

AD-A091 572

MINNESOTA UNIV MINNEAPOLIS DEPT OF COMPUTER SCIENCE

F/6 9/2

BINARY TREES AND PARALLEL SCHEDULING ALGORITHMS. (U)

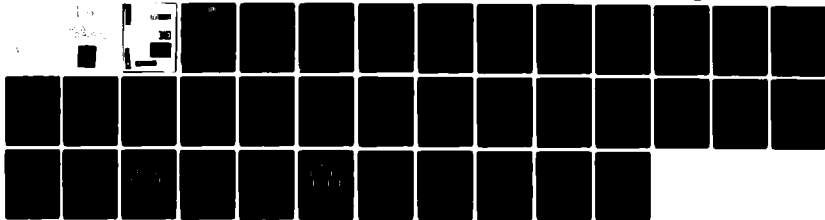
SEP 80 E DEKEL, S SAHNI

N00014-80-C-0650

UNCLASSIFIED

TR-80-19

NL



LEVEL II

12

3 Computer Science Department
136 Lind Hall
Institute of Technology
University of Minnesota
Minneapolis, Minnesota 55455

NOT J0023879

Binary Trees and
Parallel Scheduling Algorithms
by
Eliezer Dekel and Sartaj Sahni
Technical Report 80-19 ✓
September 1980

(22 14 12)

Cover design courtesy of Ruth and Jay Leavitt

DISTRIBUTION STATEMENT A
Approved for public release;
Distribution Unlimited

80 10 27 062

Binary Trees and Parallel Scheduling Algorithms*

Eliezer Dekel and Sartaj Sahni
University of Minnesota

Abstract:

^A This paper examines the use of binary trees in the design of efficient parallel algorithms. Using binary trees, we develop efficient algorithms for several scheduling problems. The shared memory model for parallel computation is used. Our success in using binary trees for parallel computations, indicates that the binary tree is an important and useful design tool for parallel algorithms.

Keywords and Phrases:

Parallel algorithms, design methodologies, complexity, scheduling, shared memory model.

Accession For	
NTIS GRA&I	<input checked="" type="checkbox"/>
DDC TAB	<input type="checkbox"/>
Unannounced	<input type="checkbox"/>
Justification	
By _____	
Distribution/ _____	
Availability Codes	
Dist	Available/ or special
A	

*This research was supported in part by the National Science Foundation under grant MCS80-005856 and in part by the Office of Naval Research under contract N00014-80-C-0650.

i. Introduction

Algorithm design techniques for single processor computers have been extensively studied. For example Horowitz and Sahni [15] extoll the virtues of such design methods as: divide-and-conquer; dynamic programming; greedy method; backtracking; and branch-and-bound. These methods generally lead to efficient sequential (i.e. single processor) algorithms for a variety of problems. These algorithms, however, are not very efficient for computers with a very large number of processors. In this paper, we propose a design method that we have found useful in the design of algorithms for computers that have many processors. The method proposed here is called the binary tree method. While this method has been used in the design of parallel algorithms earlier; here we attempt to show its broad applicability to the design of such algorithms. It is hoped that further research will bring to light some other basic design tools for parallel algorithms. One should note that trees have been used extensively in the design of efficient sequential algorithms. In fact, divide-and-conquer; backtracking; and branch-and-bound all use an underlying computation tree [15]. The use of binary trees as proposed here is quite different from the use of trees in sequential computation.

With the continuing dramatic decline in the cost of hardware, it is becoming feasible to economically build computers with thousands of processors. In fact, Batcher [5] describes a computer (MPP) with 16,384 processors that is currently being built for NASA. In coming years, one can expect to see computers with a hundred thousand or even a million processing elements. This expectation has motivated the study of parallel algorithms. Since the complexity of a parallel algorithm depends very much on the architecture of the parallel computer on which it is run, it is necessary to keep the architecture in mind when designing the algorithm. Several parallel architectures have been proposed and studied. In this paper we shall deal directly only with the single instruction stream, multiple data stream (SIMD) model. Our technique and algorithms readily adapt to the other models (eg: multiple instruction stream multiple data stream (MIMD) and data flow models). SIMD computers have the following characteristics:

- (1) They consist of p processing elements (PEs). The PEs are indexed $0, 1, \dots, p-1$ and an individual PE may be referenced as in $PE(i)$. Each PE is capable of performing the standard arithmetic and logical operations. In addition, each PE knows its index.
- (2) Each PE has some local memory.
- (3) The PEs are synchronized and operate under the control of a single instruction stream.

- (4) An enable/disable mask can be used to select a subset of the PEs that are to perform an instruction. Only the enabled PEs will perform the instruction. The remaining PEs will be idle. All enabled PEs execute the same instruction. The set of enabled PEs can change from instruction to instruction.

While several SIMD models have been proposed and studied, in this paper we wish to make a distinction between the shared memory model (SMM) and the remaining models; all of which employ an interconnection network and use no shared memory. In the shared memory model, there is a common memory available to each PE. Data may be transmitted from PE(i) to PE(j) by simply having PE(i) write the data into the common memory and then letting PE(j) read it. Thus, in this model it takes only $O(1)$ time for one PE to communicate with another PE. Two PEs are not permitted to write into the same word of common memory simultaneously. The PEs may or may not be allowed to simultaneously read the same word of common memory. If the former is the case, then we shall say that read conflicts are permitted.

Most algorithmic studies of parallel computation have been based on the SMM ([1], [7], [8], [11], [12], [13], [24], [25], [30]). This model is, however, not very realistic as it assumes that the p PEs can access any p words of memory (i word per PE) in the same time slot. In practice, however, the memory will be divided into blocks so that no two PEs can simultaneously access words in the same block. If two or more PEs wish to access words in the same memory block then the requests will get queued. Each PE will be served in a different time slot. Thus, in the worst case $O(p)$ time could be spent transferring data to the p PEs. All the papers cited earlier ignore this and take the time for a simultaneous memory access by all PEs to be $O(1)$.

SIMD computers with restricted interconnection networks appear to be more realistic. In fact, the ILLIAC IV is an example of such a machine. There are several other such machines that are currently being fabricated. The largest of these is the massively parallel processor (MPP) designed by K. Batcher. It has $p=16K$. A block diagram of a SIMD computer with an interconnection network is given in Figure 1.1. Observe that there is no shared memory in this model. Hence, PEs can communicate amongst themselves only via the interconnection network.

While several interconnection networks have been proposed (see [33]), we shall describe only three interconnection networks here. These are: mesh, cube, and perfect shuffle. The corresponding computer models are described below. Figure 1.2 shows the resulting interconnection

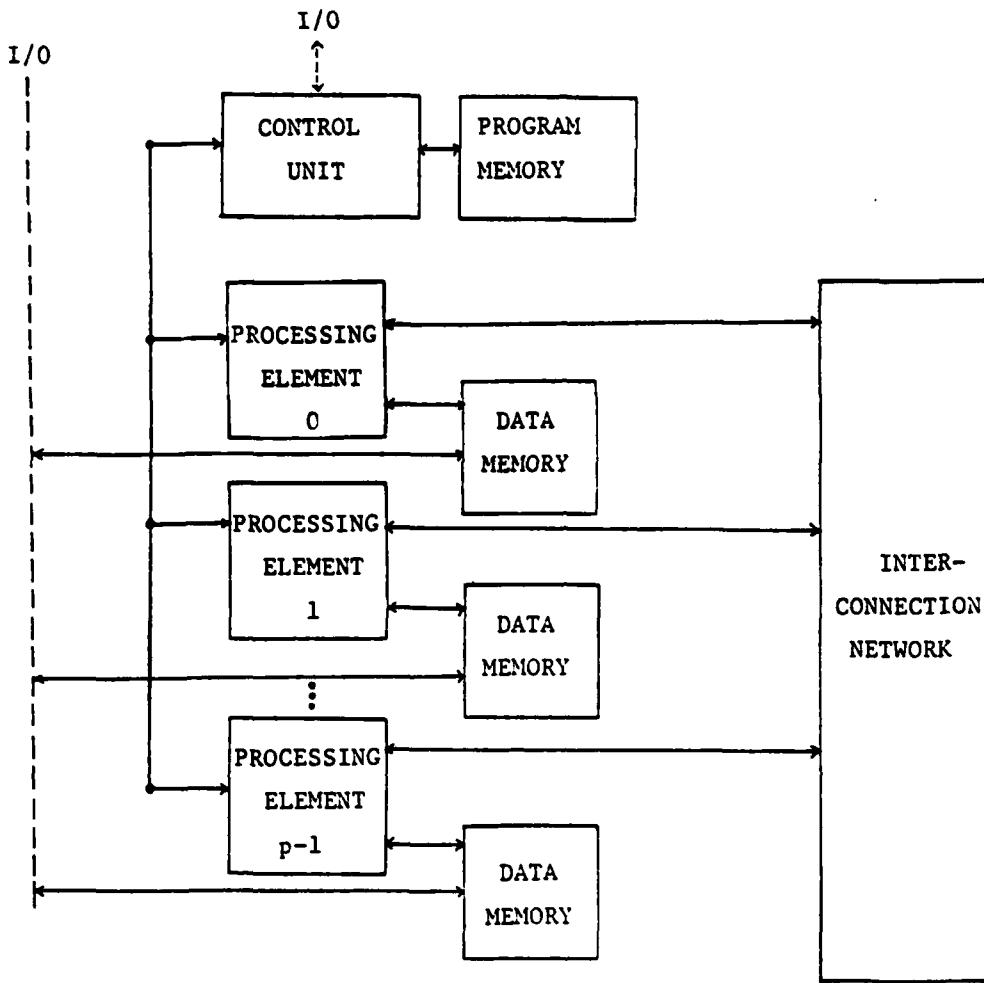


Figure 1.1 Block diagram of an SIMD computer.

patterns.

i) Mesh Connected Computer (MCC)

In this model the PEs may be thought of as logically arranged as in a k dimensional array $A(n_{k-1}, n_{k-2}, \dots, n_0)$ where n_i is the size of the i th dimension and $p = n_{k-1} * n_{k-2} * \dots * n_0$. The PE at location $A(i_{k-1}, \dots, i_0)$ is connected to the PEs at locations $A(i_{k-1}, \dots, i_{j+1}, \dots, i_0)$, $0 \leq j < k$, provided they exist. Data may be transmitted from one PE to another only via this interconnection pattern. The interconnection scheme for a 16 PE MCC with $k=2$ is given in Figure 1.2(a).

ii) Cube Connected Computer (CCC)

Assume that $p=2^q$ and let $i_{q-1} \dots i_0$ be the binary representation of i for $i \in [0, p-1]$. Let $i^{(b)}$ be the number whose binary representation is $i_{q-1} \dots i_{b+1} \bar{i}_b i_{b-1} \dots i_0$ where \bar{i}_b is the complement of i_b and $0 \leq b < q$. In the cube model, PE(i) is connected to PE($i^{(b)}$), $0 \leq b < q$. As in the mesh model, data can be transmitted from one PE to another only via the interconnection pattern. Figure 1.2(b) shows an 8 PE CCC configuration.

iii) Perfect shuffle Computer (PSC)

Let p, q, i and $i^{(b)}$ be as in the cube model. Let $i_{q-1} \dots i_0$ be the binary representation of i . Define SHUFFLE(i) and UNSHUFFLE(i) to, respectively, be the integers with binary representation $i_{q-2} i_{q-3} \dots i_0 i_{q-1}$ and $i_{q-1} \dots i_0$. In the perfect shuffle model, PE(i) is connected to PE($i^{(0)}$), PE(SHUFFLE(i)), and PE(UNSHUFFLE(i)). These three connections are, respectively, called exchange, shuffle, and unshuffle. Once again, data transmission from one PE to another is possible only via the connection scheme. An 8 PE PSC configuration is shown in Figure 1.2(c).

