  
3 if (rang == 1) coord[0]=2;  
4 MPI_Type_create_subarray(1,profil_tab,profil_sous_tab,coord,MPI_ORDER_C,MPI_INT,&motif);  
5 MPI_Type_commit(&motif);  
6 deplacement_initial=0  
7 MPI_File_set_view(descripteur,deplacement_initial,MPI_INT,motif,"native",MPI_INFO_NULL);
```

## Exemple 1 : code complet

```
1  /* read_view01 */
2  #include <mpi.h>
3  #include <stdio.h>
4  #include <stdlib.h>
5
6  int main(int argc, char *argv[]) {
7      int rang, i, coord, n1=4, n2=2, nb_valeurs=10;
8      int valeurs[nb_valeurs];
9      MPI_Datatype motif;
10     MPI_File descripteur;
11     MPI_Offset deplacement_initial;
12     MPI_Status statut;
13     MPI_Init(&argc, &argv);
14     MPI_Comm_rank(MPI_COMM_WORLD, &rang);
15     if (rang == 0) coord=0;
16     if (rang == 1) coord=2;
17     MPI_Type_create_subarray(1, &n1, &n2, &coord, MPI_ORDER_C, MPI_INT, &motif);
18     MPI_Type_commit(&motif);
19
20     MPI_File_open(MPI_COMM_WORLD, "donnees.dat", MPI_MODE_RDONLY, MPI_INFO_NULL,
21                 &descripteur);
22     deplacement_initial=0;
23     MPI_File_set_view(descripteur, deplacement_initial, MPI_INT, motif,
24                      "native", MPI_INFO_NULL);
25     MPI_File_read(descripteur, valeurs, nb_valeurs, MPI_INT, &statut);
26     printf("Lecture processus %d :", rang);
27     for(i=0; i<nb_valeurs; i++) {printf("%d ", valeurs[i]);} printf("\n");
28     MPI_File_close(&descripteur);
29     MPI_Finalize();
30 }
```

## Exemple 2 : utilisation successive de plusieurs vues

dép\_initial 0

type\_élément  MPI\_INTEGER

motif\_1 

dép\_initial 2 entiers

type\_élément  MPI\_INTEGER

motif\_2 

Figure 63 – Exemple 2 : utilisation successive de plusieurs vues

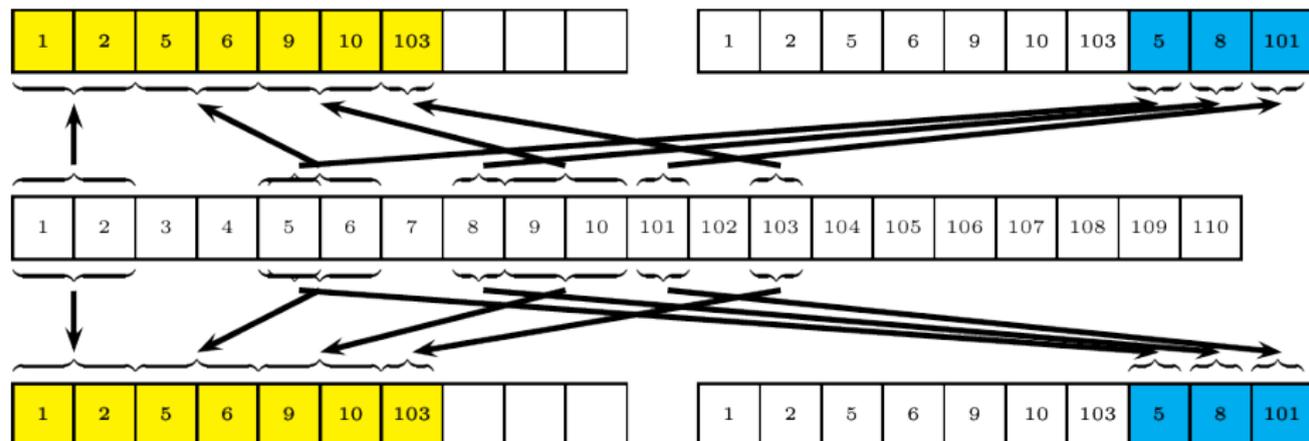
```
1 /* read_view02 */
2 #include <mpi.h>
3 #include <stdio.h>
4 #include <stdlib.h>
5
6 int main(int argc, char *argv[]) {
7     int rang, n1, n2, coord, nb_octets_entier, nb_valeurs=10, valeurs[nb_valeurs], i;
8     MPI_Datatype motif_1, motif_2;
9     MPI_File descripteur;
10    MPI_Offset deplacement_initial;
11    MPI_Status statut;
12
13    MPI_Init (&argc, &argv);
14    MPI_Comm_rank (MPI_COMM_WORLD, &rang);
```

## Exemple 2 (suite du code)

