<u>Monday, 02/06/12</u>

Ch. 10 Design and Coding of Parallel Programs:

2 examples: – Jacobi method for linear systems → next semester
            – Sorting

Example of a sorting problem: State concretely with 4 processes here!

Have 100 integers from 1, 2, ..., 100 given across 4 processes.

Want to obtain these numbers sorted and equally distributed across the 4 processes, that is 1,...,25 on Process 0, 26 to 50 on P1, etc.

This is modified from Pacheco to make clearer and to make more realistic. Think of these as indices into some list of large object (like of FEM triangles that are to distributed into 4 subdomains after refinement/coarsening for purpose of load balancing.

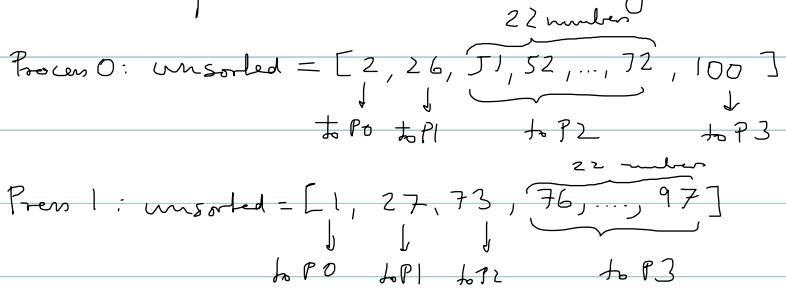Critical idea: Never store all integers on one process.

Very basic algorithm:
 – Perform local sort of the 25 local integers
 – On each process, send the integers in the range 1 to 25 to Process 0, in range 26 to 50 to Proc 1, etc.
   On each process, receive these numbers, but we do not know which process the numbers will come from.
   Notice: We do know that there are only 25 numbers to be received.
 – Notice that these numbers are not sorted, so
   need to do local sort

Example: Variable names are "unsorted" for the original 25 integers on each process and "sorted" for the 25 at the end.

Process 0: unsorted = [ 2 , 26 , $\overbrace{51 , 52 , ... , 72}^{22\text{ numbers}}$ , 100 ]
  ↓        ↓              ↓                    ↓
 to P0   to P1          to P2               to P3

Process 1: unsorted = [ 1 , 27 , 73 , $\overbrace{76 , ... , 97}^{22\text{ numbers}}$ ]
  ↓      ↓    ↓              ↓
 to P0  to P1 to P2        to P3

$\vdots$

Need communication among the 4 processes such that we receive:

on Process 0:   sorted = [ 2 , 1 , ... ]

on Process 1:   sorted = [ 26 , 27 , ... ]

on Process 2:   sorted = [ 51 , .., 72 , 73 , ... ]

on Process 3:   sorted = [ 100 , 76 , ... , 97 , ... ]

MPI_Alltoall  ( _____          for this example:
   void          * send_buffer,                  unsorted
   int             send_count,
   MPI_Datatype    send_type, _____          MPI_INT
                                                sorted
   void          * recv_buffer,
   int             recv_count,
   MPI_Datatype    recv_type, ___    ___       MPI_INT
MPI_Comm          comm)

This command is only good if you send the $\underline{\text{same}}$ number of data to each process and receive that same number from all processes !

$\Rightarrow$ There are variable version of $\underline{\underline{all}}$ MPI Commands exactly for this purpose, here MPI_Alltoall$\underline{\underline{v}}$ , but others are MPI_Scatterv$\underline{\underline{}}$ , MPI_Gatherv$\underline{\underline{v}}$
$\Rightarrow$ See Appendix A in Pacheco  —→ HW on power method

```
MPI_Alltoallv (
    void            * send_buffer,
    int             * send_counts,
    int             * send_displacements,
    MPI_Datatype      send_type,

    void            * recv_buffer,
    int             * recv_counts,
    int             * recv_displacements,

    MPI_Datatype      recv_type,

    MPI_Comm          comm)
```

array of length np, dynamically allocated

```
int    *unsorted, *send_counts, *send_displacements;
int    *sorted, *recv_counts, *recv_displacements;
Let n be total number of integers.
unsorted = (int) malloc ( (n/np) * sizeof (int ))
sorted =           — // —

send_counts = (int) malloc ( np * sizeof (int));
send_displacements =    — // —
recv_counts        =    — // —
recv_displacements =    — // —
— Perform local sort  => "unsorted" is now a sorted list of n/np int.


for (i=0; i < np; i++)   send_counts [i] = 0

for ( j=0 ; j < n/np ; j++) {
       k  = (unsorted [j]-1)/ (n/np)  /* = process id to send to */
       (send_counts [k])++
    }
```

```
send_displacement [0] = 0
for ( i = 1 ; i < np ; i++)
    send_displacements [ i ] = send_displacements [i-1] + send_counts [i-1]
```

Now, we need the recv_counts. But there is not enough info locally. Rather, the info is in all the send_counts on the other processes =>

```
MPI_Alltoall (
    send_counts ,    1,   MPI_INT ,
    recv_counts,     1,  MPI_INT,  MPI_COMM_WORLD )
recv_displacements [0] = 0
for ( i = 1 ; i < np ; i++)
    recv_displacements [i] = recv_displacements [i-1] + recv_counts [i-1]

MPI_Alltoallv (
    unsorted , send_counts, send_displacements, MPI_INT ,
    sorted  ,  recv_counts, recv_displacements, MPI_INT ,
    MPI_COMM_WORLD )
```
— Perform local sort of "sorted"

( In example of load-balancing of FEM: Now you communicate the actual, large data )

————————————————

Contrast to Pacheco example:

(1) He allows for missing and repeated numbers => Sorted might <u>not</u> have length $n/np (= 25)$ => need to do dynamic memory re-allocating based on $\sum_{i=0}^{np-1}$ recv_counts [i] = recv_displacements [np-1] + recv_counts [np-1]

(2) Memory as large as $n$ might now be needed on some process => contradicts spirit of load-balancing