

BSP ALGORITHM DESIGN: INTRODUCTION

Disclaimer: THESE NOTES DO NOT SUBSTITUTE THE TEXTBOOK FOR THIS CLASS. THE NOTES SHOULD BE USED IN CONJUNCTION WITH THE TEXTBOOK AND THE MATERIAL PRESENTED IN CLASS. IF A STATEMENT IN THESE NOTES SEEMS TO BE INCORRECT, REPORT IT TO THE INSTRUCTOR SO THAT IT BE FIXED IMMEDIATELY. THESE NOTES ARE ONLY DISTRIBUTED TO THE STUDENTS TAKING THIS CLASS WITH A. GERBESSIOTIS IN FALL 2004; DISTRIBUTION OUTSIDE THIS GROUP OF STUDENTS IS NOT ALLOWED.

Traditional vs Architecture Independent Parallel Algorithm Design

As an example of how traditional PRAM algorithm design differs from architecture independent parallel algorithm design, example algorithm for broadcasting in a parallel machine is introduced.

Problem: In a parallel machine with p processors numbered $0, \dots, p - 1$, one of them, say processor 0, holds a one-word message. The problem of *broadcasting* involves the dissemination of this message to the local memory of the remaining $p - 1$ processors.

The performance of a well-known exclusive PRAM algorithm for broadcasting is analyzed below in two ways under the assumption that no concurrent operations are allowed. One follows the traditional (PRAM) analysis that minimizes parallel running time. The other takes into consideration the issues of communication and synchronization as viewed under the BSP model. This leads to a modification of the PRAM-based algorithm to derive an architecture independent algorithm for broadcasting whose performance is consistent with observations of broadcasting operations on real parallel machines.

Broadcasting: PRAM Algorithm 1

Algorithm. Without loss of generality let us assume that p is a power of two. The message is broadcast in $\lg p$ rounds of communication by binary replication. In round $i = 1, \dots, \lg p$, each processor j with index $j < 2^{i-1}$ sends the message it currently holds to processor $j + 2^{i-1}$ (on a shared memory system, this may mean copying information into a cell read by this processor). The number of processors with the message at the end of round i is thus 2^i .

Analysis of Algorithm. Under the PRAM model the algorithm requires $\lg p$ communication rounds and so many parallel steps to complete. This cost, however, ignores synchronization which is for free, as PRAM processors work synchronously. It also ignores communication, as in the PRAM the cost of accessing the shared memory is as small as the cost of accessing local registers of the PRAM.

Under the BSP cost model each communication round is assigned a cost of $\max\{L, g \cdot 1\}$ as each processor in each round sends or receives at most one message containing the one-word message. The BSP cost of the algorithm is $\lg p \cdot \max\{L, g \cdot 1\}$, as there are $\lg p$ rounds of communication. As the communicated information by any processors is small in size, it is likely that latency issues prevail in the transmission time (ie bandwidth based cost $g \cdot 1$ is insignificant compared to the latency/synchronization reflecting term L).

In high latency machines the dominant term would be $L \lg p$ rather than $g \lg p$. Even though each communication round would last for at least L time units, only a small fraction g of it is used for actual communication. The remainder is wasted. It makes then sense to increase communication round utilization so that each processor sends the one-word message to as many processors as it can accommodate within a round.

Broadcasting: Algorithm 2

Input: p processors numbered $0 \dots p - 1$. Processor 0 holds a message of length equal to one word.

Output: The problem of *broadcasting* involves the dissemination of this message to the remaining $p - 1$ processors.

Algorithm 2. In one superstep, processor 0 sends the message to be broadcast to processors $1, \dots, p - 1$ in turn (a “sequential”-looking algorithm).

Analysis of Algorithm 2.

The communication time of Algorithm 2 is $1 \cdot \max\{L, (p - 1) \cdot g\}$ (in a single superstep, the message is replicated $p - 1$ times by processor 0).

