

Processor Allocation and Task Scheduling of Matrix Chain Products on Parallel Systems

Heejo Lee, *Member, IEEE*, Jong Kim, *Member, IEEE*,
Sung Je Hong, *Member, IEEE*, and Sunggu Lee, *Member, IEEE*

Abstract—The problem of finding an optimal product sequence for sequential multiplication of a chain of matrices (the matrix chain ordering problem, MCOP) is well-known and has been studied for a long time. In this paper, we consider the problem of finding an optimal product schedule for evaluating a chain of matrix products on a parallel computer (the matrix chain scheduling problem, MCSP). The difference between the MCSP and the MCOP is that the MCOP pertains to a product sequence for single processor systems and the MCSP pertains to a sequence of concurrent matrix products for parallel systems. The approach of parallelizing each matrix product after finding an optimal product sequence for single processor systems does not always guarantee the minimum evaluation time on parallel systems since each parallelized matrix product may use processors inefficiently. We introduce a new processor scheduling algorithm for the MCSP which reduces the evaluation time of a chain of matrix products on a parallel computer, even at the expense of a slight increase in the total number of operations. Given a chain of n matrices and a matrix product utilizing at most P/k processors in a P -processor system, the proposed algorithm approaches $k(n-1)/(n+k\log(k)-k)$ times the performance of parallel evaluation using the optimal sequence found for the MCOP. Also, experiments performed on a Fujitsu AP1000 multicomputer show that the proposed algorithm significantly decreases the time required to evaluate a chain of matrix products in parallel systems.

Index Terms—Matrix chain product, parallel matrix multiplication, matrix chain scheduling problem, processor allocation, task scheduling.

1 INTRODUCTION

1.1 Background

MATRIX multiplication is a computation intensive part of many commonly used scientific computing applications. Many algorithms in numerical and nonnumerical problems are solved efficiently using matrix-matrix multiplications [1]. Also, in the case of parallel algorithms, some problems that are not solved efficiently using sequential algorithms turn out to be solved efficiently using matrix-matrix multiplications or variations of matrix-matrix multiplication [2]. Such an application uses matrix multiplication as a basic computational kernel so that a chain of matrices is frequently required and the matrices are multiplied consecutively [2], [3], [4], [5].

In the evaluation of a chain of matrix products with n matrices, $\mathcal{M} = M_1 \times M_2 \times \cdots \times M_n$, where M_i is an $m_i \times m_{i+1}$ ($m_i \geq 1$) matrix, the product sequence greatly affects the total number of operations required to evaluate \mathcal{M} , even though the final result is the same for all product sequences due to the associative law of matrix multiplication. In the worst case, an arbitrary product sequence of matrices is $\Omega(T_{opt}^3)$, where T_{opt} is the minimum number of operations

required to evaluate a chain of matrix products [6]. The matrix chain ordering problem (MCOP) focuses on finding a product sequence for a set of matrices such that the total number of operations is minimized.

An exhaustive search to find an optimal solution for the MCOP is not a good strategy since the number of possible product sequences of a chain of matrix products with n matrices is $\Theta(4^n/n^{3/2})$, which is known as the *Catalan number* [7]. Thus, determining the optimal sequence by this method is very time consuming. There has been much work reported for solving the MCOP. The MCOP was first studied by Godbole [8] and solved using dynamic programming in $O(n^3)$ time. Chin [9] suggested an approximation algorithm, which runs in $O(n)$ time and finds a near-optimal sequence. An optimal sequential algorithm, which runs in $O(n \log(n))$ time, was given by Hu and Shing [10], [11]. This algorithm solves the MCOP by solving the equivalent problem of finding the optimal triangulation of a convex polygon. Ramanan [12] presented a simpler algorithm for the MCOP and obtained the tight lower bound of $\Omega(n \log(n))$ for a related problem.

1.2 Parallel Evaluation Time

Let us refer to the time required to find an optimal product sequence for a chain of matrices as the *ordering time* and the time required to execute the product sequence as the *evaluation time*. Many parallel algorithms aimed at reducing the *ordering time* have been studied using the dynamic programming method [13], [14], [15], [16] and the convex polygon triangulation method [17], [18], [19]. Bradford et al. [16] proposed a parallel algorithm based on dynamic programming, which runs in $O(\log^3(n))$ time with

- H. Lee is with Ahnlab, Inc., 8F V-Valley Bldg., 724 Suseo-dong, Gangnam-gu, Seoul 135-744, South Korea. E-mail: heejo@ahnlab.com.
- J. Kim and S.J. Hong are with the Department of Computer Science and Engineering, Pohang University of Science and Technology, Pohang 790-784, South Korea. E-mail: {jkim, sjhong}@postech.ac.kr.
- S. Lee is with the Department of Electrical Engineering, Pohang University of Science and Technology, Pohang 790-784, South Korea. E-mail: slee@postech.ac.kr.

Manuscript received 2 June 1999; revised 3 Apr. 2002; accepted 14 Oct. 2002. For information on obtaining reprints of this article, please send e-mail to: tpd@computer.org, and reference IEEECS Log Number 109987.

$n/\log(n)$ processors on the CRCW PRAM model. Czumaj [17] proposed an $O(\log^3(n))$ time algorithm based on the triangulation of a convex polygon, which runs with $n^2/\log^3(n)$ processors on the CREW PRAM. Also, in [18], he proposed a faster approximation algorithm that finds a near-optimal solution in $O(\log(n))$ time on the CREW PRAM and in $O(\log \log(n))$ time on the CRCW PRAM. Ramanan [19] gives an optimal algorithm that runs in $O(\log^4(n))$ time using n processors on a CREW PRAM.

Now, let us consider the *evaluation time* of a chain of matrix products. In a single processor system, the evaluation of a chain of matrix products using the optimal product sequence for the MCOP guarantees the minimum evaluation time since the sequence guarantees the minimum number of operations. However, in parallel systems, parallel computation of each matrix product using the product sequence found for the minimum number of operations does not guarantee the minimum evaluation time. This is because the evaluation time in parallel systems is affected by various factors such as dependencies among tasks, communication delays, and processor efficiency. To this date, there has been no research reported in the open literature on ways to reduce the matrix chain evaluation time in a parallel system.

1.3 New Contributions

We introduce the problem of parallel evaluation of matrix chain products. Our technical contributions are two-fold. First, we formally present the problem of finding the matrix product schedule for parallel systems (the matrix chain scheduling problem, MCSP) in order to reduce the evaluation time. The MCSP is proven to be an NP-complete problem.

Second, we propose a novel scheduling algorithm which finds a matrix product schedule that, while possibly increasing the total number of required operations, decreases the evaluation time of a chain of matrix products by finding sets of matrix products that can be executed concurrently. We show that the proposed algorithm can significantly reduce the time required to evaluate a chain of matrix products in parallel systems.

This paper is organized as follows: Section 2 presents a formal description of the processor scheduling problem for a chain of matrix products and shows that the given problem is NP-complete. In Section 3, we present a processor allocation method for multiplying a number of independent matrix products concurrently. In Section 4, we propose a matrix chain scheduling algorithm that dramatically reduces the evaluation time of a chain of matrix products by using processors efficiently in parallel systems. In Section 5, we analyze the evaluation performance of sequences found with the proposed method and sequences found for the MCOP. We also compare the proposed method with various other evaluation methods through experiments on a Fujitsu AP1000 parallel system. In Section 6, many practical issues and further extensions of the MCSP are discussed. Finally, in Section 7, we summarize and conclude the paper.

2 MATRIX CHAIN SCHEDULING PROBLEM

2.1 Notation

- P : the number of processors in a parallel system.
- \mathcal{M} : a matrix chain product of n matrices, i.e., $\mathcal{M} = M_1 \times M_2 \times \cdots \times M_n$.
- M_i : an $m_i \times m_{i+1}$ matrix ($m_i \geq 1, 1 \leq i \leq n$).
- L : a product sequence tree for a matrix chain \mathcal{M} .
- $L_{i,j}$: the sequence subtree of L for $(M_i \times \cdots \times M_j)$.
- C : the minimum number of computations for evaluating \mathcal{M} .
- ΔC : the amount of increased computation incurred by modifying the current sequence tree.
- $p_{i,j}$: the number of processors assigned for evaluating $(M_i \times \cdots \times M_j)$.
- $T_{i,j}(p_{i,j})$: the execution time for evaluating $(M_i \times \cdots \times M_j)$ on $p_{i,j}$ processors.
- (m_i, m_j, m_k) : a single matrix product for multiplying an $m_i \times m_j$ matrix by an $m_j \times m_k$ matrix.
- $\Phi(m_i, m_j, m_k, p)$: the execution time of a single matrix product (m_i, m_j, m_k) using p processors.
- $D(x)$: the set of divisors of x , i.e., $D(x) = \{d \mid d \text{ divides } x\}$.
- $LD(x, y)$: the largest divisor in $D(x)$ that is not larger than y .
- $SD(x, y)$: if $x > y$, then $SD(x, y)$ is the smallest divisor in $D(x)$ that is larger than y . Otherwise, $SD(x, y) = x$.
- m : the largest dimension among all of the matrices, i.e., $m = \max_{1 \leq i \leq n+1} (m_i)$.

