

NOTICE WARNING CONCERNING COPYRIGHT RESTRICTIONS:
The copyright law of the United States (title 17, U.S. Code) governs the making of photocopies or other reproductions of copyrighted material. Any copying of this document without permission of its author may be prohibited by law.

Two Papers on a Tree-Structured Parallel Computer¹

Jon Louis Bentley²

H. T. Kung

**Department of Computer Science
Carnegie-Mellon University
Pittsburgh, Pennsylvania 15213**

25 September 1979

DEPARTMENT
of
COMPUTER SCIENCE



Carnegie-Mellon University

Two Papers on a Tree-Structured Parallel Computer¹**Jon Louis Bentley²****H. T. Kung****Department of Computer Science
Carnegie-Mellon University
Pittsburgh, Pennsylvania 15213****25 September 1979****Abstract**

This report consists of two papers describing various aspects of a new tree-structured parallel computer. The first paper, "A tree machine for searching problems" by J. L. Bentley and H. T. Kung, describes the basic architecture of the machine. A set of N elements can be maintained on an N -processor version of the machine such that insertions, deletions, queries and updates can all be processed in $2 \lg N$ time units. The queries can be very complex, including problems arising in ordered set manipulation, data bases, and statistics. The machine is pipelined so that M successive operations can be performed in $M-1 + 2 \lg N$ time units. The paper studies both the basic machine structure and a VLSI implementation of the machine. The second paper, "A parallel algorithm for constructing minimum spanning trees" by J. L. Bentley, shows how an $(N/\lg N)$ -processor version of the machine can solve the problem of constructing minimum spanning trees in time proportional to $N \lg N$. This algorithm is an improvement over existing algorithms in several ways.

¹This research was supported in part by the Defense Advanced Research Projects Agency under Contract F33615-78-C-1551 (monitored by the Air Force Office of Scientific Research), in part by the National Science Foundation under Grant MCS 78-236-76, and in part by the Office of Naval Research under Contract N00014-76-C-0370.

²Also with the Department of Mathematics.

A Tree Machine for Searching Problems¹

Jon Louis Bentley²

H. T. Kung

Department of Computer Science
Carnegie-Mellon University
Pittsburgh, Pennsylvania 15213

30 August 1979

Abstract

In this paper we describe a new tree-structured machine (suitable for VLSI implementation) that solves a large class of searching problems. A set of N elements can be maintained on an N -processor version of this machine such that insertions, deletions, queries and updates can all be processed in $2 \lg N$ time units. The queries can be very complex, including problems arising in ordered set manipulation, data bases, and statistics. The machine is pipelined so that M successive operations can be performed in $M-1 + 2 \lg N$ time units. In this paper we will study both the basic machine structure and the actual implementation of the machine.

¹This research was supported in part by the Defense Advanced Research Projects Agency under Contract F33615-78-C-1551 (monitored by the Air Force Office of Scientific Research), in part by the National Science Foundation under Grant MCS 78-236-76, and in part by the Office of Naval Research under Contract N00014-76-C-0370.

²Also with the Department of Mathematics.

Table of Contents

1. Introduction	1
2. The Abstract Machine	2
3. An Architecture	9
4. Discussion of Implementation	16
5. Conclusions	19
Acknowledgements	21
References	22

1. Introduction

Very Large Scale Integrated circuitry (VLSI) has been increasing in speed and decreasing in size at an amazing rate over the past decade, and it promises to continue at this rate far into the next decade (see Mead and Conway [1979]). In this paper we will describe a tree-structured machine for solving searching problems that is ideally suited for implementation in VLSI. The searching problems that the machine solves arise in a number of applications areas (including ordered set manipulation, data bases and statistics), and it is able to solve all of the problems very efficiently.

Before describing this machine in detail, it is helpful to characterize its contribution in general terms. The authors believe that there is a spectrum of impacts that advances in VLSI technology will have on computer architecture. At one extreme, this technology will allow conventional architectures to be implemented as smaller, faster and cheaper machines -- this will lead to more sophisticated interconnections of conventional machines (see, for example, Swan, Fuller and Siewiorek [1977], or Sequin, Despain and Patterson [1978]). Also at this end of the spectrum will be minor (register level) architectural changes that exploit certain features of VLSI; this area has been explored by Sites [1979]. At the other extreme, VLSI architectures have been proposed that are radical departures from the von Neumann tradition (see, for example, Backus [1978], Mago [1979] or Wilner [1978]). In this paper we will investigate an approach that lies between these two extremes: a high-performance, special-purpose, non-von Neumann computing device that is designed to be used in conjunction with a conventional computer. In general, such devices should be constructed only when they solve a problem satisfying two criteria: the problem should currently consume large quantities of computer time, and the proposed special-purpose device must be much more efficient than conventional ways of solving the particular problem. When such a problem is identified it is reasonable to augment a general-purpose computing system with a special-purpose device for solving the problem; the structure of such a system is depicted in Figure 1. Many such special-purpose devices have recently been

proposed; see, for example, Kung [1979] and Kung and Leiserson [1978].

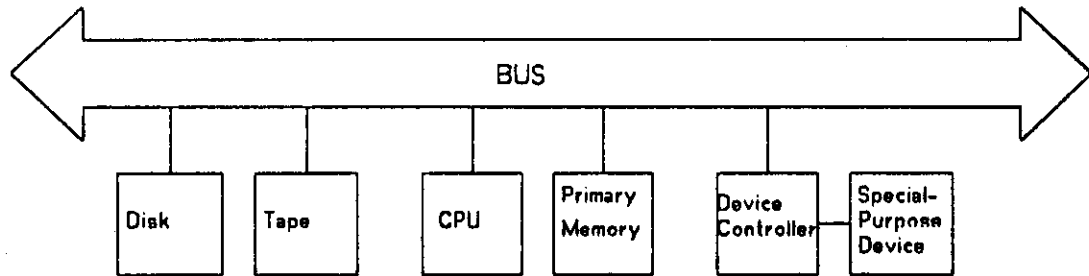


Figure 1. General system structure.

In this paper we will investigate a special-purpose machine for solving searching problems. This machine is described at an abstract level in Section 2, where we will also review some necessary background in searching problems. An architecture (that is, a user's view) of the machine is described in Section 3, and issues of implementing that architecture in VLSI are discussed in Section 4. Conclusions are then offered in Section 5.

2. The Abstract Machine

In this section we will investigate the tree-structured searching machine at an abstract level, apart from the details of architecture or implementation. The general searching problem it solves calls for maintaining a *file* of fixed-format *records*. We must be able to perform the operations of *inserting* a new record into the file, *deleting* an existing record from the file, *updating* records in the file, and *querying* the file to answer questions. Before we examine the general searching problem, we will investigate one searching problem in particular.

That particular problem is called *member searching*. In its abstract form, it involves maintaining a set of elements so we can determine if a new element is a member of the set. In concrete applications, other information is usually also requested. For example, after finding that a particular social security number is a member of a set of social security numbers, we often wish to retrieve other information (such as Year-to-Date taxes). We will now investigate how the tree

machine solves the abstract member searching problem, and then return in the next section to the complicating issues that arise in applications.