```
15  n1=4;n2=2;coord=0;
16  MPI_Type_create_subarray(1,&n1,&n2,&coord,MPI_ORDER_C,MPI_INT,&motif_1);
17  MPI_Type_commit(&motif_1);
18  n1=3;n2=1;coord=2;
19  MPI_Type_create_subarray(1,&n1,&n2,&coord,MPI_ORDER_C,MPI_INT,&motif_2);
20  MPI_Type_commit(&motif_2);
21
22  MPI_File_open(MPI_COMM_WORLD,"donnees.dat",MPI_MODE_RDONLY,MPI_INFO_NULL,
23              &descripteur);
24
25  /* Lecture en utilisant la premiere vue */
26  deplacement_initial=0;
27  MPI_File_set_view(descripteur,deplacement_initial,MPI_INT,motif_1,
28                  "native",MPI_INFO_NULL);
29  MPI_File_read(descripteur,valeurs,4,MPI_INT,&statut);
30  MPI_File_read(descripteur,&(valeurs[4]),3,MPI_INT,&statut);
31
32  /* Lecture en utilisant la seconde vue */
33  MPI_Type_size(MPI_INT,&nb_octets_entier);
34  deplacement_initial=2*nb_octets_entier;
35  MPI_File_set_view(descripteur,deplacement_initial,MPI_INT,motif_2,
36                  "native",MPI_INFO_NULL);
37  MPI_File_read(descripteur,&(valeurs[7]),3,MPI_INT,&statut);
38  printf("Lecture processus %d :",rang);
39  for(i=0;i<nb_valeurs;i++){printf("%d ",valeurs[i]);} printf("\n");
40  MPI_File_close(&descripteur);
41  MPI_Finalize();
42 }
```

# MPI-IO Vues

## Exemple 2 : illustration



```
> mpiexec -n 2 read_view02
```

```
Lecture processus 1 : 1, 2, 5, 6, 9, 10, 103, 5, 8, 101
```

```
Lecture processus 0 : 1, 2, 5, 6, 9, 10, 103, 5, 8, 101
```

# MPI-IO Vues

## Exemple 3 : gestion des trous dans les types de données

dép\_initial 0 entiers

type\_élément  MPI\_INTEGER

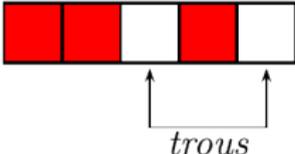
motif 

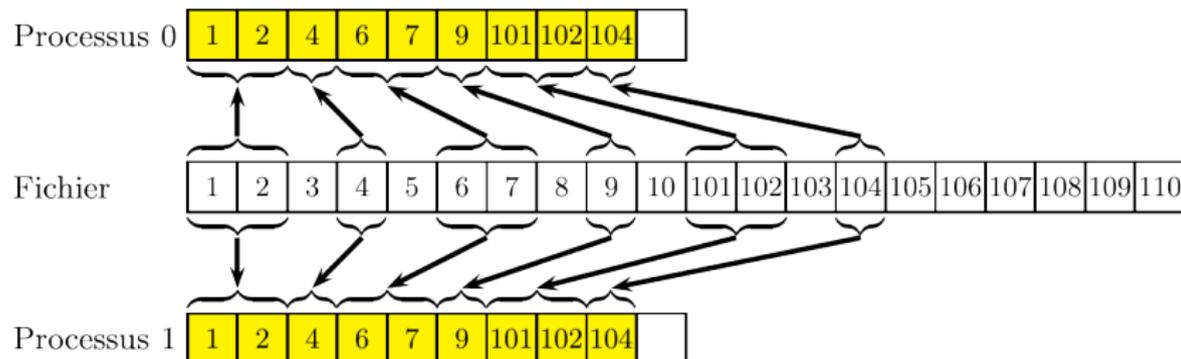
Figure 64 – Exemple 3 : gestion des trous dans les types de données

