

# Programmation Parallèle avec MPI



**Babacar Diop**

*Dpt. d'Informatique*

**UFR des Sciences Appliquées et de Technologies**  
**Université Gaston Berger de Saint-Louis**

2018/2019

# Sommaire

- ✓ Gestion de l'environnement MPI
- ✓ Communications point à point
- ✓ Communication collective
- ✓ Types de données MPI

# Message Passing Interface

- **Norme universelle** pour fournir une communication entre tous les processus dans un système de mémoire distribuée.
- La plupart des plates-formes informatiques parallèles couramment utilisées fournissent au moins une implémentation de MPI.
- Implémenté comme une collection de fonctions prédéfinies dans une bibliothèque et peuvent être appelées à partir de langages telles que C, C ++, Fortran, etc ...

# Message Passing Interface

- MPI est la première bibliothèque de transmission de messages normalisée, indépendante du vendeur.
- Avantages : portabilité, efficacité et flexibilité.
- MPI n'est pas une norme IEEE ou ISO, mais un "standard" pour l'écriture de programmes de transmission de messages sur des plates-formes HPC(High Performance Computing).

# MPI interface

MPI Library	Where?	Compilers
<b>MVAPICH</b>	Linux clusters	GNU, Intel, PGI, Clang
<b>Open MPI</b>	Linux clusters	GNU, Intel, PGI, Clang
<b>Intel MPI</b>	Linux clusters	Intel, GNU
<b>IBM BG/Q MPI</b>	BG/Q clusters	IBM, GNU
<b>IBM Spectrum MPI</b>	Coral Early Access and Sierra clusters	IBM, GNU, PGI, Clang

# MPI interface

MPI Build Scripts - Linux Clusters			
Implementation	Language	Script Name	Underlying Compiler
MVAPCH2	C	<code>mpicc</code>	C compiler for loaded compiler package
	C++	<code>mpicxx</code> <code>mpic++</code>	C++ compiler for loaded compiler package
	Fortran	<code>mpif77</code>	Fortran77 compiler for loaded compiler package. Points to mpifort.
		<code>mpif90</code>	Fortran90 compiler for loaded compiler package. Points to mpifort.
		<code>mpifort</code>	Fortran 77/90 compiler for loaded compiler package.
Open MPI	C	<code>mpicc</code>	C compiler for loaded compiler package
	C++	<code>mpiCC</code> <code>mpic++</code> <code>mpicxx</code>	C++ compiler for loaded compiler package
	Fortran	<code>mpif77</code>	Fortran77 compiler for loaded compiler package. Points to mpifort.
		<code>mpif90</code>	Fortran90 compiler for loaded compiler package. Points to mpifort.
		<code>mpifort</code>	Fortran 77/90 compiler for loaded compiler package.

# Niveau de support de threads

- Les bibliothèques MPI varient dans leur niveau de support de thread:
  - **MPI\_THREAD\_SINGLE** - **Niveau 0**: un seul thread.
  - **MPI\_THREAD\_FUNNELED** - **Niveau 1**: multi-threading avec un seul thread principal où sont dirigés tous les appels MPI.
  - **MPI\_THREAD\_SERIALIZED** - **Niveau 2**: multi-threading et sérialisation. Absence de parallélisme
  - **MPI\_THREAD\_MULTIPLE** - **Niveau 3**: Multithreading pouvant appeler MPI sans restrictions. Parallélisme absolue

# Niveau de support de threads

Exemple simple permettant de déterminer la prise en charge au niveau du thread.

```
#include "mpi.h"
#include <stdio.h>
int main( int argc, char *argv[] ) {
    int provided, claimed;
    /** Select one of the following
    MPI_Init_thread( 0, 0, MPI_THREAD_SINGLE, &provided );
    MPI_Init_thread( 0, 0, MPI_THREAD_FUNNELED, &provided );
    MPI_Init_thread( 0, 0, MPI_THREAD_SERIALIZED, &provided );
    MPI_Init_thread( 0, 0, MPI_THREAD_MULTIPLE, &provided );
    ***/
    MPI_Init_thread( 0, 0, MPI_THREAD_MULTIPLE, &provided );
    MPI_Query_thread( &claimed );
    printf( "Requete support de thread = %d Niveau de support= %d\n", claimed,
    provided );
    MPI_Finalize();
}
```



# Communicateurs et groupes : MPI

- ✓ MPI utilise des objets appelés communicateurs et groupes pour définir quel ensemble de processus peut communiquer entre eux
- ✓ La plupart des routines MPI nécessite de spécifier un communicateur comme argument
  - ✓ **MPI\_COMM\_WORLD** : communicateur prédéfini qui inclut tous vos processus MPI.

