



The Ministry of Education and Science of the Russian Federation

Lobachevsky State University of Nizhni Novgorod

Computing Mathematics and Cybernetics faculty

The competitiveness enhancement program
of the Lobachevsky State University of Nizhni Novgorod
among the world's research and education centers

Strategic initiative

“Achieving leading positions in the field of supercomputer technology
and high-performance computing”

Introduction to MPI

Lectures 8-9. Parallel Matrix Multiplication Methods

Nizhni Novgorod

2014

Lectures_8,9_. Parallel Matrix Multiplication Methods

Matrix multiplication is one of the most important problems of matrix computations. Lectures 8 and 9 review a number of parallel matrix multiplication algorithms. Two of them are based on block striped data decomposition (Lecture 8). The other method (Lecture 9) uses checkerboard data decomposition.

Lecture 8

8.1. Problem Statement

Multiplying an $m \times n$ matrix A with m rows and n columns and an $n \times l$ matrix B with n rows and l columns produces an $m \times l$ matrix C with m rows and l columns. Each element of the matrix C is calculated according to the formula:

$$c_{ij} = \sum_{k=0}^{n-1} a_{ik} \cdot b_{kj}, 0 \leq i < m, 0 \leq j < l. \quad (8.1)$$

As it can be seen in (8.1), each element of the matrix C is the result of the inner product of the corresponding row of the matrix A and column of the matrix B :

$$c_{ij} = (a_i, b_j^T) \cdot a_i = (a_{i0}, a_{i1}, \dots, a_{i(n-1)}), b_j^T = (b_{0j}, b_{1j}, \dots, b_{(n-1)j})^T. \quad (8.2)$$

This algorithm executes $m \cdot n \cdot l$ multiplications and the same number of additions of the initial matrix elements. In case of square matrices, the size of which is $n \times n$, the number of the executed operations is the order $O(n^3)$. There are also sequential matrix multiplication algorithms of smaller computational complexity (for instance, the Strassen algorithm). But studying these algorithms though requires certain efforts and for simplicity we will use the above described sequential algorithm as the basis for parallel method development in this section. We will also assume further that all matrices are square and their sizes are $n \times n$.

8.2. Sequential Algorithm

The sequential matrix multiplication algorithm includes three nested loops:

```
// Sequential matrix multiplication algorithm
double MatrixA[Size][Size];
double MatrixB[Size][Size];
double MatrixC[Size][Size];
int i, j, k;
...
for (i=0; i<Size; i++){
    for (j=0; j<Size; j++){
        MatrixC[i][j] = 0;
        for (k=0; k<Size; k++){
            MatrixC[i][j] = MatrixC[i][j] + MatrixA[i][k]*MatrixB[k][j];
        }
    }
}
```

This algorithm is an iterative procedure and calculates sequentially the rows of the matrix C . In fact, a result matrix row is computed per outer loop (loop variable i) iteration (see Figure 8.1)

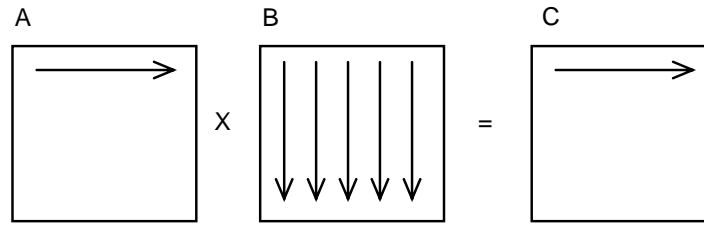


Figure 8.1 During the first iteration of loop variable i the first matrix A row and all the columns of matrix B are used to compute the elements of the first result matrix C row

As each result matrix element is a scalar product of the initial matrix A row and the initial matrix B column, it is necessary to carry out $n^2(2n-1)$ operations to compute all elements of the matrix C . As a result the time complexity of matrix multiplication is

$$T_1 = n^2 \cdot (2n - 1) \cdot \tau \quad (8.3)$$

where τ is the execution time for an elementary computational operation such as multiplication or addition.

8.3. Matrix Multiplication in Case of Block-Striped Data Decomposition

Let us consider two parallel matrix multiplication algorithms. Matrices A and B are partitioned into continuous sequences of rows or columns (*stripes*).

8.3.1. Computation Decomposition

As it is clear from the definition of matrix multiplication, all elements of the matrix C may be computed independently. As a result, a possible approach for parallelizing the matrix multiplication is to define the basic computational subtask as the problem of computing an element of the result matrix C . To carry out all the necessary computations each subtask must contain a row of the matrix A and a column of the matrix B . The total number of subtasks in case of this approach appears to be equal to n^2 (according to the number of elements of the matrix C).

One may note that the level of parallelism achieved in this approach is somewhat excessive. As a rule, in carrying out practical computations the number of the subtasks formed exceeds the number of the available processors. As a result, the aggregation stage of basic subtasks becomes inevitable. In this respect it is reasonable to aggregate the computations at the stage of selecting the basic subtasks. A possible solution is to combine all the computations related not with one, but with several elements of the result matrix C in a single subtask. For further discussion we

will define the basic computational subtask as the problem of computing all row elements of the matrix C . This approach decreases the total number of subtasks up to value n .

A row of the matrix A and all the columns of the matrix B must be available for carrying out all the necessary computations of the basic subtasks. The simple solution to the problem is duplicating the matrix B in all the subtasks, but it is unacceptable because of sizeable memory expenses needed for data storage. As a result, computations should be implemented so that subtasks contain only a part of the data needed for the computations at any given moment. The access to the other part of the data should be provided by means of data communications. Two possible ways to carry out parallel computations of this type are considered in this lecture.