BSP Algorithms

Broadcasting: Algorithm 3

Algorithm 3. Both Algorithm 1 and Algorithm 2 can be viewed as extreme cases of an Algorithm 3. The main observation is that up to L/g words can be sent in a superstep at a cost of L . Then, It makes sense for each processor to send L/g messages to other processors. Let $k - 1$ be the number of messages a processor sends to other processors in a broadcasting step. The number of processors with the message at the end of a broadcasting superstep would be k times larger than that in the start. We call k the degree of replication of the broadcast operation.

Architecture independent Algorithm 3. In each round, every processor sends the message to $k - 1$ other processors. In round $i = 0, 1, \dots$, each processor j with index $j < k^i$ sends the message to $k - 1$ distinct processors numbered $j + k^i \cdot l$, where $l = 1, \dots, k - 1$. At the end of round i (the $(i + 1)$ -st overall round), the message is broadcast to $k^i \cdot (k - 1) + k^i = k^{i+1}$ processors. The number of rounds required is the minimum integer r such that $k^r \geq p$, The number of rounds necessary for full dissemination is thus decreased to $\lg_k p$, and the total cost becomes $\lg_k p \max \{L, (k - 1)g\}$.

At the end of each superstep the number of processors possessing the message is k times more than that of the previous superstep. During each superstep each processor sends the message to exactly $k - 1$ other processors. Algorithm 3 consists of a number of rounds between 1 (and it becomes Algorithm 2) and $\lg p$ (and it becomes Algorithm 1).

```
BROADCAST (0, p, k)
1.  my_pid = bsp_pid();  mask_pid = 1;
2.  while (mask_pid < p) {
3.    if (my_pid < mask_pid)
4.      for (i = 1, j = mask_pid; i < k; i ++, j + = mask_pid) {
5.        target_pid = my_pid + j;
6.        if (target_pid < p)
7.          bsp_put(target_pid, &M, &M, 0, sizeof(M));
8.      }
9.    bsp_sync();
10.   mask_pid = mask_pid * k;
11. }
```