# Rang d'un processus

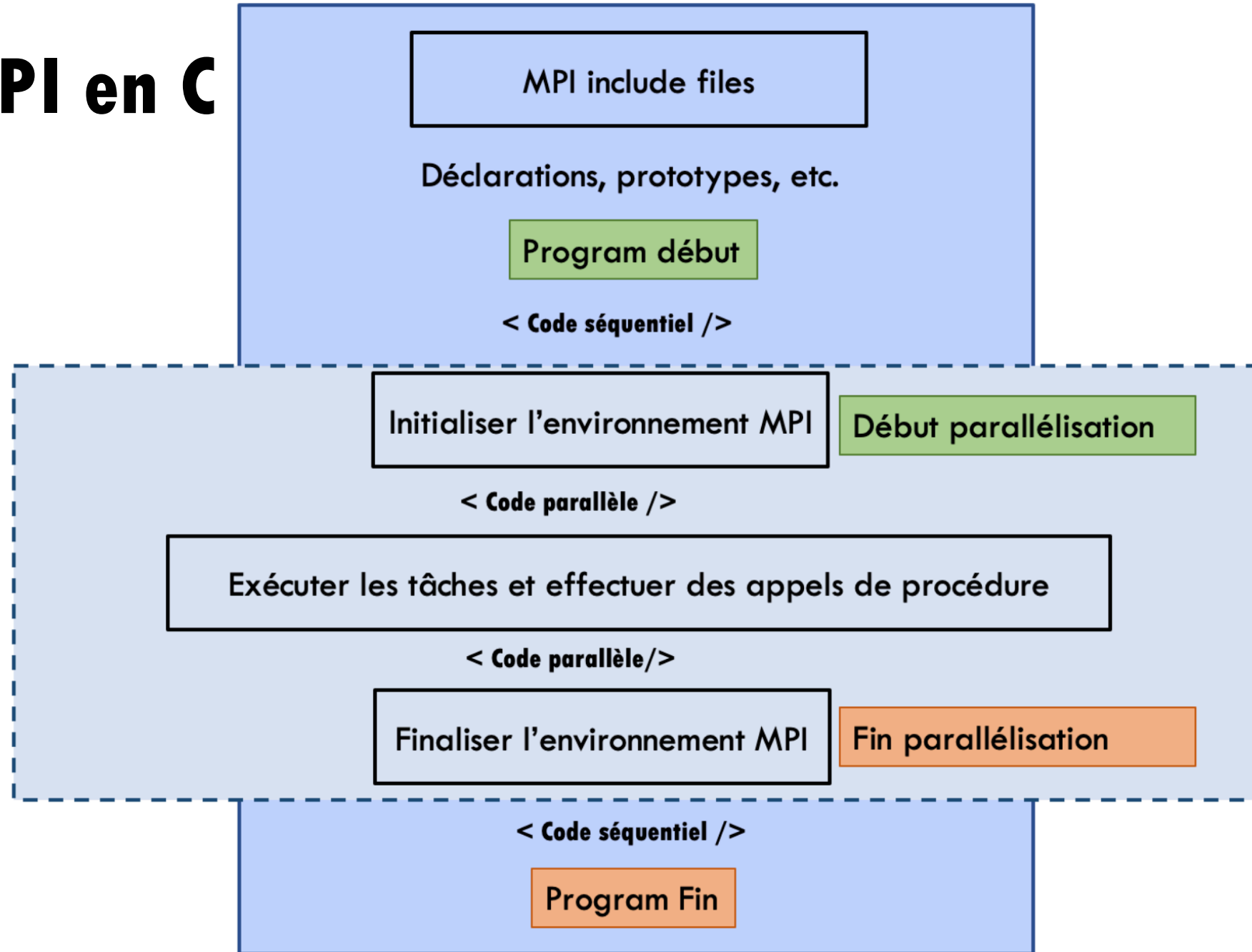
- Dans un communicateur, chaque processus a son propre identifiant, un entier unique attribué par le système lors de l'initialisation du processus.
- Un rang est parfois aussi appelé "ID du processus". Les rangs sont contigus et commencent à zéro.
- Les rangs sont utilisés par le programmeur pour spécifier la source et la destination des messages. Souvent utilisé de manière conditionnelle par l'application pour contrôler l'exécution du programme

(if rank = 0, do this / if rank = 1, do that)

# Gestion des erreurs

- La plupart des routines MPI incluent un paramètre de code retour/erreur.
- Toutefois, conformément à la norme MPI, le comportement par défaut d'un appel MPI consiste à sortir en cas d'erreur. Cela signifie que vous ne pourrez probablement pas capturer un code d'erreur autre que **MPI\_SUCCESS (zero)**.

# Structure Code MPI en C



# Initialisation: **MPI\_Init**

- **MPI\_Init**

- Initialise l'environnement d'exécution MPI.
- Cette fonction doit être appelée dans chaque programme MPI, avant toute autre fonction MPI et une seule fois dans un programme MPI.
- Il peut être utilisé pour transmettre les arguments en ligne de commande à tous les processus
  - **MPI\_Init (&argc, &argv)**
  - **MPI\_INIT (ierr)**

# Communicateur: **MPI\_Comm\_size**

- **MPI\_Comm\_size**

- Retourne le nombre total de processus MPI dans le communicateur spécifié, tel que MPI\_COMM\_WORLD. Si le communicateur est MPI\_COMM\_WORLD, il représente le nombre de processus disponibles pour votre application.

- **MPI\_Comm\_size (comm, & size)**

- **MPI\_COMM\_SIZE (comm, size, ierr)**

# Communication : `MPI_Comm_rank`

- **`MPI_Comm_rank`**

- Renvoie le rang du processus MPI appelant dans le communicateur spécifié. Initialement, chaque processus se voit attribuer un rang entier unique compris entre 0 et le nombre de tâches - 1 au sein du communicateur `MPI_COMM_WORLD`. Si un processus devient associé à d'autres communicateurs, il aura également un rang unique dans chacun d'eux.

- `MPI_Comm_rank(comm, & rank)`

- `MPI_COMM_RANK(comm, rang, ierr)`

# MPI\_Abort

- **MPI\_Abort**

- Termine tous les processus MPI associés au communicateur. Dans la plupart des implémentations MPI, il met fin à TOUS les processus, quel que soit le communicateur spécifié.
  - **MPI\_Abort (comm, errorcode)**
  - **MPI\_ABORT (comm, errorcode, ierr)**



# Nom de Processeur

- **MPI\_Get\_processor\_name**

- Renvoie le nom du processeur.
- Retourne également la longueur du nom. La mémoire tampon pour "nom" doit comporter au moins MPI\_MAX\_PROCESSOR\_NAME en caractères. Ce qui est retourné dans "name" dépend de la mise en œuvre.