8.3.2. Analysis of Information Dependencies

To compute a row of the matrix C each subtask must have a row of the matrix A and access to all columns of the matrix B . Possible ways to organize parallel computations are described below.

1. The first algorithm. The algorithm is an iterative procedure, the number of iterations is equal to the number of subtasks. Each subtask holds a row of the matrix A and a column of the matrix B at each algorithm iteration. At each iteration the scalar products of rows and columns containing in the subtasks are computed, and the corresponding elements of the result matrix C are obtained. After completing of all iteration computations the columns of matrix B must be transmitted so that subtasks should have new columns of the matrix B and new elements of the matrix C could be calculated. This transmission of columns among the subtasks must be executed in such a way that all the columns of matrix B should have appeared in each subtask sequentially.

A possible simple scheme to provide the required communications of the columns of matrix B among the subtasks is to present the topology of the information dependencies of the subtasks as a ring structure. In this case the subtask i , $0 \leq i < n$, will transmit its column of matrix B to the subtask $i+1$ at each iteration (in accordance with the ring structure subtask $n-1$ transmits its data to the subtask 0) – see Figure 8.2. After the algorithm termination the required condition will be provided, i.e. all the columns of matrix B will appear sequentially in each subtask.

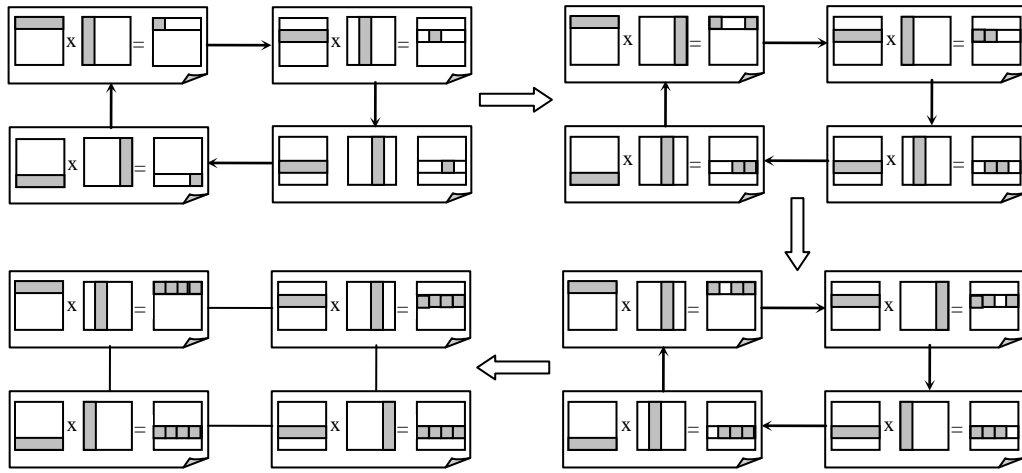


Figure 8.2 General scheme of data communications for the first parallel algorithm of matrix multiplication in case of block-stripped decomposition

Figure 8.2 presents the iterations of the matrix multiplication algorithm for the case when matrices have four rows and four columns ($n=4$). At the beginning of the computations each subtask i , $0 \leq i < n$, holds i -th row of the matrix A and i -th column of the matrix B . As a result the subtask i can compute the element c_{ii} of the result matrix C . Further each subtask transmits its column of matrix B to the following subtask in accordance with the ring structure. These actions should be repeated until all the iterations of the parallel algorithm are completed.

Program realization. Let us consider software implementation of presented method:

```
int BaseMatrixMultiplication( double* pAMatrix,
                             double* pBMatrix,
                             double* pCMatrix,
                             const int Size )
{
    MPI_Datatype ColumnType, Rtp;
    int ProcNum, ProcRank;
    MPI_Status status;
    int StripSize;
    double* Result;
    double* rbuf;
    double* vertStrip;
    int i, j, z, k;
    MPI_Comm_size ( MPI_COMM_WORLD, &ProcNum );
    MPI_Comm_rank ( MPI_COMM_WORLD, &ProcRank );
    StripSize = Size / ProcNum;

    if (StripSize == 0)
    {
        StripSize = 1;
    }
    Result = (double*) malloc (StripSize * StripSize * sizeof(double));

    rbuf = (double*) malloc (Size * StripSize * sizeof(double));
    vertStrip = (double*) malloc (Size * sizeof(double));
    MPI_Type_vector ( Size, 1, Size, MPI_DOUBLE, &ColumnType );
    MPI_Type_commit ( &ColumnType );
    MPI_Type_vector ( StripSize, StripSize, Size, MPI_DOUBLE, &Rtp );
    MPI_Type_commit ( &Rtp );
```

```

MPI_Scatter( pAMatrix, Size * StripSize, MPI_DOUBLE, rbuf,
            Size * StripSize, MPI_DOUBLE, 0, MPI_COMM_WORLD);
for (z = 0; z < ProcNum; z++)
{
    if(ProcRank == 0)
    {
        for(i = 1; i < ProcNum; i++)
        {
            for (j = 0; j < StripSize; j++)
                MPI_Send(pBMatrix + j + z * StripSize, 1, ColumnType,
                        i, 0, MPI_COMM_WORLD);
        }
    }
    else
    {
        double tmp;
        for (k = 0; k < StripSize; k++)
        {
            MPI_Recv(vertStrip, Size, MPI_DOUBLE, 0, 0,
                     MPI_COMM_WORLD, &status);

            for(i = 0; i < StripSize; i++)
            {
                tmp = 0.0;
                for(j = 0; j < Size; j++)
                {
                    tmp += rbuf[i * Size + j] * vertStrip[j];
                }
                Result[i * StripSize + k] = tmp;
            }
        }
        MPI_Send( Result, StripSize * StripSize,
                 MPI_DOUBLE, 0, 0, MPI_COMM_WORLD);
    }
    if(ProcRank == 0)
    {
        double tmp;
        for (i = 0; i < StripSize; i++)
        {
            for (j = 0; j < Size; j++)
            {
                tmp = 0.0;
                for(k = 0; k < Size; k++)
                {
                    tmp += pAMatrix[ j * Size + k ] *
                        pBMatrix[ k * Size + i + z * StripSize ];
                }
                pCMatrix[ j * Size + i + z * StripSize ] = tmp;
            }
        }
        for(i = 1; i < ProcNum; i++)
        {
            MPI_Recv(pCMatrix + i*Size*StripSize + z * StripSize, 1,
                     Rtp, i, 0, MPI_COMM_WORLD, &status);
        }
    }
}
MPI_Type_free ( &ColumnType );
MPI_Type_free ( &Rtp );