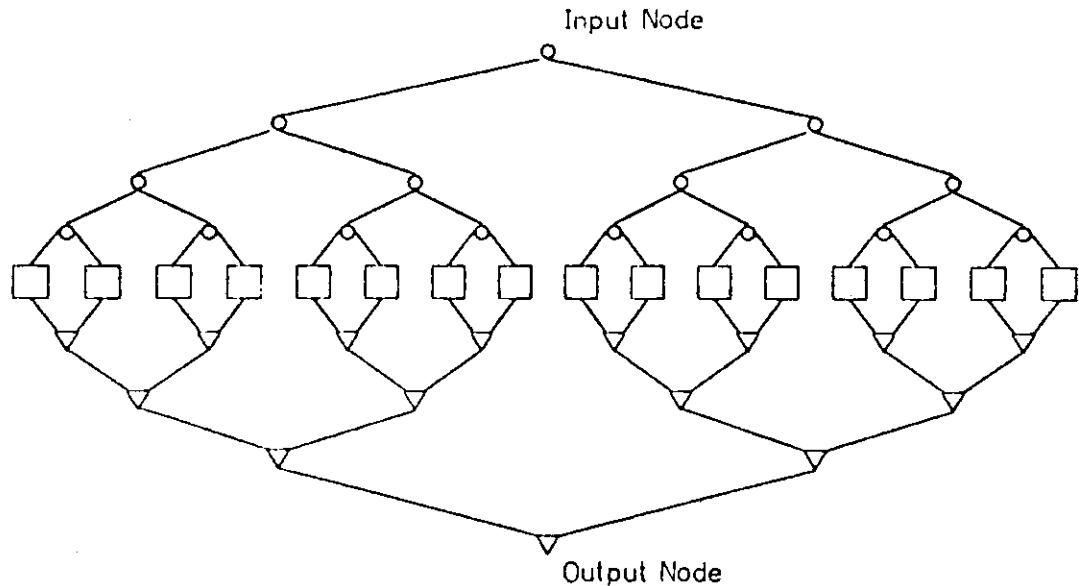


Figure 2. Structure of the tree machine.

The basic organization of the tree-structured searching machine is depicted in Figure 2. There are three kinds of nodes in the machine: circles (which broadcast data), squares (which have limited storage and computation power), and triangles (which "combine" answers to queries). A set of N elements is stored in this machine by placing each element of the set into a distinct square node of the tree. Consider now the problem of performing the member search to answer the query "Is 17 an element of the set?". We accomplish this by inserting 17 into the input node and broadcasting it down the tree -- $\lg N$ steps later the value 17 will arrive at all of the squares. This situation is illustrated in Figure 3a. At that point we compare the values stored in each square to 17 and set a bit to one if the value is equal to 17 and zero otherwise; this is shown in Figure 3b. We can now combine the bits together through the bottom portion of the network by letting each triangle compute the logical or of its two inputs, as illustrated in Figure 3c. So after a total of $2 \lg N$ time units have passed since the query was posed, a single bit emerges from the output node telling whether or not 17 is an element of the set. We have thus described a procedure for determining whether a given object is a member of the

set whose elements are stored in the square nodes.

It is important to note that the tree machine has a very regular data flow: the data moves in discrete steps in only one direction (from the input node to the output node). Thus if many successive elements are going to be tested for membership in the set stored in the square nodes, then the process of answering those queries can be pipelined. As the value of the first element to be tested is going down the tree, the next value can follow one step behind, and so on. If M successive tests are performed in this manner, exactly $M-1 + 2 \lg N$ time units pass between the entry of the first query at the top of the tree and the exit of the last of the answers at the bottom of the tree.

The tree machine is able to solve many problems besides member searching. For example, if a multiset of elements (that is, a set in which one element can appear many times) were stored in the square nodes of the tree, we might wish to count how many times a given object appears in the set. We proceed exactly as we did for member searching, first broadcasting the given element through the circles to the square nodes. We load a one into each square if its element is equal to the given object and zero otherwise, and then combine the answers by letting the triangles *sum* the values of their inputs. Another example is given by nearest neighbor searching. If we wish to find the distance to the element of the set that is closest to 17, then we do the following: broadcast 17 through the input node to all squares, subtract the value stored in the square from 17 and take the absolute value of the difference, and finally take the minimum of all those values by having the triangles return the *minimum* of their two inputs. As for member searching, for both member counting and nearest neighbor searching, we can answer a single query in $2 \lg N$ time and a series of M queries in $M-1 + 2 \lg N$ time.

In general, the tree machine can solve any problem that can be phrased as computing some function over every element in the set (such as equality or absolute value of difference) and then combining the values of those functions by some associative, commutative binary operator. For example, the rank of an element X in

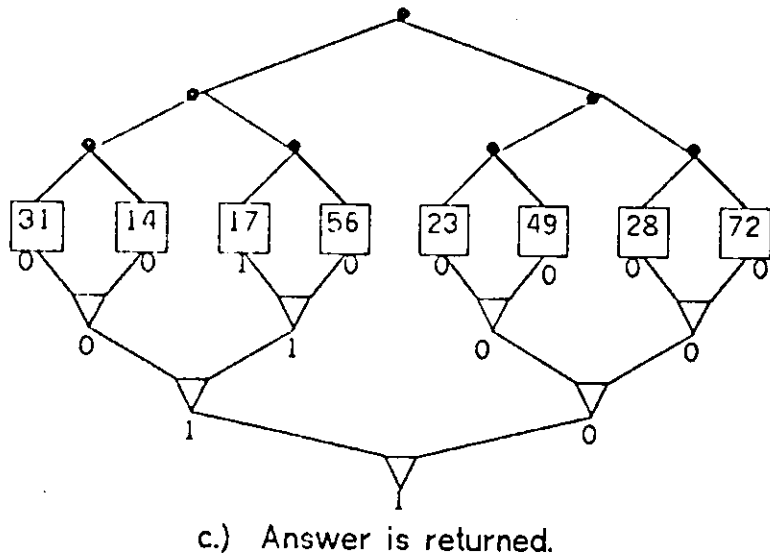
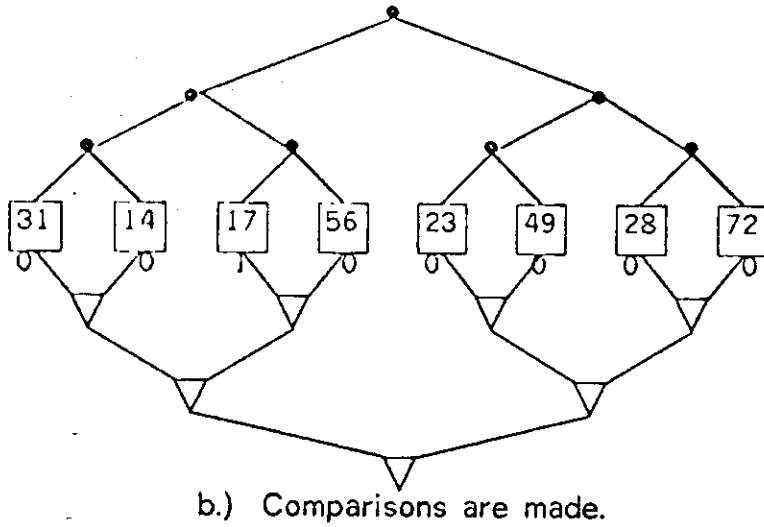
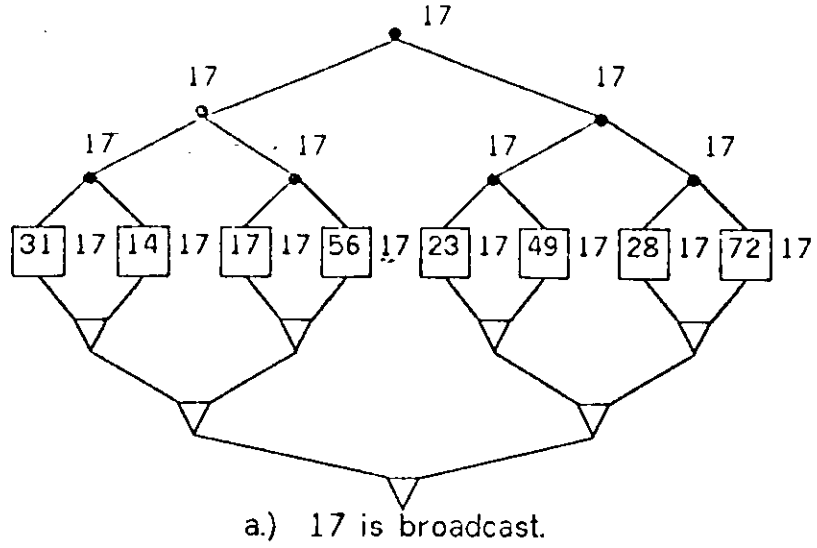


