# Communicators

Introduction

The purpose of communicators is to create subgroups on which we can carry out operations such as collective or point-to-point communications. Each subgroup will have its own communication space.
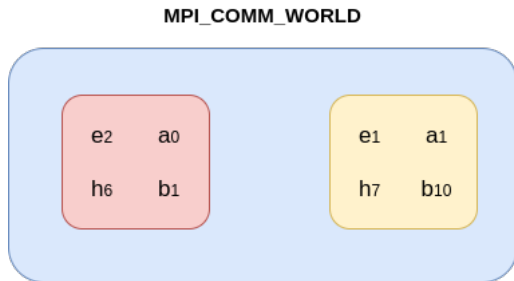
**MPI_COMM_WORLD**



Figure: Communicator partitioning

# Communicators

Example

For example, we want to broadcast a collective message to even-ranked processes and another message to odd-ranked processes.

- Looping on send/recv can be very detrimental especially if the number of processes is high. Also a test inside the loop would be compulsory in order to know if the sending process must send the message to an even or odd process rank.

- A solution is to create a communicator containing the even-ranked processes, another containing the odd-ranked processes, and initiate the collective communications inside these groups.

# Communicators

Default communicator

- A communicator can only be created from another communicator. The first one will be created from the MPI_COMM_WORLD.
- After the MPI_INIT() call, a communicator is created for the duration of the program execution.
- Its identifier MPI_COMM_WORLD is an integer value defined in the header files.
- This communicator can only be destroyed via a call to MPI_FINALIZE().
- By default, therefore, it sets the scope of collective and point-to-point communications to include all the processes of the application.

# Communicators

Groups and communicators

- A communicator consists of :
  - □ A group, which is an ordered group of processes.
  - □ A communication context put in place by calling one of the communicator construction subroutines, which allows determination of the communication space.
- The communication contexts are managed by MPI (the programmer has no action on them : It is a hidden attribute).
- In the MPI library, the following subroutines exist for the purpose of building communicators : MPI_COMM_CREATE(), MPI_COMM_DUP(), MPI_COMM_SPLIT()
- The communicator constructors are collective calls.
- Communicators created by the programmer can be destroyed by using the MPI_COMM_FREE() subroutine.

# Communicators

### Partitioning of a communicator

In order to solve the problem example :

- Partition the communicator into odd-ranked and even-ranked processes.
- Broadcast a message inside the odd-ranked processes and another message inside the even-ranked processes.
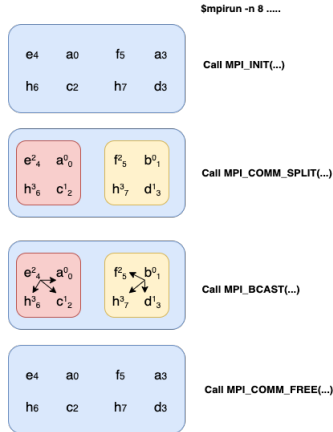


Figure: Communicator creation/destruction

# Communicators

The MPI_COMM_SPLIT() subroutine allows :

- Partitioning a given communicator into as many communicators as we want.
- Giving the same name to all these communicators : The process value will be the value of its communicator.

- Method :
    1. Define a colour value for each process, associated with its communicator number.
    2. Define a key value for ordering the processes in each communicator
    3. Create the partition where each communicator is called new_comm

```
1    MPI_COMM_SPLIT ( comm , color , key , new_comm , code )
2
3    integer , intent ( in ) :: comm , color , key
4
5    integer , intent ( out ) :: new_comm , code
```

```
1    Comm.Split ( self , int color=0 , int key=0 )
```

A process which assigns a color value equal to MPI_UNDEFINED will have the invalid communicator MPI_COMM_NULL for new_com.

# Communicators

Example

Let's look at how to proceed in order to build the communicator which will subdivide the communication space into odd-ranked and even-ranked processes via the MPI_COMM_SPLIT() constructor.

# Communicators

Partitioning of a communicator with MPI_COMM_SPLIT()
   Example :

```
1   ...
2   world_rank = comm.Get_rank(); world_size = comm.Get_size()
3
4   m = 5; a = np.zeros(m)
5
6   if world_rank==2: a[:] = 2.
7   if world_rank==5: a[:] = 5.
8
9   key = world_rank
10  if world_rank==2 or world_rank==5 :
11      key=-1
12  color = world_rank%2
13
14  newcomm = comm.Split(color, key)
15
16  row_rank = newcomm.Get_rank(); row_size = newcomm.Get_size()
17
18  commname = "Comm-"+str(color)
19  newcomm.Set_name(commname)
20
21  print("WORLD RANK/SIZE: {RANK}/{SIZE} \t ROW RANK/SIZE: {ROW_RANK}/{ROW_SIZE} \t ↩
        Comm name: {commname}".format(RANK=world_rank, SIZE=world_size, ROW_RANK=↩
        row_rank, ROW_SIZE=row_size, commname=commname))
22
23  newcomm.Bcast(a, root=0);
24
25  newcomm.Free()
```