```

```

free ( Result );
free ( rbuf );
free ( vertStrip );

return 0;
}

```

2. The second algorithm. The difference of the second algorithm from the first one is that the subtasks contain not columns but rows of matrix B . As a result, data multiplication of each subtask is the multiplication of the row elements of the matrix B by a corresponding row element of the matrix A . Therefore, a row of partial results for matrix C is obtained in each subtask.

In case of this scheme of data decomposition for matrix multiplication, it is necessary to provide sequential obtaining all rows of the matrix B by all in the subtasks, the multiplication of the row elements of the matrix B by a corresponding row element of the matrix A and summation of the new values and the previously computed ones. The ring structure of information dependencies may be also used to provide the necessary sequence of communications of the rows of the matrix B among the subtasks (see Figure 8.3).

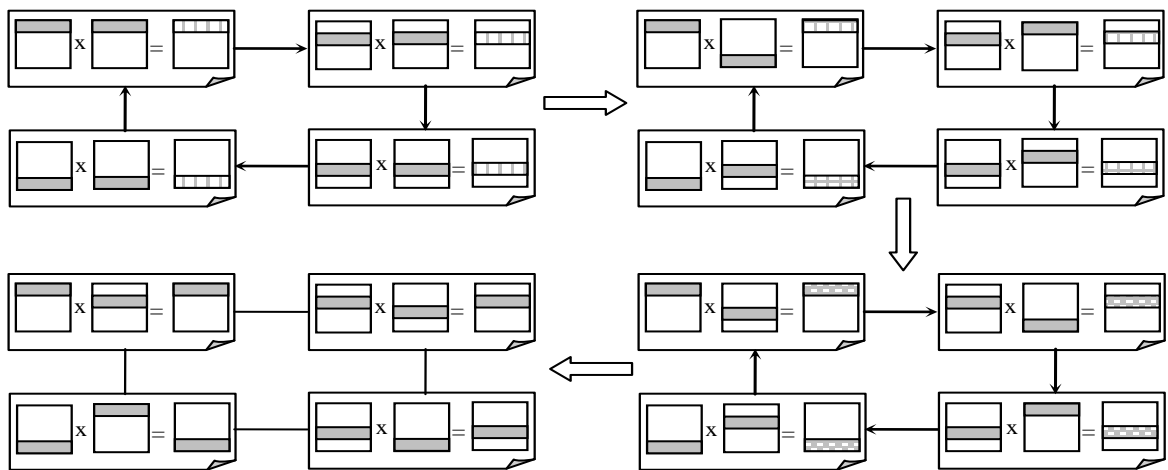


Figure 8.3 General scheme of data communications for the second parallel algorithm of matrix multiplication in case of block-striped decomposition

Figure 8.3 presents the iterations of the matrix multiplication algorithm in the case when matrices have 4 rows and 4 columns ($n=4$). At the beginning of the computations each subtask i , $0 \leq i < n$, holds i -th rows of the matrix A and the matrix B . As a result of multiplication the subtask defines i -th row of the partial results for the matrix C . Then each subtask transmits its row of the matrix B to the following subtask according to the ring structure of information dependencies. The described actions are repeated until all the iterations of the parallel algorithm are completed.

Program realization. Let us consider software implementation of presented method:

```

void MatrixMultiplication(double *pAMatrix, double *pBMatrix,