- **MPI\_Get\_processor\_name (& name, & resultlength)**
- **MPI\_GET\_PROCESSOR\_NAME (name, resultlength, ierr)**

# Version de MPI

- **MPI\_Get\_version**

- Renvoie la version du standard MPI implémentée par la bibliothèque.

- **MPI\_Get\_version (& version, & subversion)**

- **MPI\_GET\_VERSION (version, sous-version, ierr)**

# MPI initialisé ?

- **MPI\_Initialized**

- Indique si **MPI\_Init** a été appelé - retourne un état logique vrai (1) ou faux (0).
- MPI nécessite que **MPI\_Init** soit appelé une seule fois par processus. Cela peut poser problème aux modules qui souhaitent utiliser MPI et sont prêts à appeler **MPI\_Init** si nécessaire. **MPI\_Initialized** résout ce problème.
  - **MPI\_Initialized (&flag)**
  - **MPI\_INITIALIZED (indicateur, ierr)**

# Timing

- **MPI\_Wtime**

- Renvoie une durée d'horloge murale écoulée en secondes (double précision) sur le processeur appelant.

- **MPI\_Wtime ()**

- **MPI\_WTIME ()**

- **MPI\_Wtick**

- Renvoie la résolution en secondes (double précision) de MPI\_Wtime.

- **MPI\_Wtick ()**

- **MPI\_WTICK ()**

# Finalisation

- **MPI\_Finalize**
- Termine l'environnement d'exécution MPI.
  - Cette fonction doit être la dernière routine MPI appelée dans chaque programme MPI - aucune autre routine MPI ne peut être appelée après.
    - **MPI\_Finalize ()**
    - **MPI\_FINALIZE (ierr)**