2.2 Problem Description

We consider the problem of finding the schedule with minimum evaluation time for \mathcal{M} on a P processor parallel system. The number of operations needed to multiply a matrix A of size $m_i \times m_j$ by a matrix B of size $m_j \times m_k$ is $m_i m_j m_k$.¹ Many parallel algorithms for matrix multiplication have been developed for various parallel architectures [23], [24], and the execution time of matrix multiplication depends on the algorithm used and the architecture on which the algorithm runs. However, for a broader discussion, we assume that a simple parallel algorithm [24] is used and the execution time of matrix multiplication is proportional to the number of operations on a processor. Therefore, for multiplying A by B using p processors, the execution time, $\Phi(m_i, m_j, m_k, p)$, can be approximated as follows:

$$\Phi(m_i, m_j, m_k, p) \approx \begin{cases} m_i m_j m_k & \text{if } 1 \leq p \leq \frac{m_i m_j m_k}{\log(m_j)}, \\ \frac{m_i m_j m_k}{p} \log\left(\frac{p}{m_i m_k}\right) & \text{if } \frac{m_i m_j m_k}{\log(m_j)} < p \leq m_i m_j m_k. \end{cases}$$

When p_{ij} processors are allocated for evaluating $(M_i \times \cdots \times M_j)$, the evaluation time consists of two parts: the partial matrix chain evaluation time and the single matrix

1. Even though Strassen's algorithm [20], [21] and its variants perform fewer than n^3 operations for $n \times n$ matrix multiplication, these faster algorithms are regarded as inappropriate methods for matrix chain products due to more erroneous results caused by round-off errors and larger storage requirements than the usual inner-product type algorithm [22]. Therefore, we assume that a simple algorithm is used so that n^3 operations are required for $n \times n$ matrix multiplication. This is also the assumption made in other research work on the MCOP.

product execution time. The two partial matrix chains are $(M_i \cdots M_k)$ and $(M_{k+1} \cdots M_j)$ for any k ($i \leq k < j$). The evaluation time of the two matrix product chains is dependent on the evaluation method. One method is to evaluate sequentially $(M_i \cdots M_k)$ and $(M_{k+1} \cdots M_j)$ using all available processors in p_{ij} . The other method is to evaluate both $(M_i \cdots M_k)$ and $(M_{k+1} \cdots M_j)$ concurrently by partitioning p_{ij} into $p_{i,k}$ and $p_{k+1,j}$ such that $p_{i,k} + p_{k+1,j} \leq p_{ij}$. The minimum evaluation time $T_{i,j}(p_{i,j})$ of $(M_i \cdots M_j)$ on $p_{i,j}$ processors is found from the following recurrence relation:

$$T_{i,j}(p_{i,j}) \approx \min_{i \leq k < j} \left(T_{i,k}(p_{i,k}) + T_{k+1,j}(p_{k+1,j}) + \Phi(m_i, m_{k+1}, m_{j+1}, p_{i,j}), \right. \\ \left. \min_{p_{i,k} + p_{k+1,j} \leq p_{i,j}} \left(\max \left(T_{i,k}(p_{i,k}), T_{k+1,j}(p_{k+1,j}) \right) \right) + \Phi(m_i, m_{k+1}, m_{j+1}, p_{i,j}) \right).$$

The problem of finding the schedule that results in the minimum evaluation time, $T_{1,n}(P)$, is equivalent to finding the best schedule, $k_{i,j}$, for $(M_i \cdots M_j)$ with the processor allocation p_{ij} to $L_{i,j}$. Therefore, the MCSP is defined as follows:

MCSP: Given M and P , find the product sequence for evaluating M and the processor schedule for the sequence such that the evaluation time is minimized.

2.3 MCSP Complexity

The complexity of the MCSP depends on the number of processors available for the MCSP. Consider the case in which there are sufficient processors for multiplying any number of matrices concurrently. For a matrix product (m_i, m_{k+1}, m_j) , we can allocate the number of processors that guarantees the minimum execution time for the product, i.e., $\log(m_{k+1})$. Then, using dynamic programming, the problem can be solved in polynomial time according to the following recurrence relation:

$$T_{i,j} = \begin{cases} \min_{i \leq k < j} \left(\max(T_{i,k}, T_{k+1,j}) + \log(m_{k+1}) \right) & \text{if } 1 \leq i < j \leq n, \\ 0 & \text{if } i = j, 1 \leq i \leq n. \end{cases}$$

Therefore, in the case of an unlimited number of processors, the problem of finding the schedule for evaluating M with the minimum time has a polynomial time algorithm. However, in general, the number of available processors is limited and not sufficient to allocate the maximum number of processors for each product.

Now, we show that the MCSP is NP-complete using a reduction from the processor partitioning problem which is known to be NP-complete [25].

Theorem 1. *The MCSP is NP-complete.*

Proof. The decision version of the MCSP is obviously in NP. There is a nondeterministic algorithm that generates a processor schedule with a product sequence of M . Given the schedule, it can be decided in polynomial time

whether the schedule length is less than a certain value by finding the longest path in a tree graph.

The processor partitioning problem [25], denoted by *PPP*, is to find the schedule with the minimum completion time of n tasks on a partitionable P processor system ($n \leq P$). In the *PPP*, a task j is characterized as δ_j and $t_j(p_j)$, where δ_j denotes the maximum number of allocable processors for task j and $t_j(p_j)$ denotes the execution time of task j on p_j processors ($1 \leq p_j \leq \delta_j$). The *PPP* of deciding whether there exists a schedule whose completion time is less than D is NP-complete [25].

An instance of the *PPP* can be transformed to an instance of the MCSP. For the *PPP*, $t_j(p_j)$ can be a linear function of p_j , i.e., $t_j(p_j) = t_j(1)/p_j$. Let M_{2j-1} be a $1 \times t_j(1)$ matrix and M_{2j} be a $t_j(1) \times 1$ matrix for $1 \leq j \leq n$. Now, define Φ for the MCSP as follows:

$$\Phi(1, t_j(1), 1, p_j) = t_j(p_j) \quad \text{for } 1 \leq p_j \leq \delta_j.$$

Thus, the *PPP* with n tasks is transformed to an MCSP with $2n$ matrices in polynomial time.

Next, let us show that, if there is a feasible solution for the MCSP, there exists a solution for the *PPP*. In the transformed MCSP, an optimal sequence for the MCSP has the products $(M_1 \times M_2), (M_3 \times M_4), \dots, (M_{2n-1} \times M_{2n})$ because these n products satisfy the optimal sequence property in Theorem 1 of [9] and because any other sequence not only reduces the degree of concurrency, but also increases the computation. Therefore, an optimal schedule for the transformed MCSP is found in a product sequence tree having these independent n products as leaves. Since execution of $n-1$ nonleaf products with single operations requires at least $\lceil \log(n-1) \rceil$ time, finding a schedule of the transformed MCSP whose completion time is less than $D + \lceil \log(n-1) \rceil$ solves the *PPP* whose schedule length is less than D . Noting that the *PPP* is a special case of the MCSP completes the proof. \square

Since the problem of finding an optimal schedule for the MCSP is an NP-hard optimization problem, we propose a heuristic algorithm in Section 4. The algorithm enhances the evaluation performance of an n -matrix product chain on a parallel system by partitioning the parallel system and executing multiple matrix products simultaneously; this also enhances the overall efficiency of the system. A processor allocation method for executing multiple matrix products is discussed in the next section.

3 PROCESSOR ALLOCATION FOR MATRIX PRODUCTS

In order to devise an efficient matrix chain scheduling algorithm, we need to study the characteristics of matrix products. In this section, we discuss how many processors should be allocated for a single matrix product and determine the optimal processor allocation for executing independent matrix products simultaneously.

3.1 Processor Allocation for a Single Matrix Product

Many parallel matrix multiplication algorithms have been developed with various parallel architectures [21], [24]. A multiplication of two $n \times n$ matrices requires at most n^3 processors. In the best case, it takes $O(\log(n))$ time with n^3 processors by using n processors to get one element of the result matrix.² Each set of n processors executes a multiplication in one step and sums n elements within $O(\log(n))$ steps. However, in this case, we expect low utilization of processors. While summing n data for $\log(n)$ steps, some processors stay idle. Moreover, these $\log(n)$ steps are communication steps, not computation steps. Then, how many processors should be allocated for multiplying two matrices?

One distinctive feature of parallel matrix multiplication is the jerky behavior affecting the execution time, speedup, and efficiency of the computation. This is mainly due to the load imbalance, which is referred to as the “integer effect” in [27]. It is known that many parallel algorithms are based on allocating processors for computing a part, especially as a subblock, of the result matrix [5]. For executing a product (m_i, m_j, m_k) on p processors, let us consider the case where the operations for computing $m_i m_k$ elements are distributed among p processors. Then, the execution time $\Phi(m_i, m_j, m_k, p)$, the speedup $S(p)$, and the efficiency $E(p)$ are estimated as follows:

$$\Phi(m_i, m_j, m_k, p) = \left\lceil \frac{m_i m_k}{p} \right\rceil m_j,$$

$$S(p) = \frac{\Phi(m_i, m_j, m_k, 1)}{\Phi(m_i, m_j, m_k, p)} = m_i m_k / \left\lceil \frac{m_i m_k}{p} \right\rceil, \text{ and}$$

$$E(p) = \frac{S(p)}{p} = m_i m_k / \left(p \left\lceil \frac{m_i m_k}{p} \right\rceil \right).$$

Ideally, the speedup is expected to be $S(p) = p$ so that $E(p) = 1$. However, in the worst case, the integer effect causes a speedup $S(p)$ close to $p/2$, which is a factor of 2 away from what we would normally hope for. This means that half of the processors do not have any effect on the execution time. Therefore, for efficient execution, the number of processors allocated to a single matrix product should not be an arbitrary number. When multiplying an $m_i \times m_j$ matrix by an $m_j \times m_k$ matrix, the proper number of processors is one of divisors of $m_i m_k$. For multiplying a 2×4 matrix by a 4×3 matrix, the number of allocable processors could be 1, 2, 3, or 6. In this case, the load of each processor can be well balanced.

Let $D(x)$ denote the set of divisors of x , i.e., $D(x) = \{d | d \text{ divides } x\}$. The number of processors allocated to a given matrix product (m_i, m_j, m_k) should be an element in $D(m_i m_k)$. Thus,

2. On a CRCW PRAM, the matrix multiplication runs in constant time with the assumption that, when there are write conflicts by multiple processors attempting to write numbers in the same location, the sum of the numbers is written in that location [26]. Since the CRCW PRAM with an arbitrary number of processors is not realistic in practice, we are not considering that case.

$$\Phi(m_i, m_j, m_k, p) = \frac{m_i m_j m_k}{p} \text{ for } p \in D(m_i m_k).$$

Even though there may exist an algorithm using $m_i m_j m_k$ processors, we allocate at most $m_i m_k$ processors for processor efficiency.

3.2 Processor Allocation for the Concurrent Computation of Multiple Matrix Products

In this section, we discuss a processor allocation method for independently computing multiple matrix products. When executing multiple parallel tasks concurrently, one good heuristic for allocating processors to each task is “proportional allocation” [28], [29]. The proportional allocation algorithm allocates a number of processors proportional to the computation amount of each task. This algorithm tries to minimize the completion time of all tasks by balancing the execution times of the tasks. However, it assumes that the execution time of a task decreases if more processors are allocated to it and that any number of processors may be allocated. These assumptions however, do not hold for the MCSP.