```

```

double *pCMatrix, const int Size)
{
    MPI_Status Status;
    int iter;
    int NextProc;
    int PrevProc;
    int shift;
    int k, i, j;
    int ProcNum;
    int ProcRank;
    int RowSize;
    double *AMatrixRow;
    double *BMatrixRow;
    double *ProcSum;
    MPI_Comm_size(MPI_COMM_WORLD, &ProcNum);
    MPI_Comm_rank(MPI_COMM_WORLD, &ProcRank);
    RowSize = Size / ProcNum;
    if(RowSize == 0)
    {
        RowSize = 1;
    }

    assert(RowSize * ProcNum == Size);

    AMatrixRow = ( double* ) malloc ( Size * RowSize * sizeof(double) );
    BMatrixRow = ( double* ) malloc ( Size * RowSize * sizeof(double) );
    ProcSum = ( double* ) calloc ( Size * RowSize , sizeof(double) );
    MPI_Scatter(pAMatrix, Size * RowSize, MPI_DOUBLE,
               AMatrixRow, Size * RowSize, MPI_DOUBLE, 0, MPI_COMM_WORLD);
    MPI_Scatter(pBMatrix, Size * RowSize, MPI_DOUBLE,
               BMatrixRow, Size * RowSize, MPI_DOUBLE, 0, MPI_COMM_WORLD);

    for(iter = 0; iter < Size / RowSize; iter++)
    {
        shift = ( ProcRank * RowSize + iter * RowSize ) % ( Size );

        for(k = 0; k < RowSize; k++)
        {
            for (i = 0; i < RowSize; i++)
            {
                for(j = 0; j < Size; j++)
                {
                    ProcSum[k * Size + j] += (AMatrixRow + shift) [ k * Size + i ]
                    * BMatrixRow[ i * Size + j];
                }
            }
        }

        NextProc = ProcRank - 1;

        if(NextProc < 0)
        {
            NextProc = ProcNum - 1;
        }

        PrevProc = (ProcRank + 1)%ProcNum;
        MPI_Sendrecv_replace(BMatrixRow, Size * RowSize, MPI_DOUBLE,
                             NextProc, 0, PrevProc, 0, MPI_COMM_WORLD, &Status);
    }
}

```



```

MPI_Gather(ProcSum, RowSize * Size, MPI_DOUBLE, pCMatrix,
          RowSize * Size, MPI_DOUBLE, 0, MPI_COMM_WORLD);

free ( ProcSum );
free ( AMatrixRow );
free ( BMatrixRow );

}

```

8.3.3. Computational Experiment Results

The experiments were carried out on the computational cluster on the basis of processors Intel XEON 4 EM64T, 3000 Mhz and Gigabit Ethernet under OS Microsoft Windows Server 2003 Standard x64 Edition.

The results of the computational experiments are shown in Table 8.1. The experiments were performed with the use of 2, 4 and 8 processors.

Table 8.1. The results of the computational experiments for the first parallel algorithm of matrix multiplication based on the block-striped data decomposition

Matrix Size	Serial Algorithm	2 processors		4 processors		8 processors	
		Time	Speed Up	Time	Speed Up	Time	Speed Up
500	0,8752	0,3758	2,3287	0,1535	5,6982	0,0968	9,0371
1000	12,8787	5,4427	2,3662	2,2628	5,6912	0,6998	18,4014
1500	43,4731	20,9503	2,0750	11,0804	3,9234	5,1766	8,3978
2000	103,0561	45,7436	2,2529	21,6001	4,7710	9,4127	10,9485
2500	201,2915	99,5097	2,0228	56,9203	3,5363	18,3303	10,9813
3000	347,8434	171,9232	2,0232	111,9642	3,1067	45,5482	7,6368

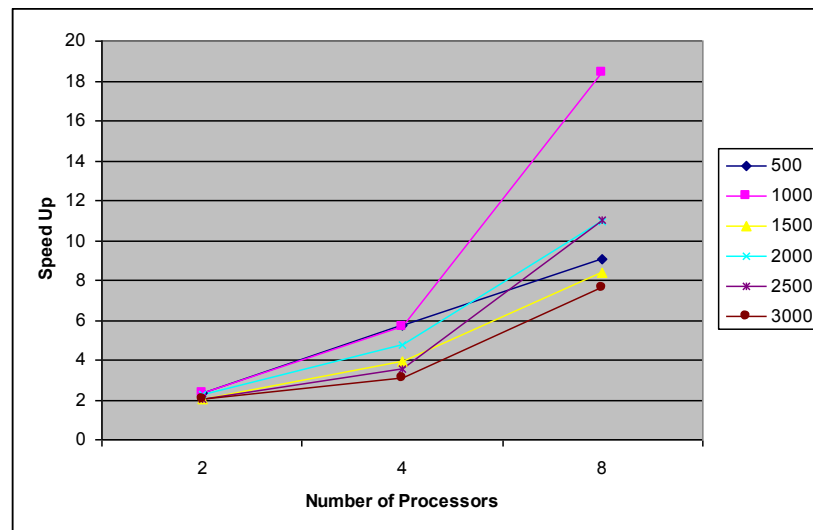


Figure 8.4 Speedup for the first parallel algorithm of matrix multiplication (block-striped matrix decomposition)

8.4. References

The problem of matrix multiplication is broadly discussed in science. As additional training materials we may recommend the works by Kumar, et al. (1994) and Quinn (2004). The problems of parallel execution of matrix multiplication are also discussed in Dongarra, et al. (1999).

Blackford, et al. (1997) may be useful for considering some aspects of parallel software development. This book describes the software library of numerical methods ScaLAPACK, which is well-known and widely used.

Lecture 9

9.1. Fox Algorithm of Matrix Multiplication in Case of Checkerboard Data Decomposition

In designing the parallel methods of matrix multiplication the checkerboard block matrix decomposition is widely used just as the block-striped matrix partitioning. Let us analyze this method of computations in detail.