```
1  /* read_view03_indexed */
2  #include <mpi.h>
3  #include <stdio.h>
4  #include <stdlib.h>
5
6  int main(int argc, char *argv[]) {
7      int rang, nb_octets_entier, nb_valeurs=9, valeurs[nb_valeurs], i;
8      int longueurs[2], deplacements[2];
9      MPI_Datatype motif_temp, motif;
10     MPI_Aint borne_inf, etendue;
11     MPI_File descripteur;
12     MPI_Offset deplacement_initial;
13     MPI_Status statut;
```

## Exemple 3 (suite du code)

```
14 MPI_Init(&argc,&argv);
15 MPI_Comm_rank(MPI_COMM_WORLD,&rang);
16
17 deplacements[0]=0; deplacements[1]=3;longueurs[0]=2;longueurs[1]=1;
18 MPI_Type_indexed(2,longueurs,deplacements,MPI_INT,&motif_temp);
19
20 /* Motif : type MPI d'etendu 5*MPI_INT */
21 MPI_Type_size(MPI_INT,&nb_octets_entier);
22 MPI_Type_get_extent(motif_temp,&borne_inf,&etendue);
23 etendue = etendue+nb_octets_entier;
24 MPI_Type_create_resized(motif_temp,borne_inf,etendue+borne_inf,&motif);
25 MPI_Type_commit(&motif);
26
27 MPI_File_open(MPI_COMM_WORLD,"donnees.dat",MPI_MODE_RDONLY,MPI_INFO_NULL,
28              &descripteur);
29
30 deplacement_initial=0;
31 MPI_File_set_view(descripteur,deplacement_initial,MPI_INT,motif,
32                 "native",MPI_INFO_NULL);
33 MPI_File_read(descripteur,valeurs,9,MPI_INT,&statut);
34
35 printf("Lecture processus %d :",rang);
36 for(i=0;i<nb_valeurs;i++) {printf("%d ",valeurs[i]);} printf("\n");
37 MPI_File_close(&descripteur);
38 MPI_Finalize();
39 }
```

## Exemple 3 : illustration



```
> mpiexec -n 2 read_view03
```

```
Lecture, processus 0 : 1, 2, 4, 6, 7, 9, 101, 102, 104
```

```
Lecture, processus 1 : 1, 2, 4, 6, 7, 9, 101, 102, 104
```

## Exemple 3 : implémentation alternative utilisant un type structure

