

Back (briefly) to MPI: more on communicators / topology

- Quick overview of communicators in MPI
- Groups
- Virtual topologies
- Cartesian Mappings
- Example: back to the 2-D matrix product

15-1

MPI communication groups

- MPI supports grouping of processes
- Can be used to organize application around tasks
- Collective communications operations can be performed within a group.
- Can create virtual communication topologies: Each group has its own communicator
- Groups and communicators are MPI objects

15-2

– MpiComm

15-2

Main MPI Group functions

- *MPI_Group_size*: returns number of processes in group
- *MPI_Group_rank*: returns rank of calling process in group
- *MPI_Group_compare*: compares group members and group order
- *MPI_Group_translate_ranks*: translates ranks of processes in one group to those in another group
- *MPI_Comm_group*: returns the group associated with a communicator
- *MPI_Group_union*: creates a group by combining two groups
- *MPI_Group_intersection*: creates a group from the intersection of two groups

15-3

– MpiComm

15-3

- *MPI_Group_difference*: creates a group from the difference between two groups
- *MPI_Group_incl*: creates a group from listed members of an existing group
- *MPI_Group_excl*: creates a group excluding listed members of an existing group
- *MPI_Group_range_incl*: creates a group according to first rank, stride, last rank
- *MPI_Group_range_excl*: creates a group by deleting according to first rank, stride, last rank
- *MPI_Group_free*: marks a group for deallocation

15-4

– MpiComm

15-4

Creating a group: Example from P. Patcheco's book

► Declarations etc.

```
main(int argc, char* argv[]) {
    int p;
    int q; /* = sqrt(p) */
    int my_rank;
    MPI_Group group_world;
    MPI_Group row_group;
    MPI_Comm row_comm;
    int* process_ranks;
    int proc;
    int test = 0;
    int sum;
    int my_rank_in_row;

    MPI_Init(&argc, &argv);
    MPI_Comm_size(MPI_COMM_WORLD, &p);
    MPI_Comm_rank(MPI_COMM_WORLD, &my_rank);
```

15-5

– MpiComm

15-5

► Setting up communicator

```
...
q = (int) sqrt((double) p);
/*--- List of processes in new communicator */
process_ranks = (int*) malloc(q*sizeof(int));
for (proc = 0; proc < q; proc++)
    process_ranks[proc] = proc;
/*--- Get group underlying MPI_COMM_WORLD */
MPI_Comm_group(MPI_COMM_WORLD, &group_world);
/*--- Create the new group */
MPI_Group_incl(group_world, q, process_ranks,
    &row_group);
/*--- Create the new communicator */
MPI_Comm_create(MPI_COMM_WORLD, row_group,
    &row_comm);
```

15-6

– MpiComm

15-6

► Testing

```
/*--- Test collective ops in row_comm */
if (my_rank < q) {
    MPI_Comm_rank(row_comm, &my_rank_in_row);
    if (my_rank_in_row == 0) test = 1;
    MPI_Bcast(&test, 1, MPI_INT, 0, row_comm);
}
MPI_Reduce(&test, &sum, 1, MPI_INT, MPI_SUM, 0,
    MPI_COMM_WORLD);
if (my_rank == 0) {
    printf("q = %d, sum = %d\n", q, sum);
}
MPI_Finalize();
}
```

15-7

– MpiComm

15-7

Explore: MPI_Comm_split

```
#include <mpi.h>
int MPI_Comm_split(MPI_Comm comm, int color,
    int key, MPI_Comm *newcomm)
```

- This function allows to partition a group associated with comm into disjoint subgroups, one for each value of color.
- Each subgroup contains processes of the same color.
- A new communicator is created for each subgroup and returned in newcomm.
- We will see something similar when considering cartesian grids



Explore this [MPI site is enough]



See how you can use it to redo previous example (2 colors)

15-8

– MpiComm

15-8

MPI Virtual topologies

- Relevant to Lab3.
- MPI allows to create 'virtual' topologies [e.g. a 3-D grid of processors]
- These topologies are virtual: no relation between architecture of the parallel machine and the process topology.
- A virtual topology is built from MPI communicators and groups
- Two types of topologies: *Cartesian* (1-D, 2-D, 3-D), *Graph*
- Uses: applications with specific communication patterns [Lab 3]
- Also: MPI may reduce communication costs by optimizing process mapping