9.2 Computation Decomposition

In case of this method of data decomposition the initial matrices A and B and the result matrix C are subdivided into sets of blocks. For simplicity the further explanations we will assume all the matrices are square of $n \times n$ size, the number of vertical blocks and the number of horizontal blocks are the same and are equal to q (i.e. the size of all block is equal to $k \times k$, $k=n/q$). In

case of this data decomposition method the multiplying matrices A and B as blocks may be represented as follows:

$$\begin{pmatrix} A_{00} & A_{01} & \dots & A_{0q-1} \\ & & & \\ & & & \\ A_{q-10} & A_{q-11} & \dots & A_{q-1q-1} \end{pmatrix} \times \begin{pmatrix} B_{00} & B_{01} & \dots & B_{0q-1} \\ & & & \\ & & & \\ B_{q-10} & B_{q-11} & \dots & B_{q-1q-1} \end{pmatrix} = \begin{pmatrix} C_{00} & C_{01} & \dots & C_{0q-1} \\ & & & \\ & & & \\ C_{q-10} & C_{q-11} & \dots & C_{q-1q-1} \end{pmatrix},$$

where each block C_{ij} of matrix C is computed in accordance with the expression:

$$C_{ij} = \sum_{s=0}^{q-1} A_{is} B_{sj}.$$

In case of the checkerboard block data decomposition it is reasonable to define the basic computational subtasks on the basis of the computations performed over the matrix blocks. As a result the basic subtask can be defined as the problem of computing of a block of the matrix C .

To perform all the necessary computations the basic subtasks should have the corresponding sets of the matrix A rows and the matrix B columns. The placement all the necessary data in each subtask will inevitably lead to duplicating and to a considerable increase of the size of memory used. As a result, the computations must be executed in such a way that the subtasks should contain only a part of the data necessary for computations at any given moment, and the access to the rest of the data should be provided by means of data communications. One of the possible approaches (*the Fox algorithm*) will be discussed further in this Lecture.

9.3 Analysis of Information Dependencies

To develop a parallel matrix multiplication method based on the checkerboard decomposition scheme it should be reminded that in this case the basic subtasks are responsible for computing the separate blocks of the matrix C . It is also required that each subtask should hold only one block of the multiplying matrices at each iteration.

To enumerate the subtasks the indices of the blocks C_{ij} contained in the subtasks can be used for enumeration. Thus, the subtask (i,j) computes the block C_{ij} . So the set of subtasks forms a square grid, which corresponds to the structure of the checkerboard block decomposition of the matrix C .

The Fox algorithm can be used to perform matrix multiplication computations under these conditions (see for instance, Fox et al. (1987), Kumar et al. (1994)).

In accordance with the Fox algorithm each basic subtasks (i,j) holds four matrix blocks:

- Block C_{ij} of matrix C , computed by the subtask;
- Block A_{ij} of matrix A , placed in the subtask before the beginning of computations;

– Blocks A'_{ij} , B'_{ij} of matrices A and B , obtained by the subtask in the course of computations.

Parallel algorithm execution includes:

- **The initialization stage.** Each subtask (i,j) obtains blocks A_{ij} , B_{ij} . All elements of blocks C_{ij} in all subtasks are set to zero;
- **The computation stage.** At this stage the following operations are carried out at each iteration l , $0 \leq l < q$:
 - For each row i , $0 \leq i < q$, the block A_{ij} of subtask (i,j) is transmitted to all the subtasks of the same processor grid row; index j , which defines the position of the subtask in the row, is computed according to the following expression:

$$j = (i + l) \bmod q,$$

where \bmod is operation of obtaining the remainder in integer division;

- Blocks A'_{ij} , B'_{ij} obtained as a result of subtask communications are multiplied and added to block C_{ij} :

$$C_{ij} = C_{ij} + A'_{ij} \times B'_{ij};$$

- Blocks B'_{ij} of each subtask (i,j) are transmitted to the subtasks, which are upper neighbors in the processor grid columns (the first row blocks are transmitted to the last row of the grid).

To illustrate these rules we show the state of blocks in each subtask in the course of executing iterations of the computation stage (for the grid of 2×2). See Figure 9.1.