Applying proportional allocation to matrix products is inappropriate. Let us consider the case of computing two matrix products $(2, 3, 8)$ and $(4, 4, 3)$ on 20 processors. Since both products involve the same amount of work, i.e., $2 \times 3 \times 8 = 4 \times 4 \times 3$, 10 processors are allocated for each product by proportional allocation. Then, the completion time of the two matrix products is $\max(2 \times 3 \times 8/8, 4 \times 4 \times 3/6) = 8$ units of time since $LD(16, 10) = 8$ and $LD(12, 10) = 6$. However, if we allocate eight processors to the first product and 12 processors to the second product, the completion time is only $\max(6, 4) = 6$ units of time. From the fact that the completion time is bounded by the longer execution time, we can reduce the completion time by allocating unused processors to the matrix product that requires a longer execution time.

The following algorithm describes the processor allocation algorithm for two independent matrix products. Given two matrix products $X = (m_x, m_{x+1}, m_{x+2})$ and $Y = (m_y, m_{y+1}, m_{y+2})$, we let $\Phi(X, p)$ and $\Phi(Y, p)$ be shorthand notation for $\Phi(m_x, m_{x+1}, m_{x+2}, p)$ and $\Phi(m_y, m_{y+1}, m_{y+2}, p)$, respectively.

Discrete Processor Allocation for Two Matrix Products (DPA)

Input: Two matrix products $X = (m_x, m_{x+1}, m_{x+2})$ and $Y = (m_y, m_{y+1}, m_{y+2})$ and a set of P processors.

Output: The number of processors allocated to the matrix products X and Y , denoted as P_x and P_y , which satisfy $1 \leq P_x, P_y \leq P$ and $P_x + P_y \leq P$.

1. $P_{prop} = P \cdot m_x m_{x+1} m_{x+2} / (m_x m_{x+1} m_{x+2} + m_y m_{y+1} m_{y+2})$
2. $d_i = LD(m_x m_{x+2}, P_{prop})$
3. $d_{i+1} = SD(m_x m_{x+2}, P_{prop})$
4. $d_j = LD(m_y m_{y+2}, P - P_{prop})$
5. $d_{j+1} = SD(m_y m_{y+2}, P - P_{prop})$
6. if $\Phi(X, d_i) \geq \Phi(Y, d_j)$ then
7. if $\max(\Phi(X, d_{i+1}), \Phi(Y, LD(P - d_{i+1}))) < \Phi(X, d_i)$ then
8. $P_x = d_{i+1}, P_y = LD(m_y m_{y+2}, P - d_{i+1})$

9. else
10. $P_x = d_i, P_y = LD(m_y m_{y+2}, P - d_i)$
11. endif
12. elseif $\max(\Phi(X, LD(P - d_{j+1})), \Phi(Y, d_{j+1})) < \Phi(Y, d_j)$
then
13. $P_x = LD(m_x m_{x+2}, P - d_{j+1}), P_y = d_{j+1}$
14. else
15. $P_x = LD(m_x m_{x+2}, P - d_j), P_y = d_j$
16. endif
17. endif

Even if we use a naive search algorithm for finding divisors, it will take $O(\max(m_i, m_j))$ time for finding both $LD(m_i m_j, p)$ and $SD(m_i m_j, p)$. We let m denote the largest dimension among all of the matrices, i.e., $m = \max_{1 \leq i \leq n+1} (m_i)$. Then, the time complexity of the discrete processor allocation (DPA) algorithm is $O(m)$. The following lemma and theorem show that DPA guarantees the minimum completion time for two matrix products.

Lemma 1. *Given two matrix products $X = (m_x, m_{x+1}, m_{x+2})$ and $Y = (m_y, m_{y+1}, m_{y+2})$ on P processors, an optimal processor allocation has at least one of four assignments: $P_x = d_i, P_x = d_{i+1}, P_y = d_j, \text{ or } P_y = d_{j+1}$.*

Proof. If an allocation does not have any of the above four assignments, then $P_x < d_i$ or $P_x > d_{i+1}$, and $P_y < d_j$ or $P_y > d_{j+1}$. In this allocation, P_x and P_y are of the following four cases: The first case is $P_x < d_i$ and $P_y < d_j$, the second is $P_x < d_i$ and $P_y > d_{j+1}$, the third is $P_x > d_{i+1}$ and $P_y < d_j$, and the fourth is $P_x > d_{i+1}$ and $P_y > d_{j+1}$. Without loss of generality, let us assume $\Phi(X, d_i) \geq \Phi(Y, d_j)$.

Case 1. $P_x < d_i$ and $P_y < d_j$ Since $\max(\Phi(X, P_x), \Phi(Y, P_y)) > \Phi(X, d_i)$, the allocation of P_x and P_y does not guarantee the minimum completion time of X and Y . This means that the completion time by this allocation is longer than that by the allocation of d_i and d_j to X and Y .

Case 2. $P_x < d_i$ and $P_y > d_{j+1}$ As in Case 1, since $\max(\Phi(X, P_x), \Phi(Y, P_y)) = \Phi(X, P_x) > \Phi(X, d_i)$, the allocation of P_x and P_y does not result in the minimum completion time.

Case 3. $P_x > d_{i+1}$ and $P_y < d_j$ In this case, the completion time with the allocation of P_x and P_y is $\max(\Phi(X, P_x), \Phi(Y, P_y)) = \Phi(Y, P_y)$. Let us compare this allocation with $P'_x = d_{i+1}, P'_y = LD(Y, P - P'_x)$. In the allocation of P'_x and P'_y , the completion time is $\max(\Phi(X, P'_x), \Phi(Y, P'_y)) = \Phi(Y, P'_y)$. Also, since $P'_x < P_x$, it is the case that $P'_y \geq P_y$. In the case of $P'_y > P_y$, the completion time $\Phi(Y, P'_y)$ with the allocation of P'_x and P'_y is shorter than the completion time $\Phi(Y, P_y)$ with the allocation of P_x and P_y so that the allocation of P_x and P_y does not guarantee the minimum completion time of X and Y . In the case of $P'_y = P_y$, the allocation of P_x and P_y does not take less time than the allocation of P'_x and P'_y , and the optimal allocation with the minimum completion time can be found with $P'_x = d_{i+1}$.

Case 4. $P_x > d_{i+1}$ and $P_y > d_{j+1}$ In this case, we cannot allocate P_x and P_y since $P_x + P_y > P$.

For all four cases, the allocation of P_x and P_y does not take less time than allocation with one of the four assignments. Therefore, the optimal allocation with the

minimum completion time is found with one of the four assignments. \square

Theorem 2. *DPA guarantees the minimum completion time for two matrix products.*

Proof. Since DPA finds the allocation with the minimum completion time among the allocations having one of the four assignments in Lemma 1, DPA guarantees the minimum completion time for computing two matrix products simultaneously. \square

In the next section, DPA is used for partitioning the processors allocated to a leaf product into two groups to run another candidate product simultaneously. Thus, we can compute multiple matrix products independently with an increased degree of concurrency.

4 MATRIX CHAIN SCHEDULING ALGORITHM

The proposed scheduling algorithm consists of three stages. First, the algorithm finds the optimal product sequence for the MCOP. Next is the top-down processor assignment stage. In this stage, processors are partitioned and proportionally assigned to each subtree according to their computation amount in order to balance the evaluation time of both the left and right partial matrix product chains. The third stage is the bottom-up execution stage that executes products independently from the leaf and tries to modify the product sequence to enhance concurrency so as to reduce the evaluation time of \mathcal{M} . This is done by finding the points that change the product sequence but do not excessively increase the total number of operations. After changing the product sequence, the processors allocated to the leaf product are relocated using the DPA.

4.1 Optimal Product Sequence by the MCOP

The product sequence of \mathcal{M} determines the number of operations to be executed in single processor systems. In parallel systems, the number of operations is not the sole factor determining the evaluation time, but it affects the evaluation time greatly nonetheless. Hence, our scheduling algorithm begins with the optimal product sequence found for the MCOP. There have been many works reported for finding the optimal product sequence that guarantees the minimum number of operations for any chain of matrix products. The optimal product sequence can be found in $O(n \log(n))$ time using a sequential algorithm [10], [11]. In addition, many parallel algorithms that run in polylog time have been studied [16], [17], [19]. Thus, by using these parallel algorithms, it is possible to find the optimal product sequence within polylog time on P processor systems.

Let us assume that the sequence and the number of operations found for the MCOP is stored in two tables named $S[n, n]$ and $W[n, n]$, respectively. $W[i, j]$ has the minimum number of operations for evaluating $L_{i,j}$ and $S[i, j]$ has the matrix index for partitioning the matrix chain $(M_i \times \cdots \times M_j)$. Note that the algorithm for the MCOP may not have computed $L_{i,j}, S[i, j]$ or $W[i, j]$ for all i, j .

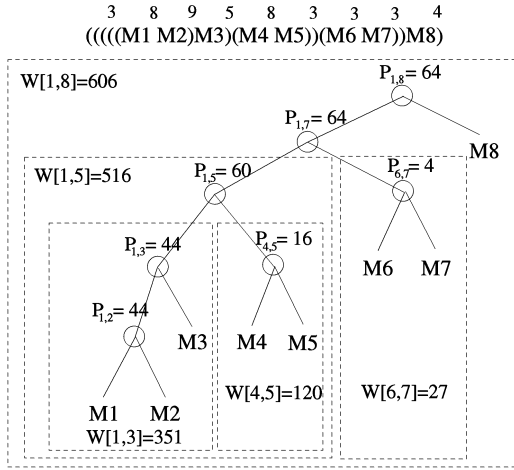


Fig. 1. Top-down processor assignment.

4.2 Top-Down Processor Assignment

In the top-down processor assignment stage, the number of processors assigned to two partial matrix chains is set to be proportional to the computation amount of a subtree. If $p_{i,j}$ processors have been assigned to $L_{i,j}$, then $p_{i,j} \times \frac{W[i,S[i,j]]}{(W[i,S[i,j]]+W[S[i,j]+1,j])}$ processors are assigned to subtree $L_{i,S[i,j]}$ and $p_{i,j} \times \frac{W[S[i,j]+1,j]}{(W[i,S[i,j]]+W[S[i,j]+1,j])}$ processors are assigned to subtree $L_{S[i,j]+1,j}$.