# Communicators

Partitioning of a communicator with MPI_COMM_SPLIT()

Results :

```
mpirun -n 6 python3 split_communicator.py

WORLD RANK/SIZE: 5/6    ROW RANK/SIZE: 0/3    Comm name: Comm-1
WORLD RANK/SIZE: 0/6    ROW RANK/SIZE: 1/3    Comm name: Comm-0
WORLD RANK/SIZE: 1/6    ROW RANK/SIZE: 1/3    Comm name: Comm-1
WORLD RANK/SIZE: 2/6    ROW RANK/SIZE: 0/3    Comm name: Comm-0
WORLD RANK/SIZE: 3/6    ROW RANK/SIZE: 2/3    Comm name: Comm-1
WORLD RANK/SIZE: 4/6    ROW RANK/SIZE: 2/3    Comm name: Comm-0
```

# Communicators

Communicator built from a group

- We can also build a communicator by defining a group of processes : Call to MPI_COMM_GROUP(), MPI_GROUP_INCL(), MPI_COMM_CREATE(), MPI_GROUP_FREE()
- This process is however far more cumbersome than using MPI_COMM_SPLIT() whenever possible.

# Communicators

Topologies

- In most applications, especially in domain decomposition methods where we match the calculation domain to the process grid, it is helpful to be able to arrange the processes according to a regular topology.

- MPI allows defining virtual cartesian or graph topologies.

  - Cartesian topologies :
    - I Each process is defined in a grid.
    - I Each process has a neighbour in the grid.
    - I The grid can be periodic or not.
    - I The processes are identified by their coordinates in the grid.

  - Graph topologies :
    - I Can be used in more complex topologies.

# Communicators

Cartesian topologies

- A Cartesian topology is defined from a given communicator named comm_old, calling the MPI_CART_CREATE() subroutine.

- We define :
  - □ An integer ndims representing the number of grid dimensions.
  - □ An integer array dims of dimension ndims showing the number of processes in each dimension.
  - □ An array of ndims logicals which shows the periodicity of each dimension.
  - □ A logical reorder which shows if the process numbering can be changed by MPI.

```
1 MPI_CART_CREATE ( comm_old , ndims , dims , periods , reorder , comm_new , code )
2
3 integer , intent ( in ) :: comm_old , ndims
4 integer , dimension ( ndims ) , intent ( in ) :: dims
5 logical , dimension ( ndims ) , intent ( in ) :: periods
6 logical , intent ( in ) :: reorganization
7
8 integer , intent ( out ) :: comm_new , code
```

```
1 Intracomm . Create_cart ( self , dims , periods = None , bool reorder = False )
```

# Communicators

## 2D Example

Example on a grid having 2 domains along x and 2 along y, periodic in y.

```python
from mpi4py import MPI

comm = MPI.COMM_WORLD
nb_procs = comm.Get_size()
rank = comm.Get_rank()

periods = tuple([False, False])
reorder = False
dims = [2,2]


cart2d = comm.Create_cart(
dims     = dims,
periods = periods,
reorder = reorder
)
```

- If reorder = .false. then the rank of the processes in the new communicator (comm_2D) is the same as in the old communicator (MPI_COMM_WORLD).
- If reorder = .true., the MPI implementation chooses the order of the processes.
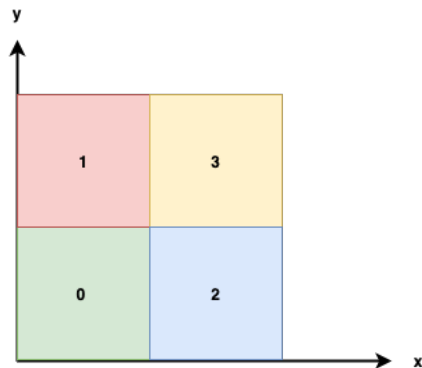
# Communicators