```
1  /* read_view03_struct */
2
3  profil_tab=3;profil_sous_tab=2;coord=0;
4  MPI_Type_create_subarray(1,&profil_tab,&profil_sous_tab,&coord,MPI_ORDER_C,MPI_INT,&temp_motif1);
5  profil_tab=2;profil_sous_tab=1;
6  MPI_Type_create_subarray(1,&profil_tab,&profil_sous_tab,&coord,MPI_ORDER_C,MPI_INT,&temp_motif2);
7  MPI_Type_size(MPI_INT,&nb_octets_entier);
8  déplacements[0]=0;déplacements[1]=3*nb_octets_entier;
9  bloc[0]=1;bloc[1]=1;type[0]=temp_motif1;type[1]=temp_motif2;
10 MPI_Type_create_struct(2,bloc,déplacements,type,&motif);
11 MPI_Type_commit(motif);
```

## Conclusion

MPI-IO offre une interface de haut niveau et un ensemble de fonctionnalités très riche. Il est possible de réaliser des opérations complexes et de tirer parti des optimisations implémentées dans la bibliothèque. MPI-IO offre aussi une bonne portabilité.

## Conseils

- L'utilisation des sous-programmes à positionnement explicite dans les fichiers doit être réservée à des cas particuliers, l'utilisation **implicite** de pointeurs individuels avec des vues offrant une interface de plus haut niveau.
- Lorsque les opérations font intervenir l'ensemble des processus (ou un sous-ensemble identifiable par un sous-communicateur MPI), il faut généralement privilégier la forme **collective** des opérations.
- Exactement comme pour le traitement des messages lorsque ceux-ci représentent une part importante de l'application, le **non-bloquant** est une voie privilégiée d'optimisation à mettre en œuvre par les programmeurs, mais ceci ne doit être implémenté qu'**après** qu'on se soit assuré du comportement correct de l'application en mode bloquant.

## Conclusion

## Conclusion

- Utiliser les communications point-à-point bloquantes, ceci avant de passer aux communications non-bloquantes. Il faudra alors essayer de faire du recouvrement calcul/communications.
- Utiliser les fonctions d'entrées-sorties bloquantes, ceci avant de passer aux entrées-sorties non-bloquantes. De même, il faudra alors faire du recouvrement calcul/entrées-sorties.
- Écrire les communications comme si les envois étaient synchrones (`MPI_Ssend()`).
- Éviter les barrières de synchronisation (`MPI_Barrier()`), surtout sur les fonctions collectives qui sont bloquantes.
