

Tu. 05/22/12:

Ch. 7 Communicators and Topologies:

MPI\_COMM\_WORLD is the only communicator that exists initially. We can create additional ones as subsets of it.

Ex.: classic example to dedicate one process to do the I/O and exclude it from the calculations. If we want to run a method like power method, then you would call it only with the non-I/O processes.

Main example of Ch. 7 is Fox's Algorithm  $\rightarrow$  Sec. 7.2 and 7.8  
Fox  $\Leftrightarrow$  matrix-matrix product  $C = AB$ ,  $A, B, C \in \mathbb{R}^{n \times n}$ ,  
with only local matrices  $l_A, l_B, l_C$  actually stored  
with a 2-D subdivision of  $A$  to  $l_A$  etc.

Assume  $p$  processes such that  $p = q^2$  for integer  $q$

ex.:  $q = 1, 2, 3, 4, 5, 6, 7, 8$

$\Leftrightarrow p = 1, 4, 9, 16, 25, 36, 49, 64$

only possible ones currently

If  $A \in \mathbb{R}^{n \times n}$ , then  $l_A \in \mathbb{R}^{l_n \times l_n}$  with  $l_n = n/q$

$\Rightarrow$  reasonable  $q$  values only even ones

Also want to use all cores on each node, i.e., want

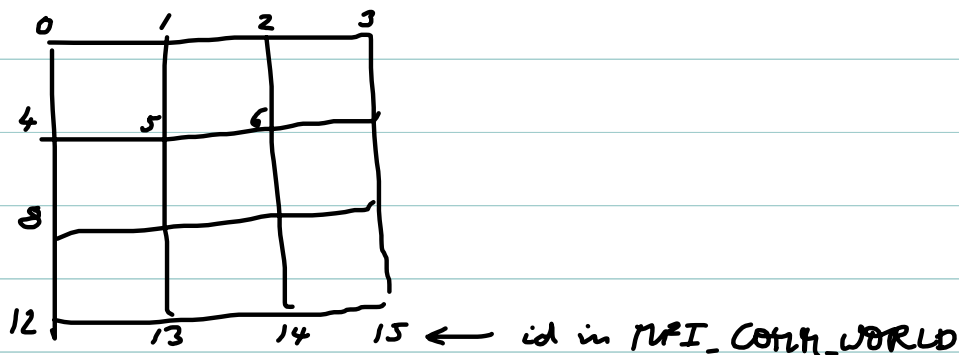
$p = 1, 4, 16, 36, 64$

$\Leftrightarrow 1, 1, 2, 9, 16$  nodes needed

Merge simply: We need to put some thought into designing the numerical experiments by choosing good  $n$  and  $p$ .

Communicators: how to create subdivisions

Idea: Think of  $p$  processes as a 2-D grid



Want to create communicators for each row and column.

Option 1: most general MPI command

Specify MPI ranks (id) to include in the communicator  
MPI\_Group group\_world

MPI\_Comm\_group (MPI\_COMM\_WORLD, &group\_world)

We cannot access a communicator directly, but we have to do it via a Group object.

Ex.: create comm for 2nd row: In example above, 2nd row is  $id = 4, 5, 6, 7$  if  $q = 4 \Leftrightarrow id = k$  for  $q \leq k < 2q$  for  $(k = 0; k < q; k++)$

$process\_rank[k] = q + k$

MPI\_Comm\_incl (group\_world, q, process\_rank, &row2\_comm\_group)

MPI\_Comm\_create (&row\_comm\_group)

→ Need to do this manually for each row and column.

option 2: Use condition to split up communicator

=> can create several communicators at once, all with same variable name

```
MPI_Comm_split(  
    MPI_Comm      old_comm,  
    int           split_key,  
    int           rank_key,  
    MPI_Comm      *new_comm)
```

All processes in `old_comm` with same `split_key` are put in one communicator, with their ranks ordered according to `rank_key`.

```
my_row = id / 9
```

```
my_col = id % 9
```

```
MPI_Comm_split ( MPI_COMM_WORLD ,  
                my_row , my_col , & row_comm )
```

```
MPI_Comm_split ( MPI_COMM_WORLD ,  
                my_col , my_row , & col_comm )
```

Notice all communicators created in one shot.