Broadcasting $n > p$ words: Algorithm 4

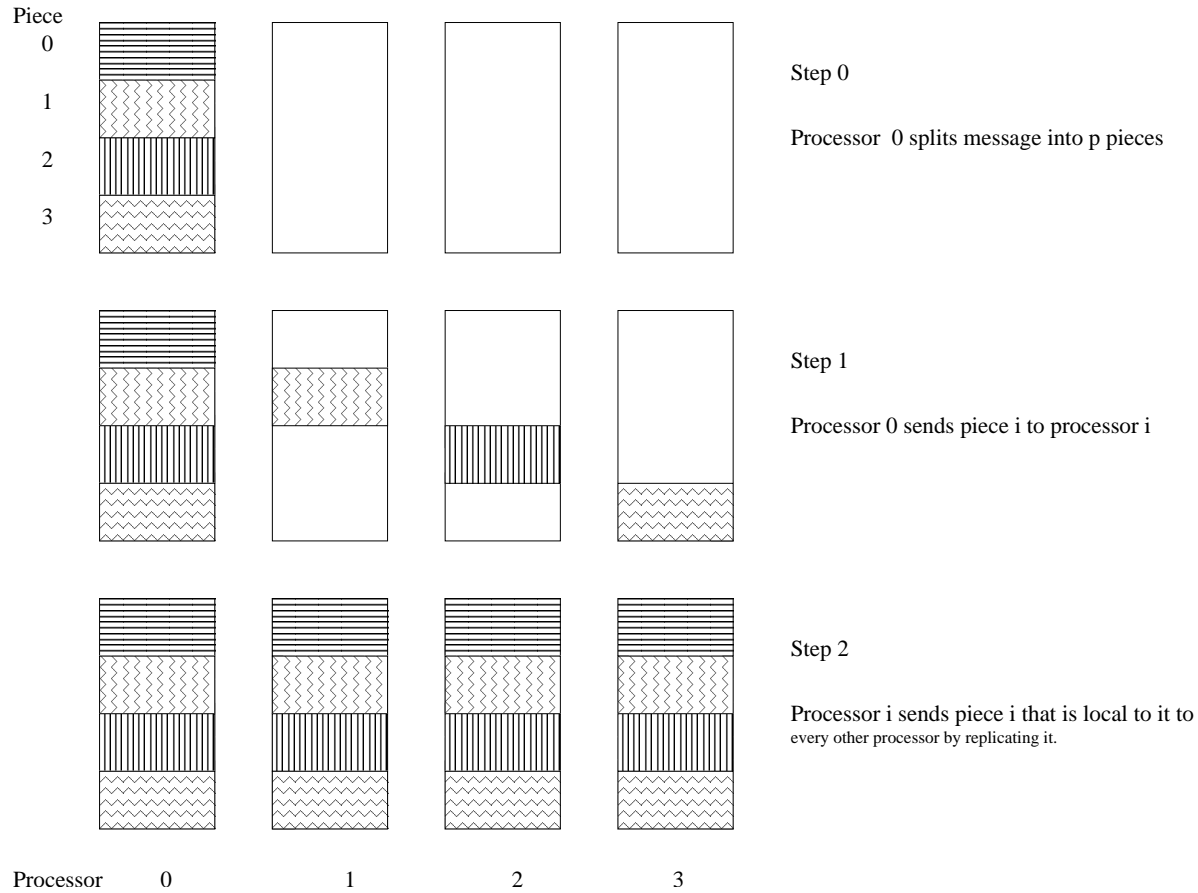
Now suppose that the message to be broadcast consists of not a single word but is of size $n > p$. Algorithm 4 may be a better choice than the previous algorithms as one of the processors sends or receives substantially more than n words of information. There is a broadcasting algorithm, call it Algorithm 4, that requires only two communication rounds and is optimal (for the communication model abstracted by L and g) in terms of the amount of information (up to a constant) each processor sends or receives.

Algorithm 4.

Two-phase broadcasting

The idea is to split the message into p pieces, have processor 0 send piece i to processor i in the first round and in the second round processor i replicates the i -th piece $p - 1$ times by sending each copy to each of the remaining $p - 1$ processors (see attached figure).

BSP Algorithms
 Broadcasting $n > p$ words: Example



BSP Algorithms

Parallel Prefix

Exercise. What can you say about parallel prefix? Analyze the BSP performance of the PRAM algorithm for parallel prefix. Can you halve its number of supersteps yet maintain the same BSP cost?

The structure of the four algorithms described for broadcasting can also be used to derive algorithms for parallel prefix that require similar number of supersteps (at most twice as many).

Algorithm 1 gives rise to a “sequential”-like parallel prefix algorithm. Algorithm 2 gives rise to the binary tree based algorithm that requires $2 \lg n$ supersteps. The corresponding PRAM algorithm, however, (that also runs on the butterfly) requires half as many supersteps and is thus more efficient on the BSP model. Algorithm 3 gives rise to the equivalents of 2 when the number of supersteps needs to be decreased.

We can generalize the prefix problem so that each processor instead of holding a single scalar value, holds a sequence/vector of scalar values n . In the case $n > p$, implementations following the counterparts of Algorithm 1,2 and 3 for broadcasting fail to provide optimal algorithms.

Algorithm 4 gives rise to a two-phase parallel prefix algorithm that is more efficient in architectures with large L for large independent prefix problems n .

BSP Algorithms

Matrix Computations

SPMD program design stipulates that processors execute a single program on different pieces of data. For matrix related computations it makes sense to distribute a matrix evenly among the p processors of a parallel computer. Such a distribution should also take into consideration the storage of the matrix by say the compiler so that locality issues are also taken into consideration (filling cache lines efficiently to speedup computation). There are various ways to divide a matrix. Some of the most common ones are described below.