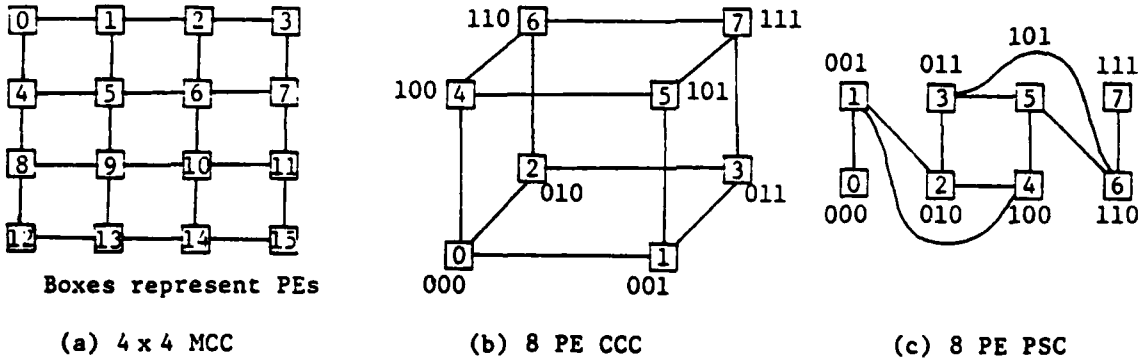


Figure 1.2

It should be noted that the MCC model requires $2k$ connections per PE, the CCC model requires $\log p$ (all logarithms in this paper are base 2) and the PSC model requires only three connections per PE. The SMM requires a large (and impractical) amount of PE to memory connections to permit simultaneous memory access by several PEs. It should also be emphasized that in any time instance, only one unit of data (say one word) can be transmitted along an interconnection line. All lines can be busy in the same time

instance though.

Each of the four models (including the SMM) described above has received much attention in the literature. Agerwala and Lint [1], Arjomandi [2], Csanky [8], Eckstein [11] and Hirschberg [12] have developed algorithms for certain matrix and graph problems using the SMM. Hirschberg [13], Muller and Preparata [24] and Preparata [30] have considered the sorting problem for SMM. The evaluation of polynomials on the SMM has been studied by Munro and Paterson [25], while arithmetic expression evaluation has been considered by Brent [7] and others. Efficient algorithms to sort and perform data permutations on an MCC can be found in Thompson and Kung [38], Nassimi and Sahni [26] and [27], and Thompson [37]. Thompson's algorithms [37] can also be used to perform permutations on a CCC and a PSC. Lang [19], Lang and Stone [20], and Stone [36] show how certain permutations may be performed using shuffles and exchanges. Nassimi and Sahni [28] develop fast sorting and permutation algorithms for a CCC and a PSC. Dekel, Nassimi, and Sahni [9] present efficient matrix multiplication and graph algorithms for CCCs and PSCs.

The algorithms considered in this paper are described explicitly only for the SMM. The algorithms are readily translated into algorithms for the other SIMD models. In some cases, it may be necessary to use the data broadcasting algorithms developed by Nassimi and Sahni [29] to accomplish this adaptation to the other models.

Throughout this paper, we assume that no read conflicts are allowed. To see the importance of this assumption, consider the partition problem. In this problem we are given n numbers a_1, a_2, \dots, a_n and we wish to determine if there is a subset S of $\{1, 2, \dots, n\}$ such that $\sum_{i \in S} a_i = \sum_{i \notin S} a_i$. This can be done in $O(\log n)$ time if read conflicts are allowed. The first phase of this algorithm uses $\lfloor n/\log n \rfloor 2^n$ PEs. PEs are divided into 2^n groups of $\lfloor n/\log n \rfloor$ PEs each. The PE groups are indexed $0, 1, \dots, 2^{n-1}$. Each PE group, i , considers the subset $S_i = \{a_j \mid \text{bit } j \text{ of } i \text{ is } 1\}$. The elements in S_i are added in $O(\log n)$ time using the $\lfloor n/\log n \rfloor$ PEs in the PE group (this is described later in this section). Next, the elements in \bar{S}_i are added. If $\sum_{j \in S_i} a_j = \sum_{j \notin S_i} a_j$ then one of the PEs in group i sets $V(i)$ to 1; otherwise $V(i)$ is set to 0. In the second phase, Valiant's [39] $O(\log \log m)$ algorithm is used to determine the maximum $V(i)$. Since there are 2^n $V(i)$'s, this takes $O(\log n)$ time. The answer to the partition problem is "yes" iff the maximum $V(i)$ is 1. The total time taken by the above algorithm is $O(\log n)$.

The procedure described above has read conflicts in two of its steps. First, when the PE groups are computing sums, many PEs will attempt to simultaneously read the same a_i . To remove these conflicts, we will need to make 2^n copies of each a_i , one copy for each PE group. This takes $O(n)$ time using no read conflicts. Second, Valiant's algorithm also has read conflicts. Removing these also takes $O(n)$ time. So, the complexity of our parallel partition problem algorithm is $O(\log n)$ if read conflicts are permitted, and is $O(n)$ if they are not.

We first illustrate the binary tree method on a very simple problem. Let us consider how we might compute the sum $\sum_{i=1}^n A(i)$, $n \geq 1$. The most frequently used sequential algorithm for this computation uses the parsing $\sum_{i=1}^n A(i) = (\dots((A(1) + A(2)) + A(3)) + \dots + A(n))$. To arrive at an efficient parallel algorithm, it is necessary to consider the parsing $\sum_{i=1}^n A(i) = (\dots((A(1) + A(2)) + (A(3) + A(4))) + ((A(5) + A(6)) + (A(7) + A(8)))) + \dots$. Computation corresponding to this parsing scheme is best described by a complete binary tree with n leaves. Figure 1.3 describes the computation for the case $n = 11$.

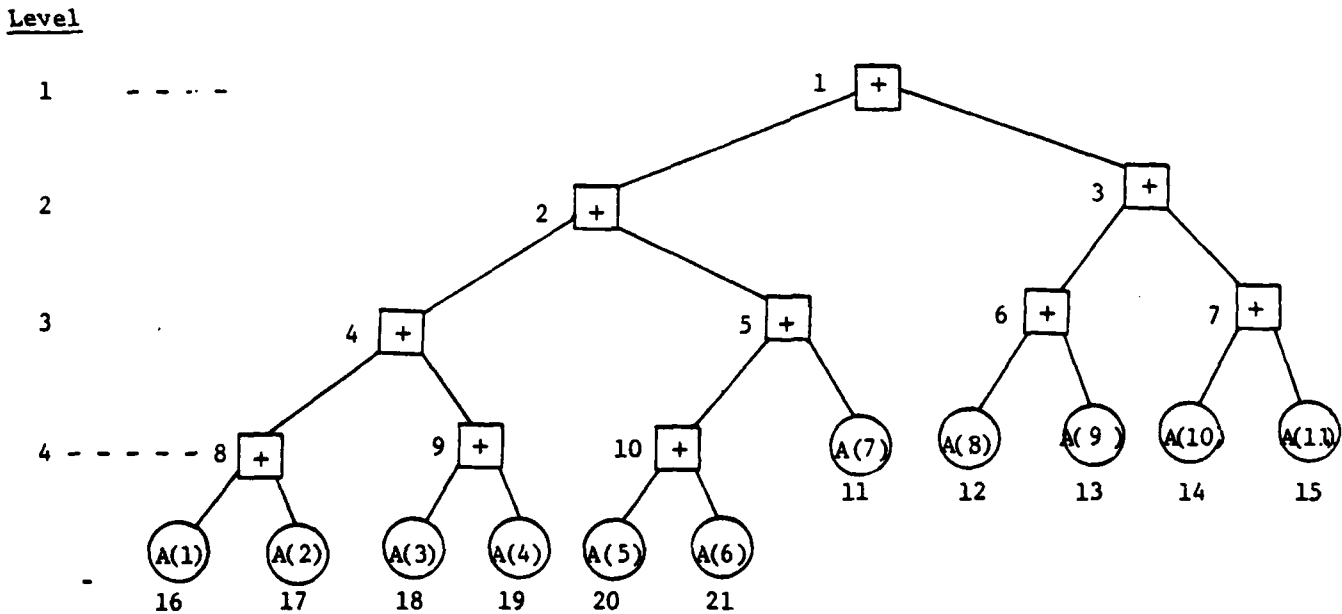


Figure 1.3 Computation tree for $\sum_{i=1}^n A(i)$

The square nodes represent nodes at which addition is to be

performed. The circular nodes represent initial data. Nodes have been numbered using the standard numbering scheme for complete binary trees. Node indices appear outside the nodes. Let $V(i)$ be the corresponding $A()$ value for node i if i denotes a circular node. $V(i)$ is initially undefined for the other nodes. Thus for Figure 1.3, $V(17) = A(2)$; $V(13) = A(9)$; $V(21) = A(6)$; etc. Using the tree of Figure 1.3, $\sum_{i=1}^{17} A(i)$ may be computed in 4 steps using 4 PEs as follows:

step 1: Use three PEs to compute, in parallel $V(8) = V(16) + V(17)$; $V(9) = V(18) + V(19)$; and $V(10) = V(20) + V(21)$

step 2: Use four PEs to compute, in parallel, $V(i) = V(2i) + V(2i + 1)$, $4 \leq i \leq 7$

step 3: Use two PEs to compute, in parallel, $V(i) = V(2i) + V(2i + 1)$, $2 \leq i \leq 3$

step 4: Use one PE to compute $V(1) = V(2) + V(3) = \sum_{i=1}^{17} A(i)$.

From the nature of a binary computation tree, it is clear that parallel addition needs at most $\lfloor n/2 \rfloor$ PEs. The parallel addition algorithm is described more formally in Figure 1.4. In lines 2 and 5, the use of a $\leq b \leq c$ means that this line is to be executed in parallel for all b satisfying the inequality. Line 2 can be performed in two steps using $\lfloor n/2 \rfloor$ PEs. Line 4 needs at most $\lfloor n/2 \rfloor$ PEs. It is clear that the complexity of procedure SUM_n is $O(\log n)$.

```

line  procedure  $SUM_n(A, n)$ 
        //compute  $\sum_{i=1}^n A(i)$  using  $\lfloor n/2 \rfloor$  PEs//
        //initialize  $v$ //
1        $k \leftarrow \lfloor \log_2 n \rfloor$ ;  $j \leftarrow 2^k$ ;  $t \leftarrow 2*(n-j)$ ;  $p \leftarrow n-i$ 
2        $V(p+i) \leftarrow A((i+t-i) \bmod n + i)$ ,  $i \leq n$ 
3       for  $i \leftarrow k$  down to 0 do //add by levels//
4          $V(j) \leftarrow V(2j) + V(2j + 1)$ ,  $2^i \leq j \leq \min\{p, 2^{i+1}-1\}$ 
5       end for
6       return  $V(1)$ 
7     end  $SUM$ 

```

Figure 1.4

In addition to analyzing the complexity of a parallel algorithm, one often (see Savage [32]) also computes the effectiveness of processor utilization (EPU). This is

defined relative to a specific problem P; the complexity of the fastest sequential algorithm known for P; and the parallel algorithm A for problem P.

$EPU(P,A) =$

$$\frac{\text{complexity of the fastest sequential algorithm for P}}{\text{number of PEs used by A} * \text{complexity of A}}$$

For the case of procedure SUMi,

$$EPU=O\left(\frac{n}{n/2 * \log n}\right) = O\left(\frac{1}{\log n}\right)$$

Note that $0 \leq EPU \leq 1$ and that an EPU close to 1 is considered 'good'. For the case of computing $\sum_{i=1}^n A(i)$, we can actually arrive at an $O(\log n)$ algorithm with an EPU of $O(1)$ (i.e., using only $\lfloor n/\log n \rfloor$ PEs) [32]. This is done by dividing the n $A(i)$ s into $\lfloor n/\log n \rfloor$ groups, each group containing at most $\lfloor \log n \rfloor$ of the $A(i)$ s. Each of these groups is assigned to a PE which sequentially computes the sum of the numbers in the group. This takes $O(\log n)$ time. Now, we need to sum up these $\lfloor n/\log n \rfloor$ group sums. Procedure SUMi can be used to compute this sum in $O(\log n)$ time.

Note that the discussion carried out so far concerning the computation of $\sum_{i=1}^n A(i)$ applies just as well to the computation of $\bigoplus_{i=1}^n A(i)$ where \bigoplus is any associative operator (for example, \max , \min , $*$, etc). Hence, $\max_{1 \leq i \leq n} \{A(i)\}$; $\min_{1 \leq i \leq n} \{A(i)\}$; $\prod_{1 \leq i \leq n} A(i)$; etc can all be computed in $O(\log n)$ time using $\lfloor n/\log n \rfloor$ PEs.