For example, given a chain of eight matrices with dimensions $\{3, 8, 9, 5, 8, 3, 3, 3, 4\}$ and a 64 processor system, processors are assigned as in Fig. 1.

4.3 Bottom-Up Concurrent Execution

After assigning processors to each subtree, the matrix products are executed concurrently and independently, starting from the leaf products. However, there are cases in which some processors stay idle. When there are idle processors in the execution of $L_{i,j}$, we try to modify the product sequence to use these idle processors by tracing the ancestors of the leaf node of $L_{i,j}$ in order to find a candidate for concurrent execution.

This upward trace continues until a suitable candidate or a sibling which is not a leaf node is found. For example, let us consider executing the sequence tree $L_{1,9}$, shown in Fig. 2, which represents

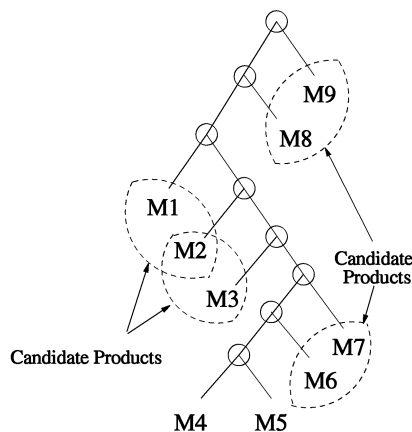


Fig. 2. Candidate products on a sequence tree.

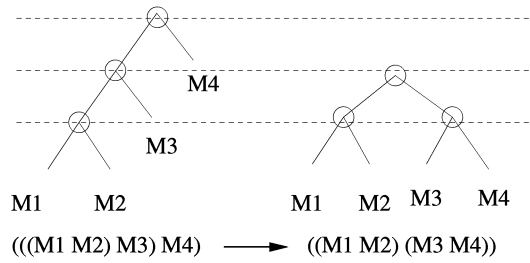


Fig. 3. A sequence modification.

$$(((M_1(M_2(M_3(((M_4M_5)M_6)M_7))))M_8)M_9).$$

In the execution of (M_4M_5) , the possible candidates for concurrent execution are (M_1M_2) , (M_2M_3) , (M_6M_7) , and (M_8M_9) . There are other types of candidate products, e.g., (M_7M_8) , which are not considered in this paper because such cases result in more modifications to the optimal sequence with no obvious benefit over other candidates.

The matrix product (M_yM_{y+1}) is a candidate of the leaf product (M_xM_{x+1}) if one of the following two conditions is satisfied.

Condition-1: $y > x + 1$ and the node associated with M_{y+1} has the left child node associated with M_y in the sequence tree;

Condition-2: $y < x - 1$ and the node associated with M_y has the right child node associated with M_{y+1} in the sequence tree.

When we modify the product sequence to execute candidate products simultaneously in the current execution phase, there is some loss due to an increase in the total number of operations. Therefore, we have to check whether the modification is beneficial or not.

In the case where the MCOP sequence $(((M_1M_2)M_3)M_4)$ is changed to $((M_1M_2)(M_3M_4))$ as shown in Fig. 3, the total number of operations changes from $C = m_1m_2m_3 + m_1m_3m_4 + m_1m_4m_5$ to

$$C' = m_1m_2m_3 + m_3m_4m_5 + m_1m_3m_5.$$

The amount of increased computation is

$$\Delta C = m_3m_4m_5 + m_1m_3m_5 - m_1m_3m_4 - m_1m_4m_5.$$

In general, when we have a product sequence, such as in Fig. 4, the amount of increased computation for multiplying $M_y \times M_{y+1}$ concurrently with $M_x \times M_{x+1}$ is as follows:

$$\Delta C = m_{y+1}m_{y+2}(m_y - m_z) + m_zm_y(m_{y+2} - m_{y+1}).$$

In this equation, m_z represents the row of the intermediate matrix (or matrix M_z itself) that is going to be multiplied

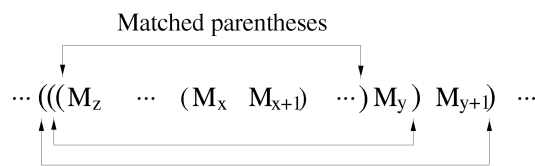


Fig. 4. A snapshot of a product sequence.

with the result of $M_y \times M_{y+1}$. In other words, there is a left parenthesis to the left of matrix M_z that matches the right parenthesis on the right side of matrix M_{y+1} . In the case of $y < z$, i.e., $(M_y(M_{y+1}(\cdots(M_x M_{x+1})\cdots M_z)))$, the amount of increased computation is

$$\Delta C = m_{y+1}m_{y+2}(m_y - m_{z+1}) + m_{z+1}m_y(m_{y+2} - m_{y+1}).$$

Finding m_z (or m_{z+1}), which is very important for the analysis of ΔC , can be done by traversing sequence tree L . If both M_y and M_{y+1} are right children, then M_z is searched by traversing the left child recursively from the parent node of M_y . Similarly, if both are left children, then M_z is searched by traversing the right child from the parent node of M_y .

In the case where $p_{i,j}$ processors are allocated to the matrix product $(M_x M_{x+1})$, but all $p_{i,j}$ processors cannot be utilized by the matrix product (i.e., $m_x m_{x+2} < p_{i,j}$), then we try to modify the product sequence L . At that time, we can decide to modify the current sequence by the following lemma.

Lemma 2. *If a leaf product $(M_x M_{x+1})$ has a candidate product $(M_y M_{y+1})$ and the DPA algorithm will allocate p_x and p_y processors to the two matrix products, respectively, then evaluation using the modified sequence reduces the evaluation time when $\Delta C < \min(\Phi(m_x, m_{x+1}, m_{x+2}, m_x m_{x+2}) \times (p_x + p_y - m_x m_{x+2}), m_y m_{y+1} m_{y+2})$.*

Proof. There are two necessary conditions for modifying a product sequence to have better performance. The first condition is that the utilization of idle processors (i.e., $p_x + p_y - m_x m_{x+2}$) should be greater than the computation increase resulting from modifying the product sequence tree. The work of idle processors can be estimated as the product of the number of utilized processors and the available time for these processors. Hence, the following condition should be satisfied:

$$\Delta C < \Phi(m_x, m_{x+1}, m_{x+2}, m_x m_{x+2}) \times (p_x + p_y - m_x m_{x+2}).$$

Also, the amount of computation given to idle processors, which is the time for multiplying $(M_y M_{y+1})$, should be more than ΔC . Therefore, the other condition to be satisfied is

$$\Delta C < m_y m_{y+1} m_{y+2}.$$

Thus, the lemma is satisfied. \square

If a candidate product $(M_y M_{y+1})$ satisfies Lemma 2, then it would be better to change product sequence $L_{i,j}$ to multiply the candidate product concurrently with $(M_x M_{x+1})$. This means that the unallocated idle processors can do more work than the increased computation required by the change in the product sequence.

When the candidate product is found, the subtree $L_{i,j}$ is modified and the processors $p_{i,j}$ are redistributed among the products in $L_{i,j}$ (including the candidate product). Also, processors are allocated proportionally to each product. This results in an enhancement of the overall system performance due to an increase in processor efficiency.

4.4 The Proposed Scheduling Algorithm

The proposed scheduling algorithm for evaluating a matrix chain product is formulated as follows:

Two-Pass Matrix Chain Scheduling Algorithm

Stage-1 MCOP

1. Find the optimal product sequence for the MCOP by using a parallel algorithm.
2. Generate the sequence tree L .

Stage-2 Top-Down Processor Assignment

1. Initialize $i = 1, j = n, p_{i,j} = P$.
2. If i is not $S[i, j]$, then allocate $p_{i,j} \times W[i, S[i, j]] / (W[i, S[i, j]] + W[S[i, j] + 1, j])$ processors to $L_{i,S[i,j]}$.
3. If j is not $S[i, j] + 1$, then allocate $p_{i,j} \times W[S[i, j] + 1, j] / (W[i, S[i, j]] + W[S[i, j] + 1, j])$ to $L_{S[i,j]+1,j}$.
4. If i is $j + 1$ or j , then finish this stage; otherwise, call this algorithm recursively, once with $i = i, j = S[i, j]$ and once with $i = S[i, j] + 1, j = j$.

Stage-3 Bottom-Up Concurrent Execution

For all leaf products, execute the following steps until there are no more unscheduled leaf products:

1. Let $(M_k M_{k+1})$ be a leaf product and $p_{k,k+1}$ be the number of processors allocated to the leaf product. If $p_{k,k+1} < m_k m_{k+2}$, then go to 5.
2. Find a candidate product by tracing ancestors of the leaf product using postorder traversal. If there is no such candidate product, go to 5.
3. Let the product $(M_l M_{l+1})$ be a candidate product found by tracing ancestors of the leaf product $(M_k M_{k+1})$. Check whether the candidate product satisfies Lemma 2. If not, go to 2.
4. Modify the sequence tree such that the candidate product $(M_l M_{l+1})$ can be executed concurrently with $(M_k M_{k+1})$. Relocate $p_{k,k+1}$ processors using the DPA algorithm and go to 1 for each leaf product of the two split subtrees.
5. Schedule the leaf product on $\min(p_{i,j}, m_k m_{k+2})$ processors. Set the parent of the leaf product as a new leaf product.

The scheduling algorithm starts from the sequence for the MCOP and then tries to modify the sequence to increase the degree of concurrency. The evaluation time of a matrix chain product is affected by the amount of computation required and the degree of concurrency. The amount of computation required is minimized by using the MCOP sequence and the degree of concurrency is maximized with a complete binary tree. The optimal product sequence with the minimum evaluation time has a form that is somewhere in between that of the MCOP sequence and the complete binary tree. The proposed scheduling algorithm moves from the MCOP sequence to a near-optimal sequence.

For purposes of efficiency, the scheduling algorithm modifies the current product sequence when the candidate product satisfies Lemma 2. Even though we can select the most suitable candidate among a number of candidates satisfying Lemma 2 by traversing the sequence tree, the scheduling algorithm uses the first candidate that satisfies Lemma 2 in order to minimize the scheduling time.

