

Monday, 01/23/12:

## Ch. 6 Grouping Data for Communications:

Motivating example:

trapezoidal rule:  $a, b \in \mathbb{R}$ ,  $n \in \mathbb{Z}$ , want to broadcast them, but in one MPI call, as opposed to 3:

```
MPI_Bcast(&a, 1, MPI_DOUBLE, ...)
```

```
MPI_Bcast(&b, 1, MPI_DOUBLE, ...)
```

```
MPI_Bcast(&n, 1, MPI_INT, ...)
```

One idea to combine data might be to define a vector, but this would only work for homogeneous datatypes, like  $a$  and  $b$  together.

How to use this idea to compute min and max in one call.

```
Have MPI_Reduce(&l_x, &x_min, 1, MPI_DOUBLE, MPI_MIN, ... )  
MPI_Reduce(&l_x, &x_max, ... , MPI_MAX, ... )
```

replace this by:

```
double v[2], vminmax[2]
```

```
v[0] = l_x, v[1] = -l_x
```

```
MPI_Reduce(v, vminmax, 2, MPI_DOUBLE, MPI_MIN, ... )
```

```
x_min = vminmax[0]
```

```
x_max = -vminmax[1]
```



negative sign

We need to handle combinations of any datatypes  $\Rightarrow$  Done by

### Derived datatypes

Idea of chptr 6: from general to specific!

MPI has intrinsic datatypes like MPI\_INT, MPI\_DOUBLE, etc.

Based on these, we can construct new types that are then created at runtime, potentially optimized. ( $\rightarrow$  Ch. 7).

What is the essence of a datatype?

→ type signature = Sequence of pairs of (type, displacement)

$\{ (t_0, d_0), (t_1, d_1), (t_2, d_2), \dots, (t_{n-1}, d_{n-1}) \}$

$t_i$  = type of the  $i$ th data (can be derived itself  $\Rightarrow$  doubly derived datatypes possible)

$d_i$  = displacement in memory location relative to starting address of  $i$ -th data

Most general MPI command:

```
MPI_Type_struct (
    int      count
    int      *bl      (= block lengths)
    MPI_Aint *dl      (= displacements)
    MPI_Datatype *t1
    MPI_Datatype *new_mpi_t ) ;
```

return of length count

Ex.: Pacheco pp. 92-94

C convention for new datatype

count = 3

$bl[0] = bl[1] = bl[2] = 1$  ( $\Leftrightarrow$  a, b, n scalars)

$t1[0] = t1[1] = MPI\_DOUBLE$  (for a and b)

$t1[2] = MPI\_INT$  (for n)

$d1[0] = 0$  (displacement of 0th data from itself)

~~/\* displacement of &b from &a : \*/~~

$MPI\_Address (&a, &start\_address)$

$MPI\_Address (&b, &address)$

$MPI\_Aint$

$d1[1] = address - start\_address$  (= displacement in units Aint)

~~/\* displacement of &n from &a : \*/~~

$MPI\_Address (&n, &address)$

$d1[2] = address - start\_address$

$MPI\_Type\_struct (count, bl, t1, d1, &new\_mpi\_t)$

$MPI\_Type\_commit (&new\_mpi\_t)$

$\Rightarrow$  necessary to actually create the derived type;

only needed for final one in case of doubly derived types

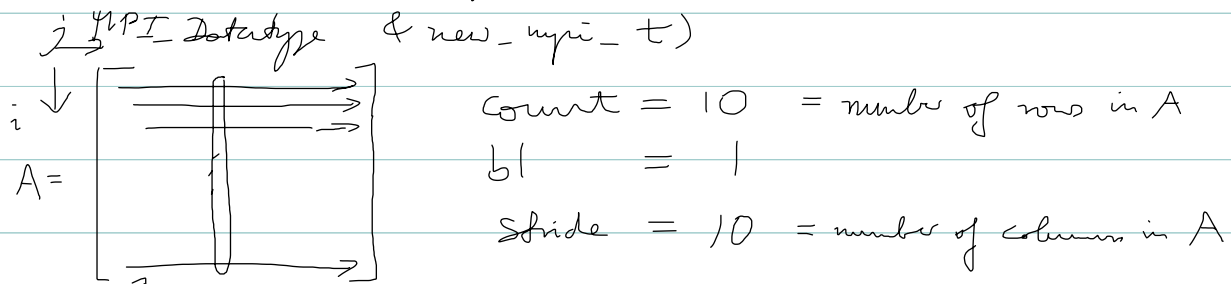
There are several earlier derived types for arrays of homogeneous data

Following Pacheco, consider example of `double A[10][10]` assumed to be consecutive in memory.