Figure 3. A member search.

a set (that is, the number of elements in the set less than X) can be calculated by

storing in each square a one if the element is less than X and zero otherwise; the final answer is then computed by having the triangles add their inputs. Other problems defined on totally ordered sets that can be solved by the tree machine include predecessor (what is the greatest element less than the given?), successor (what is the least element greater than the given?), and minimum (what is the least element in the set?). In general, the tree machine can solve all of the "Decomposable Searching Problems" defined by Bentley and Saxe [1979]. That reference contains both an algebraic definition of the class and a list of over twenty particular searching problems in the class.

The tree machine is also able to answer much more complicated kinds of queries (of the form that arise in data base applications, for instance). Suppose, for example, that every square node of the tree contains a record with ten keys. We might want to know how many records there are in the file with first key equal to a given value, the second key at least as great as the third key, the fourth key in a certain range, and so on. This type of query is easily answered: we merely broadcast each of the conditions down to the square nodes, keeping track in each node of whether it has satisfied all the conditions shipped so far. We load a one if all conditions have been satisfied and a zero otherwise, and combine by having the triangles sum their inputs. Many applications call for a list of the satisfying records instead of merely their count, and this can be accomplished by letting the triangles compute the *union* of their inputs. This can be viewed intuitively by observing each triangle independently, and imagining a person "tapping" the entire machine at each time step. As each triangle is tapped, there are three cases to consider: if it has no items in its inputs, it reports that; if it has one item, it returns it; and if it has two items, it returns only one (delaying the other until the next tap). This "tapping" process continues as long as there are elements that have yet to be reported. (Note that to compute unions in this manner, the pipelining must be carefully designed to ensure that no "overflow" occurs.)

Having discussed searching at some length, we will now turn to the issues of maintaining the set of elements stored in the square nodes. A tree machine with N

square nodes (where N is a power of two) can store up to N records. A new record can be inserted into the set by placing it in any unused square. We find such a square by having each circle keep track of the number of unused square descendants of each of his two sons. When a request comes to the root for a new (unused) position, he passes the request to one of his sons with unused square leaves, and so on. Mechanically, this is accomplished by turning off all of the squares except the one finally chosen as the holder of the new record; this square is then loaded with the desired data. Note that a single record can be inserted in $\lg N$ steps, and a set of M records can be inserted in $M-1 + \lg N$ steps.

Another maintenance operation is that of updating a set of records: this can be easily accomplished by broadcasting the conditions that the changed records must meet, turning off all processors that do not meet the conditions, and then making the desired changes. (Although the update set will often have just one element, an example of a "mass update" might be processed on the first of the month: for all salesmen with Month-Of-Starting-Employment equal to This-Month, add one to Years-Of-Service.) To delete a single record we set a flag in its square node saying that it is unused and then adjust the counts in all of the circles above it. This can be accomplished either by pushing information "backward" to the top of the tree (adding one to each counter as you go), or by doing a dummy reinsertion of that element, and modifying the counters on the way down. The time for either of these operations is proportional to $\lg N$. Notice that after a set of elements in squares have been identified for deletion, they can be deleted in parallel (in a single step) and all counters can be reset (by pushing the information up the tree) in $\lg N$ steps. Although having information go up the tree is handy for deletion, it does complicate the basic design severely; this feature might therefore not be implemented.

So far in our discussion each machine has represented but one set. In some applications, however, a given user might wish to represent many sets, or many users might want to use the machine independently for their respective sets. Either of these can be accomplished so long as the sum of the sizes of the sets is less than N , the number of square nodes. Although we could "slice" the machine into

sections to accomplish this, there is a much more elegant solution. Namely, a fixed portion of each record is dedicated to a "set identification field", or "SetID". To process an operation on Set 56 (or a set belonging to user 56), we have as a prelude to the operation the sequence "check SetID for equality with 56 and turn off the processor if not equal". (Notice that we are not requiring that all records in all sets be of the same format, but just that they have one field in common.) In an environment with much sharing, this prelude will occur so often that it might be advantageous to provide a single instruction that accomplishes its purpose.

Although so far we have used the tree machine to solve only searching problems, it can be applied to many other problems as well. For instance, it can be used to sort a set of M elements in time proportional to M (as long as M is between $\lg N$ and N , where N is the number of square nodes in the tree machine). This is accomplished by making two passes through the M elements: the first inserts the elements into the machine, and the second counts for each element the number of elements less than it (that is, it computes the element's rank, as we saw before). This tells precisely where each element occurs in sorted order (the output is a permutation vector), and it is then trivial to arrange the elements into sorted order. By use of pipelining, both steps run in time linear in M . Note that it was critical to phrase sorting as a counting problem, rather than as extracting the minimum, to make use of pipelining in the second step -- this algorithm essentially implements an N^2 algorithm in N time by using all N processors in parallel. There are many other examples of such speedups for problems that are not *prima facie* searching problems. Two such examples are computing all nearest-neighbor pairs in a k -dimensional point set (which arises in data analysis) and reporting all pairwise design rule violations in a VLSI mask (a design automation task). The application of this machine to the problem of constructing minimum spanning trees has been discussed by Bentley [1979]--he shows how an $N/\lg N$ -processor version of the tree machine can construct the minimum spanning tree of an N -node graph in $O(N \lg N)$ time, which is optimal for complete graphs. Other applications of tree-structured machines have been studied by Browning [1979].

This concludes our discussion of the machine at an abstract level, and we can now state the properties that a concrete embodiment of the machine must possess. There must be three kinds of nodes in such a machine: circles, squares and triangles. The circles must broadcast data and have a small amount of state (namely, to remember how many unused squares are descendants of each of their sons). The only processing required of a circle is incrementing or decrementing by one. The squares, however, must have substantial memory and computation power. Each square must have enough processing capability to handle the most difficult kinds of queries and updates desired and (usually) enough memory to store the largest record in most applications. The triangles must be able to combine answers. Most of the "combinators" we desire are very simple to implement; these are *and*, *or*, *min*, *max*, and *plus*. The only complicated combinator is *union*, and we are willing to "turn off" pipelining in the presence of that operator.

3. An Architecture

In Section 2 we described the tree-structured searching machine at an abstract level, ignoring many issues of implementation. In this section we will move one step closer to an implementation, and describe a particular architecture (that is, a user's view of the machine) realizing the abstract machine. It is essential that the reader understand that the architecture we will investigate is not proposed as the best possible architecture realizing the abstract machine of the last section. Rather, it is put forth only as evidence that there is *at least one* reasonably efficient architecture for the machine. In Section 4 we will discuss how this architecture can be implemented in VLSI.