Suppose that instead of computing just $\sum_{i=1}^n A(i)$, we wish to compute $S_j = \sum_{i=1}^j A(i)$, $1 \leq j \leq n$. We shall refer to this problem as the partial sums problem. When computing S_n using the sequential algorithm, we obtain S_i , $1 \leq i \leq n$ as a by-product and so, in this case, no additional effort need be expended. In the case of procedure SUMi (and its refinement to the case of $\lfloor n/\log n \rfloor$ PEs), however, all the S_i s are not computed during the computation of S_n . Following the computation of S_n , the remaining S_i s can be obtained by making one pass down the binary tree. In this pass each node transmits to its children the sum of the values to the left of the child.

Let $A(i:i) = (1, 1, 2, 3, 1, 2, 1, 2, 3, 4, 2)$. The computation tree of Figure 1.3 is redrawn in Figure 1.5.

The index of each node appears outside it. Inside each node there are two numbers. The upper number is V as defined for procedure SUM1. The lower number in each node is L ; where for any node i , L is defined as:

$$L(i) = \begin{cases} 0 & i=1 \\ L(i/2) & i \text{ is even} \\ L(i/2)+V(i-1) & i \text{ is odd} \end{cases}$$

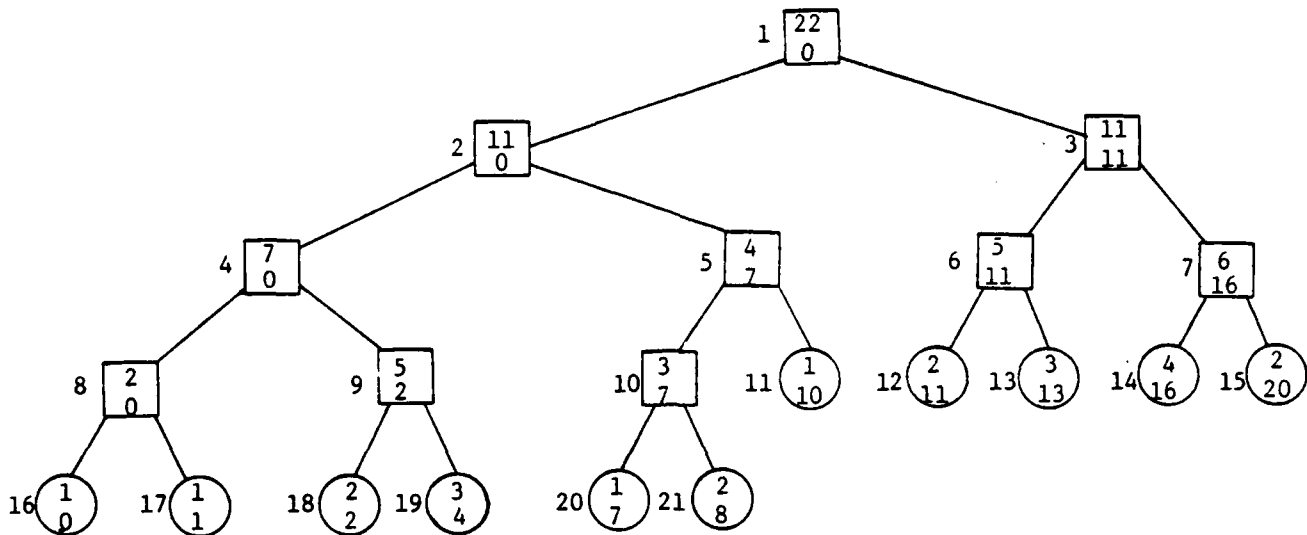


Figure 1.5

One may easily verify that if i is a circular node representing $A(j)$, then $L(i) = \sum_{i < p < j} A(p)$. Hence, from the L values of the circular nodes, one can easily obtain all the partial sums. Our first algorithm for the partial sums problem is PSUM1 (Figure 1.6). This algorithm simply computes the $V(i)$ s in the first pass and the $L(i)$ s in the second. Finally, the S values are computed.

As in the case of SUM1, the parallelism of lines 4 and 8 requires only $n/2$ PEs. Using $n/2$ PEs, line 2 can be done in two steps. Actually, procedure PSUM1 can be run in $O(\log n)$ time using only $\lfloor n/\log n \rfloor$ PEs. The idea here, is the same as that for SUM1.

The perfect shuffle connection scheme seems to be well suited to the binary tree method as it contains an underlying complete binary tree. If we let the PEs represent nodes in a complete binary tree, then the left child of PE i is PE $2i$, and the right child is PE $2i + 1$. Since $2i =$

```

line  procedure PSUMi (A,n,S)
        //compute S(i) =  $\sum_{j=i}^n A(j)$ ,  $i \leq n$ //
1       k <- |log2 n|; j <- 2k; t <- 2*(n-j); p <- n-i
2       V(p+i) <- A((i+t-1) mod n + 1), i < n
3       for i <- k down to 0 do //add by levels//
4         V(j) <- V(2j) + V(2j+1), 2i ≤ j ≤ min{p, 2i+1-1}
5       end for
        //compute Ls//
6       L(i) <- 0
7       for i <- 1 to k+i do //compute L by levels//
8         L(j) <- if j even then L(j/2)
                    else L(j/2)+V(j-i)
                    endif
                    2i ≤ j ≤ min{n, 2i+1-1}
9       end for
10      S((i+t-1) mod n + 1) <- L(p+i)+V(p+i), i < n
11      end PSUMi

```

Figure 1.6

SHUFFLE(i), and $2i + 1 = \text{EXCHANGE}(\text{SHUFFLE}(i))$; the downward pass is easily carried out. Also, $\text{PARENT}(i) = \text{UNSHUFFLE}(i)$, i even and $\text{PARENT}(i) = \text{UNSHUFFLE}(\text{EXCHANGE}(i))$, i odd. So the complexity analysis for SUMi, and PSUMi hold even when a PSC is used. For a binary tree with n leaves, a PSC with $n-1$ PEs is needed, however.

By using a slightly different computation tree and rearranging the order of computation, one can arrive at a one pass algorithm for the partial sums problem. Let $A(0:n-1)$ be the n numbers to be added. Let $S(0:n-1)$ denote the partial sum array. A 2^k -block of array elements consists of all array elements whose indices differ only in the least significant k bits. The 2^1 -blocks of $A(0:10)$ are $[0,1], [2,3], [4,5], [6,7], [8,9]$, and $[10]$; the 2^2 -blocks are $[0,1,2,3], [4,5,6,7]$, and $[8,9,10]$; etc. Two 2^k -blocks are sibling blocks iff their union is a 2^{k+1} -block. Thus, $[0,1]$ and $[2,3]$ are sibling blocks; so also are $[0,1,2,3]$ and $[4,5,6,7]$. However, $[2,3]$ and $[4,5]$ are not sibling blocks. The one pass algorithm computes S by first computing the partial sums for all 2^0 -blocks of A . In this case, $S(i)=A(i)$. Next, S is computed for all 2^1 -blocks; then for all 2^2 -blocks; ...; and finally for the single 2^q -block where $q=|\log_2 n|$.

Let X and Y be two sibling 2^k -blocks. Let X be the block containing all elements with bit k equal to 0. The union of X and Y is a 2^{k+1} -block. Relative to this 2^{k+1} -block, the S values for elements of X are the same as with respect to the corresponding 2^k -block. The S values for elements in Y however change by the sum of the A elements corresponding to the 2^k -block X . Figure 1.7 gives the S values and 2^k -blocks when S values are computed by blocks as

described above. Blocks are enclosed in brackets.

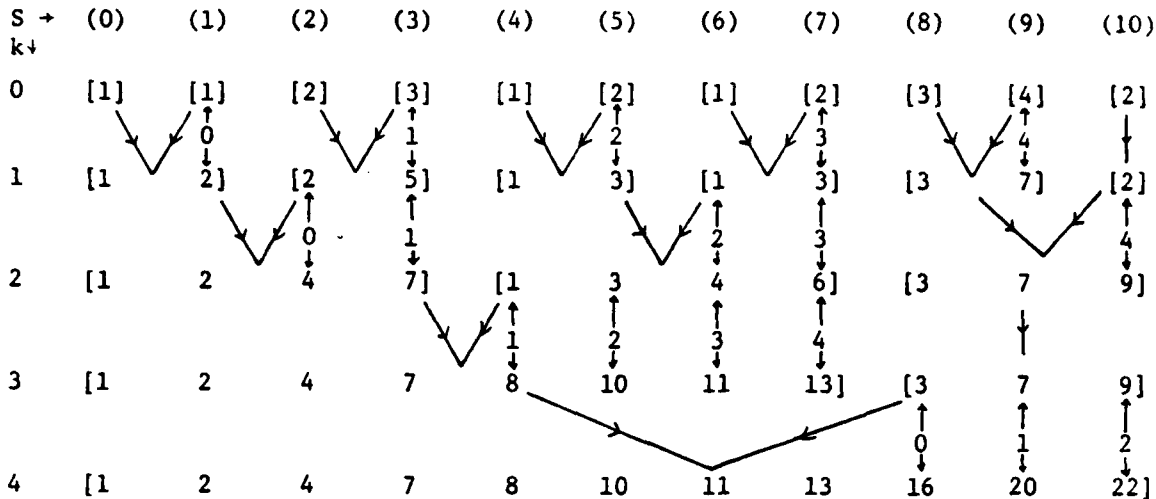


Figure 1.7 Computing S by blocks

The updating of S when going from one block size to the next is easily performed if we keep track of the sum of the A(i)s in each 2^k -block. For this purpose, we use an auxiliary array T. T(i) for i in a given 2^k -block (except possibly the rightmost 2^k -block) is the sum of all the A(i)s in that block. Before we can formally specify the partial sums algorithm, we need a processor assignment scheme. Figure 1.7 shows a processor assignment scheme for our example. Processors are assigned only to compute the S values that change. Thus, when $k=0$, PE(0) computes S(0); PE(1) computes S(1); PE(2) computes S(2); PE(3) computes S(3); PE(4) computes S(4); PE(5) computes S(5); PE(6) computes S(6); PE(7) computes S(7); PE(8) computes S(8); PE(9) computes S(9); and PE(10) computes S(10). When $k=1$, PE(0) computes S(1); PE(1) computes S(2); PE(2) computes S(3); PE(3) computes S(4); PE(4) computes S(5); PE(5) computes S(6); PE(6) computes S(7); PE(7) computes S(8); PE(8) computes S(9); and PE(9) computes S(10). PEs 10 and 11 are idle when $k=1$. Let ... $i_3 i_2 i_1 i_0$ be the binary representation of i. The PE assignment rule is obtained by defining the function $f(i, j) = \dots i_{j+1} i_j \dots i_{j-1} \dots i_0$. For any k, PE(i) computes $S(f(i, k) + 2^k)$ provided that this index of S is no more than $n-1$. The one pass partial sums algorithm is stated as procedure PSUM2 (Figure 1.8). PSUM2 uses $\lfloor n/2 \rfloor$ PEs indexed 0 through $\lfloor n/2 \rfloor - 1$.

It should be easy to see that our earlier ideas regarding the use of only $\lfloor n/\log n \rfloor$ PEs carry over to the case of PSUM2. So, PSUM2 can be modified to obtain an $O(\log n)$ one pass algorithm using only $\lfloor n/\log n \rfloor$ PEs. For the modified algorithm, $EPU=O(1)$.

```
line procedure PSUM2 (A,S,n)
  //one pass partial sums//
  1 declare A(0:n-1),S(0:n-1), T(0:n-1)
  2 for each PE(i) do in parallel
    //initialize S and T for 20-blocks//
  3   j<- f(i,0)
  4   S(j)<- T(j)<- A(j)
  5   S(j+1)<- T(j+1)<- A(j+1)
  6   for k<- 0 to  $\lfloor \log_2 n \rfloor$  -1 do
    //combine 2k-blocks//
  7     j<- f(i,k)
  8     if j+2k<n then
  9       S(j+2k)<- S(j+2k)+T(j)
10      T(j+2k)<- T(j+2k)+T(j)
11      T(j)<- T(j+2k)
12     endif
13   end for
14 end for
15 end PSUM2
```

Figure 1.8 One pass partial sums algorithm

2. Parallel Scheduling Algorithms

In this section, we develop fast parallel algorithms for a variety of scheduling problems. Each of these algorithms is arrived at using the binary tree method of section 1. We shall refrain from providing explicit formal statements such as those of Figures 1.4, 1.6, and 1.8, of these algorithms. Instead, we shall describe the algorithms informally and illustrate them with an example. One should note that we are interested in both the complexity as well as the EPU of the algorithms developed.