One way to distribute a matrix is by using block distributions. Split an array into blocks of size $n/p_1 \times n/p_2$ so that $p = p_1 \times p_2$ and assign the i -th block to processor i . This distribution is suitable for matrices as long as the amount of work for different elements of the matrix is the same.

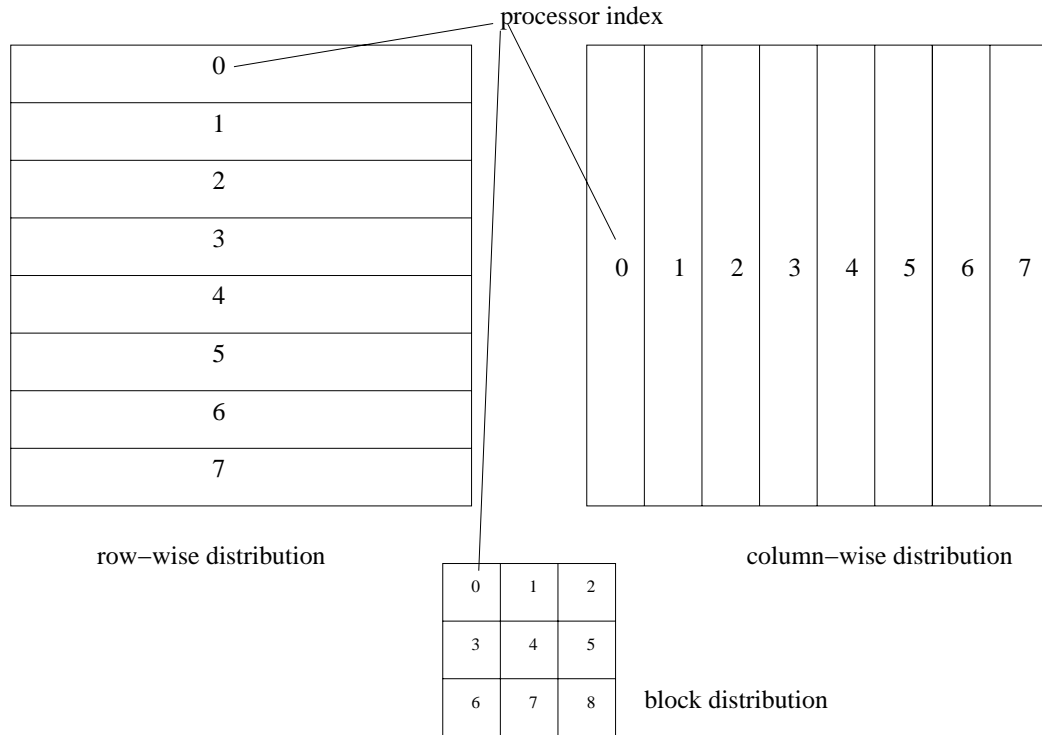
The most common block distributions are.

- **column-wise** (block) distribution. Split matrix into p column stripes so that n/p consecutive columns form the i -th stripe that will be stored in processor i . This is $p_1 = 1$ and $p_2 = p$.
- **row-wise** (block) distribution. Split matrix into p row stripes so that n/p consecutive rows form the i -th stripe that will be stored in processor i . This is $p_1 = p$ and $p_2 = 1$.
- **block** or **square** distribution. This is the case $p_1 = p_2 = \sqrt{p}$, i.e. the blocks are of size $n/\sqrt{p} \times n/\sqrt{p}$ and store block i to processor i .

There are certain cases (eg. LU decomposition, Cholesky factorization), where the amount of work differs for different elements of a matrix. For these cases block distributions are not suitable.