- La programmation mixte MPI/OpenMP peut apporter des gains d'extensibilité, mais pour que cette approche fonctionne bien, il est évidemment nécessaire d'avoir de bonnes performances OpenMP à l'intérieur de chaque processus MPI. Un cours est dispensé à l'IDRIS (<https://cours.idris.fr>).

## Travaux pratiques MPI – Exercice 8 : Équation de Poisson

On considère l'équation de Poisson suivante :

$$\begin{cases} \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} = f(x, y) & \text{dans } [0, 1] \times [0, 1] \\ u(x, y) = 0 & \text{sur les frontières} \\ f(x, y) = 2 \cdot (x^2 - x + y^2 - y) \end{cases}$$

On va résoudre cette équation avec une méthode de décomposition de domaine :

- L'équation est discretisée sur le domaine via la méthode des différences finies.
- Le système obtenu est résolu avec un solveur suivant la méthode de Jacobi.
- Le domaine global est découpé en sous domaines.

La solution exacte est connue et est  $u_{exacte}(x, y) = xy(x - 1)(y - 1)$

## Travaux pratiques MPI – Exercice 8 : Équation de Poisson

Pour discrétiser l'équation, on définit une grille constituée d'un ensemble de points  $(x_i, y_j)$

$$x_i = i h_x \quad \text{pour } i = 0, \dots, ntx + 1$$

$$y_j = j h_y \quad \text{pour } j = 0, \dots, nty + 1$$

$$h_x = \frac{1}{(ntx + 1)}$$

$$h_y = \frac{1}{(nty + 1)}$$

$h_x$  : pas suivant  $x$

$h_y$  : pas suivant  $y$

$ntx$  : nombre de points intérieurs suivant  $x$

$nty$  : nombre de points intérieurs suivant  $y$

Il y a au total  $ntx+2$  points suivant  $x$  et  $nty+2$  points suivant  $y$

## Travaux pratiques MPI – Exercice 8 : Équation de Poisson

- Soit  $u_{ij}$  l'estimation de la solution à la position  $x_i = ih_x$  et  $y_j = jh_y$ .
- La méthode de Jacobi consiste à calculer :

$$u_{ij}^{n+1} = c_0(c_1(u_{i+1j}^n + u_{i-1j}^n) + c_2(u_{ij+1}^n + u_{ij-1}^n) - f_{ij})$$

$$\text{avec : } c_0 = \frac{1}{2} \frac{h_x^2 h_y^2}{h_x^2 + h_y^2}$$

$$c_1 = \frac{1}{h_x^2}$$

$$c_2 = \frac{1}{h_y^2}$$

## Travaux pratiques MPI – Exercice 8 : Équation de Poisson

- En parallèle, les valeurs aux interfaces des sous-domaines doivent être échangées entre les voisins.
- On utilise des cellules fantômes, ces cellules servent de buffer de réception pour les échanges entre voisins.