All the scheduling problems to be discussed assume that n jobs have to be scheduled on m identical machines. Associated with job i is a four-tuple (r_i, d_i, p_i, w_i) where r_i is its release time; d_i is its due time; p_i is its processing requirement; and w_i is its weight, $1 \leq i \leq n$. The processing of no job can commence until its release time. No job can be scheduled for processing on more than one machine at any time instant. Job i is completed after it has been processed for p_i time units. If a job does not complete by its due time, it is tardy. In a nonpreemptive schedule, job i is scheduled to process on a single machine from some start time s_i to the completion time $s_i + p_i$, $1 \leq i \leq n$. In a preemptive schedule it is permissible to split the processing of jobs over machines as well as over non-adjacent time intervals.

2.1 Minimizing Maximum Lateness

Let S be a schedule for the n jobs (r_i, d_i, p_i, w_i) . Let c_i be the completion time of job i . The lateness of job i is defined to be $c_i - d_i$. The maximum lateness, L_{\max} , is $\max\{c_i - d_i\}$. We wish to obtain an m machine nonpreemptive schedule that minimizes L_{\max} . This problem is known to be NP-hard [22]. So, we shall consider only special cases of this problem, i.e., cases for which a polynomial time sequential algorithm is known. Specifically, we shall consider the following cases: (i) $p_i = 1$, $1 \leq i \leq n$ and all release times are integer; (ii) $m = 1$ (i.e., the number of machines is 1) and preemption is allowed; and (iii) cases (i) and (ii) with precedence constraints. These three cases are considered in sections 2.1.1, 2.1.2, and 2.1.3 respectively. Since the weights w_i play no part in the L_{\max} problem, we shall only consider triples (r_i, d_i, p_i) in these subsections.

2.1.1 $p_i = 1$, $1 \leq i \leq n$ and all release times are integer.

Jackson [16] has shown that when $m = 1$ and all jobs have the same release time, L_{\max} is minimized by scheduling the jobs in nondecreasing order of due times. Horn [14] and Baker and Sue [3] have generalized this method to the case when $m = 1$ and all jobs do not have the same release time. An optimal one machine schedule is now obtained by assigning jobs to time slots, one slot at a time starting at time 0. When we are considering the time slot $[i, i+1]$, we select a job with least due time from among the set of available jobs. (The set of available jobs consists of all jobs not yet selected that have a release time less than or equal to i .) If this set is empty, then this slot is left idle. This strategy can be implemented to run in $O(n \log n)$ time on a single processor computer. Blazewicz [6] has extended this idea to the general case, $m > 1$. His algorithm also schedules by time slots. Let J be the set of jobs available when slot $[i, i+1]$ is to be scheduled. If $|J| \leq m$ then all the available jobs are processed in $[i, i+1]$. If $|J| > m$, then we select m jobs with least due times.

In developing the parallel algorithm, we first consider the case $m = 1$. The algorithm of Horn is readily seen to be highly sequential. No decision concerning time slot $[i, i+1]$ can be made unless we know the jobs that are available at this time. This of course depends on which jobs were selected for the earlier time slots. So, a straightforward adaptation of Horn's algorithm would need n steps (one for each time slot). The overall complexity of the resulting parallel algorithm would be $O(n)$. This is not very good. We are really interested in algorithms with complexity $O(\log^k n)$ for some k .

Despite the highly sequential nature of Horn's method, his idea can be used to arrive at a parallel algorithm with complexity $O(\log^2 n)$. This is accomplished using the binary tree method. It is helpful to consider an example. Suppose we have 14 jobs with r_i , and d_i as specified in Figure 2.1(a). The first step in our proposed parallel algorithm is to sort the jobs by release times (into nondecreasing order). Jobs with the same release time are sorted into nondecreasing order of due time. Let R_1, R_2, \dots, R_k be the k distinct release times of the n jobs ($R_1 < R_2 < \dots < R_k$). Let $R_{k+1} = \infty$. For our example, the sorted sequence of jobs is shown in Figure 2.1(b); $k=4$; and $R_1=2, R_2=5, R_3=6, R_4=9$, and $R_5=\infty$.

i	1	2	3	4	5	6	7	8	9	10	11	12	13	14
r_i	5	2	2	5	2	2	6	2	5	6	9	9	9	9
d_i	8	7	7	10	3	6	10	5	12	12	16	14	15	16

(a)

i	5	8	6	2	3	1	4	9	7	10	12	13	11	14
r_i	2	2	2	2	2	5	5	5	6	6	9	9	9	9
d_i	3	5	6	7	7	8	10	12	7	17	11	15	16	16

(b)

Figure 2.1

Next, a binary computation tree is associated with the problem. The tree used is the unique complete binary tree with k leaves. With each node in this tree, we associate a time interval (t_L, t_R) . Assume that the leaf nodes are numbered i through k , left to right. The i th leaf node has associated with it the interval (R_i, R_{i+1}) , $i < i \leq k$. The interval (t_L, t_R) associated with a nonleaf node, N , is obtained from the intervals associated with the two children of this node. $t_L(N) = t_L$ (left child of N) and $t_R(N) = t_R$ (right child of N). For our example, the binary computation tree together with time intervals is shown in Figure 2.2.

A schedule that minimizes L_{\max} is now obtained by making two passes over this computation tree. The first pass is made level by level towards the root; the second is made level by level from the root to the leaves. Let P be any node in the computation tree. Let the interval associated with P be (t_L, t_R) . The set of available jobs, $A(P)$ for P

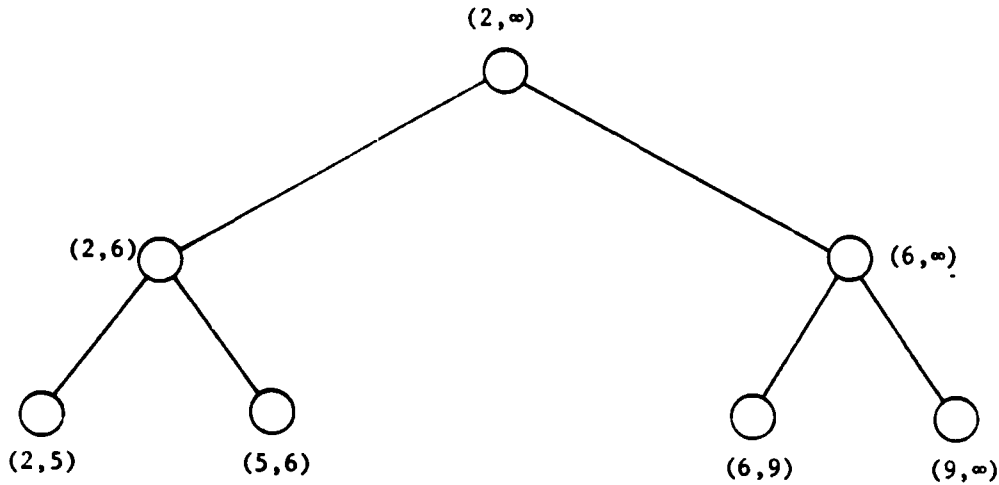


Figure 2.2 Computation tree for the example of Figure 2.1.

consists exactly of those jobs that have a release time r_i such that $t_L < r_i < t_R$. This set of jobs may be partitioned into two subsets, respectively called the used set and the transferred set. The set of used jobs consists exactly of those available jobs that will be scheduled between t_L and t_R for the L_{max} problem defined by the job set $A(P)$. The remaining jobs in $A(P)$ make up the transferred set. For our example, the set of available jobs for the node representing the interval $(2, 6)$ is $\{5, 8, 6, 2, 3, 1, 4, 9\}$. If Horn's algorithm is used on this set of jobs, then jobs 5, 8, 6, and 2 will get scheduled in the interval from 2 to 6. Hence, the used set is $\{5, 8, 6, 2\}$ and the transferred set is $\{3, 1, 4, 9\}$.

In the first of the two passes mentioned above, the used and transferred sets for each of the nodes in the computation tree are determined. For a leaf node the used and transferred sets are determined by directly using Jackson's rule. If P is a leaf node for the interval (t_L, t_R) , then the used set is obtained by selecting jobs from the available job set $A(P)$ for P in nondecreasing order of due times. Since jobs with the same release time have already been sorted by due times, the used set consists of the first $\min\{|A(P)|, t_R - t_L\}$ jobs in $A(p)$. The remaining jobs form the transferred set. For our example, for the interval $(2, 5)$, the set of used jobs is $\{5, 8, 6\}$ while the set of transferred jobs is $\{2, 3\}$; for the interval $\{5, 6\}$, the used set is $\{1\}$ and the transferred set is $\{4, 9\}$; etc. Figure 2.3 shows the used and transferred sets for each of the leaf nodes for our example. The solid vertical line separates the used jobs from the transferred jobs.

For a nonleaf node, the used and transferred sets may be computed from the used and transferred sets of its children. Let P be a nonleaf node and let U_L , U_R , T_L , and T_R be the used and transferred sets for its left and right children respectively. Let (t_L, t_P) , (t_L, t_R) , and (t_L, t_R) be the intervals, respectively, associated with node P , its left child, and its right child. Clearly, $t_L = t_L$; $t_R = t_R$; and $t_P = t_P$. It should be clear that if Horn's algorithm is used to schedule the available jobs $A(P)$ then the jobs in U_L will be the ones scheduled from t_L to t_P . The set of jobs scheduled from t_P to t_R will be some subset of $T_L \cup U_R$. Let Q denote the $\min\{|T_L \cup U_R|, t_R - t_P\}$ jobs of $T_L \cup U_R$ that have least due times. It is not too difficult to see that Q is the subset of $A(P)$ that is scheduled by Horn's algorithm in the interval t_P to t_R . Hence, the used set for P is $U_L \cup Q$ and the transferred set is $T_R \cup A(P) - Q$. Observe that if U_L , U_R , T_L , and T_R are in nondecreasing order of deadlines, then the set Q can be obtained by merging together U_L and T_L and selecting the first $\min\{|T_L \cup U_R|, t_R - t_P\}$ jobs from the merged list. Q can next be merged with U_R to obtain the used set in nondecreasing order of due times. Another merge yields the transferred set in nondecreasing order of due times. Figure 2.3 gives the used and transferred sets in nondecreasing order of due times for all nodes in our example computation tree.

In the second pass, the used sets are updated so that the used set for a node representing the interval (t_L, t_R) is precisely the subset of jobs (from amongst all n jobs) that is scheduled in this interval by Horn's algorithm when solving the L_{\max} problem for the entire job set. This is done by working down the computation tree level by level starting with the root. The used set for the root node is unchanged in this pass. If P is a node whose used set been updated then the used sets for the left child and the right child of P are obtained in the following way. Let the interval associated with P be (t_L, t_R) and let the interval associated with its left child be (t_L, t_P) . Let V be the subset of the used set of P consisting solely of jobs with a release time less than t_L . Let U be the current used set (i.e. the one computed in the first pass) for the left child of P . Let W be the set obtained by merging U and V (note that U and V are disjoint and that both are ordered by due times). The new used set, U^1 , for the left child of P consists of the first $\min\{|W|, t_P - t_L\}$ jobs in W . The used set for the right child of P consists of all jobs in the used set for P that are not included in U^1 .