BSP Algorithms

Block distributions



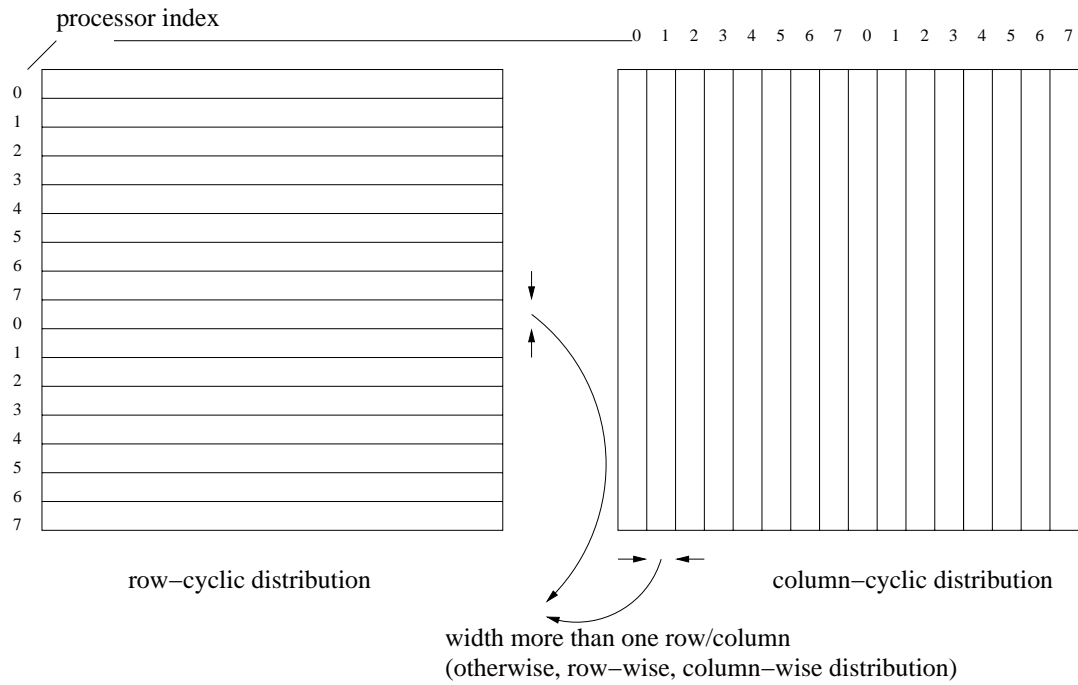
Matrix Distributions : Block cyclic

In block cyclic distributions the rows (similarly for columns) are split into q groups of n/q consecutive rows per group, where potentially $q > p$, and the i -th group is assigned to a processor in a cyclic fashion.

- **column-cyclic** distribution. This is an one-dimensional cyclic distribution. Split matrix into q column stripes so that n/q consecutive columns form the i -th stripe that will be stored in processor $i \% p$. The symbol $\%$ is the mod (remainder of the division) operator. Usually $q > p$. Sometimes the term **wrapped-around column** distribution is used for the case where $n/q = 1$, i.e. $q = n$.
- **row-cyclic** distribution. This is an one-dimensional cyclic distribution. Split matrix into q row stripes so that n/q consecutive rows form the i -th stripe that will be stored in processor $i \% p$. The symbol $\%$ is the mod (remainder of the division) operator. Usually $q > p$. Sometimes the term **wrapped-around row** distribution is used for the case where $n/q = 1$, i.e. $q = n$.
- **scattered** distribution. Let $p = q_i \cdot q_j$ processors be divided into q_j groups each group P_j consisting of q_i processors. Particularly, $P_j = \{jq_i + l \mid 0 \leq l \leq q_i - 1\}$. Processor $jq_i + l$ is called the l -th processor of group P_j . This way matrix element (i, j) , $0 \leq i, j < n$, is assigned to the $(i \bmod q_i)$ -th processor of group $P_{(j \bmod q_j)}$. A scattered distribution refers to the special case $q_i = q_j = \sqrt{p}$.