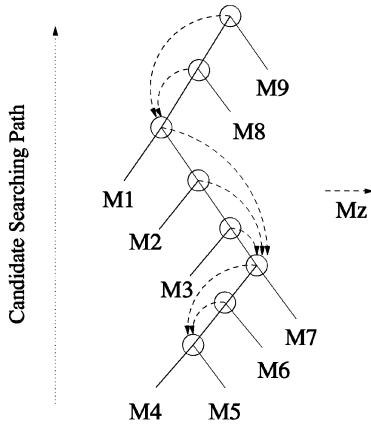


Fig. 5. Candidate searching and passing the information of the skewed point M_z .

4.5 Algorithm Complexity

The time complexity of the proposed algorithm is analyzed as follows: Stage-1 and Stage-2 can be done within $O(n)$ time. In Stage-3, to reduce the time for checking Lemma 2, we pass the information of the skewed point (M_z for ΔC) to the next parent product when we are tracing the ancestors from a leaf product, as shown in Fig. 5.

Then, we do not need to traverse the children of a candidate product to find M_z since M_z is passed from the previously traced child. This allows Step 3 of Stage-3 to be completed in constant time. The maximum number of products traced to check concurrent execution is $(n-3)$. The total number of products that may be traced in Stage-3 is $(n-3) + (n-4) + \dots + 1 = (n-2)(n-3)/2 = O(n^2)$. Also, in Step 4 of Stage-3, the number of sequence modifications is at most $(n-4)$. Since the DPA algorithm for two matrix products takes $O(m)$, the time complexity for Step 4 of Stage-3 is $O((n-4)m)$. Therefore, the time complexity of the proposed algorithm is $O(n^2 + nm)$.

4.6 Scheduling Example

The following simple example illustrates the proposed scheduling algorithm and compares the expected evaluation time of the product sequence by the proposed algorithm with that of the optimal product sequence for the MCOP.

In a system with 50 processors, let us consider the case of evaluating a chain of matrix products with five matrices. Given five matrices,

$$M_1: 6 \times 2, M_2: 2 \times 7, M_3: 7 \times 5, M_4: 5 \times 7, M_5: 7 \times 8,$$

Stage-1 finds the product sequence with the minimum number of operations for the MCOP as

$$(M_1(((M_2M_3)M_4)M_5)).$$

The MCOP sequence tree is represented as the left tree of Fig. 6. In Stage-2, we assign 50 processors to each matrix product. In Stage-3, since the leaf product (M_2M_3) cannot utilize the 50 allocated processors, we try to modify the product sequence. The product (M_4M_5) is found to be a candidate product. By checking Lemma 2, we get $p_x = 10$, $p_y = 40$ using the DPA algorithm,

$$\Delta C = 7 \times 8(5-2) + 2 \times 5(8-7) = 178,$$

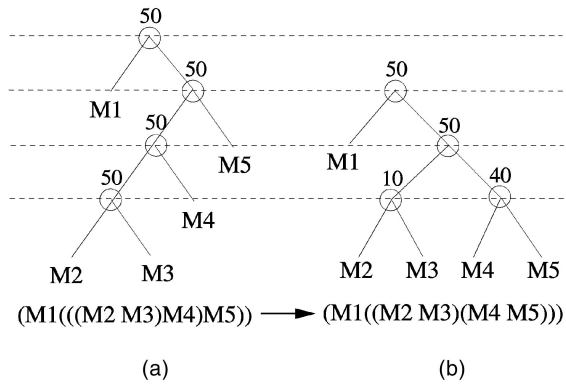


Fig. 6. The sequence trees for (a) the MCOP and (b) the proposed scheduling algorithm.

and $\Phi(2, 7, 5, 10) = 7$. Since

$$\Delta C = 178 < \min(7 \times (10 + 40 - 10), 5 \times 7 \times 8) = 280,$$

the product sequence is modified to $(M_1((M_2M_3)(M_4M_5)))$ as shown in the right tree of Fig. 6.

Let us compare the evaluation time of the optimal MCOP sequence with that of the product sequence found by the proposed scheduling algorithm. When we evaluate the matrix chain by the optimal MCOP sequence, it takes

$$(2 \times 7 \times 5) / \min(50, 2 \times 5) + (2 \times 5 \times 7) / \min(50, 2 \times 7) + (2 \times 7 \times 8) / \min(50, 2 \times 8) + (6 \times 2 \times 8) / \min(50, 6 \times 8) = 21$$

units of time. The evaluation time of the product sequence by the proposed scheduling algorithm is

$$\max(2 \times 7 \times 5 / \min(10, 2 \times 5), 5 \times 7 \times 8 / \min(40, 5 \times 8)) + 2 \times 5 \times 7 / \min(50, 2 \times 7) + 6 \times 2 \times 8 / \min(50, 6 \times 8) = 7 + 5 + 2 = 14$$

units of time. The product sequence by the proposed scheduling algorithm requires 526 operations, which is 178 operations more than the MCOP sequence having the minimum number of operations (348). However, the proposed algorithm requires less time than the MCOP sequence. This is due to the concurrent execution of multiple matrix products which increases system efficiency and reduces the total evaluation time.

5 PERFORMANCE COMPARISON

5.1 Expected Performance

In this section, the performance of the evaluation sequence found by the proposed algorithm is analyzed and compared with the optimal sequence for the MCOP. Evaluation time, speedup, and efficiency are used as the performance metrics. The upper bound and the lower bound of performance obtained by the proposed method are compared with those by the MCOP sequence in terms of these metrics.

If a schedule is given for \mathcal{M} , the evaluation sequence of \mathcal{M} can be represented as a tree with $n-1$ tasks. Let us denote task i as v_i , the number of operations in v_i as w_i , and the

maximum number of processors allocated to v_i as p_i^{max} . We let $w_i = m_{i_1} m_{i_2} m_{i_3}$ and $p_i^{max} = m_{i_1} m_{i_2} m_{i_3}$ for $v_i = (m_{i_1}, m_{i_2}, m_{i_3})$.

The evaluation time of \mathcal{M} in a single processor system is the same as the number of operations. Therefore, the evaluation time, represented as $T_{SEQ}(1)$, is bounded by the minimum number of operations of a MCOP sequence. If the number of operations of the MCOP sequence is W , $T_{SEQ}(1)$ is found as follows:

$$T_{SEQ}(1) = \sum_{i=1}^{n-1} \Phi(m_{i_1}, m_{i_2}, m_{i_3}, 1) = \sum_{i=1}^{n-1} w_i = W.$$

On a P -processor parallel system, the evaluation time of the MCOP sequence in which each matrix product is parallelized one at a time, denoted by $T_{MCOP}(P)$, is calculated as follows:

$$\begin{aligned} T_{MCOP}(P) &= \sum_{i=1}^{n-1} \Phi(m_{i_1}, m_{i_2}, m_{i_3}, \min(P, p_i^{max})) \\ &= \sum_{i=1}^{n-1} \frac{w_i}{\min(P, p_i^{max})}. \end{aligned}$$

In the case of $p_i^{max} \geq P$ for all $1 \leq i \leq n-1$, $T_{MCOP}(P) = W/P$. Under ordinary circumstances, $p_i^{max} \leq P$ for some i and $T_{MCOP}(P)$ is larger than W/P .

The proposed two-pass scheduling method for the MCSP allocates processors to the MCOP sequence tree by top-down proportional assignment and modifies the MCOP sequence when it results in a reduction of the evaluation time. Therefore, the evaluation time of the proposed method, represented as $T_{MCSP}(P)$, is not more than the evaluation time of the MCOP sequence with the top-down processor assignment. When $p_{i,j}$ processors are allocated to an $L_{i,j}$ subtree in the MCOP sequence tree L , the evaluation time $T_{MCSP}(P)$ is found as $T_{1,n}(P)$ using the following recurrence relation:

$$\begin{aligned} T_{i,j}(P_{i,j}) &= \max_{L_{i,k} \in L} (T_{i,k}(p_{i,k}), T_{k+1,j}(p_{k+1,j})) \\ &\quad + \Phi(m_i, m_{k+1}, m_{j+1}, p_{i,j}). \end{aligned}$$

Without loss of generality, $p_{i,j}^{max} = m_i m_{j+1}$. Then, $T_{MCSP}(P) = \frac{W}{P}$ for the case of $p_{i,j}^{max} \geq p_{i,j}$ for all $L_{i,j}$ in L . Let the root node task of subtree $L_{i,j}$ be v_r such that $v_r = (m_i, m_{k+1}, m_{j+1})$ and $p_r^{max} = p_{i,j}^{max} = m_i m_{j+1}$. The condition $p_{i,j}^{max} \leq p_{i,j}$ is sufficient to make $p_r^{max} \leq P$ valid. The upper bound of $T_{MCSP}(P)$ is $T_{MCOP}(P)$, thus

$$T_{MCSP}(P) \leq T_{MCOP}(P).$$

Now, we measure the maximum performance gain of the proposed method with respect to the MCOP sequence. In the best case scenario, the MCOP sequence has no concurrency (i.e., sequential evaluation), but the proposed method changes it to a complete binary tree without increased operations. One obvious case is " n " square matrix chain products. When each task v_i for $1 \leq i \leq n-1$ can utilize P/k processors in a P -processor system, i.e., $p_i^{max} = P/k$, the performance of the MCOP sequence is

$$T_{MCOP}(P) = \frac{kW}{P}, \quad S_{MCOP}(P) = \frac{P}{k}, \quad E_{MCOP}(P) = \frac{1}{k}.$$

For the complete binary tree with n nodes, the height is $\log(n)$. The tasks whose depths are not more than $\log(k)$ will run on P/k processors. Other tasks with depth i ($\log(k) < i \leq \log(n)$) will run on $(P/k)(1/2)^{i-\log(k)}$ processors. Therefore, $T_{MCSP}(P)$ is transformed as follows:

$$T_{MCSP}(P) = \sum_{i=1}^{\log(k)} \frac{w_i}{P/k} + \sum_{i=\log(k)+1}^{\log(n)} \frac{w_i}{(P/k) \cdot 1/2^{i-\log(k)}}.$$