The basic structure of the architecture we will investigate is that studied in Section 2 (illustrated in Figure 2). The flow of instructions and data in the machine is exclusively from the input node (at the top of the figure) to the output node (at the bottom) -- we will not have deletions that employ any "backwards flow". The machine is based on 16-bit instructions and 32-bit data words (which are interpreted either as integers in two's-complement or as 32-bit vectors). The top

data paths in the machine (the son links from circles in Figure 2) are 16 bits wide; the bottom data paths (links to triangles) are 80 bits wide. The entire machine operates synchronously; an operation is (perhaps) performed at each node and data is transmitted from the node to its sons on each *major cycle*. Having described the machine at this gross level, we will now examine the circles, squares and triangles individually.

The primary function of the circle on each major cycle is to broadcast what it just received to its sons. In only three contexts must it perform a more sophisticated operation. As a new element is being inserted, it must decide which way to direct the insertion (to one that has unused square leaves) and then decrement the appropriate counter by one; it then ships a "no-op" to the other son. The no-op is effected by having one bit in the instruction turned off as the 16-bit instruction is passed to the "other" son. To accomplish a deletion we insert an instruction packet of three 16-bit instructions at the root node. The first instruction is the deletion and the next two 16-bit words contain the binary address of the node to be deleted. The circles can tell by looking at appropriate bits of the address whether they should increment one of their counters as they see this instruction. The final capability the circles must have is that of passing data to the squares, without interpreting that as an instruction to them; we will return to this issue as we discuss the squares.

While the circles have the simplest architectures of the three units we will see, the squares have the most complex. The abstract machine requires that the squares be able to store data and to perform enough calculations to answer queries and perform updates. This architecture will accomplish both these tasks by shipping combinations of instructions and data to the machine. We now have to make a fundamental design decision: should the individual squares be special-purpose devices (honed for a particular view of the tree machine's task), or should they be (in some limited sense) general-purpose computing devices? We will choose the latter course, and make each square a "baby" von Neumann computer; it is important, however, to emphasize that this is merely a design decision and not an

inherent property of the abstract machine.

Each square will be a small von Neumann-like processor that receives its instructions and data from an external, 16-bit stream. An individual processor contains sixteen 32-bit words of memory, two 32-bit registers, and a vector of eight single-bit data flags ($F[0]$, $F[1]$, ..., $F[7]$). The processor also contains an eight-bit Set Identification number (SetID), and an Instruction Register. The first bit of the F vector ($F[0]$) is used as the "Active" bit of the processor; a special "Enable" command turns on all processors (by setting $F[0]$ to one), and a processor can conditionally turn itself off by storing a zero in $F[0]$. The basic layout of the machine is shown in Figure 4 (notice that because the machine is rotated 90° , the data flows from right to left rather than from top to bottom).

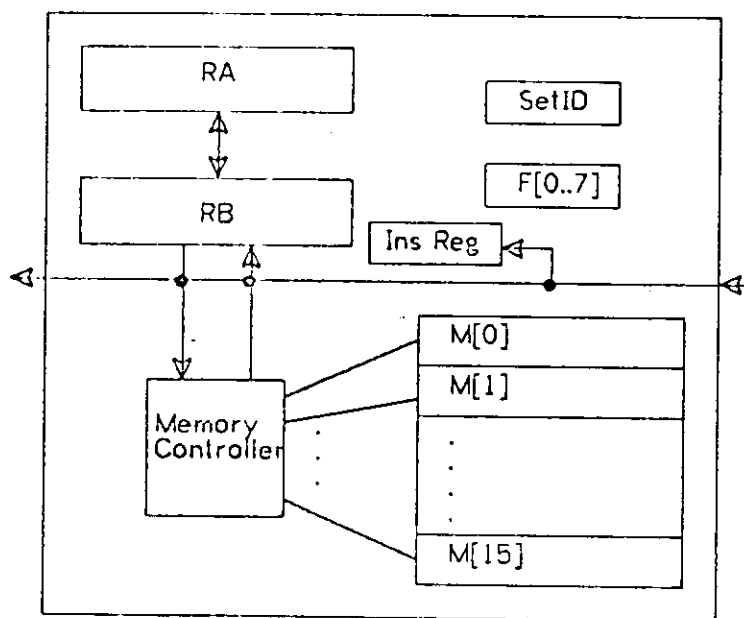


Figure 4. Components of the square.

The 16-bit instruction format for the square processor is shown in Figure 5. The first bit of an instruction processed by the squares is always zero; a one in that bit signifies an instruction that is ignored by the squares but passed on to the triangles. The two Fam bits specify one of the four families to which an instruction can belong (Arithmetic-Logical, Load-Store, Bit or Special), and the Code gives the opcode of the instruction. There is a one bit flag (Flag) in each instruction, and arguments to

the instruction are either two four-bit addresses (A1 and A2), an 8-bit string (Name) or a five-bit integer (Num). The actual instructions are described by group in an ISP-like language in Table 1. All of the arithmetic-logical instructions are zero-address instructions, combining registers RA and RB and storing the result in either RA or RB (usually RA). The load-store instructions specify one of 16 memory addresses as their operand; the data movement is then between that address and the register RB. The bit operations generally have two addresses: they combine the first and the second operands, storing the result in the first. The exceptions to this pattern are the unary not operator and the compare (comp) operation; the latter compares RA with RB and stores in the first bit (F[A1]) whether or not the values are equal and tells which inequality in the second bit (F[A2]) -- this is just a straightforward encoding of three states into two bits.

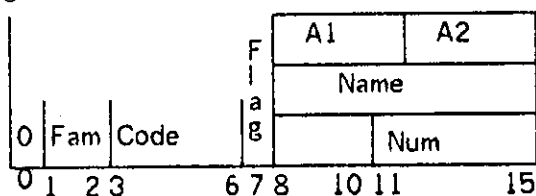


Figure 5. Instruction format.

The only instructions that are not entirely obvious are the special instructions. The *enable* instruction turns on all processors in the tree. The *ins* (insert) instruction turns on precisely one processor, turning off the rest (and decrementing the counters in the circles). The *del* (delete) instruction has no effect on the processors; it only increments the appropriate counters in the circles (the squares must ignore the two following instructions packets, though -- they are just the processor address). The *ship* instruction allows data to enter the RB register from the data/instruction stream. The Flag bit tells whether the next one or two 16-bit packets should be loaded into RB; the data can then be processed as desired. The *chksid* and *setsid* instructions are for manipulating the 8-bit SetID register; the former turns off the processor if SetID is not equal to Name, and the latter loads the SetID field from Name.

Arithmetic-Logical

add	→	RA ← RA + RB
sub	→	RA ← RA - RB
neg	→	RA ← - RA
rand	→	RA ← RA ∧ RB
ror	→	RA ← RA ∨ RB
rxor	→	RA ← RA ⊕ RB
rnot	→	RA ← ~ RA
shift Num	→	RA ← RA left shifted by Num
tab	→	RA ← RB
tba	→	RB ← RA
swap	→	RA ↔ RB

Load-Store

ldb Num	→	RB ← M[Num]
stb Num	→	M[Num] ← RB

Bit

band A1,A2	→	F[A1] ← F[A1] ∧ F[A2]
bor A1,A2	→	F[A1] ← F[A1] ∨ F[A2]
bxor A1,A2	→	F[A1] ← F[A1] ⊕ F[A2]
bnot A1	→	F[A1] ← ~ F[A1]
comp A1,A2	→	F[A1] ← RA=RB; F[A2] ← RA < RB

Special

enable	→	F[0] ← 1
ins	→	F[0] ← this processor selected
del	→	(defined in text)
ship Flag	→	(defined in text)
chksid Name	→	F[0] ← SetID = Name
setsid Name	→	SetID ← Name

Table 1. Instruction set for squares.