Let us now go through this second pass on our example. Let P be the root node. $(t_L, t_R) = (2, 00)$ and $V = \emptyset$. Hence, the new used set for the left child of P is simply its old used set. The used set for the right child of P becomes $\{3, 7, i, 4, i2, 9, i3, ii, i4, i0\}$. Now, let P be the right child of the root. $(t_L, t_R) = (6, 00)$; $V = \{3, i, 4, 9\}$; $W = \{3,$

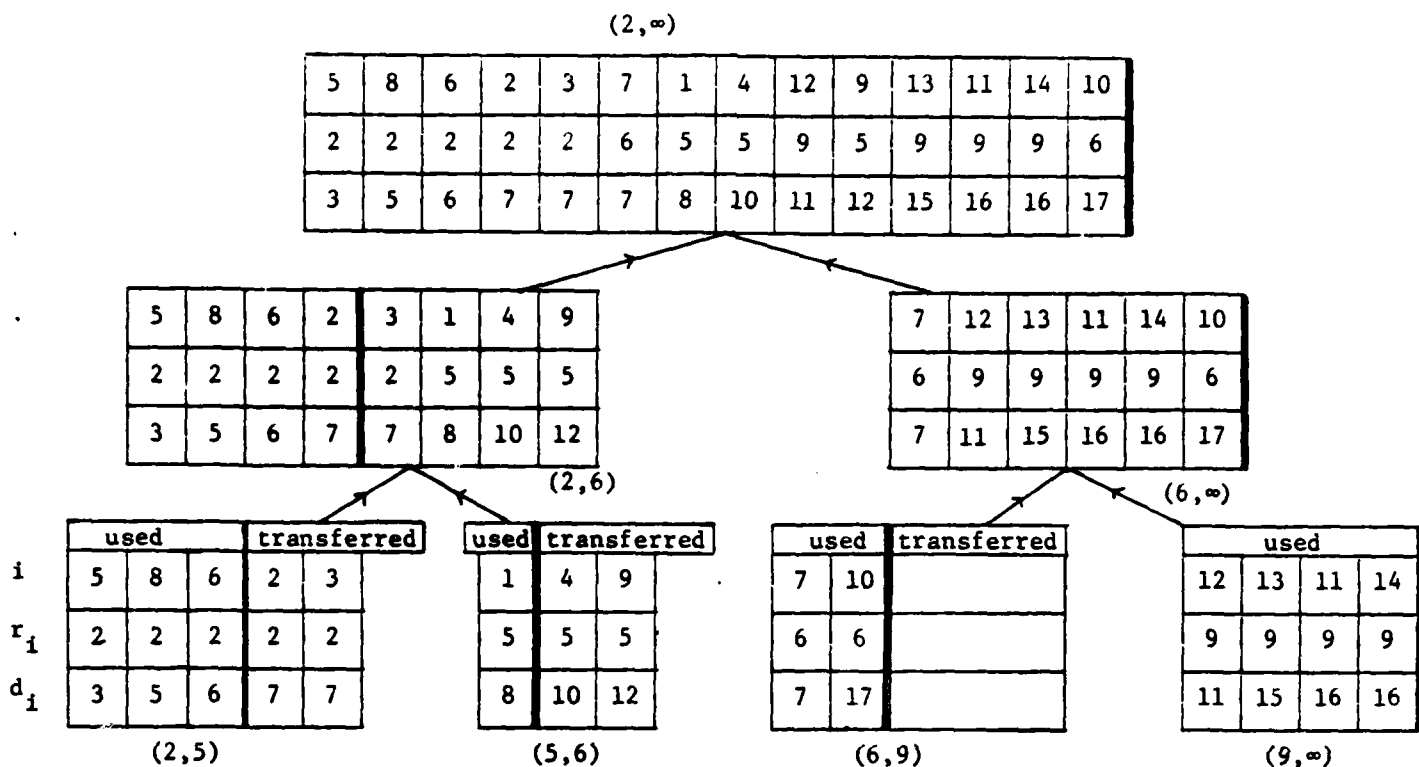


Figure 2.3 First pass of the L_{\max} algorithm

$7, 1, 4, 9, 10$ }. The new used set for the left child of P is $\{3, 7, 1\}$. The new used set for the right child of P is $\{4, 12, 9, 13, 11, 14, 10\}$. Figure 2.4 shows the new used sets for all the nodes in the computation tree.

From the definition of an updated used set, it follows that the schedule defined by the leaf nodes (for our example, this is: job 5 at time 2, job 8 at time 3, job 6 at time 4, job 2 at time 5, etc.) minimizes L_{\max} . The correctness of the node updating procedure is easily seen. If P is the root node, then it represents the interval (R, ∞) . All jobs are necessarily scheduled in this interval by Horn's algorithm. Hence, the updated used set for this node consists of all n jobs. Now, let P be any nonleaf node for which we have obtained the updated used set. Assume that this is in fact the correct updated used set, i.e., it consists exactly of those jobs scheduled by Horn's algorithm in that interval. We shall show that the updating procedure gives the correct used sets for the left and right child of P . Let t_L, t_R, t_P, V, W, U , and U^1 be as defined in the updating procedure. Let X be the used set for P . From the way the first pass works, it follows that only jobs from $N = U \cup U^1 \cup V$ are candidates for scheduling by Horn's algorithm, in

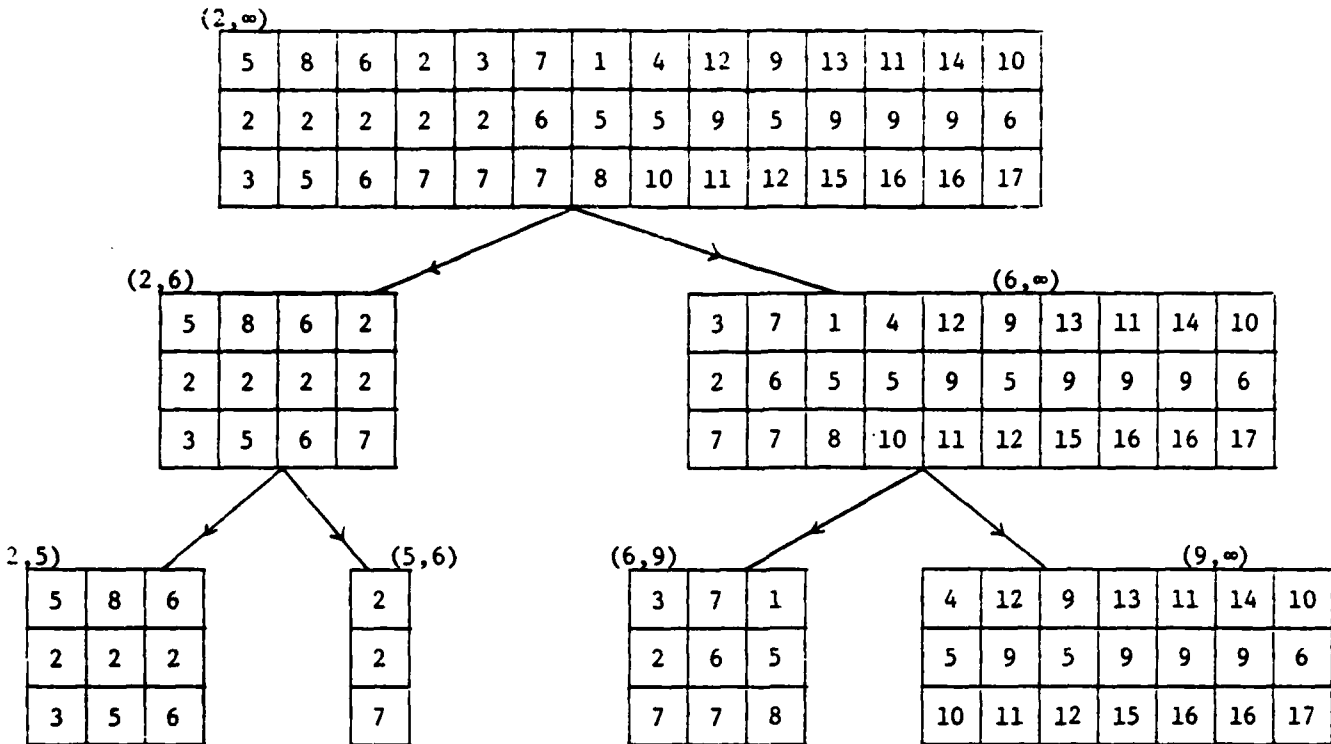


Figure 2.4 Results of second pass.

the interval (t_L, t_R^i) . It is a simple matter to see that only $\min\{|W|, t_L^i, t_R^i\}$ of these can be scheduled in this interval; further these jobs are selected in nondecreasing order of due times. Hence, U^1 is correctly computed. From this it follows that the used set for the right child must be $X - U^1$.

Having established the correctness of our parallel procedure, we are ready to determine its complexity as well as the required number of PEs. The first step₂ consists of sorting the jobs. This can be done in $O(\log^2 n)$ time using $\lfloor n/2 \rfloor$ PEs [4]. In both the first and second passes over the computation tree we are essentially performing a fixed number of merges of ordered sets at each node. Using Batcher's bitonic merge scheme ([4], [18]), a p element ordered set can be merged with a q element ordered set using $\lfloor (p+q)/2 \rfloor$ PEs in $O(\log(p+q))$ time. Hence, the overall complexity of our parallel L_{\max} algorithm is $O(\log^2 n)$. The number of PEs used is $\lfloor n/2 \rfloor$. The EPU of this algorithm is $O(n \log n / (n/2 \log^2 n)) = O(1/\log n)$.

Our parallel L_{\max} algorithm for the case $m=i$ easily generalizes to the case $m>i$. The two passes over the computation tree are changed so that all uses of $t_R - t_L$ and $t_R^1 - t_L^1$ are replaced by $m(t_R - t_L)$ and $m(t_R^1 - t_L^1)$ respectively. The schedule is obtained from the updated used sets of the leaf nodes. The i th job in this used set is assigned to the $i \bmod m + 1$ th machine.

2.1.2 $m=i$ and preemptions permitted

Horn's [14] algorithm for this problem is quite similar to the sequential algorithm for the case discussed in section 2.1.1 and also has a sequential complexity that is $O(n \log n)$. A schedule with minimum L_{\max} is obtained by starting at the first release time and considering an available job, i , with least due time. Let the processing time of this job be p . Let the time to the next release time be t and let the current time be T . Job i is scheduled from T to $T + \min\{p, t\}$. The current time changes from T to $T + \min\{p, t\}$ and the remaining processing time for job i becomes $p - \min\{p, t\}$. Next, from the available job set at the current time T a job with minimum due time is selected for processing, and so on.

The parallel algorithm of section 2.1.1 can be adapted to this case. Jobs are sorted as before and two passes are made over the tree. In the first pass, used and transferred sets are computed for each node. In the second pass, the used sets are updated. For the first pass, the used and transferred sets for the leaf nodes are obtained by computing the partial sum sequence for the ordered set of available jobs for each leaf (see the algorithm of Figure 1.8). Next, for each leaf we determine the first partial sum, j , (if any) that exceeds the value of $t_R - t_L$ for that node. If there is no such partial sum, then all the available jobs are used. If there is, then the used set consists of jobs $1, 2, \dots, j-1$ together with a fraction, f , of job j . This fraction is chosen such that the sum of the processing times of jobs $1, 2, \dots, j-1$ and f times that of job j equals $t_R - t_L$. The transferred set consists of $1-f$ of job j together with the remaining jobs.

For nonleaf nodes, the used and transferred sets are computed from the corresponding sets for the left and right children. Let P be a nonleaf node. Let Q and S be its left and right children respectively. The used set for P is obtained by merging (according to due times) the transferred set of Q with the used set of S , to obtain W . The partial sums for W are computed and W is partitioned into W_1 and W_2 such that the sum of the processing times for the jobs in W_1 equals $\min\{\text{sum of processing times in } W, t_R^2 - t_L^2\}$ where (t_L^2, t_R^2) is the interval associated with node S . Observe that

this partitioning of W may require us to split one of the jobs in W in the same way as was done for leaf nodes. The used set for P is obtained by merging together W_1 and the used set for Q . The transferred set for P is obtained by merging together W_2 and the transferred set for S .

The updating of the second pass is also carried out in a manner similar to that used in section 2.1.1. The updated used set for the root node consists of all n jobs. Let P be a node for which the updated used set has been computed. Let (t_L, t_R) be the interval associated with P . Let Q and S , respectively, be the left and right children of P . Let the interval associated with Q be (t_L, t_R^1) . Define V to be the set of all jobs in the used set of P that have a release time less than t_L . Merge V and the current used set of Q together. Let the resulting ordered set be W . Compute the partial sums for W and partition W into W_1 and W_2 as was done in the first pass. Once again, it may be necessary to split a job into two to accomplish this. The used set for Q is W_1 . The remaining jobs in the used set of P (including possibly a remaining fraction of a job that went into W_1) constitute the used set for S .

Once the updated used sets for the leaves have been computed, a schedule minimizing L_{\max} is obtained by scheduling the used sets of the leaves in the intervals associated with them. For each such interval, the scheduling is in non-decreasing order of due time.

The correctness of the algorithm described above follows from the correctness of Horn's algorithm and the discussion in section 2.1.1. The algorithm can be run in $O(\log^2 n)$ time using at most $3n/2$ PEs. Note that because jobs may split, we may at some level have a total of $n+2k$ jobs (or job parts). Recall that k denotes the number of distinct release times and that at each node at most one additional job split can occur. Because of the effective increase in number of jobs, more than $\lfloor n/2 \rfloor$ PEs are needed here, while only $\lfloor n/2 \rfloor$ were needed in section 2.1.1. The EPU is still $O(1/\log n)$.