## Travaux pratiques MPI – Exercice 8 : Équation de Poisson

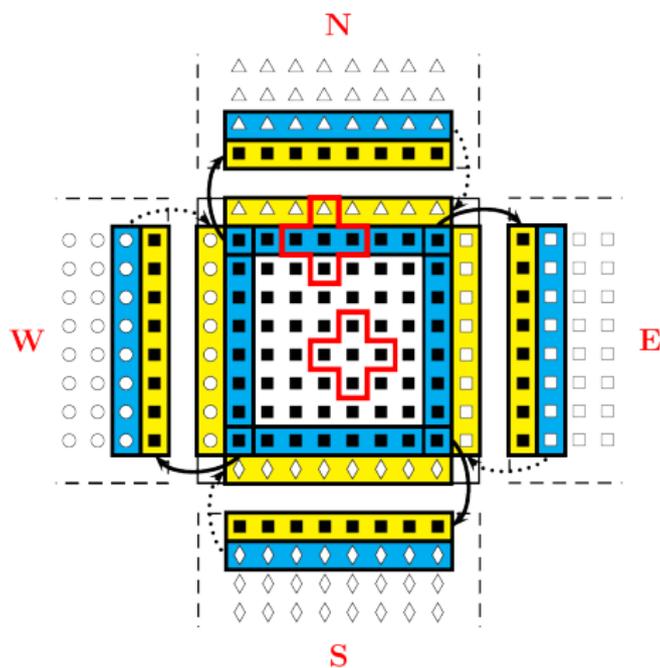


Figure 65 – Échange de points aux interfaces

## Travaux pratiques MPI – Exercice 8 : Équation de Poisson

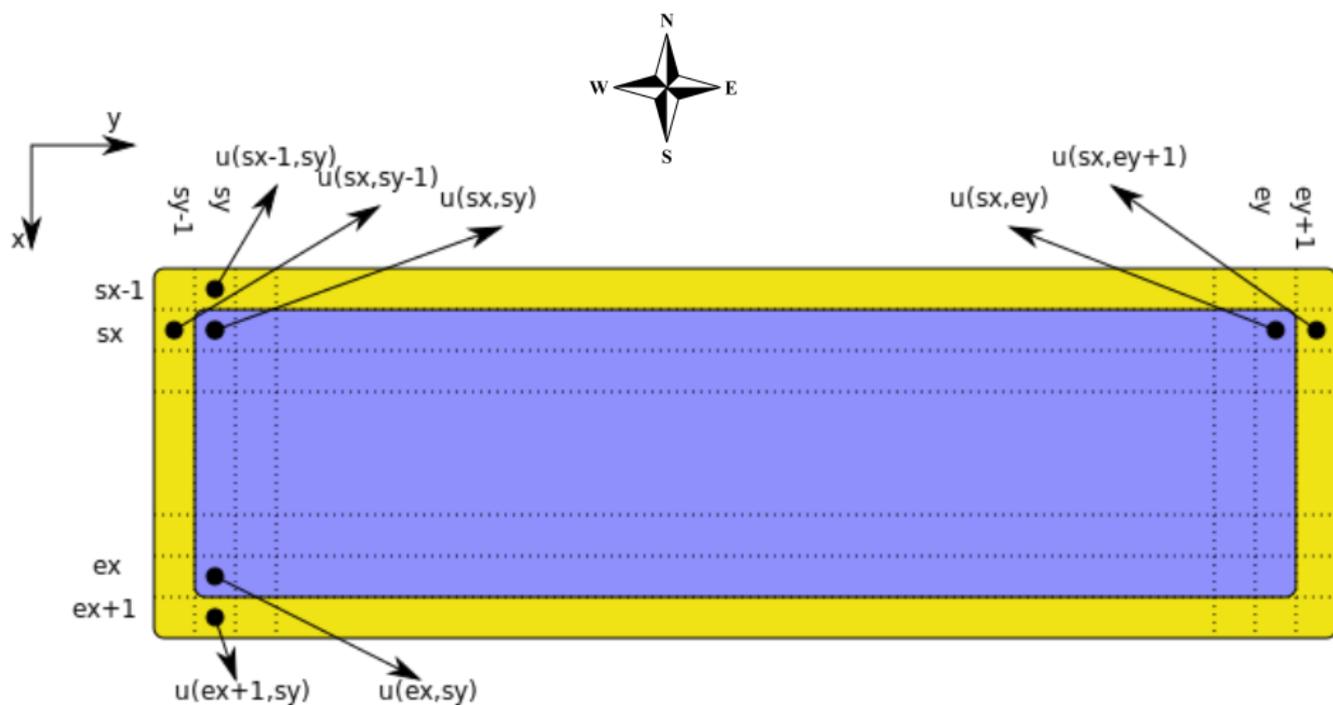


Figure 66 – Numérotation des points dans les différents sous-domaines

## Travaux pratiques MPI – Exercice 8 : Équation de Poisson

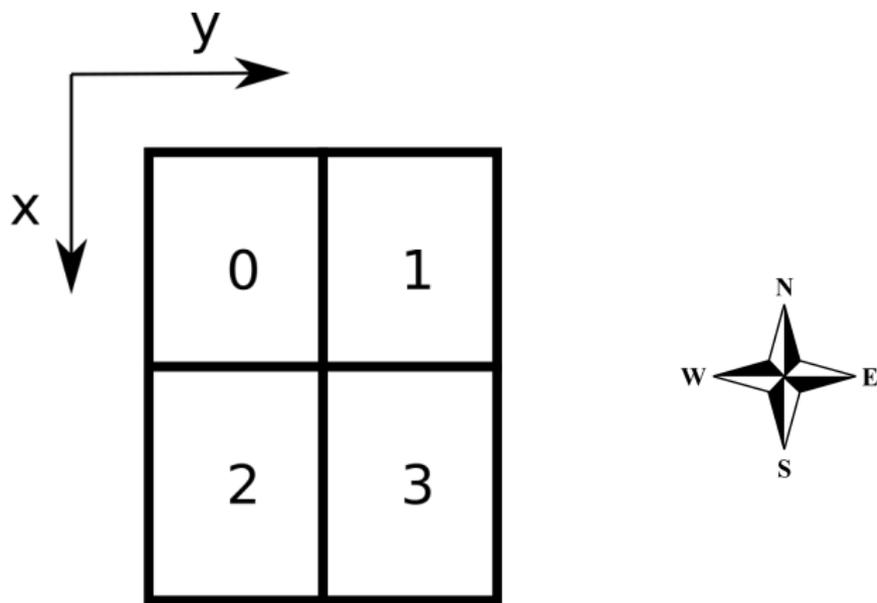
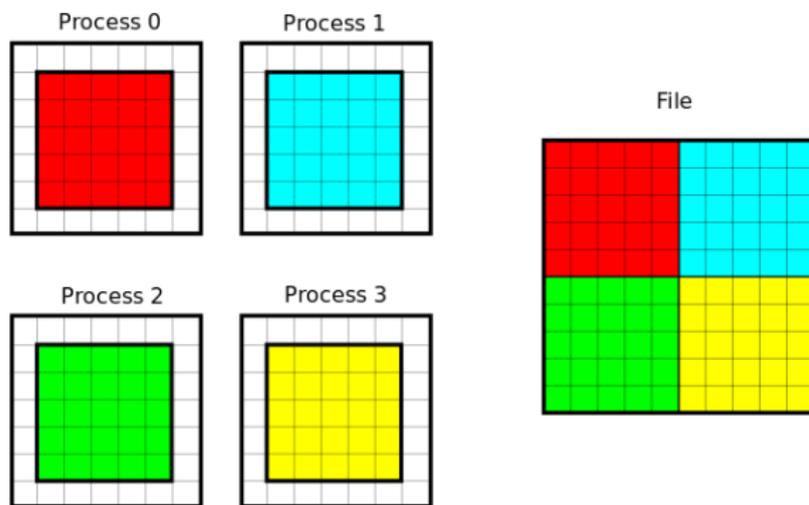


Figure 67 – Rang correspondant aux différents sous-domaines

## Travaux pratiques MPI – Exercice 8 : Équation de Poisson



**Figure 68** – Ecriture de la matrice  $u$  globale dans un fichier

Il s'agit de définir :

- Une vue, pour ne voir dans le fichier que la partie de la matrice  $u$  globale que l'on possède ;
- Un type afin d'écrire la matrice  $u$  locale (sans les cellules fantômes) ;
- Appliquer la vue au fichier ;
- Faire l'écriture en une fois.

## Travaux pratiques MPI – Exercice 8 : Équation de Poisson

- initialiser l'environnement MPI ;
- créer la topologie cartésienne 2D ;
- déterminer les indices de tableau pour chaque sous-domaine ;
- déterminer les 4 processus voisins d'un processus traitant un sous-domaine donné ;
- créer deux types dérivés *type\_ligne* et *type\_colonne* ;
- échanger les valeurs aux interfaces avec les autres sous-domaines ;
- calculer l'erreur globale. Lorsque l'erreur globale sera inférieure à une valeur donnée (précision machine par exemple), alors on considérera qu'on a atteint la solution ;
- reformer la matrice *u* globale (identique à celle obtenue avec la version monoprocesseur) dans un fichier `donnees.dat`.

## Travaux pratiques MPI – Exercice 8 : Équation de Poisson

- Un squelette de la version parallèle est proposé : il s'agit d'un programme principal (`poisson.c`) et de plusieurs sous-programmes. Les modifications sont à effectuer dans le fichier `parallel.c`.
- Pour compiler utilisez la commande `make`, pour exécuter le code utilisez la commande `make exe`. Pour vérifier les résultats, utilisez la commande `make verification` qui exécute un programme de relecture du fichier `donnees.dat` puis le compare avec la version monoprocasseur.