15-9

– MpiComm

15-9

MPI Cartesian topologies

- A few of the commands used for working with MPI Cartesian topologies:
 - *MPI_Dims_create*: Create N-Dimensional arrangement of PEs in cartesian grid.
 - *MPI_Cart_create*: Create N-Dimensional virtual topology/cartesian grid.
 - *MPI_Cart_coords*: Get local PE coordinates in new cartesian grid
 - *MPI_Cart_sub*: Partitions a communicator into subgroups which form lower-dimensional cartesian subgrids.
 - *MPI_Cart_shift*: Used to find processor neighbors. Returns the shifted source and destination ranks, given a shift direction and amount.

15-10

– MpiComm

15-10

Illustration (examples from Lab3)

```
#include <mpi.h>
int MPI_Cart_create(MPI_Comm comm_old, int ndims,
    const int dims[], const int periods[],
    int reorder, MPI_Comm *comm_cart)
```

"Makes a new communicator to which Cartesian topology information has been attached."

Parameter	i/o	Definition
comm_old:	in	input communicator (handle)
ndims:	in	number of dimensions of cartesian grid
dims:	in	integer array of size ndims specifying the number of processes in each dim.
periods:	in	logical array of size ndims specifying for each dimension
reorder:	in	ranking may be reordered (true) or not (false) (logical)
comm_cart:	out	communicator with new cartesian topology (handle)

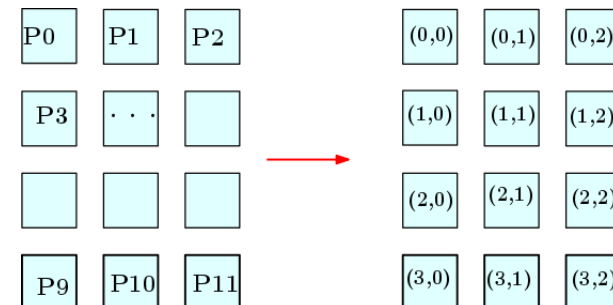
15-11

– MpiComm

15-11

- One note: *reorder* tells whether or not to allow reordering of processes so to achieve a good embedding with respect to physical architecture.
- PEs are ordered from (0:np-1); numbered/ranked in row-major (C lang) order.

Create a (4,3) grid from a 12 process group:



15-12

– MpiComm

15-12

- In Lab3:
- `MPI_Comm comm2D` declares a communicator named `comm2D`.
- We then create a cartesian communicator with a $qi \times qj$ grid (2-D) topology:

```
/*----- CREATE 2-D communicator*/
dims[0] = qi;
dims[1] = qj;
periods[0]=periods[1]=1;
/*-----2D topology created*/
MPI_Cart_create(MPI_COMM_WORLD,2,dims,periods,0,&
comm2D);
/*----- my rank in 2-D form */
MPI_Cart_coords(comm2D,myid,2,coords);
```

15-13

– MpiComm

15-13

```
/*----- rank of West neighbor */
coords[0]=qi;
coords[1]=qj-1;
MPI_Cart_rank(comm2D,coords,&westNB);
```

- Finds the rank of neighbor to the west. `westNB` can now be used to communicate with node $(qi, qj - 1)$ in grid.
- This is very convenient as it simplifies communication.

15-14

– MpiComm

15-14

An important function `MPI_Cart_sub`

```
#include <mpi.h>
int MPI_Cart_sub(MPI_Comm comm,
const int remain_dims[],
MPI_Comm *comm_new)
```

- Partitions a communicator into subgroups which form lower-dimensional cartesian subgrids
- `remain_dims`: This is a logical vector the i th entry of which specifies whether the i th dimension is kept in the subgrid (true) or is dropped (false)
- `comm_new`: (out) communicator containing the subgrid that includes the calling process (handle)
- Recall `MPI_Comm_split` - this is similar.



Show how you can use to implement the 2-D matrix multiply [page 12-17 of Lecture notes set 12]

15-15

– MpiComm

15-15