Example 2.1: Figure 2.5 gives an example job set. Since the jobs are already in the order desired, we may begin directly with the first pass over the computation tree. Figure 2.6 gives the result of the first pass. Figure 2.7 gives the result of the second pass. []

Sorted input

i	1	2	3	4	5	6	7	8	9	10
P _i	2	2	6	6	6	6	11	20	20	20
d _i	4	25	8	14	22	25	14	21	24	30
P _i	2	1	2	8	4	1	2	1	3	8

Figure 2.5

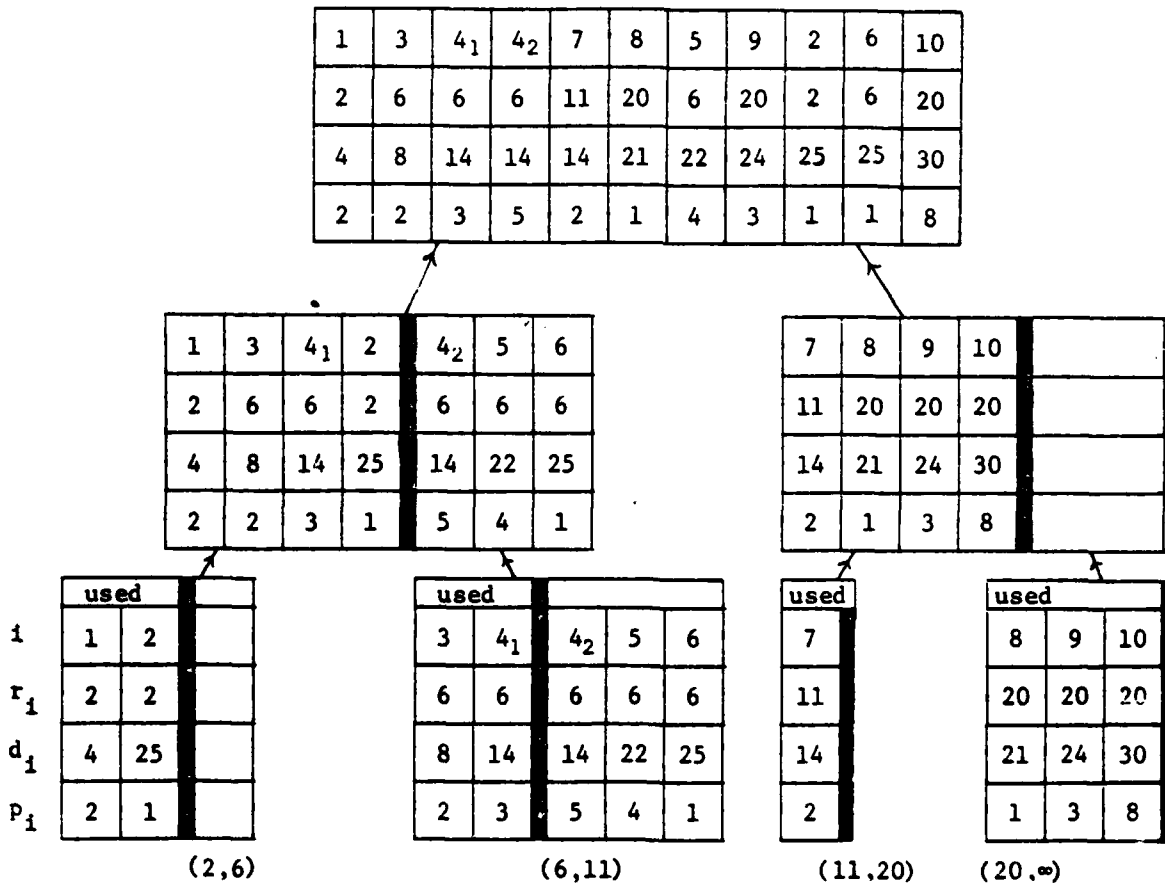


Figure 2.6

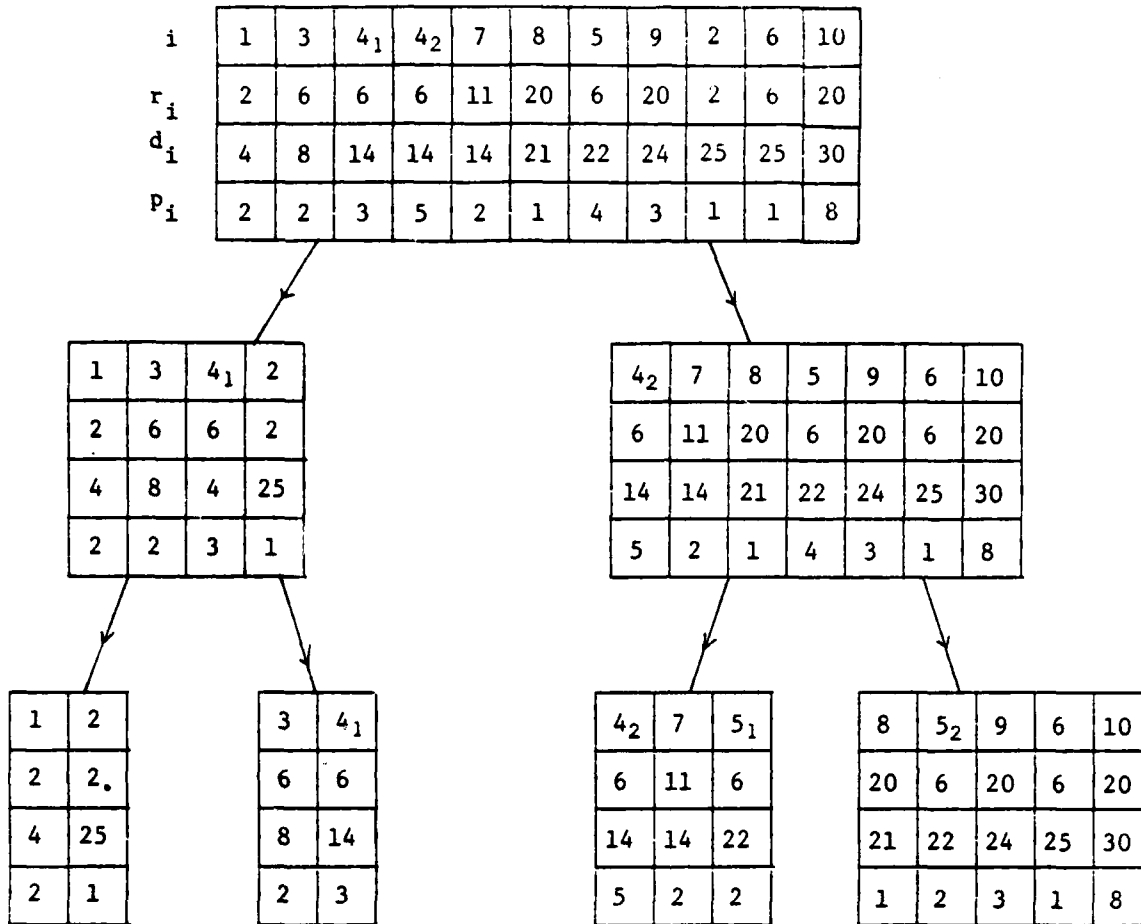


Figure 2.7

2.1.3 Precedence Constraints

Suppose that the set of jobs to be scheduled defines a partial order $<$. $i < j$ means that the processing of job j cannot commence until the processing of job i has been completed. Let (r_i, d_i, p_i) be the release, processing, and due times of job i . Modify the release and due times as below:

$$r_i^i = \max\{r_i, \max_{i < j} \{r_j + p_j\}\}$$

$$d_i^i = \max\{d_i, \max_{i < j} \{d_j - p_j\}\}$$

Rinooy Kan [31] has observed that a schedule minimizing L_{\max} when $p_i=1$, the r_i s are integer, and $<$ is a partial order can be obtained by simply using Horn's algorithm (cf. section

2.1.1) on the jobs $(r_i^i, p_i=i, d_i^i)$, $1 \leq i \leq n$ with no precedence constraints. Since the modified release and due times can be computed in $O(\log^2 n)$ time using the critical path algorithm of [9], a schedule minimizing L_{\max} in the presence of precedence constraints can be obtained in $O(\log^2 n)$ time ($m=i, p_i=i$). The number of PEs needed by the algorithm of [9] is $n^3/\log_2 n$, so the EPU of the resulting algorithm is $O(n^2 \log n / (n^3 \log_2 n)) = O(1/(n \log n))$.

When $m=i$, a partial order $<$ is specified, and preemptions are allowed, a schedule minimizing L_{\max} can be obtained by computing modified release and due times as above and then using the algorithm of section 2.1.2 on the modified jobs. The resulting algorithm has complexity $O(\log^2 n)$; uses $O(n^3/\log n)$ PEs; and has an EPU that is $O(1/(n \log n))$.

2.2 Minimizing Total Costs

Let (r_i, d_i, p_i, w_i) , $1 \leq i \leq n$ define n jobs. Let S be any one machine schedule for these jobs. The completion time c_i of job i is the time at which it completes processing. Job i is tardy iff $c_i > d_i$. The tardiness T_i of job i is $\max\{0, c_i - d_i\}$. When $p_i=i$, Horns [14] algorithm described in section 2.1.2 also finds a schedule that minimizes $\sum T_i$.

A schedule that minimizes $\sum w_i c_i$ when $p_i=i$, $1 \leq i \leq n$ and $m=i$ can be obtained by extending Smith's rule (see Rinnooy Kan [31]). Smith's rule [35] minimizes $\sum w_i c_i$ when $r_i=0$, $1 \leq i \leq n$. It essentially schedules jobs in nondecreasing order of \bar{p}_i/w_i . The extension to the case when $p_i=i$, $1 \leq i \leq n$ and the r_i 's may be different (but integer) works in following way. Scheduling is done time slot by time slot. From the set of available jobs for any slot, a job with least $1/w_i$ (or equivalently, maximum w_i) is selected and scheduled in this slot. This procedure is quite similar to that used for the L_{\max} problem with $p_i=i$ (see section 2.1.1). The only difference is that Smith's rule replaces the use of Jackson's rule 2.1.1. The used and transferred sets are now kept in nonincreasing order of weights.

Since the preemptive schedule obtained by the algorithm of section 2.1.2 also minimizes $\sum T_i$, this problem is easily solved in parallel. When $\sum c_i$ is to be minimized, $m=i$, and preemptions are permitted, the algorithm of section 2.1.2 can still be used. This time, however, the used and transferred sets are maintained in nondecreasing order of p_i rather than d_i [31].

Number of Tardy Jobs

Now, let us consider the problem of minimizing the number of tardy jobs when $m=1$ and all jobs have the same release time. Without loss of generality, we may assume that all jobs have a release time $r_i=0$. The fastest sequential algorithm for this problem is due to Hodgson and Moore [23]. It consists of the following three steps:

Step 1: Sort the n jobs into nondecreasing order of due times. Initialize the set R of tardy jobs to be empty.

Step 2: If there is no tardy job in the current sorted sequence, then append the jobs in R to this sequence. This yields the desired schedule. Stop.

Step 3: Find the first tardy job in the current sorted sequence. Let this be in position j . Find the job with the largest processing time from amongst the first j jobs in this sequence. Remove this job from the sequence and add it to R . Go to step 2.

The time complexity of the Hodgson and Moore algorithm is $O(n \log n)$. As in the case of the Hodgson and Moore algorithm, our parallel algorithm for this problem begins by sorting the jobs into nondecreasing order of due times. Within due times, jobs are sorted by p_i . Let D_1, D_2, \dots, D_k ($D_1 < D_2 < \dots < D_k$) be the k distinct due times associated with the n jobs. Let $D_0=0$. We next consider the unique complete binary tree that has exactly k leaves. If the leaf nodes of this tree are considered from left to right, then with the i th leaf we associate the interval (D_{i-1}, D_i) . The interval associated with a nonleaf node is (t_1, t_2) iff there exists t_3 such that (t_1, t_3) and (t_3, t_2) are the intervals, respectively, associated with its left and right children. If the interval (t_1, t_2) is associated with some node P , then all jobs with a due time d such that, $t_1 < d \leq t_2$ are associated with that node.