# Exemple illustratif

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

int main(int argc, char *argv[]) {
    int numtasks, rank, len, rc;
    char hostname[MPI_MAX_PROCESSOR_NAME];
    // initialize MPI
    MPI_Init(&argc,&argv); // get number of tasks
    MPI_Comm_size(MPI_COMM_WORLD,&numtasks);
    MPI_Comm_rank(MPI_COMM_WORLD,&rank);
    MPI_Get_processor_name(hostname, &len);
    printf ("No de taches = %d Mon Rang= %d exécutant sur %s\n",
numtasks,rank,hostname);
    // do some work with message passing
    // done with MPI
    MPI_Finalize();
}
```

# Communication point à point

- **Types d'opérations**

- Opérations impliquant le transfert de messages entre deux tâches MPI différentes
- Une tâche effectue l'envoi et l'autre la réception correspondante
  - Envoi synchrone
  - Blocage de l'envoi/blocage de la réception
  - Envoi non bloquant / réception non bloquante
  - Envoi tamponné
  - Envoi / réception combinés
  - "Prêt" à envoyer
- Toute fonction d'envoi peut être associée à n'importe quelle fonction de réception.

# Buffering

- **Rôle**

- Mise en mémoire tampon

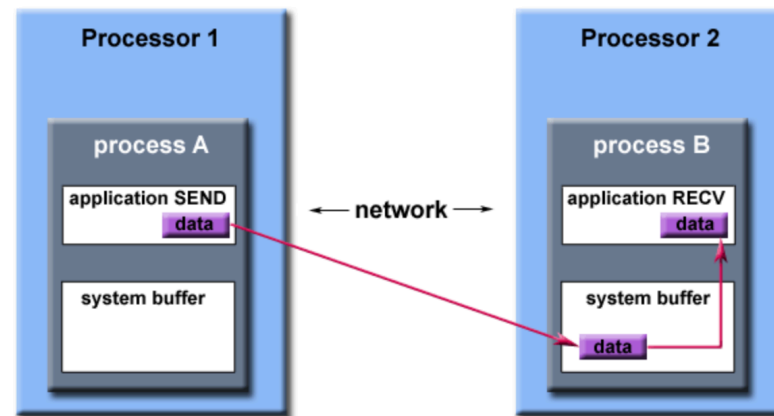
- **Objectif**

- Gérer le stockage des données lorsque deux tâches ne sont pas synchronisées



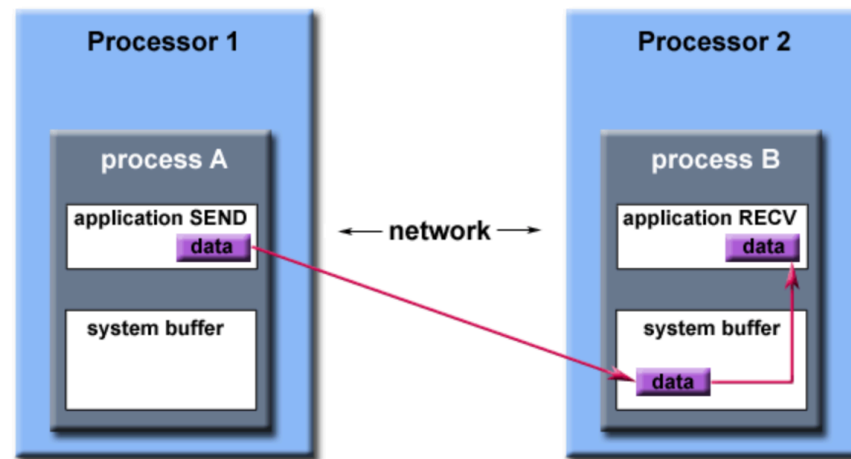
# Buffering

- Considérons les deux cas suivants:
  - Une opération d'envoi se produit 5 secondes avant que la réception soit prête. Où se trouve le message pendant que la réception est en attente?
  - Plusieurs envois arrivent à la même tâche de réception qui ne peut accepter qu'un seul envoi à la fois. Qu'advient-il des messages en cours de "sauvegarde"?



# Buffering

- L'implémentation MPI décide de ce qu'il advient des données dans ces types de cas.
- En règle générale, une zone de mémoire tampon appelée **buffer** est réservée aux données en transit.
- Par exemple:



# Buffering

- La mémoire tampon du système est:
  - Opaque pour le programmeur et entièrement gérée par la bibliothèque MPI
  - Une ressource finie qui peut être facile à épuiser
  - Capable d'exister du côté de l'envoi, du côté de la réception ou des deux
  - Un élément susceptible d'améliorer les performances du programme car il permet aux opérations d'envoi-réception d'être asynchrones.

# **Communication Point-à-Point sur MPI**

# Envoie/Réception bloquants

- Un envoi bloquant ne « retournera » qu'après la possibilité de modifier le tampon d'application (vos données d'envoi) en vue de sa réutilisation.
