

Communicators and Topologies: Matrix Multiplication Example

Ned Nedialkov

Dept. of Computing and Software
McMaster University, Canada
nedialk@mcmaster.ca

January 2012

Fox's algorithm

A and B are $n \times n$ matrices

Compute $C = AB$ in parallel

Let $q = \sqrt{p}$ be an integer such that it divides n , where p is the number of processes

Create a Cartesian topology with processes (i, j) , $i, j = 0, \dots, q - 1$

Denote $m = n/q$

Distribute A and B by blocks on p processes such that A_{ij} and B_{ij} are $m \times m$ blocks stored on process (i, j)

On process (i, j) , we want to compute

$$\begin{aligned} C_{i,j} = \sum_{k=0}^{q-1} A_{i,k} B_{k,j} &= A_{i,0} B_{0,j} + A_{i,1} B_{1,j} + \cdots + A_{i,i-1} B_{i-1,j} \\ &\quad + A_{i,i} B_{i,j} + A_{i,i+1} B_{i+1,j} + \cdots + A_{i,q-1} B_{q-1,j} \end{aligned}$$

Rewrite this as

stage	compute
0	$C_{i,j} = A_{i,i} B_{i,j}$
1	$C_{i,j} = C_{i,j} + A_{i,i+1} B_{i+1,j}$
\vdots	\vdots
	$C_{i,j} = C_{i,j} + A_{i,q-1} B_{q-1,j}$
	$C_{i,j} = C_{i,j} + A_{i,0} B_{0,j}$
	$C_{i,j} = C_{i,j} + A_{i,1} B_{1,j}$
\vdots	\vdots
$q - 1$	$C_{i,j} = C_{i,j} + A_{i,i-1} B_{i-1,j}$

Each process computes in stages

stage 0

- ▶ process (i, j) has $A_{i,j}, B_{i,j}$ but needs $A_{i,i}$
- ▶ process (i, i) broadcasts $A_{i,i}$ across processor row i
- ▶ process (i, j) computes $C_{i,j} = A_{i,i}B_{i,j}$

stage 1

- ▶ process (i, j) has $A_{i,j}, B_{i,j}$, but needs $A_{i,i+1}, B_{i+1,j}$
 - ▶ shift the j th block column of B by one block up
(block 0 goes to block $q - 1$)
 - ▶ process $(i, i + 1)$ broadcasts $A_{i,i+1}$ across processor row i
- ▶ process (i, j) computes $C_{i,j} = C_{i,j} + A_{i,i+1}B_{i+1,j}$

Similarly on next stages

Algorithm outline

On process (i,j) :

$$C_{ij} = 0$$

for $s = 0$ to $q - 1$

$$k = (i + s) \bmod q$$

broadcast $A_{i,k}$ across process row i

$$C_{i,j} = C_{i,j} + A_{i,k}B_{k,j}$$

if $s \neq q - 1$

send $B_{k,j}$ to $((i - 1) \bmod q, j)$

receive $B_{k+1,j}$ from $((i + 1) \bmod q, j)$

Example

Consider multiplying two 3×3 block matrices:

$$\begin{bmatrix} C_{00} & C_{01} & C_{02} \\ C_{10} & C_{11} & C_{12} \\ C_{20} & C_{21} & C_{22} \end{bmatrix} = \begin{bmatrix} A_{00} & A_{01} & A_{02} \\ A_{10} & A_{11} & A_{12} \\ A_{20} & A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} B_{00} & B_{01} & B_{02} \\ B_{10} & B_{11} & B_{12} \\ B_{20} & B_{21} & B_{22} \end{bmatrix}$$

Process (i, j) stores A_{ij} , B_{ij} and computes C_{ij}

Stage 0

process (i,i mod 3)	broadcasts along row i	A_{00}, B_{00}	A_{00}, B_{01}	A_{00}, B_{02}
(0, 0)	A_{00}	A_{11}, B_{10}	A_{11}, B_{11}	A_{11}, B_{12}
(1, 1)	A_{11}	A_{22}, B_{20}	A_{22}, B_{21}	A_{22}, B_{22}
(2, 2)	A_{22}			

Process (i,j) computes

$$\begin{array}{lll} C_{00} = A_{00}B_{00} & C_{01} = A_{00}B_{01} & C_{02} = A_{00}B_{02} \\ C_{10} = A_{11}B_{10} & C_{11} = A_{11}B_{11} & C_{12} = A_{11}B_{12} \\ C_{20} = A_{22}B_{20} & C_{21} = A_{22}B_{21} & C_{12} = A_{22}B_{22} \end{array}$$

Shift-rotate on the columns of B :

$$\begin{array}{lll} A_{00}, B_{10} & A_{00}, B_{11} & A_{00}, B_{12} \\ A_{11}, B_{20} & A_{11}, B_{21} & A_{11}, B_{22} \\ A_{22}, B_{00} & A_{22}, B_{01} & A_{22}, B_{02} \end{array}$$

Stage 1

process $(i, (i+1) \bmod 3)$	broadcasts along row i		A_{01}, B_{10}	A_{01}, B_{11}	A_{01}, B_{12}
$(0, 1)$	A_{01}	\rightarrow	A_{12}, B_{20}	A_{12}, B_{21}	A_{12}, B_{22}
$(1, 2)$	A_{12}		A_{20}, B_{00}	A_{20}, B_{01}	A_{20}, B_{02}
$(2, 0)$	A_{20}				