The set $J(P)$ of jobs associated with any node P may be partitioned into two sets $S(P)$ and $R(P)$. $S(P)$ and $R(P)$ are defined in the following way. Consider the problem of obtaining a schedule that minimizes the number of tardy jobs for $J(P)$ assuming that all jobs in $J(P)$ have a release time t_i ((t_1, t_2) is the interval associated with P). $S(P)$ is the set of non-tardy jobs in this schedule while $R(P)$ is the set of tardy jobs. It is well known [16] that if all jobs in $S(P)$ are scheduled in nondecreasing order of due times then no job in $S(P)$ will be tardy. From the definition of S and R , it is clear that $S(\text{root})$ defines the set of non-tardy jobs in a schedule for all n jobs that minimizes the number of tardy jobs. These jobs may be scheduled at the front of

the schedule in nondecreasing order of due times. The remaining jobs can be scheduled, in any order, after the jobs in $S(\text{root})$.

For a leaf node P , $S(P)$ and $R(P)$ are easily computed. First the partial sum sequence for $J(P)$ is obtained (recall that the jobs associated with P are in nondecreasing order of p_i). Let the interval associated with P be (t_1, t_2) . All jobs with a partial sum that is less than or equal to $t_2 - t_1$ are in $S(P)$. The remainder are in $R(P)$.

Let us consider an example. Figure 2.8(a) shows a set of 10 jobs. In Figure 2.8(b), these jobs have been ordered by due times and within due times by p_i . There are four distinct due times, and we have $D(\emptyset:4) = (0, 8, 15, 17, 25)$. Figure 2.9 shows the complete binary tree with four leaves. The interval associated with each node is also given. The S and R sets for each of the leaf nodes are also shown.

i	1	2	3	4	5	6	7	8	9	10
p_i	4	3	5	6	4	3	3	4	3	3
d_i	15	25	8	8	15	25	17	25	8	25

(a)

i	9	3	4	1	5	7	2	6	10	8
p_i	3	5	6	4	4	3	3	3	3	4
d_i	8	8	8	15	15	17	25	25	25	25

(b)

Figure 2.8

The computation of S and R for a nonleaf node P is done using the S and R sets of its left child Q and its right child T . Let the interval associated with Q and T , respectively, be (t_L, t_L^i) and (t_R, t_R^i) . It is clear that $S(T) \subset S(P)$ and that $R(Q) \subset R(P)$. To get the remaining jobs in $S(P)$, we merge together the jobs in $S(Q)$ and $R(T)$. Let the resulting ordered set be W . The partial sum sequence of the processing times of the jobs in W is next computed. Let V be the subset of W consisting of jobs that have a partial sum sequence no more than $t_R^i - t_L^i$. Let $X = W - V$. Clearly, $V \subset S(P)$. However, $V \cup S(T)$ may not equal $S(P)$ as it is possible for (at most) one of the jobs in X to also be in $S(P)$. To determine this job, we first determine for each due time D_i , $t_R^i \leq D_i < t_R^i$, a job in X that has least processing time

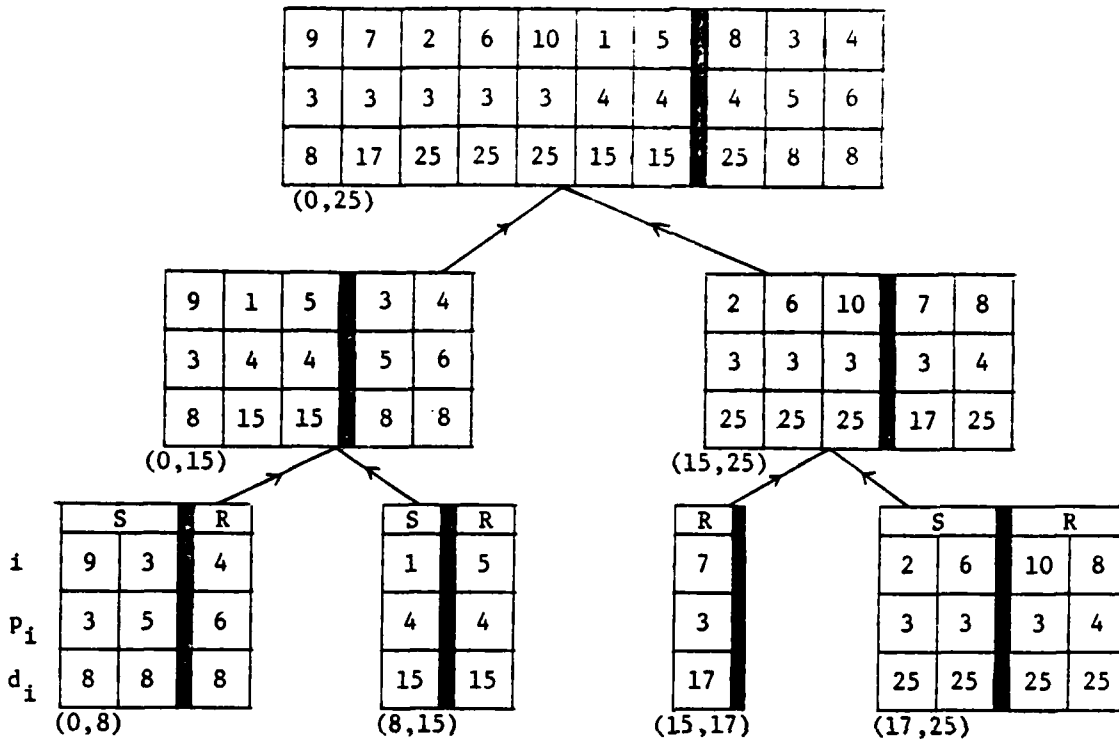


Figure 2.9

amongst all jobs in X with due time D_i . If there are no jobs in X with a certain due time D_i , then no job corresponding to this due time is selected. Let the set of jobs determined in this way be $U = \{J_1, J_2, \dots, J_q\}$. Let $\Delta = t_R^i - t_L^i - \sum p_i$. For each due time D_i , $t_R^i \leq D_i < t_R^i$, determine the sum of the processing times of all jobs in $S(T)$ with due times no more than D_i . Let this sum be Y_i . Let $\Delta_i = D_i - Y_i - t_R^i$. Now, compute $\gamma_i = \min_{j \geq i} \{\Delta_j\}$. It can be seen that the job (if any) in U with due time D_i can be in $S(P)$ only if its processing time is less than or equal to $\Delta + \gamma_i$. This information is used to remove from U those jobs that cannot possibly be in $S(P)$. From the remaining jobs, the job r with minimum processing time is selected and added to $S(P)$. $R(P) = R(Q) \cup (X - \{r\})$. The S and R sets for all nonleaf nodes in our example are specified in Figure 2.9.

The sets U and $\{\Delta_i\}$ can be computed in $O(\log n)$ time using $O(n)$ PEs if S and R are available in nondecreasing order of due times (so it is necessary to keep two copies of each S and R ; one ordered by processing times and one by due times). The γ_i s may be computed in $O(\log n)$ time using $O(n/\log n)$ PEs using a modified version of the partial sums

algorithm. Merging $S(Q)$ and $R(T)$ by processing times or by due times requires $O(\log n)$ time and $n/2$ PEs. So, all the work needed to be done at any level can be accomplished in $O(\log n)$ time with $O(n)$ PEs. The overall complexity of our parallel algorithm is therefore $O(\log^2 n)$ and its EPU is $O(i/\log n)$.

Job Sequencing With Deadlines

The problem of minimizing the sum of the weights of the tardy jobs is commonly referred to as the job sequencing with deadlines problem [15]. It is assumed that $r_i = 0$, and $p_i = 1$, $1 \leq i \leq n$. When the assumption $p_i = 1$ is not made, the problem is known to be NP-hard [17]. We shall now proceed to show how the binary tree method leads to an efficient parallel algorithm for this problem. We shall explicitly consider only the case $m=1$. When $m>1$, the problem can be transformed into an equivalent $m=1$ problem. Further, all the d_i 's are assumed to be integers.

An $O(n \log n)$ sequential algorithm for this problem appears in [15]. This algorithm builds an optimal schedule by first determining the set of jobs that are to be completed by their due times. This is done by considering the jobs in nonincreasing order of weights. The job currently being considered is added to the set of selected jobs iff it is possible to schedule this job and all previously selected jobs in such a way that all of them complete by their respective due times.

In our parallel algorithm, we begin by sorting the jobs by due times. Jobs with the same due time are sorted into nonincreasing order of weight. Figure 2.10(a) shows an example job set. Figure 2.10(b) shows that result of sorting this job set. Let the distinct due times be D_1, D_2, \dots, D_k ($D_1 < D_2 < \dots < D_k$). Let $D_0 = 0$. The computation tree to use is the unique complete binary tree with k leaves. Consider these leaves left to right. With leaf i , we associate the interval (D_{i-1}, D_i) , $1 \leq i \leq k$. Let P be a nonleaf node. Let the intervals associated with its left and right children, respectively be (t_L, t_R) and (t'_L, t'_R) . The interval associated with P is (t_L, t'_R) . The interval associated with the root is therefore $(0, D_k)$. Figure 2.11 shows the computation tree for our example. The interval associated with each node is also shown.

The set $J(P)$ of jobs associated with node P consists precisely of those jobs that have a due time d_i such that $t_L < d_i \leq t'_R$, where (t_L, t'_R) is the interval associated with P . With each node P , we may also associate two sets of jobs, $S(P)$ and $R(P)$. Consider the job sequencing with deadlines problem defined by the job set $J(P)$. Assume that all jobs

i	1	2	3	4	5	6	7	8	9	10	11	12	13	14
d _i	3	6	6	3	6	1	6	6	7	3	7	3	7	6
w _i	50	55	65	40	70	20	60	80	75	60	85	30	50	10

(a)

i	6	10	1	4	12	8	5	3	7	2	14	11	9	13
d _i	1	3	3	3	3	6	6	6	6	6	6	7	7	7
w _i	20	60	50	40	30	80	70	65	60	55	10	85	75	50

(b)

Figure 2.10

have a release time t_i . $S(P)$ consists exactly of those jobs in $J(P)$ that will be scheduled to finish by their due times in an optimal schedule for $J(P)$. $R(P)$ consists of the remaining jobs in $J(P)$. Once $S(\text{root node})$ is known, the optimal schedule for the overall job sequencing problem is also known.

For the leaf nodes, $S(P)$ and $R(P)$ are easily obtained. For each leaf node P , $S(P)$ consists of the $t_P - t_L$ jobs of $J(P)$ with largest weight (see Figure 2.11). If P is a non-leaf node, $S(P)$ and $R(P)$ are computed from the S and R sets of its children. Let Q and T , respectively, be the left and right children of P . Let the intervals associated with Q and T , respectively, be (t_L, t_Q) and (t_T, t_P) . Let $W = S(Q) \cup R(T)$ and let V be the set consisting of the $\min\{|W|, t_P - t_L\}$ jobs of W with largest weights. It is not too difficult to see that $S(P) = V \cup S(T)$. Hence $R(P) = J(P) - S(P) = R(Q) \cup (W - S(P))$. The S and R sets for each of the nodes in our example are also given in Figure 2.11.

Once the S and R sets have been computed, the optimal schedule can be obtained by sorting $S(\text{root})$ by due times and appending the jobs in $R(\text{root})$ to the end. For our example, the optimal schedule is 10, 8, 5, 3, 7, 11, 9, 2, 1, 13, 4, 12, 6, 14. The sum of the weights of the tardy jobs is 255.