- Un envoi bloquant peut être synchrone, ce qui signifie qu'un échange de données est en cours avec la tâche de réception pour confirmer un envoi sécurisé.
- Un envoi bloquant peut être asynchrone si un tampon système est utilisé pour conserver les données en vue de leur remise éventuelle à la réception.
- Une réception bloquante ne « retournera » que lorsque les données sont arrivées et sont prêtes à être utilisées par le programme.

# Envoie/Réception non bloquants

- Les routines d'envoi et de réception non bloquantes n'attendent pas que des événements de communication soient terminés, tels que la copie de messages de la mémoire de l'utilisateur dans l'espace tampon du système ou l'arrivée effective du message.
- Les opérations non bloquantes "demandent" simplement à la bibliothèque MPI d'effectuer l'opération quand elle le peut. L'utilisateur ne peut pas prédire quand cela se produira.

# Envoie/Réception non bloquants

- Il est dangereux de modifier le tampon d'application (votre espace variable) jusqu'à ce que vous sachiez que l'opération de demandée a bien été effectuée. Il y a des routines "d'attente" utilisées pour cela.
- Les communications non bloquantes sont principalement utilisées pour superposer les calculs aux communications et exploiter les gains de performances possibles.

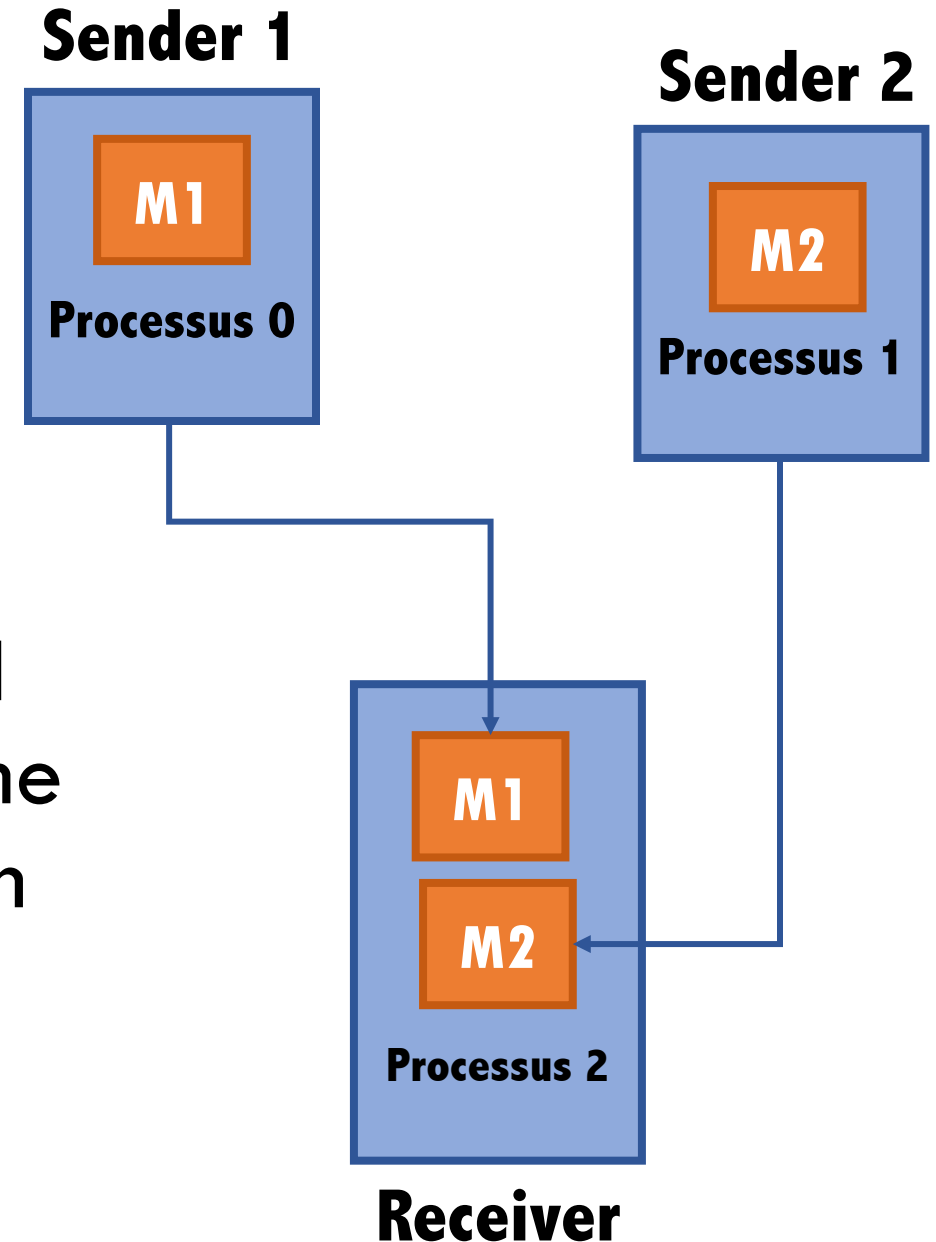
# Example

Blocking Send	Non-blocking Send
<pre>myvar = 0;  for (i=1; i&lt;ntasks; i++) {     task = i;     MPI_Send (&amp;myvar ... .. task ...);     myvar = myvar + 2      /* do some work */ }</pre>	<pre>myvar = 0;  for (i=1; i&lt;ntasks; i++) {     task = i;     MPI_Isend (&amp;myvar ... .. task ...);     myvar = myvar + 2;      /* do some work */      MPI_Wait (...); }</pre>
Safe. Why?	Unsafe. Why?



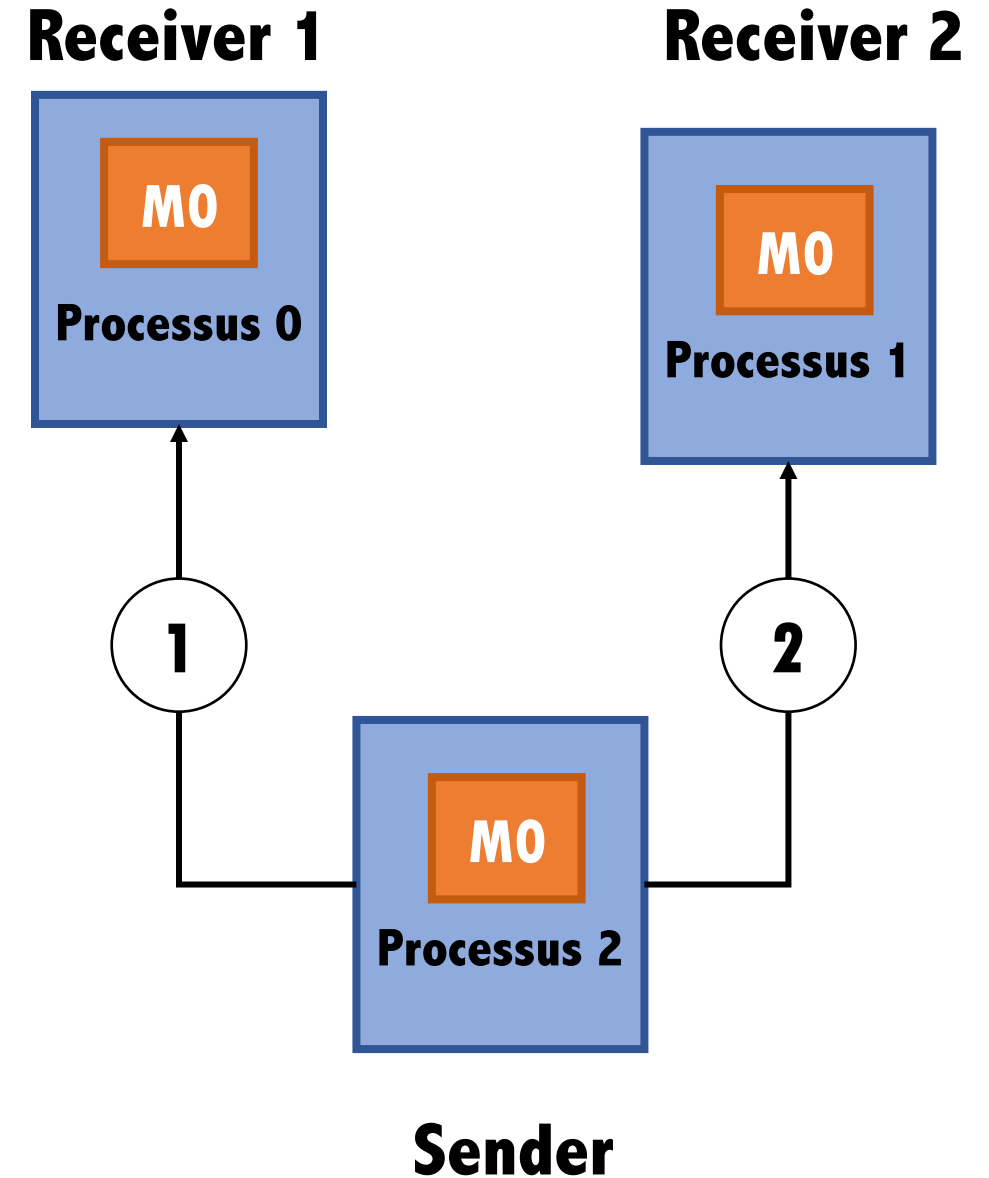
# Ordre d'envoi

- MPI garantit l'ordre des messages.
- Si 2 messages sont envoyés (M1 et M2) successivement à la même destination, l'ordre de réception sera M1 ensuite M2



# Ordre de réception

- Si 2 destinataires postent une opération de réception successivement (R1, R2) depuis le même message, alors R1 sera reçu R2

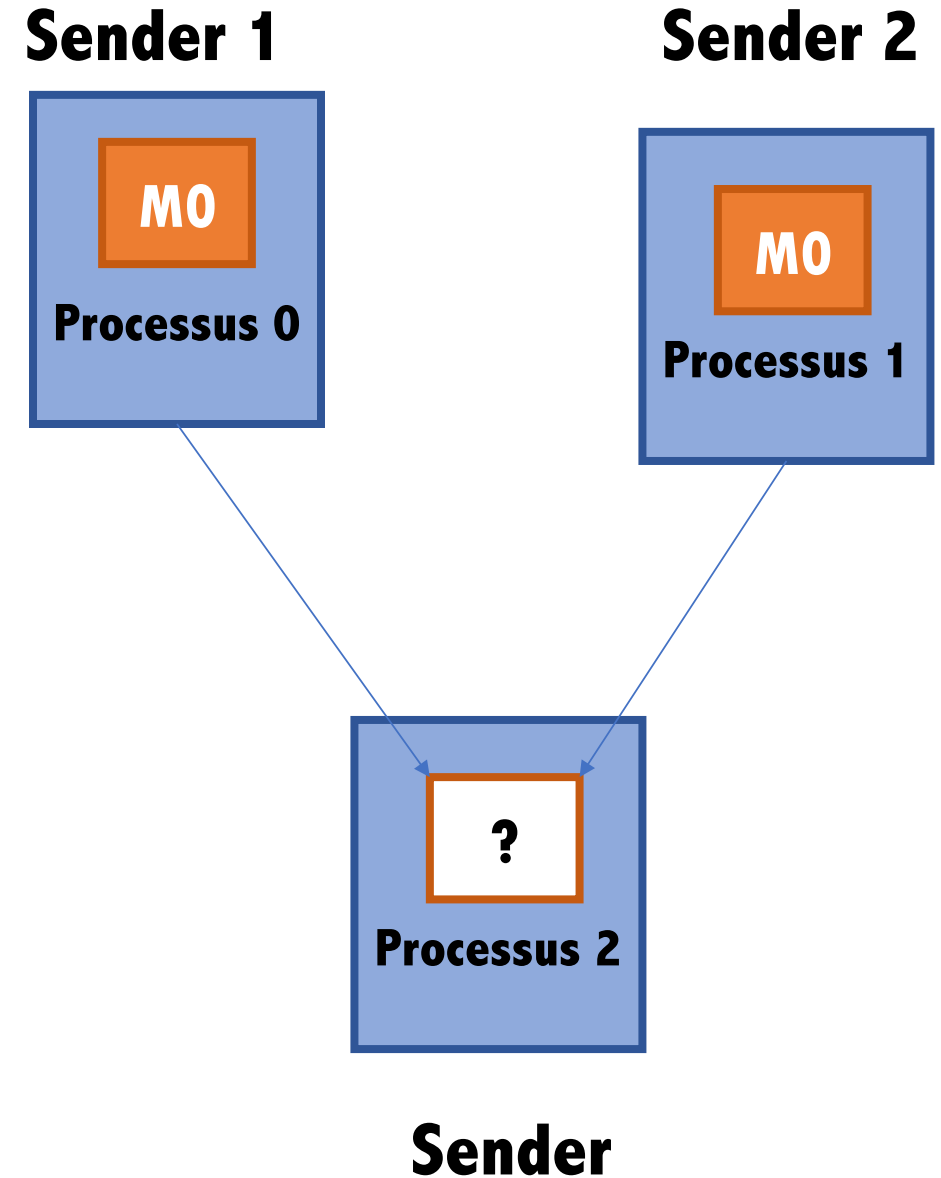


# Ordre de réception

**MPI ne garantit pas tous les cas**

Exemple: la tâche 0 envoie un message à la tâche 2. La tâche 1 envoie un message concurrent à la tâche 2.

Un seul des envois sera effectué.  
Lequel ? Juste un des 2



# Fonctions de passage de message

- Les fonctions de communication point à point de MPI ont généralement une liste d'arguments prenant l'un des formats suivants:

<b>Envoi bloquant</b>	<code>MPI_Send(buffer, count, type, dest, tag, comm)</code>
<b>Envoi non bloquant</b>	<code>MPI_Isend(buffer, count, type, dest, tag, comm, request)</code>
<b>Réception bloquant</b>	<code>MPI_Recv(buffer, count, type, source, tag, comm, status)</code>
<b>Réception non bloquante</b>	<code>MPI_Irecv(buffer, count, type, source, tag, comm, request)</code>

# Arguments de fonctions

## buffer

- Espace mémoire référençant les données à envoyer ou à recevoir.
- Dans la plupart des cas, il s'agit simplement du nom de la variable à envoyer/recevoir.
- Pour les programmes C, cet argument est passé par référence

## count

- Indique le nombre d 'éléments du type à envoyer

# Arguments de fonctions

## source, dest

- Rang du processus émetteur/récepteur

## tag

- Entier positif aléatoire attribué par le programmeur pour identifier un message de manière unique. Les opérations d'envoi et de réception doivent correspondre au tag de message.
  - `MPI_ANY_TAG` permet de recevoir n'importe quel message, quelle que soit le tag.

## comm

- Contexte de communication
- Exemple: `MPI_COMM_WORLD` est généralement utilisé

# Arguments de fonctions

## status

- Pointeur sur une structure prédéfinie **MPI\_Status**

## request

- Utilisé par les opérations d'envoi et de réception non bloquantes
- Pointeur sur une structure prédéfinie **MPI\_Request**

# Arguments de fonctions

## Types de données MPI

C Data Types		
<b>MPI_CHAR</b>	char	<b>MPI_CHARACTER</b>
<b>MPI_WCHAR</b>	wchar_t - wide character	
<b>MPI_SHORT</b>	signed short int	
<b>MPI_INT</b>	signed int	<b>MPI_INTEGER</b> <b>MPI_INTEGER1</b> <b>MPI_INTEGER2</b> <b>MPI_INTEGER4</b>
<b>MPI_LONG</b>	signed long int	
<b>MPI_LONG_LONG_INT</b> <b>MPI_LONG_LONG</b>	signed long long int	
<b>MPI_SIGNED_CHAR</b>	signed char	
<b>MPI_UNSIGNED_CHAR</b>	unsigned char	
<b>MPI_UNSIGNED_SHORT</b>	unsigned short int	
<b>MPI_UNSIGNED</b>	unsigned int	
<b>MPI_UNSIGNED_LONG</b>	unsigned long int	
<b>MPI_UNSIGNED_LONG_LONG</b>	unsigned long long int	



# Arguments de fonctions

## Types de données MPI

<b>MPI_FLOAT</b>	float	<b>MPI_REAL</b> <b>MPI_REAL2</b> <b>MPI_REAL4</b> <b>MPI_REAL8</b>
<b>MPI_DOUBLE</b>	double	<b>MPI_DOUBLE_PRECISION</b>
<b>MPI_LONG_DOUBLE</b>	long double	
<b>MPI_C_COMPLEX</b> <b>MPI_C_FLOAT_COMPLEX</b>	float _Complex	<b>MPI_COMPLEX</b>
<b>MPI_C_DOUBLE_COMPLEX</b>	double _Complex	<b>MPI_DOUBLE_COMPLEX</b>
<b>MPI_C_LONG_DOUBLE_COMPLEX</b>	long double _Complex	