Process (i, j) computes

$$\begin{array}{lll} C_{00}+ = A_{01}B_{10} & C_{01}+ = A_{01}B_{11} & C_{02}+ = A_{01}B_{12} \\ C_{10}+ = A_{12}B_{20} & C_{11}+ = A_{12}B_{21} & C_{12}+ = A_{12}B_{22} \\ C_{20}+ = A_{20}B_{00} & C_{21}+ = A_{20}B_{01} & C_{22}+ = A_{20}B_{02} \end{array}$$

Shift-rotate on columns of B :

$$\begin{array}{lll} A_{01}, B_{20} & A_{01}, B_{21} & A_{01}, B_{22} \\ A_{10}, B_{00} & A_{10}, B_{01} & A_{10}, B_{02} \\ A_{21}, B_{10} & A_{21}, B_{11} & A_{21}, B_{12} \end{array}$$

Stage 2

process $(i, (i+2) \bmod 3)$	broadcasts along row i			
(0, 2)	A_{02}	A_{02}, B_{20}	A_{02}, B_{21}	A_{02}, B_{22}
(1, 0)	A_{10}	A_{10}, B_{00}	A_{10}, B_{01}	A_{10}, B_{02}
(2, 1)	A_{21}	A_{21}, B_{10}	A_{21}, B_{11}	A_{21}, B_{12}

Process (i, j) computes

$$\begin{array}{lll}
 C_{00}+ = A_{02}B_{20} & C_{01}+ = A_{02}B_{21} & C_{02}+ = A_{02}B_{22} \\
 C_{10}+ = A_{10}B_{00} & C_{11}+ = A_{10}B_{01} & C_{12}+ = A_{10}B_{02} \\
 C_{20}+ = A_{21}B_{10} & C_{21}+ = A_{21}B_{11} & C_{22}+ = A_{21}B_{12}
 \end{array}$$

Implementation

Code adapted from P. Pacheco, Parallel Programming with MPI,
<http://www.cs.usfca.edu/~peter/ppmpi/>

```
1 void Setup_grid( GRID_INFO_T* grid )
2 {
3     int old_rank;
4     int dimensions[2];
5     int wrap_around[2];
6     int coordinates[2];
7     int free_coords[2];
8
9     /* Set up Global Grid Information */
10    MPI_Comm_size(MPI_COMM_WORLD, &(grid->p));
11    MPI_Comm_rank(MPI_COMM_WORLD, &old_rank);
```

```
13  /* We assume p is a perfect square */
14  grid->q = (int) sqrt((double) grid->p);
15  dimensions[0] = dimensions[1] = grid->q;
16
17  /* We want a circular shift in second dimension. */
18  /* Don't care about first */
19  wrap_around[0] = wrap_around[1] = 1;
20  MPI_Cart_create(MPI_COMM_WORLD, 2, dimensions,
21                  wrap_around, 1, &(grid->comm));
22
23  MPI_Comm_rank(grid->comm, &(grid->my_rank));
24  MPI_Cart_coords(grid->comm, grid->my_rank, 2,
25                   coordinates);
26
27  grid->my_row = coordinates[0];
28  grid->my_col = coordinates[1];
```

```
29
30  /* Set up row communicators */
31  free_coords[0] = 0;
32  free_coords[1] = 1;
33  MPI_Cart_sub(grid->comm, free_coords,
34                &(grid->row_comm));
35
36  /* Set up column communicators */
37  free_coords[0] = 1;
38  free_coords[1] = 0;
39  MPI_Cart_sub(grid->comm, free_coords,
40                &(grid->col_comm));
41 }
```

```
1 void Fox( int n, GRID_INFO_T* grid ,
2             LOCAL_MATRIX_T* local_A ,
3             LOCAL_MATRIX_T* local_B ,
4             LOCAL_MATRIX_T* local_C )
5 {
6     LOCAL_MATRIX_T* temp_A ;
7     int stage ;
8     int bcast_root ;
9     int n_bar; /* n/sqrt(p) */
10    int source ;
11    int dest ;
12    MPI_Status status ;
13
14    n_bar = n / grid->q ;
15    Set_to_zero( local_C );
```

```
17
18 /* Calculate addresses for circular shift of B */
19 source = (grid->my_row + 1) % grid->q;
20 dest = (grid->my_row + grid->q - 1) % grid->q;
21
22 /* Set aside storage for the broadcast block of A */
23 temp_A = Local_matrix_allocate(n_bar);
```

```
25  for (stage = 0; stage < grid->q; stage++) {  
26      bcast_root = (grid->my_row + stage) % grid->q;  
27      if (bcast_root == grid->my_col) {  
28          MPI_Bcast(local_A, 1, local_matrix_mpi_t,  
29                      bcast_root, grid->row_comm);  
30          Local_matrix_multiply(local_A, local_B, local_C);  
31      }  
32      else {  
33          MPI_Bcast(temp_A, 1, local_matrix_mpi_t,  
34                      bcast_root, grid->row_comm);  
35          Local_matrix_multiply(temp_A, local_B, local_C);  
36      }  
37      MPI_Sendrecv_replace(local_B, 1, local_matrix_mpi_t,  
38                              dest, 0, source, 0,  
39                              grid->col_comm, &status);  
40  }  
41 }
```