To illustrate the operation of the processors we will study two program segments for performing searches. The first segment is for member searching.

```

chksid  ThisSet  // Turn off undesired processors
ship    Two      // The next two packets hold the comparand
dataL
dataR
tab                    // Put comparand in RA
ldb     KeyAd    // Put key in RB
comp    1,2      // Answer is in F[1]

```

The search key enters the RB register from the data stream and is then transferred to the RA register. The program then loads the key field of the record into the RB register (KeyAd is an integer identifying which of the 16 memory words holds the key), and makes the comparison. F[1] is then one if and only if the record's key field is equal to the data shipped in the stream. At this point, the answer can be combined in the triangle network.

The next program that we will examine arises in "nearest neighbor" searching; it computes the distance between the data and the key field of the record. Since we desire the absolute value of the difference of the key and the data, we must have a conditional step in our program.

```

chksid  ThisSet
ship    Two      // RA ← Data
dataL
dataR
tab
ldb     KeyAd    // RB ← Key
comp    2,0      // If Data ≤ Key, leave processor on
swap
chksid  ThisSet  // Turn all processors back on
sub                    // RA ← |Key-Data|

```

The crucial step of this program is the comp instruction: if Data is less than Key then a one is stored in F[0], which leaves the processor on; the swap then interchanges key and data. The next instruction (chksid) turns all appropriate processors back on, and the subtract correctly computes a positive value. The triangles can then be instructed to return the minimum of these values.

The two code segments that we have just seen illustrate many of the aspects of

coding the tree machine. Many other examples have been coded, and all of them appear to be fairly efficient. More quantitatively, the ratio of tree machine instructions to "critical" operations in the task clusters very closely around 2.5. This statistic is evidence for the vindication of our design decision to make the squares general-purpose machines, rather than special devices tailored to the searching task domain. (Pursuing that alternative remains an interesting open problem.)

Before ending our discussion of the squares, it is interesting to compare the design of the processor with a more typical von Neumann processor. In some ways, we faced exactly the same problems: the choices of data representation, instruction formatting, operation set, and addressing were all taken from the von Neumann design space as discussed by Blaauw and Brooks [1979]. On the other hand, we avoided many of the issues faced by designers of typical machines; these include instruction sequencing, interrupt handling, and input/output control.

Before we discuss the architecture of the triangle, we must settle one more point about what we want it to do. In most applications that compute the minimum of a set (for instance), we want to know not only what the value of the minimum is but also what element has that value. We therefore have three objects associated with computing the minimum: the operation (minimum), the value, and the name (which is a 32-bit word associated with the value; its address or "key" in many applications). When combining two such objects, we take the value as the minimum of the two values, and the name from the name of the smaller value. The name is thus *inherited* from the minimum. We will also associate names with other binary operators: the name of *maximum* is inherited from the node with greater value; for *plus*, from a nonzero element; for *or*, from a nonzero bit vector (arbitrary if both are zero); and for *and* from a zero bit vector.

Having defined the concepts of value, name and inheritance, it is straightforward to describe the architecture of the triangles. They will operate on 80-bit packets: 16 bits of instruction, and 32 bits each of value and name. Computing *min*, *max*,

plus, *and*, and *or* are all simple. Union is a bit more detailed, but also conceptually straightforward. One aspect that we have not mentioned is the interface between the squares and the triangles; we must include instructions for transferring the contents of the RB register to the name or value field of the triangle immediately beneath it (these could be included in the load-store family). This allows us to give complete programs for answering queries. After computing the answers (as illustrated in the two segments shown above), we load them into the desired fields of the triangles, and combine them as desired.

It is important to emphasize that the architecture we have just seen is not the architecture that the ultimate user of the machine will see. Rather, there will be a hierarchy of functions available to him. At the highest level, he will be able to perform operations on sets (load a set, erase a set, for each element in the set, and so forth); at an intermediate level there are record-handling operations (defining queries or inserting, deleting and updating records); and at the lowest level there are the machine instructions themselves. At the lowest level the user can make very efficient code by knowing the details of the machine; at the higher levels he sacrifices efficiency for clean and easy code.

An important part of the implementation of this architecture is that there be a fairly sophisticated device controller for the tree machine (such as an off-the-shelf microprocessor). This controller will implement the hierarchy of functions mentioned above. This will also reduce the bus activity substantially by having the controller fetch items from main memory and issue instructions to the tree machine; it appears that having the CPU itself perform these tasks would lead to a substantial degradation in overall system performance.

4. Discussion of Implementation

In this section we will discuss one implementation of the architecture of Section 3 in VLSI technology. The fundamental description of the implementation is that it is bit-serial. There are two motivations for this: one, to exploit the shift-register technology of VLSI, and two, to use very few pins on packages.

The implementation of both the circles and the triangles described in the last section is straightforward. The squares are also easy to implement bit-serially. The 16-word memory is in fact a parallel shift register, 16 bits wide and 32 bits long. The two registers RA and RB are also shift registers. To load or store a word, RB and the memory shift register are shifted in parallel, and the memory controller of Figure 4 is just a multiplexor (decoding a 4-bit address to one of 16 lines). All of the arithmetic-logical operations are accomplished by putting a single-bit function box between the RA and RB registers, and then shifting the pair through it (all operations require at most one bit of memory). Notice that we have assumed that the squares have 32 minor cycles during each major cycle of the machine. The bit operations are straightforward to implement if the Flag array is just a small RAM. Estimates by experienced VLSI designers indicate that the chip area for the functionality in the square is about equal to the chip area required for the 512-bit memory. Using current technology, it is easy to imagine putting 16 squares on a single chip.

Now that we know how we will implement the individual processing elements (circles, squares and triangles), we must describe how to place them on a chip. The first simplification we will make is to consider them as standard binary trees rather than the "mirrored" binary tree of Figure 2; the unmirroring process is illustrated in Figure 6. We now face the problem of laying out a binary tree on a chip. This problem has been studied by Mead and Rem [1979], who suggest the space-economical layout illustrated in Figure 7. The amount of space used in that layout is proportional to the number of processors on the chip. Note that each edge in that layout is realized by two "wires" on the chip -- one for data going to the squares, and one for data coming from the squares.

Since only some fixed number of the processors in a tree machine will fit on a single chip, it is important that we discuss the packaging of the chips. The packaging strategy we propose is illustrated in Figure 8. There are two kinds of chips in that figure: the *leaf* chips and the *internal* chips. The leaf chips contain

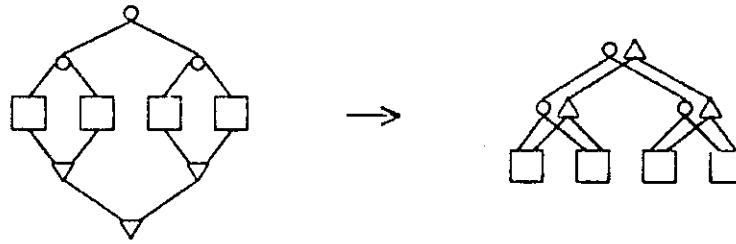


Figure 6. "Unmirroring" the tree machine.

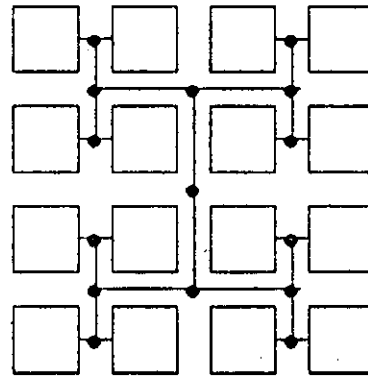


Figure 7. Tree layout on a chip.