<b>MPI_C_BOOL</b>	_Bool	<b>MPI_LOGICAL</b>
<b>MPI_INT8_T</b> <b>MPI_INT16_T</b> <b>MPI_INT32_T</b> <b>MPI_INT64_T</b>	int8_t int16_t int32_t int64_t	
<b>MPI_UINT8_T</b> <b>MPI_UINT16_T</b> <b>MPI_UINT32_T</b> <b>MPI_UINT64_T</b>	uint8_t uint16_t uint32_t uint64_t	
<b>MPI_BYTE</b>	8 binary digits	<b>MPI_BYTE</b>
<b>MPI_PACKED</b>	data packed or unpacked with MPI_Pack()/ MPI_Unpack	<b>MPI_PACKED</b>

# Comm bloquante

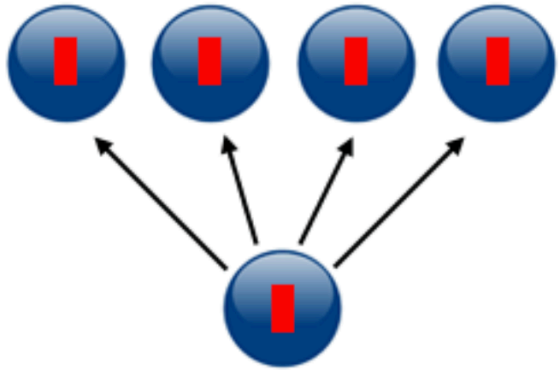
- MPI\_Send (&buf,count,datatype,dest,tag,comm)
- MPI\_Recv (&buf,count,datatype,source,tag,comm,&status)
- MPI\_Ssend (&buf,count,datatype,dest,tag,comm)
- MPI\_Sendrecv(&sendbuf,sendcount,sendtype,dest,sendtag,&recvbuf,recvcount,recvtype,source,recvtag,comm,&status)
- MPI\_Probe (source,tag,comm,&status)
- MPI\_Get\_count (&status,datatype,&count)
- MPI\_Get\_count (&status,datatype,&count)

# Comm non bloquante

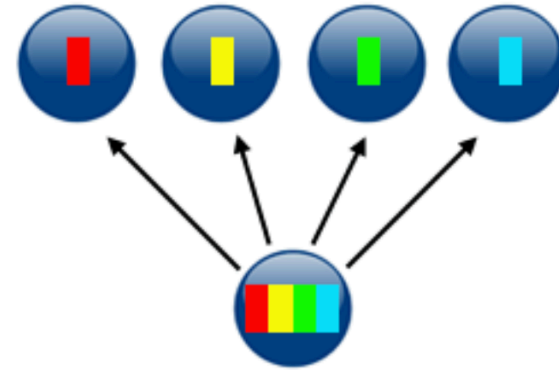
- MPI\_Isend (&buf,count,datatype,dest,tag,comm,&request)
  - MPI\_Irecv (&buf,count,datatype,source,tag,comm,&request)
  - MPI\_Ssend (&buf,count,datatype,dest,tag,comm,&request)
  - MPI\_Iprobe (source,tag,comm,&flag,&status)
  - MPI\_Test (&request,&flag,&status)
- MPI\_Testany (count,&array\_of\_requests,&index,&flag,&status)
- MPI\_Testall (count,&array\_of\_requests,&flag,&array\_of\_statuses)
- MPI\_Testsome (incount,&array\_of\_requests,&outcount, &array\_of\_offsets, &array\_of\_statuses)

# **Communication collective sur MPI**

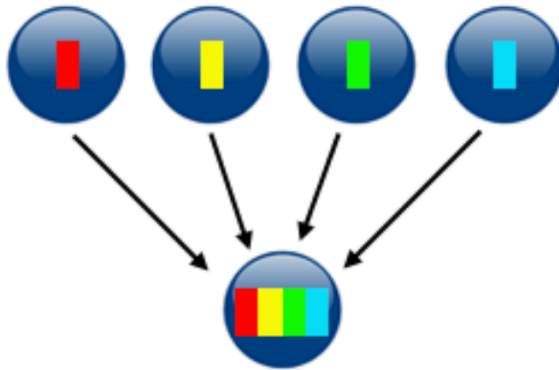
# Types de communication collective



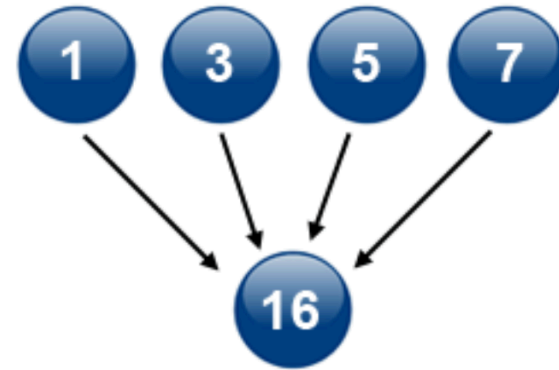
**Diffusion**



**Diffusion ciblée**



**Regroupement**



**Aggrégation**

# Caractéristiques

- Les routines de communication collectives ne prennent pas de tag
- Lors d'une opération collective, les processus sont partitionnés en de nouveaux groupes, puis reliés à de nouveaux communicateurs
- Utilisable uniquement avec les types de données prédéfinis MPI - pas avec les types de données dérivées MPI.

# Fonctions de comm collective

- MPI\_Barrier (comm)
- MPI\_Bcast (&buffer,count,datatype,root,comm)
- MPI\_Scatter (&sendbuf,sendcnt,sendtype,&recvbuf, recvcnt,recvtype,root,comm)
- MPI\_Gather (&sendbuf,sendcnt,sendtype,&recvbuf, recvcount,recvtype,root,comm)
- MPI\_Allgather (&sendbuf,sendcount,sendtype,&recvbuf, recvcount,recvtype,comm)
- MPI\_Reduce (&sendbuf,&recvbuf,count,datatype,op,root,comm)
- MPI\_Allreduce (&sendbuf,&recvbuf,count,datatype,op,comm)