Let us denote the number of operations scheduled by the proposed method as W_{MCSP} . Since $w_i = W_{MCSP}/(n-1)$ for a task v_i ,

$$\begin{aligned} T_{MCSP}(P) &= \frac{k}{P} \cdot \frac{\log(k)W_{MCSP}k}{(n-1)} + \frac{W_{MCSP}}{(n-1)P} (2^{\log(n)-\log(k)} - 1) \\ &= \frac{k \log(k)W_{MCSP}}{(n-1)P} + \frac{k(\frac{n}{k} - 1)W_{MCSP}}{(n-1)P} \\ &= \frac{(n + k \log(k) - k)W_{MCSP}}{(n-1)P}. \end{aligned}$$

Thus,

$$\begin{aligned} S_{MCSP}(P) &= \frac{(n-1)P}{(n + k \log(k) - k)}, \\ E_{MCSP}(P) &= \frac{(n-1)}{(n + k \log(k) - k)}. \end{aligned}$$

In the best case, W_{MCSP} is equal to the minimum W since $\Delta C = 0$. Thus, the performance gain of the MCSP sequence, denoted by $\Psi(P)$, becomes

$$\begin{aligned} \Psi(P) &= \frac{T_{MCOP}(P)}{T_{MCSP}(P)} = \frac{kW/P}{(n + k \log(k) - k)W/(n-1)P} \\ &= \frac{k(n-1)}{n + k \log(k) - k}. \end{aligned} \quad (1)$$

When each task can only utilize 25 percent of the processors in a system ($k = 4$) and the number of matrices in a chain is 100 ($n = 100$), $\Psi(P) = 3.8$. The proposed method evaluates matrix chain products 3.8 times faster than the method obtained by parallelizing the MCOP sequence. When $k = 1$, the performance of the proposed method is bounded by the evaluation of the MCOP sequence such that $\Psi(P) = 1$. As k or n increases, $\Psi(P)$ increases. As n increases to infinity, the proposed scheme can run almost k times faster than the method using the MCOP sequence.

From the above discussion, when the number of matrices becomes larger and larger and/or when the size of the matrices becomes smaller and smaller so that only part of the system is required to execute one product, the proposed method outperforms the MCOP-based evaluation method.

5.2 Experiment Setup

5.2.1 Parallel System

To compare the performance of the proposed method with those of other evaluation methods, we experimented on the Fujitsu AP1000 parallel system, which is a distributed-memory MIMD machine with 512 cells. Each cell is a SPARC processor with 16MB of memory. The AP1000 system has three independent networks: the B-net for broadcasting, the T-net for point-to-point communication,

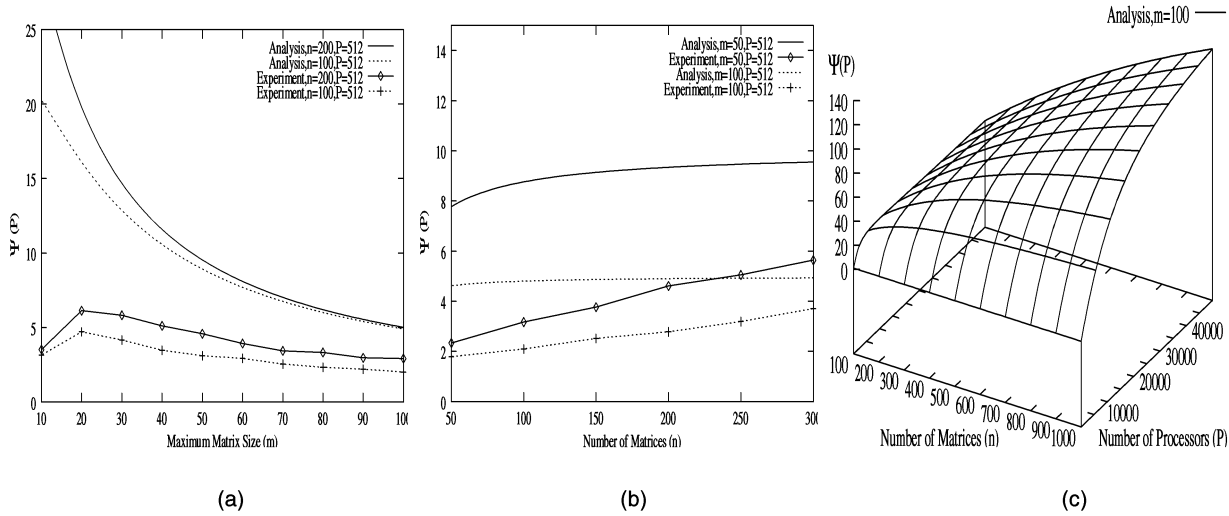


Fig. 7. Comparison of experimental results with analysis. (a) $\Psi(P)$ as a function of m with $P = 512$ and $n = 100$ or 200 . (b) $\Psi(P)$ as a function of n with $P = 512$ and $m = 50$ or 100 . (c) Expected behavior of $\Psi(P)$ as a function of P and n .

and the S-net for synchronization. The processors are interconnected by the T-net, a two-dimensional torus network with 25Mbytes/sec/port. The host computer and the processors are connected by the B-net, a broadcasting network with 50Mbytes/sec, and by the S-net for synchronization.

5.2.2 Generating Matrix Chains

The evaluation times of randomly generated matrix product chains are measured for each scheduling method. Since the initial matrix loading times are highly dependent on system characteristics such as the communication link speed and the interconnection network, the loading times are excluded in the statistics of the evaluation time. In fact, the proposed algorithm can spend less time than the sequential evaluation methods for distributing matrices to processors since several matrices can be loaded at the same time.³ The results shown in the following sections are the averages of 100 experiments.

5.2.3 Evaluation Methods

- **Linear:** evaluate by parallelizing each matrix product from the first product (M_1M_2) to the last one sequentially.
- **MCOP-Seq:** evaluate by parallelizing each matrix product using the MCOP sequence sequentially.
- **MCOP-Con:** evaluate by parallelizing each matrix product using the MCOP sequence, but execute independent matrix products concurrently by allocating the maximum number of processors.
- **MCSP-BT:** evaluate by the MCOP-Con method, but, when there are idle processors during execution, try to modify the sequence using Lemma 2.
- **MCSP-TP:** evaluate by the proposed scheduling algorithm.

3. Since some parallel computers, such as the Fujitsu AP1000, support collective communication schemes, including scatter and gather, we can reduce the matrix loading time by using these schemes.

5.3 Comparison of Theoretical Results with Experimental Results

Experimental results are compared with the upper bound of the performance computed by (1). In Figs. 7a and 7b, the lower two lines show the speedups of MCSP-TP to MCOP-Seq and the upper two lines show the upper bounds of $\Psi(P)$ with $k = P/m$. The computed performance differs from the measured $\Psi(P)$ due to a few reasons. The main reason is that the analysis is done with the best case scenario in estimating the upper bounds of $\Psi(P)$, but the cases used for experiments are not always the best case since they are generated randomly. The others come from the environmental variables simplified in the analysis.

Nevertheless, the upper bounds of $\Psi(P)$ shown in Fig. 7c inform us about the trend of the performance gain. It implies that we may gain more performance improvements by the proposed method as the number of matrices increases and the number of processors increases. The next section supports such results.

5.4 Experimental Results

Fig. 8 shows the evaluation times $T(P)$ of the five evaluation methods in milliseconds. In Fig. 8a, a chain of matrix products is generated randomly for matrices varying in size from 1 to 100 ($m = 100$) and executed on a system with 512 processors ($P = 512$). The upper two lines represent the evaluation times of sequential evaluation using Linear and MCOP-Seq and the lower three lines represent the evaluation times of the schedule sequences found by using MCOP-Con, MCSP-BT, and MCSP-TP. As the number of matrices in a matrix chain increases, the proposed MCSP-TP algorithm shows a greater performance gain. From the comparison of $T(P)$ of MCOP-Seq with that of Linear, it can be seen that reducing the amount of computation can decrease $T(P)$. However, when we allow concurrent execution, we can further decrease $T(P)$. Therefore, we confirm that the evaluation time $T(P)$ of a chain of matrix products is greatly affected by task scheduling.

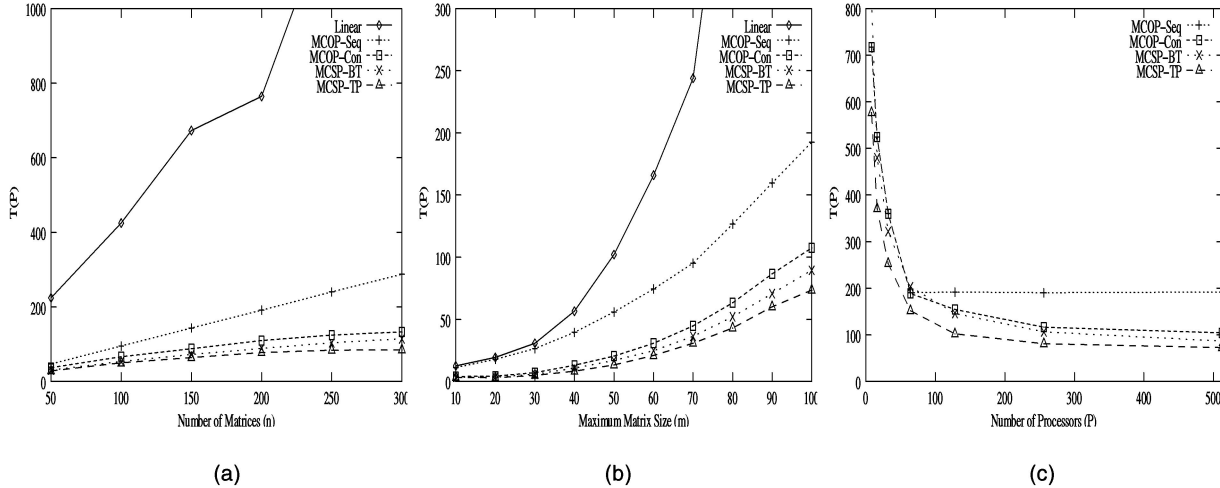


Fig. 8. Evaluation time $T(P)$ comparison. (a) $T(P)$ as a function of n with $m = 100$ and $P = 512$. (b) $T(P)$ as a function of m with $n = 200$ and $P = 512$. (c) $T(P)$ as a function of P with $n = 200$ and $m = 100$.