BSP Algorithms

Block cyclic distributions



BSP Algorithms
Scattered Distribution

0	4	8	12	0	4	8	12
1	5	9	13	1
2	6	10	14	2			
3	7	11	15	3			
0	4	8	12	...			
1	5	9	13				
2	6	10	14				
3	7	11	15				

BSP Algorithms

Matrix Multiplication

The BSP algorithm for matrix multiplication presented below was presented in the seminal work of Valiant. It works for $p \leq n^2$. Each processor is assigned the task of computing an $n/\sqrt{p} \times n/\sqrt{p}$ submatrix of the product $A \times B$. The input matrices A and B are divided into p block-submatrices, each one of dimension $m \times m$, where $m = n/\sqrt{p}$. We call this distribution of the input among the processors *block* distribution. This way, element $A(i, j)$, $0 \leq i < n, 0 \leq j < n$, belongs to the $(j/m) * \sqrt{p} + (i/m)$ -th block that is subsequently assigned to the memory of the same-numbered processor. Let A_i (respectively, B_i) denote the i -th block of A (respectively, B) stored in processor i . With these conventions the algorithm can be described in Figure 1. The following Proposition describes the performance of the aforementioned algorithm.

```
begin MULT_A ( $C, A, B, n, p$ )
1. Let  $m = n/\sqrt{p}$  ;
   Each processor is also assigned a unique processor number  $q$ ;
2. Let  $p_i = q \bmod \sqrt{p}$  ;  $p_j = q/\sqrt{p}$  ;  $C_q = 0$ ;
3.  $a_l \leftarrow A_{p_i+l*\sqrt{p}}$ ,  $0 \leq l < \sqrt{p}$ ;
4.  $b_l \leftarrow B_{p_j*\sqrt{p}+l}$ ,  $0 \leq l < \sqrt{p}$ ;
5. for  $0 \leq l < \sqrt{p}$  do
    $C_q = C_q + a_l \times b_l$ ;
end MULT_A
```

Figure 1: Procedure MULT_A.

Proposition 1 *Algorithm MULT_A for multiplying two $n \times n$ matrices A and B stored according to the block distribution requires, for any $p \leq n^2$, computation time $C_{mul}(n)$ that is given by*

$$C_{mul}(n) = \max \left\{ L, \frac{(2n-1)n^2}{p} \right\},$$

and communication time $M_{mul}(n)$ that is given by the expression

$$M_{mul}(n) = \max \left\{ L, g \frac{2n^2}{\sqrt{p}} \right\}.$$

One immediately realizes that algorithm MULT_A is not memory efficient since it requires more local memory per processor – by a factor of \sqrt{p} – than the required one. Algorithm MULT_B shown in Figure 2 is the memory efficient variant of MULT_A. It is not synchronization efficient though since its number of supersteps is not constant any more; it has been increased by a factor of \sqrt{p} . The performance of algorithm MULT_B is summarized in Proposition 2.

```

begin MULT_B (C,A,B,n,p)
1.  Let  $m = n/\sqrt{p}$  ;
    Each processor is also assigned a unique processor number  $q$ ;
2.  Let  $p_i = q \bmod \sqrt{p}$  ;  $p_j = q/\sqrt{p}$  ;  $C_q = 0$ ;
3.  for  $0 \leq l < \sqrt{p}$  do
    begin
4.     $a \leftarrow A_{((p_i+p_j+l) \bmod \sqrt{p}) * \sqrt{p} + p_i}$ ;
5.     $b \leftarrow B_{((p_i+p_j+l) \bmod \sqrt{p}) + p_j * \sqrt{p}}$ ;
6.     $C_q = C_q + a \times b$ ;
    end
end MULT_B

```

Figure 2: Procedure MULT_B.