(say) 16 square nodes and 15 circle and triangle nodes. All the communication to a leaf chip is through two wires, so the chip needs only two communications pins (besides power, ground and timing synchronization pins). Notice that this implies that with technological advances in VLSI, we will be able to place many more processors on a square chip; we are not bound by pin limitations. The internal chips would probably be constructed with seven circles and triangles on them; this implies that there is one input-output pair of wires at the top of the chip and eight pairs at the bottom. The total number of pins for this chip is therefore eighteen (plus miscellaneous pins). This chip is therefore pinbound even in today's fabrication technology; unless there are unexpected advances in packaging technology, the internal chips will probably continue to have seven or at most fifteen pairs of circles and triangles.

To get a better feeling for the size of the tree machine, we will briefly consider how one might be built today. Suppose that we put sixteen square nodes on each leaf chip, and seven circle-triangle pairs on each internal chip (both of these are

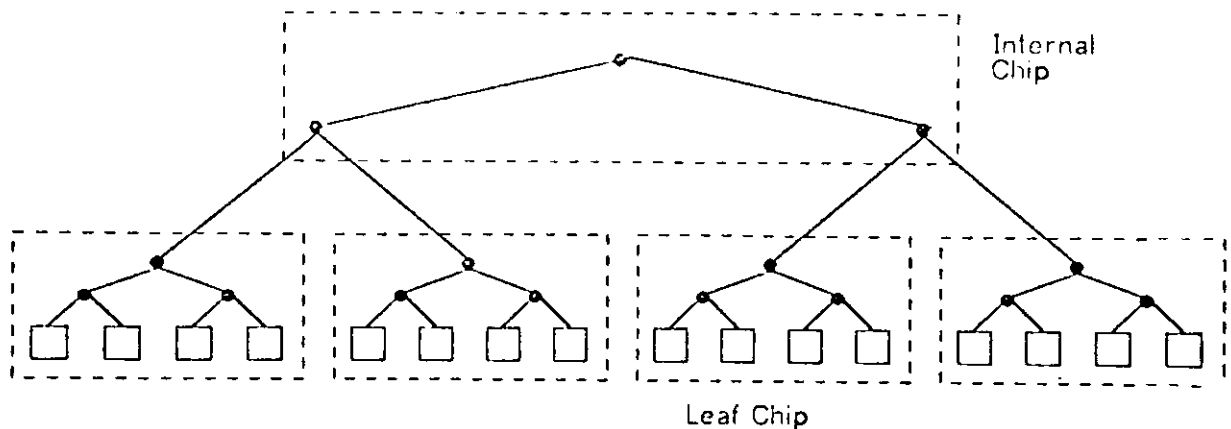


Figure 8. Two kinds of chips.

easily accomplished in today's technology). We will now put 64 leaf chips and nine internal chips on a board; this gives us 1024 square nodes. We can then put sixteen of these boards in a small cabinet, giving a tree machine of more than sixteen thousand square nodes, each holding a 512-bit record. If we assume that technology continues to double the number of components on a chip every two years, this implies that we can expect a tree machine of one million records to fit in about a cubic foot of space by the end of the 1980's.

These rough (but fairly conservative) estimates indicate that the tree machine might be one reasonable way to exploit the processing power that VLSI will give us. Before we can assert this with confidence, however, we must show that the tree machine is a wiser way to invest resources than other structures for searching. For example, might it be better to put the same resources into a large RAM memory rather than a tree machine? The authors' preliminary investigations strongly suggest that the excess cost of the tree machine compared to a RAM is very small compared to the functionality purchased, but the detailed comparison of this architecture to the RAM and its other competitors remains an open problem.

5. Conclusions

In this paper we have investigated the tree machine for searching problems on several levels. In Section 2 we studied it in an abstract setting and showed that it

can rapidly solve many searching problems, as well as some other problems that do not immediately appear to be searching problems. In Section 3 we saw an architecture (that is, a user's view) of the machine, and in Section 4 we saw that that architecture can be efficiently implemented in VLSI technology. Having studied the machine at these various levels, we will now spend a few moments summarizing the contributions of this work.

This machine can be compared with many other architectures. It is similar to an associative memory in many aspects, but it can perform many more operations than even the most powerful associative memories considered to date (see, for example, Lamb and Vanderslice [1978]). One might consider the square processors as forming a Single-Instruction, Multiple-Data stream (SIMD) computer, but each square is considerably simpler than most SIMD machines proposed to date. The tree machine is also superficially similar to the CASSM computer of Su *et al* [1979], but there are fundamental differences in the two machines at both the architectural and implementation levels. Two other machines to which it might be compared are the tree-structured machines of Mago [1979] and Sequin, Despain and Patterson [1978]. Both of these machines, however, are put forward as general-purpose computing devices, while our machine is much more specialized to the particular problem of searching.

Although we explored only one design path in this paper, it is important to remember that there are many variants of the tree machine. For example, in the unmirrored tree machine of Figure 6, the circle-triangle nodes could be made more powerful so that they could interact with passing data in more sophisticated ways, thereby substantially enhancing the machine's capability. So far we have investigated only binary trees; in certain applications, other branching factors may prove superior. Other interesting variants of the machine come from changing the amount of memory in a square processor; might it be reasonable, for instance, to have thousands of memory words in each square? Many other design paths remain unexplored -- in this paper we have only attempted to describe the fundamental concepts of the machine. _

An interesting aspect of the tree machine is what we might call its "computational structure", which is illustrated in Figure 9. That diagram has three interpretations. First, it illustrates the tree machine itself: very small input and output channels, with massive computation going on in between. Second, it describes the searching problem: a small question is asked about a large set, giving a small answer. And finally, the figure illustrates the constraints of working with pinbound VLSI: the number of pins on a chip is very small compared to the number of functional components. The fact that the abstract structure of both the searching problem and the tree machine's solution to it closely model the medium of VLSI indicates that this approach might be very successful.

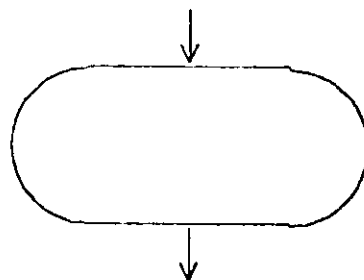


Figure 9. A computational structure.

To summarize the tree machine, the authors feel that this work has three contributions. The first is the abstract tree machine: it gives a number of nice "theoretical" solutions to a large set of problems. The second contribution is the architecture and implementation we have proposed; they indicate that this machine might be a reasonable device to build as further advances in VLSI technology occur. Finally, we feel that the "computational structure" we just investigated provides an example of the kind of argument that will justify special-purpose architectures proposed for implementation in VLSI.

Acknowledgements

The authors would like to acknowledge the careful comments of Dorothea Haken, Dave McKeown, Jim Saxe, Dick Sites and Siang Song and the helpful criticisms

received when presenting these ideas in a seminar at Xerox PARC.