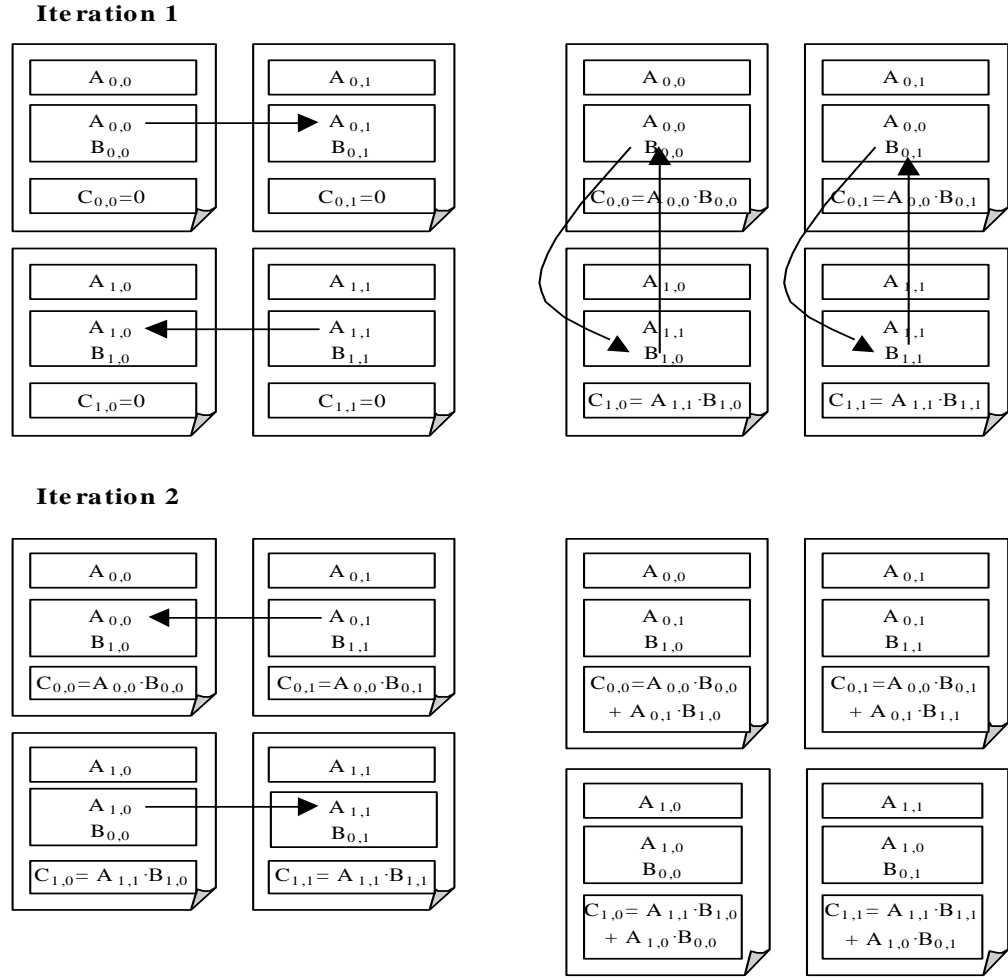


Figure 9.1 Block distribution among subtasks on iterations of the Fox algorithm

9.4 Scaling and Distributing Subtasks among Processors

The number of blocks at the checkerboard decomposition scheme can be regulated by variation of matrix block sizes. These sizes may be chosen so that the total number of the basic subtasks coincides with the number of processors p . Thus, for instance, in the simplest case when the number of processors is equal to $p = \delta^2$, the size of the block grid may be chosen equal to δ (i.e. $q = \delta$). This way to define the number of blocks makes the amount of computations in each subtask the same and, thus, uniform balancing of the computational load is achieved. In a more general case, when the number of processors and the sizes of matrices are arbitrary, computational load may not be distributed among processors equally but proper setting the sizes of the matrix blocks can provide uniform load balancing with adequate accuracy.

To execute the Fox algorithm efficiently, when the basic subtasks form a square grid and data communications consist in block transmission along rows and columns of the subtask grid, the network topology should be also a square grid. In this case it is possible to map easily the set of

subtasks onto the set of processors by placing the basic subtasks (i,j) on processors $P_{i,j}$. The required structure of the data communication network may be provided at the physical level, if the network topology is a grid or a complete graph.

9.5 Software Implementation

Here we discuss possible software implementation of the Fox algorithm for matrix multiplication in case of the checkerboard block data decomposition. The given program code contains the basic modules of the parallel program. The absence of some auxiliary functions will not hinder the process of understanding of this parallel computation scheme.

1. The main function. The main function implements the computational method scheme by sequential calling out the necessary subprograms.

```
// The Fox algorithm of matrix multiplication - checkerboard decomposition
// Program execution conditions:
//   all matrices and their blocks are square,
//   matrix blocks and processors form square grids of the same size

int ProcNum = 0;      // Number of available processes
int ProcRank = 0;     // Rank of current process
int GridSize;        // Size of virtual processor grid
int GridCoords[2];    // Coordinates of current processor in grid
MPI_Comm GridComm;    // Grid communicator
MPI_Comm ColComm;     // Column communicator
MPI_Comm RowComm;     // Row communicator

void main ( int argc, char * argv[] ) {
    double* pAMatrix;  // The first argument of matrix multiplication
    double* pBMatrix;  // The second argument of matrix multiplication
    double* pCMatrix;  // The result matrix
    int Size;          // Size of matrices
    int BlockSize;     // Sizes of matrix blocks on current process
    double *pAblock;   // Initial block of matrix A on current process
    double *pBblock;   // Initial block of matrix B on current process
    double *pCblock;   // Block of result matrix C on current process
    double *pMatrixAblock;
    double Start, Finish, Duration;

    setvbuf(stdout, 0, _IONBF, 0);

    MPI_Init(&argc, &argv);
    MPI_Comm_size(MPI_COMM_WORLD, &ProcNum);
    MPI_Comm_rank(MPI_COMM_WORLD, &ProcRank);

    GridSize = sqrt((double)ProcNum);
    if (ProcNum != GridSize*GridSize) {
        if (ProcRank == 0) {
            printf ("Number of processes must be a perfect square \n");
        }
    }
    else {
        if (ProcRank == 0)
            printf("Parallel matrix multiplication program\n");
    }
}
```

```

// Creating the cartesian grid, row and column communicators
CreateGridCommunicators();

// Memory allocation and initialization of matrix elements
ProcessInitialization ( pAMatrix, pBMatrix, pCMatrix, pAblock, pBblock,
    pCblock, pMatrixAblock, Size, BlockSize );

DataDistribution(pAMatrix, pBMatrix, pMatrixAblock, pBblock, Size,
    BlockSize);

// Execution of Fox method
ParallelResultCalculation(pAblock, pMatrixAblock, pBblock,
    pCblock, BlockSize);

ResultCollection(pCMatrix, pCblock, Size, BlockSize);
TestResult(pAMatrix, pBMatrix, pCMatrix, Size);

// Process Termination
ProcessTermination (pAMatrix, pBMatrix, pCMatrix, pAblock, pBblock,
    pCblock, pMatrixAblock);
}

MPI_Finalize();
}

```

2. The function *CreateGridCommunicators*. This function creates a communicator as a two-dimensional square grid, determines the coordinates of each process in the grid and creates communicators for each row and each column separately.

The grid is created by the function *MPI_Cart_create* (the vector *Periodic* defines the permissibility of data communications among the bordering processes of the grid columns and rows). After the grid has been created, each parallel program process will have its coordinates in the grid. The coordinates may be obtained by means of the function *MPI_Cart_coords*.