Figure: A 2D non-periodic Cartesian topology

# Communicators

## 3D Example

Example on a 3D grid having 2 domains along x, 2 along y and 2 along z, non periodic.

```python
from mpi4py import MPI

comm = MPI.COMM_WORLD
nb_procs = comm.Get_size()
rank = comm.Get_rank()

periods = tuple([False, False, False])
reorder = False
dims = [2,2,2]


cart3 = comm.Create_cart(
dims    = dims,
periods = periods,
reorder = reorder
)
```
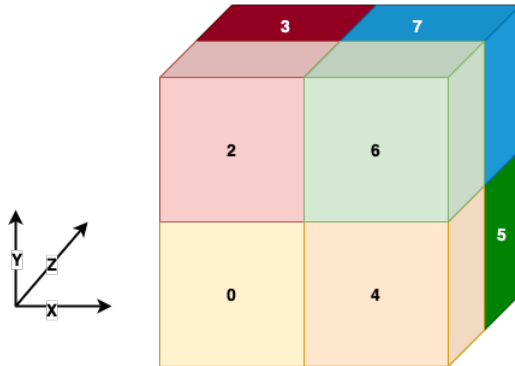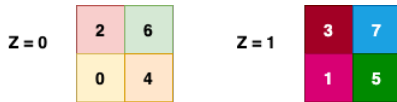
Figure: A 3D non-periodic Cartesian topology

# Communicators
Process distribution

The MPI_DIMS_CREATE() subroutine returns the number of processes in each dimension of the grid according to the total number of processes.

```
1 MPI_DIMS_CREATE(nb_procs, ndims, dims, code)
2
3 integer, intent(in) :: nb_procs, ndims
4 integer, dimension(ndims), intent(inout) :: dims
5
6 integer, intent(out) :: code
```

Remark : If the values of dims in entry are all 0, then we leave to MPI the choice of the number of processes in each direction according to the total number of processes.

# Communicators

Rank od a process

In a Cartesian topology, the MPI_CART_RANK() subroutine returns the rank of the associated process to the coordinates in the grid.

```
1 MPI_CART_RANK(comm, coords, rank, code)
2
3 integer, intent(in) :: comm
4 integer, dimension(ndims), intent(in) :: coords
5
6 integer, intent(out) :: rank, code
```

```
1    Cartcomm.Get_cart_rank(self, coords)
```

# Communicators
Coordinates of a process

In a cartesian topology, the MPI_CART_COORDS() subroutine returns the coordinates of a process of a given rank in the grid.

```
1 MPI_CART_COORDS(comm, rank, ndims, coords, code)
2
3 integer, intent(in) :: comm, rank, ndims
4 integer, dimension(ndims),intent(out) :: coords
5
6 integer, intent(out) :: code
```

```
1   Cartcomm.Get_coords(self, int rank)
```

# Communicators

Figure: A 2D non-periodic Cartesian topology
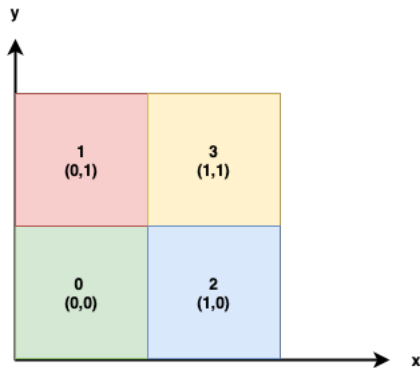
# Communicators

```python
1   from mpi4py import MPI
2
3   comm = MPI.COMM_WORLD
4   rank = comm.Get_rank()
5
6   periods = tuple([False, False])
7   reorder = False
8   dims = [2,2]
9
10  cart2d = comm.Create_cart(
11  dims    = dims,
12  periods = periods,
13  reorder = reorder
14  )
15
16  coord2d = cart2d.Get_coords(rank)
17
18  print("I'm rank", rank, "my 2d coords are", coord2d)
```

```
mpirun -n 4 python3 coordinate_2d_cart.py

I'm rank 0 my 2d coords are [0, 0]
I'm rank 1 my 2d coords are [0, 1]
I'm rank 2 my 2d coords are [1, 0]
I'm rank 3 my 2d coords are [1, 1]
```

# Communicators

Rank of neighbours

In a Cartesian topology, a process that calls the MPI_CART_SHIFT() subroutine can obtain the rank of a neighboring process in a given direction.

```
1 MPI_CART_SHIFT( comm, direction, step, rank_previous, rank_next, code)
2
3 integer, intent(in) :: comm, direction, step
4 integer, intent(out) :: rank_previous, rank_next
5
6 integer, intent(out) :: code
```

```
1 Cartcomm.Shift(self, int direction, int disp)
```