References

- Backus, J. [1978]. "Can programming be liberated from the von Neumann style? A functional style and its algebra of programs," *Communications of the ACM* 21, 8 (August 1978), pp. 613-641.
- Bentley, J. L. [1979]. "A parallel algorithm for constructing minimum spanning trees," to appear in the Seventeenth Annual Allerton Conference on Communication, Control and Computing, October 1979.
- Bentley, J. L. and J. B. Saxe [1979]. "Decomposable searching problems," in preparation. (Preliminary version by J. L. Bentley, *Information Processing Letters* 8, 5 (June 1979), pp. 244-251.)
- Blaauw, G. A. and F. P. Brooks, Jr. [1979]. *Computer Architecture*, unpublished draft.
- Browning, S. [1979]. "Computation on a tree of processors," Caltech Internal Memorandum.
- Kung, H. T. [1979]. "Let's design algorithms for VLSI", Caltech Conference on Very Large Scale Integration: Architecture, Design, Fabrication, (January 1979).
- Kung, H. T. and C. Leiserson [1979]. "Systolic arrays (for VLSI)," *Carnegie-Mellon University Computer Science Research Review, 1977-78*, pp. 37-57. To appear in Mead and Conway [1979].
- Lamb, S. M. and R. Vanderslice [1978]. "Recognition memory: low cost content-addressable parallel processor for speech data manipulation," presented at the Acoustical Society of America and Acoustical Society of Japan Joint Meeting, Session BB, Honolulu (29 November 1978).
- Mago, G. [1979]. "A network of microprocessors to execute reduction languages," to appear in *International Journal of Computer and Information Sciences*.
- Mead, C. A. and L. A. Conway [1979]. *Introduction to VLSI Systems*, to appear.
- Mead, C. A. and M. Rem [1979]. "Cost and performance of VLSI Computing

Structures," *IEEE Journal of Solid State Circuits SC-14*, 2, April 1979, pp. 455-462.

Sequin, C. H., A. M. Despain, and D. A. Patterson [1978]. "Communication in X-tree, a modular multiprocessor system," *ACM 78 Proceedings*.

Sites, R. L. [1979]. "How to use 1000 registers," Caltech Conference on Very Large Scale Integration: Architecture, Design, Fabrication, (January 1979).

Su, S. Y. W., L. H. Nguyen, A. Emam, and G. J. Lipovski [1979]. "The architectural features and implementation techniques of the multicell casom," *IEEE Trans. Comp. C-28*, 6, (June 1979), pp. 430-445.

Swan, R. J., S. H. Fuller and D. P. Siewiorek [1977]. "Cm*: A Modular Multiprocessor," *AFIPS Conf. Proc. 46*, pp. 637-644.

Wilner, W. [1978]. "Recursive machines," Xerox PARC SSL Internal Memorandum (January 21, 1978).

A Parallel Algorithm for Constructing Minimum Spanning Trees¹

Jon Louis Bentley

Departments of Computer Science and Mathematics
Carnegie-Mellon University
Pittsburgh, Pennsylvania 15213

30 August 1979

Abstract

The construction of minimum spanning trees (MSTs) of weighted graphs is a problem that arises in many applications. In this paper we will study a new parallel algorithm that constructs the MST of an N -node graph in time proportional to $N \lg N$, on an $N/(\lg N)$ -processor computing system. The primary theoretical contribution of this paper is the new algorithm, which is an improvement over Sollin's parallel MST algorithm in several ways. On a more practical level, this algorithm is appropriate for implementation in VLSI technology.

¹This research was supported in part by the Defense Advanced Research Projects Agency under Contract F33615-78-C-1551 (monitored by the Air Force Office of Scientific Research) and in part by the Office of Naval Research under Contract N00014-76-C-0370.

A Parallel Algorithm for Constructing Minimum Spanning Trees¹

Jon Louis Bentley

Departments of Computer Science and Mathematics,
Carnegie-Mellon University
Pittsburgh, Pennsylvania 15213

30 August 1979

Abstract

The construction of minimum spanning trees (MSTs) of weighted graphs is a problem that arises in many applications. In this paper we will study a new parallel algorithm that constructs the MST of an N -node graph in time proportional to $N \lg N$, on an $N/(\lg N)$ -processor computing system. The primary theoretical contribution of this paper is the new algorithm, which is an improvement over Sollin's parallel MST algorithm in several ways. On a more practical level, this algorithm is appropriate for implementation in VLSI technology.

¹This research was supported in part by the Defense Advanced Research Projects Agency under Contract F33615-78-C-1551 (monitored by the Air Force Office of Scientific Research) and in part by the Office of Naval Research under Contract N00014-76-C-0370.

30 August 1979

A Parallel MST Algorithm

Table of Contents

1. Introduction
 2. Background on the Tree Machine
 3. The Prim-Dijkstra MST Algorithm
 4. The Algorithm on the Tree Machine
 5. Conclusions
- Acknowledgements**
- References**

1. Introduction

A minimum spanning tree (MST) of a weighted graph is defined as a subset of edges of the graph that connects all nodes and whose total edge weight is minimum. The problem of computing MSTs arises in many applications. In a communications system, for example, we might wish to link together a set of cities using minimum wire length; if only wires from one city to another are allowed, this calls precisely for finding the minimum spanning tree. (Indeed, this method is used for billing purposes by the Bell System!) The problem of computing MSTs also arises in such diverse applications areas as data analysis, operations research, and transportation networks. In this paper we will investigate a new method for rapidly computing MSTs on a parallel computer.

Much previous work has been devoted to the minimum spanning tree problem. In its most general setting, the problem is defined on graphs; Cheriton and Tarjan [1976] provide a thorough treatment of algorithms for the general problem. A special case of the MST problem (and one of great practical interest) is when the input graph is the complete graph induced by a set of points in Euclidean k -space. (That is, the vertices are the points in the set, and the edge weight between two vertices is their distance in space.) Shamos [1978] gives an optimal worst-case algorithm for computing the MST of a planar point set, and Yao [1977] describes the best-known worst-case algorithm for point sets in k -space. Fast expected-time algorithms are considered by Bentley, Weide and Yao [1978].

In this paper we will study a parallel algorithm for computing MSTs of graphs in which the edge weights are implicitly given by a distance function that returns the (nonnegative) distance between two nodes. The most common example of such graphs is the Euclidean case that we just saw; there are, however, other examples (for instance, the vertices might be bit strings and the distances their Hamming distance). Assuming that the distance function can be computed in constant time, the best known uniprocessor algorithm for this problem runs in (optimal) time proportional to N^2 on a graph of N vertices. The purpose of this paper is to present

an algorithm for solving this problem on an $N/(\lg N)$ -processor parallel machine in time proportional to $N \lg N$. (A similar result was achieved by Sollin [1977]; we will compare the new algorithm to his later in the paper.)

This paper is organized into three primary sections. In Section 2 we study a particular tree-structured parallel computer, and in Section 3 we will review a well-known algorithm for constructing MSTs. These sections provide the background for Section 4, in which we will see how the algorithm can be implemented very efficiently on the tree machine. Conclusions are then offered in Section 5.

2. Background on the Tree Machine

In this section we will examine the tree-structured searching machine described by Bentley and Kung [1979]. Although this discussion will provide the necessary background for applying the tree machine to the MST problem, it will be rather brief; the interested reader can refer to Bentley and Kung [1979] for additional details.

The basic structure of the tree machine is illustrated in Figure 1. It is a "mirrored" binary tree, containing three types of nodes: circles (which broadcast data), squares (which store data and compute), and triangles (which combine their inputs). To illustrate the operation of the tree machine, let us suppose that each of N square processors holds an element of set S , and we wish to determine whether 59 is a member of S . We enter 59 into the input node of the tree, and $\lg N$ steps later it reaches all the square nodes. We then perform N comparisons (in parallel), giving N bits that tell whether the element stored in each processor is 59. We can now combine those bits in $\lg N$ steps by letting the triangles compute the logical *or* of their two inputs. We therefore receive the answer to our query $2 \lg N$ steps after it was posed.