Proposition 2 Algorithm MULT_B for multiplying two $n \times n$ matrices A and B stored according to the block distribution requires, for any $p \leq n^2$, computation time $C_{mul}(n)$ that is given by

$$C_{mul}(n) = \sqrt{p} \max \left\{ L, \frac{(2n-1)n^2}{p^{3/2}} \right\}$$

and communication time $M_{mul}(n)$ that is given by the expression

$$M_{mul}(n) = \sqrt{p} \max \left\{ L, g \frac{2n^2}{p} \right\}$$

BSP Algorithms

Experimental Results

In order to show the efficiency of algorithm design on the BSP model we present some experimental results for matrix multiplication on Cray T3D; additional results can be found in the author's Web page. Algorithm MULTT_B is a variation of MULT_B where in order to multiply A with B , matrix A is first transposed and the loop for matrix multiplication is changed accordingly. This way the access patterns for both A and B are the same (column - column as opposed to row - column) thus improving locality (cache usage), and subsequently program performance.

<i>Algorithm MULT_B</i>									
	$p = 1$		$p = 4$		$p = 16$		$p = 64$		
n	Time (sec)	Mfl rate	Time (sec)	Mfl rate	Time (sec)	Mfl rate	Time (sec)	Mfl rate	
256	4.1	7.9	1.1	7.8	0.28	7.4	0.03	13.9	
512	34.0	7.8	8.4	7.9	2.1	7.7	0.56	7.4	
1024	289.8	7.4	68.4	7.8	16.9	7.9	4.3	7.7	
2048	-	-	-	-	136.8	7.8	33.8	7.9	

Table 1: Execution time for MULT_B on the Cray T3D

<i>Algorithm MULTT_B</i>									
	$p = 1$		$p = 4$		$p = 16$		$p = 64$		
n	Time (sec)	Mfl rate	Time (sec)	Mfl rate	Time (sec)	Mfl rate	Time (sec)	Mfl rate	
256	2.3	14.3	0.58	14.4	0.15	13.7	0.03	15.1	
512	20.7	12.9	4.7	14.1	1.16	14.4	0.30	13.5	
1024	202.7	10.5	41.7	12.8	9.4	14.1	2.3	14.3	
2048	-	-	-	-	83.5	12.8	19.0	14.1	

Table 2: Execution time for MULTT_B on the Cray T3D

Matrix Multiplication: Algorithm C

Finally, we outline a matrix multiplication algorithm that is computation, communication and synchronization efficient. It fails, however, to be memory efficient, as its memory requirements are a multiplicative factor $p^{1/3}$ from the optimal. Algorithm MULT_C is outlined in the remainder of this section.

In MULT_C matrices A and B (and the result C) are split two ways into submatrices. Each matrix (A , B and the result C) is split into p “physical” block-submatrices, as in the previous algorithms, each of size $n/p^{1/2} \times n/p^{1/2}$. A “physical” block-submatrix indicates the part of the matrix stored in a single physical (processor) location (i.e. block-submatrix A_i is stored in processor i). At the same time, each of the three matrices is split into $p^{2/3}$ “virtual” block-submatrices each of size $n/p^{1/3} \times n/p^{1/3}$. A “virtual” block-submatrix indicates the block geometry that will be used in the matrix multiplication algorithm to be outlined below. The elements of a “virtual” block-submatrix may be stored in more than one physical processors.

Whereas in the first two algorithms “physical” and “virtual” block-submatrices coincided in number and dimension, in this communication efficient algorithm are clearly distinguished.

Let the “virtual” block-submatrices be identified as A_{ij} , B_{ij} and C_{ij} . Matrix multiplication will thus require the computation of all $C_{ij} = \sum_{k=1}^{p^{1/3}} C_{ijk} = \sum_{k=1}^{p^{1/3}} A_{ik}B_{kj}$, where $C_{ijk} = A_{ik}B_{kj}$.

BSP Algorithms

Algorithm C: Description

The algorithm consists of the following steps. We name the processors (i, j, k) the way we did in the matrix multiplication algorithm on the hypercubic networks.