The experimental results with varying distributions of matrix sizes (m) is shown in Fig. 8b. The evaluation times are measured with $P = 512$ for a chain of $n = 200$. As shown in the figure, as the variance in the matrix size gets larger, MCOP-Seq has better performance than Linear. This is caused by the fact that the amount of computation is reduced greatly by the MCOP sequence and that there are small numbers of idle processors during their execution when evaluating a chain of large matrices. However, we notice that the proposed MCSP-TP method still outperforms all other methods. These experiments imply that the proposed MCSP-TP method is still effective for larger matrices due to performance improvement through concurrent execution (with more concurrency than the other methods), even though there are rare exceptions.

In Fig. 8c, the evaluation time $T(P)$ between different numbers of processors in a system is compared. $T(P)$ of Linear is too high to be appeared in this figure. As the number of processors increases, $T(P)$ of MCOP-Seq becomes nondecreasing. This implies that the reduction of computation has a limited effect on reducing $T(P)$. However, $T(P)$'s of MCSP-BT and MCSP-TP decrease more than that of MCOP-Con with larger P . This implies that sequence modification to increase the degree of concurrency improves the performance by utilizing processors efficiently. Another aspect we can see from this result is that the number of processors does not affect the performance significantly, but that the processor scheduling policy is more important than computational reduction in improving the performance of evaluating a chain of matrix products.

In Table 1, we measured the amount of computation W , evaluation time $T(P)$, and efficiency $E(P)$ for each evaluation method and the performance improvement $\Psi(P)$ of MCSP-TP. Table 1 shows the results when $n = 200$, $m = 100$, and $P = 512$. Even though the amount of computation using MCOP-Seq is 50 times less than that of Linear, the evaluation time is reduced by only 4.4 times. We observe that $E(P)$ by MCOP-Seq (0.12) is significantly lower than that of Linear (0.76). Evaluation by MCSP-BT and MCSP-TP requires more computation than that of

MCOP-Con, but $T(P)$ decreases due to the efficient use of processors. Note that MCSP-BT utilizes even more processors than MCSP-TP and still has larger $T(P)$ than that of MCSP-TP. This result implies that the proposed MCSP-TP algorithm uses processors more efficiently and effectively than MCSP-BT.

Fig. 9a shows $E(P)$ as a function of P when $n = 200$ and $m = 100$. We cannot induce any correlation between $E(P)$ and $T(P)$ from the above results. It is clear that high $E(P)$ does not have high performance. Hence, $E(P)$ is not a good performance metric for evaluating a chain of matrices. As P increases, we can see that MCSP-BT and MCSP-TP sustain higher $E(P)$ than MCOP-Seq and MCOP-Con.

Fig. 9b shows $\Psi(P)$ as a function of P for the four cases of n and m . As the number of processors increases, $\Psi(P)$ increases. Also, we can see that as n increases and m decreases, the proposed method outperforms the evaluation by the optimal MCOP sequence.

5.5 Summary of Results

From the above performance evaluation, we can conclude the following:

- When evaluating a chain of matrices on a parallel system, simply reducing the number of required operations does not greatly decrease the evaluation time.

TABLE 1
Comparison with Respect to W , $T(P)$, $E(P)$, $\Psi(P)$
When $n = 200$, $m = 100$, and $P = 512$

Schedule	W	$T(P)$	$E(P)$	$\Psi(P)$
Linear	24254K	839	0.76	-
MCOP-Seq	473K	192	0.12	-
MCOP-Con	473K	104	0.23	-
MCSP-BT	811K	87	0.36	-
MCSP-TP	536K	72	0.29	2.98

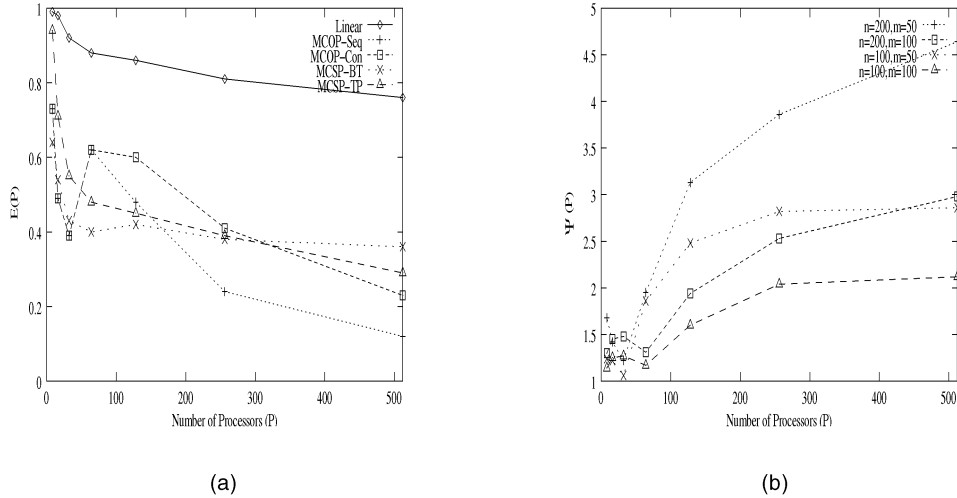


Fig. 9. Efficiency and performance gain with respect to the number of processors P . (a) $E(P)$ as a function of P . (b) $\Psi(P)$ as a function of P .

- Concurrent execution of independent multiple matrix products compensates for inefficient processor usage with parallel processing and increases the system efficiency so that performance improves greatly.
- When the number of processors becomes larger or when the number of matrices in a matrix chain increases, evaluation by the proposed scheduling algorithm MCSP-TP progressively outperforms other methods such as Linear, MCOP-Seq, MCOP-Con, and MCSP-BT.
- Even when the size of the matrices is quite large so that there are no idle processors during their evaluation, MCSP-TP is still effective due to the concurrent execution of independent matrix products with the top-down processor assignment.
- When evaluating a chain with small matrices on many processors, sequence modification to increase system efficiency greatly reduces the evaluation time.
- For matrix chain products, efficient scheduling is more effective than only increasing the number of processors.

6 EXTENSION OF THE MCSP

In this section, we discuss practical considerations for the practical use of the MSCP method and extensions of the MCSP approach.

6.1 Practical Considerations

Depending on the computing environment for evaluating \mathcal{M} , the parameters used in the proposed method are easily adjustable. The parameters are the execution time function, the maximum number of processors allocated to a product, and the minimum number of processors allocated to a product. For ease of discussion, let us denote these parameters as $\Phi(v_i, p_i)$, p_i^{max} , and p_i^{min} , respectively, for a given product v_i .

- **Matrix multiplication algorithm:** Many parallel matrix multiplication algorithms have been suggested for various architectures [24], [23]. Depending

on which algorithm is used for the MCSP, the execution time function $\Phi(v_i, p_i)$ can be specified in detail. Thus, $\Phi(v_i, p_i)$ can be changed for taking the parallel processing overheads into account. The parallel processing overheads include load imbalance, communication cost, and synchronization.

- **Parallel system architecture:** The parallel system used for evaluating \mathcal{M} also affects the parameters. For example, on a shared memory multiprocessor system, the communication cost is negligible, but, on a message-passing distributed system, the communication cost should be considered in $\Phi(v_i, p_i)$. Consequently, the parameter p_i^{max} will be more restricted in a distributed system.
- **Input matrix size:** One limitation in using the proposed algorithm is that the method is not very effective with a chain of large matrices. However, a more elaborate description of $\Phi(v_i, p_i)$ from the above parallel algorithm-architecture parameters will compensate for this limitation by using a reduced value of p_i^{max} . Also, even in the case where any one product is large enough to utilize the whole system, the top-down processor assignment leads to better performance than sequential evaluation.

Another parameter to be considered in the MCSP is the minimum number of processors (p_i^{min}) to be allocated to a product v_i . Basically, we assumed that any product can be executable on one processor, i.e., $p_i^{min} = 1$. However, depending on the size of the matrices, the minimum number of processors to be allocated to a product v_i can be restricted to be $p_i^{min} > 1$. This is mainly due to the memory capacity of a single processor. We can consider the parameter p_i^{min} in the top-down processor assignment stage.

- **Matrix distribution cost:** In the proposed method, the cost for the initial distribution of matrices among the processors was excluded, even though this cost is not negligible in a distributed-memory system.⁴ Of course, the cost of the proposed method is definitely less than sequential evaluation by the MCOP sequence. This is mainly due to the fact that the

4. The cost is negligible in a shared-memory multiprocessor system.

proposed method runs on fewer processors than sequential evaluation. The effect of this cost for scheduling the MCSP sequence is considered for future work.

6.2 MCSP Extensions

We can extend the MCSP to a few sparse matrix problems. The MCSP has characteristics such that it can be represented as a tree precedence task graph with many equivalent tree graphs and the task graph determines the number of required operations and the degree of concurrency. The proposed method for the MCSP can be applied to evaluate a chain of square matrices with sparsity and to factorize a large sparse matrix in parallel systems.

- **Sparse matrix chain products:** Sparse matrix multiplications require a different number of operations with respect to the sparsity structure. For multiplying an $n \times n$ matrix A with nonzero density d_1 by an $n \times n$ matrix B with nonzero density d_2 , the number of required operations is $d_1 d_2 n^3$ [30]. Also, for multiplying a tridiagonal matrix with any type of matrix, the required operations vary depending on the type of matrices [31]. Such a sparse matrix multiplication has the same characteristics as the MCSP. A different evaluation sequence results in a different number of required operations. We can extend the MCSP to the evaluation of a chain of square matrices with sparsity.
- **Sparse matrix factorization:** Elimination trees are used extensively in sparse matrix factorization because they present the sequence of computation and parallelism [32]. Since there are many equivalent elimination trees with different parallelization structures [33], we can follow the MCSP approach for this sparse factorization problem.

The above items have been studied extensively for a long time and each problem can be approached as a separate research subject.

7 SUMMARY AND CONCLUSION

In this paper, we introduced the matrix chain scheduling problem (MCSP) and proposed a heuristic scheduling algorithm for the MCSP. The proposed algorithm schedules matrix products to processors with the objective of enhancing concurrency at the expense of a slight increase in the required number of operations when compared to the optimal product sequence found for the matrix chain ordering problem (MCOP). We have shown that performance is significantly enhanced by the proposed algorithm using experiments on the Fujitsu AP1000 parallel system. As a result, we can confirm that efficient processor scheduling is much more important than simply reducing the total number of operations when evaluating a matrix chain product in parallel systems. Given a parallel system with a large number of processors or a matrix product chain involving many matrices, evaluation by the proposed method, greatly outperforms the parallel evaluation method which uses the optimal product sequence found for the MCOP. The main contribution of this work is the

formalization of the MCSP and the introduction of a processor allocation and task scheduling algorithm that results in a significant performance improvement when evaluating matrix chain products in parallel systems. We are currently working on extending this algorithm to evaluate a chain of square matrices in the form of sparse matrices or band matrices and to scheduling of parallel matrix factorization using elimination trees. Also, we plan to study generalizing the MCSP to scalable task scheduling on parallel systems.