Have A on process 0 and want to send one column to another process

```
MPI_Type_vector (
    int    count
    int    bl
    int    stride
```

`MPI_Datatype` `el_type` → of the original array



Notice: `A[10][10]` is in row-oriented storage in C!

```
MPI_Type_vector(10, 1, 10, MPI_DOUBLE, &col_type)
```

```
MPI_Type_commit (&col_type)
```

Usage in Pacheco:

```
if (id == 0)
    MPI_Send (&A[0][j], 1, col_type, j, 0, MPI_COMM_WORLD)
```

← destination  
address of first element in  $j$ th column of `A`

```
if (id == j)
    MPI_Recv (&A[0][j], 1, col_type, 0, 0, MPI_COMM_WORLD, &status)
```

← source

Much better example: receive not into `A[10][10]` on Proc  $j$ , but receive exactly the data sent (10 `DOUBLE`) into a vector

```
if (id == j)
    double z[10]
    MPI_Recv (z, 10, MPI_DOUBLE, 0, 0, MPI_COMM_WORLD, &status)
```

The only criterion for whether a Send and Recv work together is that the type signatures of sent and received data must match.

Some other commands:

Ex.: derived type for row of A  $A = \left[ \begin{array}{c} \text{---} \rightarrow \\ \text{---} \rightarrow \\ \text{---} \rightarrow \end{array} \right] \leftarrow \text{one row}$   
 $\text{MPI\_Type\_vector}(10, 1, 1, \text{MPI\_DOUBLE}, \&\text{row\_type})$

or:

$\text{MPI\_Type\_vector}(1, 10, 1, \text{MPI\_DOUBLE}, \&\text{row\_type})$

→ For consecutive data, can also use

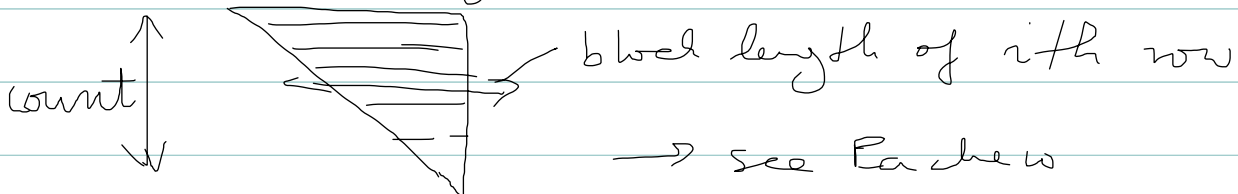
$\text{MPI\_Type\_contiguous}(\text{int count}, \text{MPI\_Datatype old\_type}, \text{MPI\_Datatype} * \text{new\_type})$

Here:  $\text{MPI\_Type\_contiguous}(10, \text{MPI\_DOUBLE}, \&\text{row\_type})$

Generalize int b1 into vector

$\text{MPI\_Type\_indexed}(\text{int count}, \text{int} * \text{b1}, \text{int} * \text{d1}, \text{MPI\_Datatype old\_type}, \text{MPI\_Datatype} * \text{new\_type})$

can be used for triangular matrix:



Finish Ch. 6:

Special intrinsic type is `MPI_PACKED`

created by `MPI_Pack`, disassembled  
by `MPI_Unpack`

→ Packed example for  $a, b, n$

Only situation for me:

if size of a variable is part of the  
message.

Pseudo code:

- Pack size, then pack array of data.
- Send whole package
- Unpack size  $\Rightarrow$  allocate space  
for array of necessary size,  
then unpack data into this array.