### Option 3: Topologies of communicator explicitly part of MPI

Command: Cartesian grids

MPI Comm      grid-comm  
int              ndims, dimsize [2], wrap-around [2],  
                  rcoords

ndims = 2

dim-size [0] = dimsize [1] = 9

wrap-around [0] = 0 ( $\Rightarrow$ ) do not wrap around

wrap-around [1] = 1 ( $\Rightarrow$ ) wrap around in 2nd dimension

rcoords = 1 ( $\Rightarrow$ ) allow for ordering of the MPI process at run

MPI\_Cart\_create ( MPI\_COMM\_WORLD, ndims, Time  
                  dim\_size, wrap-around, rcoords, & grid-comm )

MPI\_Cart\_coords (  
    MPI Comm              grid-comm ,  
    int                    my-grid-rank ,  
input  $\rightarrow$  int              ndims ,  
vector  $\rightarrow$  int              \* coords )  
and output

Ex.:  $g=4$ ,  $np=16$ ,  $ndims=2$ ,  $my\_grid\_rank=5 \Rightarrow coords=\{1,1\}$

MPI\_Cart\_rank (  
    MPI Comm              grid-comm ,  
vector input  $\rightarrow$  int        \* coords  
Scalar output  $\rightarrow$  int       \* grid-rank )

MPI\_Cart\_sub (

MPI\_Comm

int

MPI\_Comm

grid\_comm

\* free\_coords

\* new\_comm)

if free\_coords[i] = 0, then this coordinate is not free  
⇔ if remain\_coords[i] = 1, then this coordinate remains  
in the new communicator

⇒  
row\_comm by  $\begin{cases} \text{remain\_coords}[0] = 0 \\ [1] = 1 \end{cases}$   
col\_comm by  $\begin{cases} [0] = 1 \\ [1] = 0 \end{cases}$

→ pp. 125-126 setxy-grid

---

## Sec. 7.2 Fox's Algorithm

$$A = \begin{bmatrix} A_{00} & A_{01} & A_{02} & \dots \\ A_{10} & A_{11} & A_{12} & \dots \\ \vdots & & & \\ A_{g-10} & \dots & & \end{bmatrix}$$

$$AB = \begin{bmatrix} A_{00} & A_{01} \\ A_{10} & A_{11} \end{bmatrix} \begin{bmatrix} B_{00} & B_{01} \\ B_{10} & B_{11} \end{bmatrix}$$

$$= \begin{bmatrix} A_{00}B_{00} + A_{01}B_{10} & A_{00}B_{01} + A_{01}B_{11} \\ A_{10}B_{00} + A_{11}B_{10} & \dots \end{bmatrix}$$

$$= C = \begin{bmatrix} C_{00} & C_{01} \\ C_{10} & C_{11} \end{bmatrix}$$

$\Rightarrow$  Message: block matrix algebra works!  
(Pacheco shows this with scalar blocks.)

$\Leftrightarrow$  For each block in  $C$ :

$$C_{ij} = A_{i0} B_{0j} + A_{i1} B_{1j} + \dots + A_{ii} B_{ij} + \dots + A_{i, g-1} B_{(g-1)j}$$

$$= A_{ii} B_{ij} + A_{i, i+1} B_{(i+1)j} + \dots + A_{i, g-1} B_{(g-1)j} + A_{i0} B_{0j} + \dots$$

This shows the idea of Fox:

$(i, j)$  is the 2-D process ID, have  $B_{ij}$  on the local process, so start accumulating the sum by the  $A_{ii} B_{ij}$  term. Then communication will be needed to get the right pieces of  $A$  and  $B$  for each step.

Want as result  $l\_C = C_{ij}$  on Process  $(i, j)$ .

by accumulating local products:  $l\_C += l\_A * l\_B$

Initially  $l\_B = B_{ij}$ , so no communication needed,

but we need to communicate along row to  $l\_A = A_{ii}$

Code:

$l\_A\_original = l\_A$

$l\_C = 0$  initialize

for ( $i\_step = 0$ ;  $i\_step < q$ ;  $i\_step++$ ) {

$MPI\_Bcast(l\_A, l\_n \times l\_n, MPI\_DOUBLE, (i+i\_step) \% q,$   
 $row\_comm)$

$l\_C += l\_A * l\_B$

$MPI\_Sendrecv\_replace(l\_B, l\_n \times l\_n, MPI\_DOUBLE,$   
 $dest = (i-1+q) \% q, source = (i+1) \% q, col\_comm)$

$l\_A = l\_A\_original$

}

This uses some ideas from the ringsend in Ch. 3

Pacheco does not have  $l\_A\_original$ , but he

uses  $if (i == (i+ident) \% q) MPI\_Bcast(l\_A, \dots)$

else  $MPI\_Bcast(l\_D, \dots)$

then  $l\_C += l\_D * l\_B$

$MPI\_Sendrecv\_replace$  needs an auxiliary buffer  $\Rightarrow$

This code needs memory for 5 matrices of size  $(l\_n)^2$

( $l\_A, l\_B, l\_C, l\_A\_original$ , plus aux. buffer)

BLAS = basic linear algebra subprograms

level 1 ( $\Rightarrow$ ) 1 for loop like in a dot product

2 ( $\Rightarrow$ ) 2 " " " matrix-vector product

3 ( $\Rightarrow$ ) 3 " " " matrix-matrix product

BLAS3 dgemm = matrix - matrix product

function uses a block algorithm like Fox

and is much more efficient than any ordering of for loops can be.

When doing code development, think carefully about each sub-problem, e.g., here the  $L-C += L-A * L-B \rightarrow$  realize BLAS3 is better than any naive code, so use it!