Since the S and R sets are maintained in nonincreasing order of weights, the merging required at each node to compute S and R can be carried out using a parallel bitonic merge. Hence, all the computation needed at each level of the computation tree can be performed in $O(\log n)$ time using $n/2$ PEs. The overall complexity for our job sequencing with

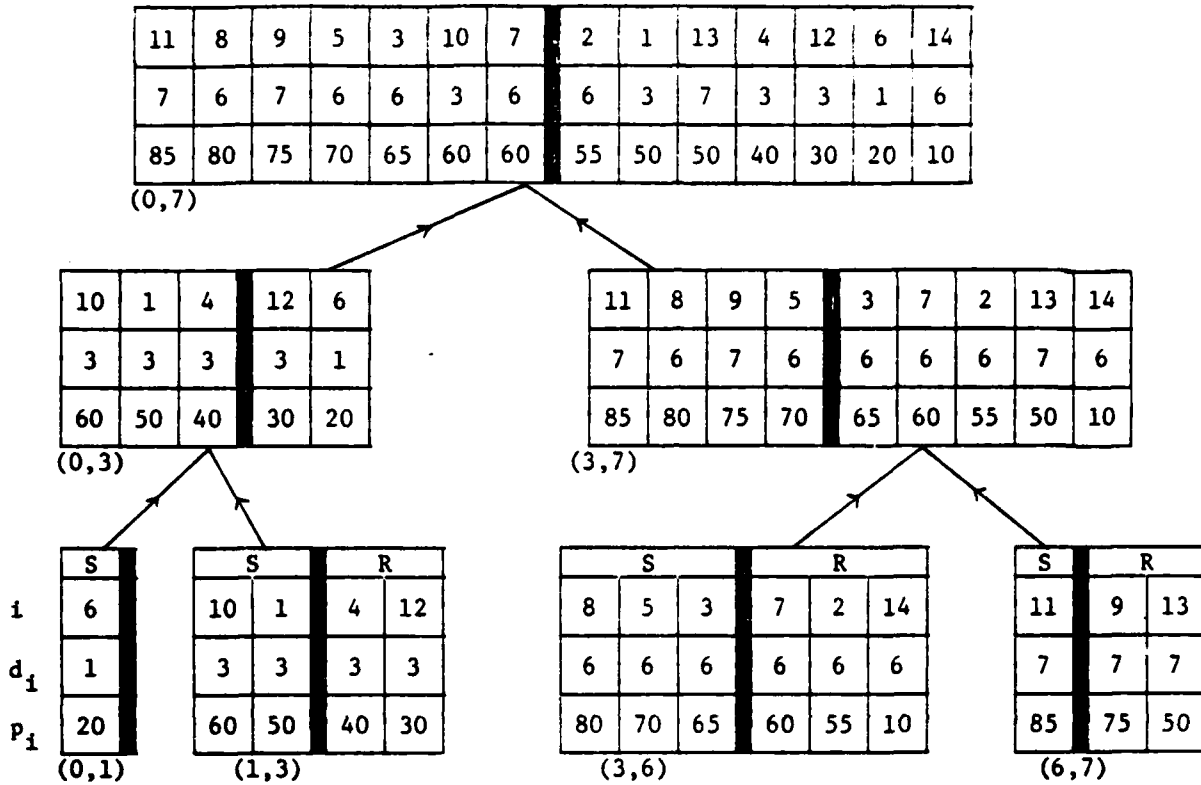


Figure 2.11

deadlines algorithm is $O(\log^2 n)$ and the EPU is $O(i/\log n)$. (In [10] Dekel and Sahni show how to solve the job sequencing problem in $O(\log n)$ time. This algorithm does not use the binary tree method and has an EPU which is considerably inferior to that of the algorithm developed here.)

Finally, we note that the parallel algorithm developed to minimize the number of tardy jobs when $m=1$ and $r_i=0$, can be adapted to obtain a one machine schedule that minimizes the sum of the weights of the tardy jobs provided that all jobs have agreeable weights. (All jobs have agreeable weights iff $p_i < p_j$ implies $w_i \geq w_j$ for all i and j .) The sequential algorithm for this problem is an extension of the Hodgson-Moore algorithm to minimize the number of tardy jobs. This extension is due to Lawler[21]. Also, Sidney's [34] extension which takes into account jobs that must necessarily be completed by their due times can also be solved by a modified version of our algorithm.

3. Conclusions

We have demonstrated that the binary computation tree is a very important tool in the design of efficient parallel algorithms. The binary tree method is closely related to the divide-and-conquer approach used to obtain many efficient sequential algorithms [15]. While divide-and-conquer algorithms do use an underlying computation structure that is a tree, the use of this tree is implicit. Further, only one pass over this tree can be made as partial results computed in the various nodes are not saved for use in further passes. In this respect, the binary tree method is more general than divide-and-conquer. The single pass algorithms discussed in this paper can, however, be just as well viewed as divide-and-conquer algorithms.

While all the parallel algorithms discussed in this paper have assumed that as many PEs as needed are available, they can be run quite easily using fewer PEs. The complexity of course will increase by a factor of q/k where k is the number of PEs available and q is the number assumed in the paper.

References

1. Agerwala, T. and Lint, B., "Communication in Parallel Algorithms for Boolean Matrix Multiplication," Proc. 1978 Int. Conf. on Parallel Processing, IEEE pp. 146-153, 1978.
2. Arjomandi, E., "A study of parallelism in graph theory," Ph.D. thesis, Computer Science department, University of Toronto, December 1975.
3. Baker, K. R. and Su, Z.-S., "Sequencing with due-dates and early start times to minimize maximum tardiness," Naval Res. Logist. Quart. vol. 21, pp. 171-176, 1974.
4. Batcher, K. E., "Sorting networks and their applications," Proc. AFIPS 1968 SJCC, Vol. 32., AFIPS Press, Montvale, NJ, pp. 307-314.
5. Batcher, K. E., "MPP- a massively parallel processor," Proc. 1979 Int. Conf. on Parallel Processing, IEEE, p 249, 1979
6. Blazewicz, J., "Simple algorithm for multiprocessor scheduling to meet deadlines," Info. Processing Letters, Vol. 6, Number 5, pp. 162-164, 1977.
7. Brent, R. P., "The parallel evaluation of general arithmetic expressions," JACM, Vol. 21, No. 2, April, 1974, pp. 201-206.
8. Csanky, L., "Fast parallel matrix inversion algorithms," Proc. 6th IEEE Symp. on Found. of Computer Science, October 1975, pp. 11-12.
9. Dekel, E., Nassimi, D., and Sahni S., "Parallel matrix and graph algorithms," University of Minnesota, TR 79-10, 1979.
10. Dekel, E. and Sahni, S., "Parallel scheduling algorithms," University of Minnesota, Technical Report, 1980. to be published.
11. Eckstein, D., "Parallel graph processing using depth-first search and breadth first search," Ph.D. Thesis, University of Iowa, 1977.
12. Hirschberg, D. S., "Parallel algorithms for the transitive closure and the connected component problems," Proc. 8th ACM Symp. on Theo. of Comput., May 1976, pp. 55-57.
13. Hirschberg, D. S., "Fast parallel sorting algorithms," CACM, Vol. 21, No. 8, August 1978, pp. 657-661.
14. Horn, W. A., "Some simple scheduling algorithms," Naval Res. Logist. Quart., Vol. 21, pp. 177-185, 1974.
15. Horowitz, E. and Sahni, S., "Fundamentals of computer algorithms," Computer Science Press, Potomac, MD, 1978.
16. Jackson, J. K., "Scheduling a production line to minimize tardiness," Research report 43, Management Science Research Project, University of California, Los Angeles, 1955.
17. Karp, R. M., "Reducibility among combinatorial problems," In: Miller, R. E., and Thatcher, J. W. (eds.), Complexity of computer computations, Plenum Press, New

- York, 1972.
18. Knuth, D. E., "The Art of Computer Programming Vol. 3: Sorting and Searching," Addison-Wesley, Reading, Mass., 1973.
 19. Lang, T., "Interconnections between processors and memory modules using the shuffle-exchange network." IEEE Trans. on Computers, C-25, No. 5, May, 1976, pp. 496-503.
 20. Lang, T. and Stone, H., "A shuffle exchange network with simplified control," IEEE Trans. on Computers, C-25, No. 1, January, 1976, pp. 55-65.
 21. Lawler, E. L., "Sequencing to minimize the weighted number of tardy jobs," Rev. Francaise Automat. Informat. Recherche Operationnelle, 10.5, Suppl. 27-33, 1976.
 22. Lenstra, J. K., "Sequencing by enumerative methods," Mathematical Centre Tract 69, Mathematisch Centrum, Amsterdam, 1977.
 23. Moore, J. M., "An n job, one machine sequencing algorithm for minimizing the number of late jobs," Management Sci. 15, pp. 102-109, 1968.
 24. Muller, D. E., and Preparata, F. P., "Bounds to complexities of networks for sorting and for switching," JACM, Vol. 22, No. 2, April 1975, pp. 195-201.
 25. Munro, I. and Paterson, M., "Optimal algorithms for parallel polynomial evaluation," JCSS, Vol. 7, 1973, pp. 189-198.
 26. Nassimi, D. and Sahni, S., "Bitonic sort on a mesh-connected parallel computer," IEEE Trans. on Computers, C-28, No. 1, January 1979, pp. 2-7.
 27. Nassimi, D. and Sahni, S., "An optimal routing algorithm for mesh connected parallel computer," JACM 27, 1, pp. 6-29, 1980.
 28. Nassimi, D. and Sahni, S., "Parallel permutation and sorting algorithms and a new generalized connection network," JACM, to appear.
 29. Nassimi, D. and Sahni, S., "Data broadcasting in SIMD Computers," Proc. 1980 Int. Conf. on Parallel Processing, IEEE, to appear in IEEE Trans on computers.
 30. Preparata, F. P., "New parallel-sorting schemes," IEEE Trans. on Computers, C-27, No. 7, July 1978, pp. 669-673.
 31. Rinnooy, Kan, A. H. G., "Machine scheduling problems: classification, complexity, and computations," Nijhoff, The Hague, 1976.
 32. Savage, C., "Parallel algorithms for graph theoretic problems," Ph.D. Thesis, University of Illinois, Urbana, August 1978.
 33. Siegal, H., "A model of SIMD machines and a comparison of various interconnection networks," Proc. IEEE Trans on Computers, C-28, 1979, pp. 907-917.
 34. Sidney, J. B., "An extension of Moore's due date algorithm," In S. E. Elmagharby (ed) 1973, Symp. on the Theory of Scheduling and its applications, Lecture Notes in Economics and Mathematical Systems, 86

- Springer, Berlin, pp. 393-398, 1973.
35. Smith, W. E., "Various optimizers for single-stage production," Naval Res. Logist. Quart., Vol. 3, pp. 59-66, 1956.
 36. Stone, H., "Parallel processing with the perfect shuffle," IEEE Trans. on Computers, C-20, 1971, pp. 153-161.
 37. Thompson, C. D., "Generalized connection networks for parallel processor interconnection," IEEE Trans. on Computers, C-27, No. 12, December, 1978, pp. 1119-1125.
 38. Thompson, C. D., and Kung, H. T., "Sorting on a mesh connected parallel computer," CACM, vol. 20, No. 4, 1977, pp. 263-271.
 39. Valiant, G. L., "Parallelism in comparison problems," SIAM. J. Comput. vol. 4, no. 3, September 1975.

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

REPORT DOCUMENTATION PAGE		READ INSTRUCTIONS BEFORE COMPLETING FORM
1. REPORT NUMBER (14) TR-80-29 2	2. GOVT ACCESSION NO. AD-A091	3. RECIPIENT'S CATALOG NUMBER 372
4. TITLE (and Subtitle) (6) Binary Trees and Parallel Scheduling Algorithms,		5. TYPE OF REPORT & PERIOD COVERED (1) Technical Report July 1980 - September 1980
7. AUTHOR(s) (10) Eliezer/Dekel and Sartaj/Sahni		6. PERFORMING ORG. REPORT NUMBER
9. PERFORMING ORGANIZATION NAME AND ADDRESS 3 Dept Computer Science Department 412026 University of Minnesota 136 Lind Hall, 207 Church St. S.E., Mpls., MN.		8. CONTRACT OR GRANT NUMBER(s) (15) N00014-80-C-0650 VNSR-M2580-005856
11. CONTROLLING OFFICE NAME AND ADDRESS Department of the Navy Office of Naval Research Arlington, Virginia 22217		10. PROGRAM ELEMENT PROJECT, TASK AREA & WORK UNIT NUMBERS (12) 372
14. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office)		12. REPORT DATE September 1980
		13. NUMBER OF PAGES 34
		15. SECURITY CLASS. (of this report) UNCLASSIFIED
16. DISTRIBUTION STATEMENT (of this Report) Unlimited		15a. DECLASSIFICATION/DOWNGRADING SCHEDULE
17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)		
18. SUPPLEMENTARY NOTES		
19. KEY WORDS (Continue on reverse side if necessary and identify by block number) Parallel algorithms, design methodologies, complexity, scheduling, shared memory model.		
20. ABSTRACT (Continue on reverse side if necessary and identify by block number) This paper examines the use of binary trees in the design of efficient parallel algorithms. Using binary trees, we develop efficient algorithms for several scheduling problems. The shared memory model for parallel computation is used. Our success in using binary trees for parallel computations, indicates that the binary tree is an important and useful design tool for parallel algorithms.		

DD FORM 1 JAN 73 1473

EDITION OF 1 NOV 68 IS OBSOLETE
S/N 0102-LF-014-6601

Unclassified 412026
SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)