

Wednesday 02/08/12:

Ch. 7, Sec. 7.2, 7.8 Fox's Algorithm

Fox:  $C = AB$ ,  $A, B, C \in \mathbb{R}^{n \times n}$

with all matrices in block-cyclic storage = block matrices

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

Each block is  $l_n \times l_n$   
 $l_n = \frac{n}{q}$ ,  $p = q^2$   
parallel processes

Goal: Never assemble any matrix larger than  $l_n \times l_n$  on any process.

Example  $n=4$ ,  $q=2$

$$AB = \begin{bmatrix} a_{00} & a_{01} & a_{02} & a_{03} \\ a_{10} & a_{11} & \dots & \dots \\ a_{20} & a_{21} & \dots & \dots \\ a_{30} & a_{31} & \dots & \dots \end{bmatrix} \begin{bmatrix} b_{00} & b_{01} \\ b_{10} & b_{11} \\ \dots & \dots \end{bmatrix} = \begin{bmatrix} A_{00} & A_{01} \\ A_{10} & A_{11} \end{bmatrix} \begin{bmatrix} B_{00} & B_{01} \\ B_{10} & B_{11} \end{bmatrix}$$
$$= \begin{bmatrix} a_{00}b_{00} + a_{01}b_{10} + a_{02}b_{20} + a_{03}b_{30} & \dots \\ a_{10}b_{00} + a_{11}b_{10} + a_{12}b_{20} + a_{13}b_{30} & \dots \end{bmatrix}$$

$$= \begin{bmatrix} A_{00}B_{00} + A_{01}B_{10} & \dots \\ \dots & \dots \end{bmatrix}$$
$$\begin{bmatrix} a_{00}b_{00} + a_{01}b_{10} & \dots \\ \dots & \dots \end{bmatrix} \quad \begin{bmatrix} a_{02}b_{20} + a_{03}b_{30} & \dots \\ \dots & \dots \end{bmatrix}$$

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

=> Message: block-matrix arithmetic works as expected.

=> Viewpoint is that each parallel process in a Cartesian process grid holds block matrices  $A_{ij}, B_{ij}, C_{ij}$

of size  $l_n \times l_n$ .

Idea of Fox: Compute  $C_{ij}$  in place on the local process by moving/copying needed data of  $A_{ik}, B_{kj}, 0 \leq k < n'$  with  $n' = n/l_n \Leftrightarrow A, B, C$  have  $n' \times n'$  blocks of size  $l_n \times l_n$

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

Idea is to compute each of these terms on Process  $(i, j)$  in the Cartesian process grid. We start with the term  $A_{ii}B_{ij}$ , that is, we start such that  $B_{ij}$  is available on the process and then proceed like

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

Here,  $\exists$  used  $(i, j)$  to denote Process  $id$  in the Cartesian grid.

Moreover, in code,  $l\_A, l\_B, l\_C$  that denote the data stored on Process  $(i, j)$ .

These  $l\_A, l\_B$  are  $A_{ij}, B_{ij}$  initially and also at the end of the algorithm, but in the intermediate stages, we will override with  $A_{ik}, B_{kj}$   $0 \leq k < n'$ .

Assume we use BLAS function  $dgemm$  for the local products  $l\_C += l\_A * l\_B$  with the accumulation in " $+=$ "  $\rightarrow$  notice that  $dgemm$  accumulates, so we will use that.

At start,  $l\_A = A_{ij}$ ,  $l\_B = B_{ij}$ .

We need  $A_{ii}$  on Process  $(i,j)$  to start the first local product  $l\_C += l\_A * l\_B \Leftrightarrow C_{ij} = A_{ii} B_{ij}$ .  
Notice that  $A_{ii}$  is used in this first stage on all processes in the  $i$ -th row of the process grid.

$\Rightarrow$  this is a broadcast within each row communicator. The root process is  $i = \text{row id}$  in the Cartesian process grid.

`MPI_Bcast(l_A, (l_n * l_n), MPI_DOUBLE, i, row_comm)`  
Then perform local product

`dgemv('N', 'N', l_n, l_n, l_n, 1.0, l_A, l_n, l_B, l_n, 1.0, l_C, l_n);`  $\Leftrightarrow$

$$l\_C = \alpha l\_A l\_B + \beta l\_C \quad \text{with } \alpha = \beta = 1.0$$

We realize now that there needs to be some index for the stage of the algorithm, so that we can express  $A_{ik}$  and  $B_{kj}$  somehow.

$$C_{ij} = A_{i+i_0} B_{i_0j} + A_{i+i_1} B_{i_1j} + \dots + A_{i+i_s} B_{i_sj} + \dots$$

produce these from mod operation  
 $\% g$

$$\Rightarrow \text{root of Bcast} = (i+s) \% g$$

Now  $s = \text{stage index}$  with  $\text{for}(s=0; s < q; s++)$   
we look at  $B_{ij}$  communication. In stage  $s$ ,  
Process  $(i, j)$  needs  $B_{i+s, j}$ .  $\Rightarrow$  All these  $B_{kj}$   
live in the same column communicator.

$\Rightarrow$  Realize Process  $(i, j)$  needs  $B_{i+s, j}$  which is  
'below' it in the same column, and the process  
'above' needs the current data in  $l\_B$ .

$\Rightarrow$  So, throughout the  $q$  stages, we need  
to 'shift'  $l\_B$  up within the column  
communicator and wrap around from top  
to bottom (like in the ring send on HW)

Most convenient MPI command for this  
is  $\text{MPI\_Sendrecv\_recv}(\text{dest}, (\text{r\_n} * \text{r\_n}),$   
 $\text{MPI\_DOUBLE}, \text{dest}, \text{source}, \text{col\_comm})$

$$(i-1+q)\%q$$

$$(i+1)\%q$$

$\uparrow$   
Since  $(i, j)$  needs  $B_{i+1, j}, B_{i+2, j}, \text{etc.}$

---

Looking back at  $B_{\text{cast}}$  of  $l\_A$ , we notice  
that Stage 0  $B_{\text{cast}}$  of  $A_{ii}$  destroys the data in all  
other  $l\_A$  that we need in future  $\Rightarrow$  So, we  
need to save original  $l\_A$  and recover it in each stage

$l\_A\_{original} = l\_A$

$l\_C = 0$  (initialized)

for ( $s=0; s < g; s++$ ) {

$MPI\_Bcast(l\_A, 1, 4 * l_n, MPI\_DOUBLE, (i+s) \% g, row\_comm)$

$dgemv(\dots, 1.0, l\_A, \dots, l\_B, \dots, \underline{1.0}, l\_C, \dots) \Leftrightarrow l\_C += l\_A l\_B$

$MPI\_Sendrecv\_replace(l\_B, \dots, (i-1+g) \% g, (i+1) \% g, col\_comm)$

$l\_A = l\_A\_{original}$  /\* recovers original value \*/  
}

Since we do  $g$  moves of  $l\_B$  upwards in  $col\_comm$ , with a wraparound from top to bottom,  $B$  will be stored as before the algorithm.

Note: Pacheco on p. 127 does not have  $l\_A\_{original}$ , but rather uses a temporary variable to program  $Bcast$  in an if statement:

if ( $j = \underbrace{(i+s) \% g}_{\text{root of } Bcast}$ )  $\Leftrightarrow$  if I am the sending process

$MPI\_Bcast(l\_A, \dots)$

else  $\Leftrightarrow$  if I am receiving data

$MPI\_Bcast(l\_A\_{tmp}, \dots)$

$dgemv(\dots, \underline{l\_A\_{tmp}}, \dots)$