- The direction parameter corresponds to the displacement axis (xyz).
- The step parameter corresponds to the displacement step.
- If a rank does not have a neighbor before (or after) in the requested direction, then the value of the previous (or following) rank will be MPI_PROC_NULL.
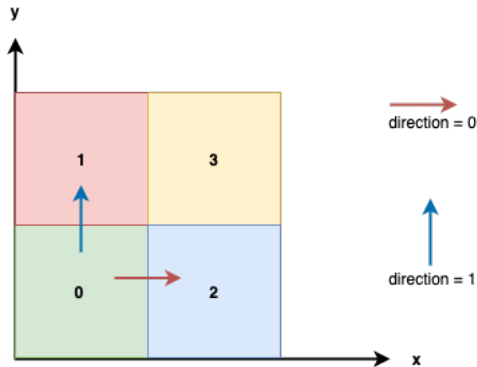
# Communicators

Example : MPI_CART_SHIFT()



Figure: Call of the MPI_CART_SHIFT() subroutine

# Communicators

Example : MPI_CART_SHIFT()

```python
from mpi4py import MPI

comm = MPI.COMM_WORLD
rank = comm.Get_rank()

periods = tuple([False,False])
reorder = False
dims = [2,2]

cart2d = comm.Create_cart(
dims     = dims,
periods = periods,
reorder = reorder
)

left,right = cart2d.Shift(direction = 0, disp=1)
low,high   = cart2d.Shift(direction = 1, disp=1)

print("I'm rank", rank, "my (left,right) neighbours are",( left, right),
"my (low,high) neighbours are",( low, high))
```

# Communicators

Example : MPI_CART_SHIFT()

```
mpirun -n 4 python3 neighbours_2d_cart.py

I'm rank 0
my (left,right) neighbours are (-2, 2) my (low,high) neighbours are (-2, 1)

I'm rank 1
my (left,right) neighbours are (-2, 3) my (low,high) neighbours are (0, -2)

I'm rank 2
my (left,right) neighbours are (0, -2) my (low,high) neighbours are (-2, 3)

I'm rank 3
my (left,right) neighbours are (1, -2) my (low,high) neighbours are (2, -2)
```
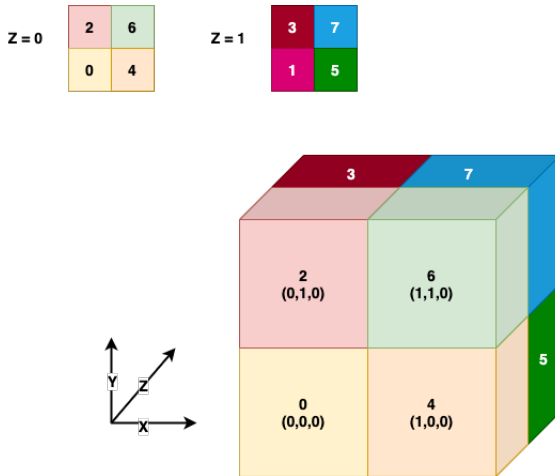
# Communicators

3D Example : coordinates and neighbours



Figure: 3D Coordinates and Neighbours

# Communicators

3D Example : coordinates and neighbours

```python
1   from mpi4py import MPI
2
3   comm = MPI.COMM_WORLD
4   rank = comm.Get_rank()
5
6   periods = tuple([False,False,False])
7   reorder = False
8   dims = [2,2,2]
9
10  cart3d = comm.Create_cart(
11  dims    = dims,
12  periods = periods,
13  reorder = reorder
14  )
15
16  coord3d        = cart3d.Get_coords(rank)
17  left,right     = cart3d.Shift(direction = 0, disp=1)
18  low,high       = cart3d.Shift(direction = 1, disp=1)
19  ahead,before   = cart3d.Shift(direction = 2, disp=1)
20
21  print("I'm rank", rank, "my 3d coords are", coord3d, "my (left,right) neighbours ↩
           are",( left, right), "my (low,high) neighbours are", (low, high), "my (ahead↩
           ,before) neighbours are", (ahead,before))
```

# Communicators

3D Example : coordinates and neighbours

```
mpirun -n 8 --oversubscribe python3 create_3d_cart.py

I'm rank 0 my 3d coords are [0, 0, 0]
my (left,right) neighbours are (-2, 4)
my (low,high) neighbours are (-2, 2)
my (ahead,before) neighbours are (-2, 1)

I'm rank 1 my 3d coords are [0, 0, 1]
my (left,right) neighbours are (-2, 5)
my (low,high) neighbours are (-2, 3)
my (ahead,before) neighbours are (0, -2)

I'm rank 2 my 3d coords are [0, 1, 0]
my (left,right) neighbours are (-2, 6)
my (low,high) neighbours are (0, -2)
my (ahead,before) neighbours are (-2, 3)

...
```