ACKNOWLEDGMENTS

The authors would like to thank Professor Chul E. Kim for many fruitful suggestions on the complexity analysis. Also, they give special thanks to Dr. Min-Ho Kyung for very helpful discussions on this work. This research was supported in part by the Ministry of Education of Korea through its BK21 program toward the Electrical and Computer Engineering Division at POSTECH.

REFERENCES

- [1] E. Horowitz and S. Sahni, *Fundamentals of Data Structures*. Potomac, Md.: Computer Science Press, 1976.
- [2] E. Dekel, D. Nassimi, and S. Sahni, "Parallel Matrix and Graph Algorithms," *SIAM J. Computing*, vol. 10, pp. 657-675, Nov. 1981.
- [3] S.-T. Yau and Y.Y. Lu, "Reducing the Symmetric Matrix Eigenvalue Problem to Matrix Multiplications," *SIAM J. Scientific Computing*, vol. 14, no. 1, pp. 121-136, 1993.
- [4] S.-S. Lin, "A Chained-Matrices Approach for Parallel Computation of Continued Fractions and Its Applications," *J. Scientific Computing*, vol. 9, no. 1, pp. 65-80, 1994.
- [5] V. Kumar, A. Grama, A. Gupta, and G. Karypis, *Introduction to Parallel Computing: Design and Analysis of Parallel Algorithms*. Benjamin/Cummings, 1994.
- [6] A. Chandra, "Computing Matrix Chain Products in Near-Optimal Time," technical report, IBM T.J. Watson Research Center, Yorktown Heights, N.Y., IBM Research Report RC 5625(#24393), 1975.
- [7] H. Gould, *Bell and Catalan Numbers*. Research Inst., Morgantown, W.Va., June 1977.
- [8] S. Godbole, "An Efficient Computation of Matrix Chain Products," *IEEE Trans. Computers*, pp. 864-866, vol. 22, no. 9, Sept. 1973.
- [9] F. Chin, "An $O(n)$ Algorithm for Determining a Near-Optimal Computation Order of Matrix Chain Product," *Comm. ACM*, pp. 544-549, 1978.
- [10] T. Hu and M. Shing, "Computation of Matrix Chain Products. Part I," *SIAM J. Computing*, vol. 11, pp. 362-373, May 1982.
- [11] T. Hu and M. Shing, "Computation of Matrix Chain Products. Part II," *SIAM J. Computing*, vol. 13, pp. 228-251, May 1984.
- [12] P. Ramanan, "A New Lower Bound Technique and Its Application: Tight Lower Bound for a Polygon Triangulation Problem," *SIAM J. Computing*, vol. 23, pp. 834-851, Aug. 1994.
- [13] L. Valiant, S. Skyum, S. Berkowitz, and C. Rackoff, "Fast Parallel Computation of Polynomials Using Few Processors," *SIAM J. Computing*, vol. 12, pp. 641-644, 1983.
- [14] W. Rytter, "Note on Efficient Parallel Computations for Some Dynamic Programming Problems," *Theoretical Computer Science*, vol. 59, pp. 297-307, 1988.
- [15] S.-H.S. Huang, H. Liu, and V. Viswanathan, "Parallel Dynamic Programming," *IEEE Trans. Parallel and Distributed Systems*, vol. 5, no. 3, pp. 326-328, Mar. 1994.
- [16] P.G. Bradford, G.J. Rawlins, and G.E. Shannon, "Efficient Matrix Chain Ordering in Polylog Time," *SIAM J. Computing*, vol. 27, no. 2, pp. 466-490, 1998.
- [17] A. Czumaj, "Parallel Algorithm for the Matrix Chain Product and the Optimal Triangulation Problems," *Proc. Symp. Theoretical Aspects of Computer Science*, pp. 294-305, 1993.

- [18] A. Czumaj, "Very Fast Approximation of the Matrix Chain Product Problem," *J. Algorithms*, vol. 21, no. 1, pp. 71-79, 1996.
- [19] P. Ramanan, "An Efficient Parallel Algorithm for the Matrix Chain Product Problem," *SIAM J. Computing*, vol. 25, pp. 874-893, Aug. 1996.
- [20] V. Strassen, "Gaussian Elimination Is Not Optimal," *Numerical Math.*, vol. 13, pp. 354-356, 1969.
- [21] G.H. Golub and C.F.V. Loan, *Matrix Computations*. second ed., Baltimore: Johns Hopkins Univ. Press, 1989.
- [22] N.-K. Tsao, "Error Complexity Analysis of Algorithms for Matrix Multiplication and Matrix Chain Product," *IEEE Trans. Computers*, vol. 30, no. 10, pp. 758-771, Oct. 1981.
- [23] C. Puglisi, "Parallel Algorithms and Architectures for Matrix Multiplication," *Computer Math. Applications*, vol. 17, no. 12, pp. 1567-1572, 1989.
- [24] A. Gupta and V. Kumar, "Scalability of Parallel Algorithms for Matrix Multiplication," *Proc. Int'l Conf. Parallel Processing*, pp. 115-123, 1993.
- [25] R. Krishnamurti and E. Ma, "An Approximation Algorithm for Scheduling Tasks on Varying Partition Sizes in Partitionable Multiprocessor Systems," *IEEE Trans. Computers*, vol. 41, no. 12, pp. 1572-1579, Dec. 1992.
- [26] S.G. Akl, *The Design and Analysis of Parallel Algorithms*. Englewood Cliffs, N.J.: Prentice Hall, Inc., 1989.
- [27] V. Sarkar, *Partitioning and Scheduling Parallel Programs for Multiprocessors*. Cambridge, Mass.: MIT Press, 1989.
- [28] C.D. Polychronopoulos and U. Banerjee, "Speedup Bounds and Processor Allocation for Parallel Programs on Multiprocessors," *Proc. Int'l Conf. Parallel Processing*, pp. 961-968, 1986.
- [29] C.D. Polychronopoulos and U. Banerjee, "Processor Allocation for Horizontal and Vertical Parallelism and Related Speedup Bounds," *IEEE Trans. Computers*, vol. 36, no. 4, pp. 410-420, Apr. 1987.
- [30] A. Schoor, "Fast Algorithm for Sparse Matrix Multiplication," *Information Processing Letters*, vol. 15, no. 2, pp. 87-89, 1982.
- [31] J. Takche, "Complexities of Special Matrix Multiplication Problems," *Computer Math. Applications*, vol. 15, no. 12, pp. 977-989, 1988.
- [32] J.W.H. Liu, "The Role of Elimination Trees in Sparse Factorization," *SIAM J. Matrix Analysis and Applications*, vol. 11, pp. 134-172, Jan. 1990.
- [33] J.W.H. Liu, "Equivalent Sparse Matrix Reordering by Elimination Tree Rotations," *SIAM J. Scientific and Statistical Computing*, vol. 9, pp. 424-444, May 1988.



Heejo Lee received the BS, MS, and PhD degrees in computer science and engineering from Pohang University of Science and Technology (POSTECH), Korea, in 1993, 1995, and 2000, respectively. Since 2001, he has been working at Ahnlab, Inc. as Chief Technology Officer (CTO) and director of the Technology Planning Department. From 2000 to 2001, he was a postdoctoral research associate at the Network Systems Lab of the Department of

Computer Sciences, Purdue University, and at the Center for Education and Research in Information Assurance and Security (CERIAS), Purdue University. His research interests include network security, parallel scientific computing, and fault-tolerant computing. He is a member of the IEEE and the IEEE Computer Society.



Jong Kim received the BS degree in electronic engineering from Hanyang University, Seoul, Korea, in 1981, the MS degree in computer science from the Korea Advanced Institute of Science and Technology, Seoul, Korea, in 1983, and the PhD degree in computer engineering from Pennsylvania State University in 1991. Since 1992, he has been an assistant professor and then an associate professor in the Department of Computer Science and Engineering, Pohang University of Science and Technology (POSTECH), Pohang, Korea. From 1991 to 1992, he was a research fellow in the Real-Time Computing Laboratory of the Department of Electrical Engineering and Computer Science, University of Michigan. His major areas of interest are fault-tolerant computing, parallel and distributed computing, and security. He is a member of the IEEE and the IEEE Computer Society.



Sung Je Hong received the BS degree in electronics engineering from Seoul National University, Korea, in 1973, the MS degree in computer science from Iowa State University, Ames, in 1979, and the PhD degree in computer science from the University of Illinois, Urbana-Champaign, Illinois, in 1983. Since 1989, he has been a professor in the Department of Computer Science and Engineering, Pohang University of Science and Technology, Pohang, Korea. From 1983 to 1989, he was a staff member of Corporate Research and Development, General Electric Company, Schenectady, New York. From 1975 to 1976, he was with Oriental Computer Engineering, Korea, as a logic design engineer. From 1973 to 1975, he served the army as a system analyst at the Army Finance Center, Korea. His current research interests include VLSI design, CAD algorithms, testing, and parallel processing. He is a member of the IEEE and the IEEE Computer Society.



Sunggu Lee received the BSEE degree with highest distinction from the University of Kansas, Lawrence, in 1985 and the MSE and PhD degrees from the University of Michigan, Ann Arbor, in 1987 and 1990, respectively. He is currently an associate professor in the Department of Electronic and Electrical Engineering at the Pohang University of Science and Technology (POSTECH), Pohang, Korea. Prior to this appointment, he was an assistant professor in the Department of Electrical Engineering at the University of Delaware in Newark, Delaware. From June 1997 to June 1998, he spent one year as a visiting scientist at the IBM T.J. Watson Research Center. His research interests are in parallel computing using clusters, fault-tolerant computing, and real-time computing. He is a member of the IEEE and the IEEE Computer Society.

► For more information on this or any other computing topic, please visit our Digital Library at <http://computer.org/publications/dlib>.