Then in addition to the grid topology a set of communicators for each grid column and row separately is created by the function *MPI_Cart_sub*.

```

void CreateGridCommunicators() {
    int DimSize[2]; // Number of processes in each dimension of the grid
    int Periodic[2]; // =1, if the grid dimension should be periodic
    int Subdims[2]; // =1, if the grid dimension should be fixed

    DimSize[0] = GridSize;
    DimSize[1] = GridSize;
    Periodic[0] = 0;
    Periodic[1] = 0;

    // Creation of the Cartesian communicator
    MPI_Cart_create(MPI_COMM_WORLD, 2, DimSize, Periodic, 1, &GridComm);

    // Determination of the cartesian coordinates for every process
    MPI_Cart_coords(GridComm, ProcRank, 2, GridCoords);

    // Creating communicators for rows
    Subdims[0] = 0; // Dimensionality fixing
    Subdims[1] = 1; // The presence of the given dimension in the subgrid
    MPI_Cart_sub(GridComm, Subdims, &RowComm);
}

```

```

// Creating communicators for columns
Subdims[0] = 1;
Subdims[1] = 0;
MPI_Cart_sub(GridComm, Subdims, &ColComm);
}

```

3. The function *ProcessInitialization*. This function sets the matrix sizes and allocates memory for storing the initial matrices and their blocks, initializes all the original problem data. In order to determine the elements of the initial matrices we will use the functions *DummyDataInitialization* and *RandomDataInitialization*.

```

// Function for memory allocation and data initialization
void ProcessInitialization (double* &pAMatrix, double* &pBMatrix,
double* &pCMatrix, double* &pAblock, double* &pBblock, double* &pCblock,
double* &pTemporaryAblock, int &Size, int &BlockSize) {
if (ProcRank == 0) {
do {
printf("\nEnter size of the initial objects: ");
scanf("%d", &Size);

if (Size%GridSize != 0) {
printf ("Size of matrices must be divisible by the grid size! \n");
}
}
while (Size%GridSize != 0);
}
MPI_Bcast(&Size, 1, MPI_INT, 0, MPI_COMM_WORLD);

BlockSize = Size/GridSize;

pAblock = new double [BlockSize*BlockSize];
pBblock = new double [BlockSize*BlockSize];
pCblock = new double [BlockSize*BlockSize];
pTemporaryAblock = new double [BlockSize*BlockSize];

for (int i=0; i<BlockSize*BlockSize; i++) {
pCblock[i] = 0;
}
if (ProcRank == 0) {
pAMatrix = new double [Size*Size];
pBMatrix = new double [Size*Size];
pCMatrix = new double [Size*Size];
//DummyDataInitialization(pAMatrix, pBMatrix, Size);
RandomDataInitialization(pAMatrix, pBMatrix, Size);
}
}
}

```

4. The function *ParallelResultCalculation*. The function *ParallelResultCalculation* executes the parallel Fox algorithm of matrix multiplication. The matrix blocks and their sizes must be given to the function as its arguments.

According to the scheme of parallel computations described in Exercise 3, it is necessary to carry out *GridSize* iterations in order to execute matrix multiplication with the use of Fox algorithm. Each of the iterations consists of the execution of the following operations:

- The broadcast of the matrix A block along the processor grid row (to execute the step we should develop the function *ABlockCommunication*),
- The multiplication of matrix blocks (to carry out the multiplication of matrix blocks we may use the function *SerialResultCalculation*, which was implemented in the course of the development of the serial matrix multiplication program),

The cyclic shift of the matrix B blocks along the column of the processor grid (the function *BBlockCommunication*).

```
void ParallelResultCalculation(double* pAblock, double* pMatrixAblock,
double* pBblock, double* pCblock, int BlockSize) {
    for (int iter = 0; iter < GridSize; iter++) {
        // Sending blocks of matrix A to the process grid rows
        ABlockCommunication (iter, pAblock, pMatrixAblock, BlockSize);
        // Block multiplication
        BlockMultiplication(pAblock, pBblock, pCblock, BlockSize);
        // Cyclic shift of blocks of matrix B in process grid columns
        BblockCommunication(pBblock, BlockSize);
    }
}
```

5. The function *ABlockCommunication*. The function broadcasts matrix A blocks to the process grid rows. The leading process *Pivot* that responsible for sending is chosen in each row of the grid. For broadcasting the pivot processes are used their blocks *pMatrixAblock* (let us to remind that these blocks transmitted to the processes at the moment of the initial data distribution). The required communications are executed by means of the function *MPI_Bcast*. It should be noted that the operation is collective, and its localization in separate process grid rows is provided by the communicators *RowComm*, which are created for the set of processes of each row separately.

```
// Broadcasting matrix A blocks to process grid rows
void ABlockCommunication (int iter, double *pAblock, double* pMatrixAblock,
int BlockSize) {

    // Defining the leading process of the process grid row
    int Pivot = (GridCoords[0] + iter) % GridSize;

    // Copying the transmitted block in a separate memory buffer
    if (GridCoords[1] == Pivot) {
        for (int i=0; i<BlockSize*BlockSize; i++)
            pAblock[i] = pMatrixAblock[i];
    }

    // Block broadcasting
    MPI_Bcast(pAblock, BlockSize*BlockSize, MPI_DOUBLE, Pivot, RowComm);
}
```

6. The function *BlockMultiplication*. The function executes block multiplication of the matrices A and B . The easiest way to perform this multiplication is to use the serial matrix multiplication algorithm. It should be noted that we provide the simplest variant of the function implementation for better understanding of the program. These calculations may be optimized to de-

crease the computation time. This optimization may be aimed, for instance, at increasing the efficiency of the processor cache, vectorizing the executed operations etc.