Step 1. Processor (i, j, k) gets A_{ik} and B_{kj} . Note that each of these two “virtual” block-submatrices may originate from more than one processors. Each processor sends at most $2n^2/p$ elements (but each one replicated $p^{1/3}$ times) and receives at most $2n^2/p^{2/3}$ elements. The communication cost of Step 1 is $\max\{L, 2gn^2/p^{2/3}\}$. Subsequently, the two submatrices are multiplied as in the sequential case a step requiring at most $\max\{L, 2n^3/p\}$ time. Partial-submatrix C_{ijk} is thus computed on processor (i, j, k) . Each element of such a submatrix is a partial sum of an element c_{lm} of the result matrix C .

Step 2. Each element of C_{ijk} is transmitted from (i, j, k) to that physical processor that stores the “physical” block-submatrix of C whose elements will be formed as sums of the receiving elements (partial sums) of C_{ijk} . Note that each (i, j, k) processor may send its elements to more than one physical processors. At the completion of this step, each of the p processors storing a block-submatrix of C of dimension $n/p^{1/2} \times n/p^{1/2}$ receives at most $p^{1/3} \cdot n^2/p$ such elements (partial sums). The complex communication performed in this step requires time $\max\{L, gn^2/p^{2/3}\}$.

Step 3. The received partial sums are added. $p^{1/3}$ partial sums are summed to give an element of C stored at a physical processor, for a total of n^2/p such elements (of a “physical” block-submatrix). The total computation time performed is $\max\{L, n^2/p^{2/3}\}$.

Proposition 3 *Algorithm MULT_C for multiplying two $n \times n$ matrices A and B stored according to the block distribution requires, for any $p \leq n^2$, computation time $C_{mul}(n)$ that is given by*

$$C_{mul}(n) \leq \max\{L, 2n^3/p\} + \max\{L, n^2/p^{2/3}\},$$

and communication time $M_{mul}(n)$ that is given by the expression

$$M_{mul}(n) = \max\left\{L, 2g\frac{n^2}{p^{2/3}}\right\} + \max\left\{L, g\frac{n^2}{p^{2/3}}\right\}.$$

BSP Algorithms

Algorithm C: Optimality

The optimality in communication of the algorithm is established by the following result.

Theorem 1 *On a model of computation that allows the operations $\{+, *\}$ only, if any processor reads s elements of A and B and computes at most s partial sums of C , it can compute at most $O(s^{3/2})$ multiplicative terms for these sums.*

This way, if a processor reads at most s elements of A and B it can compute at most $O(s^{3/2})$ multiplicative terms of C . Combined, all p processors can compute $p O(s^{3/2})$ such terms which must be $\Omega(n^3)$. Therefore $s = \Omega(n^2/p^{2/3})$ and thus algorithm MULT_C is communication optimal.

How can one prove the Theorem? It suffices to show that if A has s 1's in arbitrary position (all other positions are 0) and so has B , then the product $C = A \times B$ requires at most $O(s^{3/2})$ non-trivial multiplications (i.e. multiplications where both terms are non-zero). This can be proved by considering two sets of rows for A , the small ones having at most \sqrt{s} ones on each such row and the large one. Call A_s the submatrix of A formed by these small rows, and A_l the submatrix consisting of the large rows having at least \sqrt{s} ones. We can only have at most \sqrt{s} rows in A_l . Consider $A_l \times B$. The results can contribute at most $s^{3/2}$ non-trivial multiplications in C . The reason for that is that each row of A_l can contribute s ones when multiplied with B since B has only s ones. There are at most \sqrt{s} in A_l , i.e. claim follows. For $A_s \times B$ just note that each row of A_s has at most \sqrt{s} ones. Since only s terms are computed in C , and a term involves a row of A_s that has at most \sqrt{s} elements, this $A_s \times B$ product can only involve at most $s^{3/2}$ multiplications.