- MPI\_Reduce\_scatter (&sendbuf,&recvbuf,recvcount,datatype, op,comm)
- MPI\_Alltoall (&sendbuf,sendcount,sendtype,&recvbuf, recvcnt,recvtype,comm)
- MPI\_Scan (&sendbuf,&recvbuf,count,datatype,op,comm)

<https://www.mpich.org/static/docs/v3.2/www3/index.htm>

# Aggrégations avec MPI-Reduce

MPI Reduction Operation	
<b>MPI_MAX</b>	maximum
<b>MPI_MIN</b>	minimum
<b>MPI_SUM</b>	sum
<b>MPI_PROD</b>	product
<b>MPI_LAND</b>	logical AND
<b>MPI_BAND</b>	bit-wise AND
<b>MPI_LOR</b>	logical OR
<b>MPI_BOR</b>	bit-wise OR
<b>MPI_LXOR</b>	logical XOR
<b>MPI_BXOR</b>	bit-wise XOR
<b>MPI_MAXLOC</b>	max value and location
<b>MPI_MINLOC</b>	min value and location



# **Types natifs et définition de types dérivés sur MPI**

# Types natifs MPI

C Data Types	
<b>MPI_CHAR</b>	<b>MPI_C_COMPLEX</b>
<b>MPI_WCHAR</b>	<b>MPI_C_FLOAT_COMPLEX</b>
<b>MPI_SHORT</b>	<b>MPI_C_DOUBLE_COMPLEX</b>
<b>MPI_INT</b>	<b>MPI_C_LONG_DOUBLE_COMPLEX</b>
<b>MPI_LONG</b>	<b>MPI_C_BOOL</b>
<b>MPI_LONG_LONG_INT</b>	<b>MPI_LOGICAL</b>
<b>MPI_LONG_LONG</b>	<b>MPI_C_LONG_DOUBLE_COMPLEX</b>
<b>MPI_SIGNED_CHAR</b>	<b>MPI_INT8_T</b>
<b>MPI_UNSIGNED_CHAR</b>	<b>MPI_INT16_T</b>
<b>MPI_UNSIGNED_SHORT</b>	<b>MPI_INT32_T</b>
<b>MPI_UNSIGNED_LONG</b>	<b>MPI_INT64_T</b>
<b>MPI_UNSIGNED</b>	<b>MPI_UINT8_T</b>
<b>MPI_FLOAT</b>	<b>MPI_UINT16_T</b>
<b>MPI_DOUBLE</b>	<b>MPI_UINT32_T</b>
<b>MPI_LONG_DOUBLE</b>	<b>MPI_UINT64_T</b>
	<b>MPI_BYTE</b>
	<b>MPI_PACKED</b>

# Types dérivés

- MPI fournit plusieurs méthodes pour construire des types de données dérivés:
  - Contigu
  - Vecteur
  - Indexé
  - Struct

# Fonctions de créations de types

- **MPI\_Type\_contiguous** (**count**, **oldtype**, &**newtype**)
  - Constructeur produisant un nouveau type de données en faisant **count** copies de **oldtype**
- **MPI\_Type\_vector** (**count**, **blocklength**, **stride**, **oldtype**, &**newtype**)
  - Semblable à contiguë, mais avec des écarts réguliers (**stride**) dans les déplacements.
- **MPI\_Type\_struct**(**count**, **blocklens**[], **offsets**[], **old\_types**, &**newtype**)
  - Le nouveau type de données est formé à partir de types de données composant

# Fonctions de créations de types

- **MPI\_Type\_extent (datatype,&extent)**
  - Renvoie la taille en octets du type de données spécifié.
- **MPI\_Type\_commit (&datatype)**
  - Valide un nouveau type de données au système.
- **MPI\_Type\_free (&datatype)**
  - Libère l'objet spécifié, important pour éviter l'épuisement de la mémoire si de nombreux objets sont créés