7. The function *BblockCommunication*. The function performs the cyclic shift of blocks of the matrix B in the process grid columns. Each process transmits its block to the upper neighboring process *NextProc* in the process column and receives the block transmitted from the process *PrevProc*, which stands below it in the grid column. Data transmission is executed by means of the function *MPI_SendRecv_replace*, which provides all the necessary block transmissions using the same memory buffer *pBblock*. Besides, this function prevents possible deadlocks, which happen when data transmission begins to be performed simultaneously by several processes in the ring network topology.

```
// Cyclic shift of matrix B blocks in the process grid columns
void BblockCommunication (double *pBblock, int BlockSize) {
    MPI_Status Status;
    int NextProc = GridCoords[0] + 1;
    if ( GridCoords[0] == GridSize-1 ) NextProc = 0;
    int PrevProc = GridCoords[0] - 1;
    if ( GridCoords[0] == 0 ) PrevProc = GridSize-1;

    MPI_Sendrecv_replace( pBblock, BlockSize*BlockSize, MPI_DOUBLE,
        NextProc, 0, PrevProc, 0, ColComm, &Status);
}
```

9.6 Computational Experiment Results

The results of the experiments with the use of 4 and 9 processors are given in Table 9.1.

Table 9.1 The Results of the computational experiments for estimating the Fox parallel algorithm efficiency

Matrix Size	Serial Algorithm	Parallel Algorithm			
		4 processors		9 processors	
		Time	Speed Up	Time	Speed Up
500	0,8527	0,2190	3,8925	0,1468	5,8079
1000	12,8787	3,0910	4,1664	2,1565	5,9719
1500	43,4731	10,8678	4,0001	7,2502	5,9960
2000	103,0561	24,1421	4,2687	21,4157	4,8121
2500	201,2915	51,4735	3,9105	41,2159	4,8838
3000	347,8434	87,0538	3,9957	58,2022	5,9764

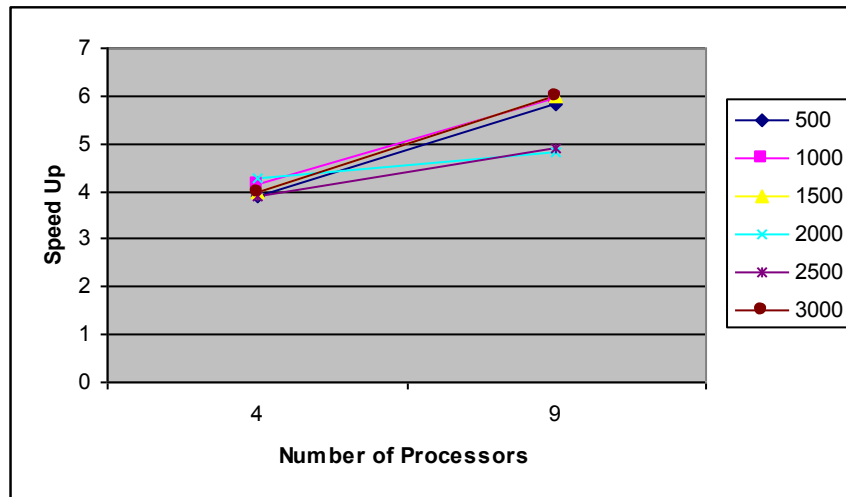


Figure 9.1 Speedup of the Fox Parallel Algorithm with Respect to Number of Processors

Test questions

1. What is the statement of the matrix multiplication problem?
2. Give the examples of the problems, which make use of the matrix multiplication operations.
3. Give the examples of various sequential algorithms of matrix multiplication operations. Is the complexity various in case of different algorithms?
4. What methods of data distribution are used in developing parallel algorithms of matrix multiplication?
5. Analyze and compute the efficiency of the block-striped algorithm for horizontal partitioning of the multiplied matrices.
6. What information communications are carried out for the algorithms in case of the block-striped data decomposition?

Practice

1. Implement two block striped algorithms of matrix multiplication. Compare runtimes.
2. Implement the Fox algorithm. Perform computational experiments. Compare experimental results with previous implementations.

References

1. **Dongarra, J.J., Duff, L.S., Sorensen, D.C., Vorst, H.A.V. (1999).** Numerical Linear Algebra for High Performance Computers (Software, Environments, Tools). Soc for Industrial & Applied Math.

2. **Blackford**, L. S., Choi, J., Cleary, A., D'Azevedo, E., Demmel, J., Dhillon, I., Dongarra, J. J., Hammarling, S., Henry, G., Petitet, A., Stanley, D. Walker, R.C. Whaley, K. (1997). Sca-lapack Users' Guide (Software, Environments, Tools). Soc for Industrial & Applied Math.
3. **Foster**, I. (1995). Designing and Building Parallel Programs: Concepts and Tools for Soft-ware Engineering. Reading, MA: Addison-Wesley.
4. **Andrews**, G. R. (2000). Foundations of Multithreaded, Parallel, and Distributed Program-ming.. – Reading, MA: Addison-Wesley.
5. **Kahaner**, D., Moler, C., Nash, S. (1988). Numerical Methods and Software. – Prentice Hall.
6. **Quinn**, M. J. (2004). Parallel Programming in C with MPI and OpenMP. – New York, NY: McGraw-Hill.
7. **Wilkinson**, B., Allen, M. (1999). Parallel programming. – Prenrice Hall.