The above example illustrates the important aspects of the tree machine. The primary function of the circles is to broadcast data (they also decide where a new record is to be inserted into the tree). The triangles' purpose is to combine their inputs; the operations they perform include *min*, *max*, *and*, *or*, and *plus*. In addition

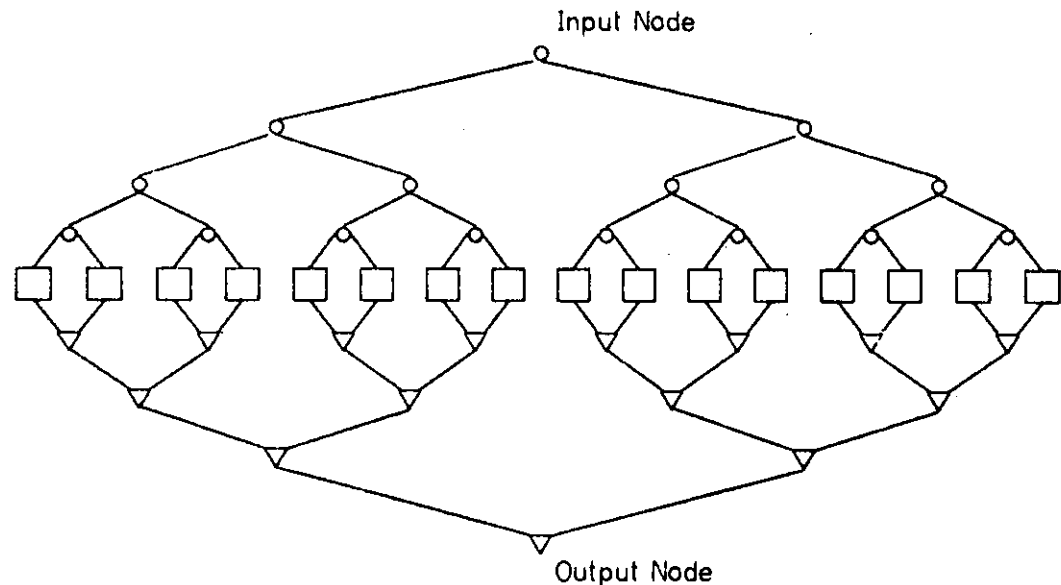


Figure 1. Structure of the tree machine.

to returning the minimum value (for instance), simple bookkeeping enables them to return also the name of the object achieving it. The squares are the most complicated nodes. Each is a "baby" von Neumann computer with a small amount of main memory. They receive their instructions and data in a single stream transmitted by the circles. (The stream for the above member search would be something like "the following is data, constant 59, load into register A, load register B from the memory word holding the key, and compare.") An important property of the tree machine is that M successive member searches can be "pipelined" (that is, the second goes down the tree immediately following the first, and so on) to run in $M-1 + 2 \lg N$ time.

Note that the tree machine can be simulated to within a constant time factor by a set of N processors operating on a large shared memory with constant access time. Such a machine is a more typical theoretical model of parallel computation (see, for example, Borodin and Munro [1975] or Kung [1979]), and the result of Section 4 holds in this model as well as for the tree machine.

Although this concludes our discussion of the tree machine on a "theoretical" level, there are two important points of practical interest that should be made. The

first is the fact that the tree machine is efficiently implementable in VLSI technology; this issue has been discussed at length by Bentley and Kung [1979]. The second point is that the algorithm we will see in Section 4 can be implemented on the concrete tree machine of Bentley and Kung [1979] simply by coding a program; no modification of the basic machine is required.

3. The Prim-Dijkstra MST Algorithm

In this section we will study an MST algorithm due to Prim [1957] and its computer implementation as given by Dijkstra [1959]. Because the algorithm has been examined at length in those pioneering papers (as well as in numerous texts), we will review it at a fairly intuitive level. Proofs of correctness and details of implementation abound in the literature.

Prim's algorithm works by starting with a *fragment* (subtree) that consists of exactly one vertex of the graph and then successively expanding the fragment until a spanning tree is obtained. Because the fragment is expanded at each step under a very local condition, this is referred to as a "greedy" algorithm. To specify that condition, we must give some definitions. The *distance* of a vertex to the fragment is defined as the weight of the shortest edge connecting that vertex to the fragment. The *nearest neighbor of the fragment* is then defined as that vertex not in the fragment with minimum distance to the fragment (with ties broken arbitrarily). It is now easy to describe Prim's algorithm. The fragment is initialized to contain an arbitrary vertex. We then expand the fragment at each stage by adding to it the nearest neighbor of the fragment. When applying this algorithm to an N -node graph, after $N-1$ stages we will have emptied the set of nonfragment nodes and thereby have built a spanning tree of the graph. Prim gives an elegant proof that the resulting spanning tree is indeed minimum (the proof of this particular algorithm is based only on Principle 2 of his paper).

If Prim's algorithm is implemented in a computer program by keeping track of the "nearest nonfragment neighbor" of every fragment node, then the resulting algorithm has $O(N^3)$ running time. Dijkstra [1959] makes the brilliant observation that an

$O(N^2)$ algorithm can be achieved by the dual strategy of keeping track of the "nearest fragment neighbor" of each nonfragment node. (Note that quadratic time is optimal because all $\binom{N}{2}$ edge weights must be examined in the general case.) The primary data structure of Dijkstra's algorithm contains, for each nonfragment node, the name of its nearest neighbor in the fragment and the distance to it. We initialize the program by choosing an arbitrary node as the single member of the fragment and assigning that node as the nearest neighbor of all other (nonfragment) nodes. In each successive stage of the algorithm, we choose the nonfragment node with minimum distance to nearest neighbor for insertion into the fragment and then update every nonfragment node to see if that chosen node is perhaps its nearest neighbor (we check if the distance to the chosen node is less than the distance to its current nearest neighbor). The informal method we have just seen is described more precisely as Algorithm 1.

```

TreeEdges  $\leftarrow \phi$ 
NonFragment  $\leftarrow \phi$ 
for  $i \leftarrow 2$  to  $N$  do
    New(a)
    a.Node  $\leftarrow i$ 
    a.NNinFrag  $\leftarrow 1$ 
    a.NNDist  $\leftarrow$  Distance(1,i)
    Insert a into NonFragment
while NonFragment  $\neq \phi$  do
    Assign b s.t. b.NNDist =  $\min_{c \in \text{NonFragment}} c.\text{NNDist}$ 
    Insert (b.Node, B.NNinFrag) into TreeEdges
    Delete b from NonFragment
    foreach a in NonFragment do
        ThisDist  $\leftarrow$  Distance(a.Node, b.Node)
        if ThisDist  $<$  a.NNDist then
            a.NNinFrag  $\leftarrow$  b.Node
            a.NNDist  $\leftarrow$  ThisDist

```

Algorithm 1. The Prim-Dijkstra MST algorithm.

Algorithm 1 assumes that the input vertices are named 1 through N and that the function Distance(i,j) returns the distance from node i to node j . The set TreeEdges

contains pairs of integers, and its final value gives the MST. The set NonFragment contains the nonfragment nodes, each of which is represented as a Pascal record with three fields: Node (an integer node number), NNinFrag (an integer name of the node's nearest neighbor in the fragment), and NNDist (the distance to the node's nearest neighbor). Whitney [1972] gives a Fortran implementation of Algorithm 1 that implements the NonFragment set as an array; the running time of his program is proportional to N^2 .

4. The Algorithm on the Tree Machine

In this section we will see how Algorithm 1 can be efficiently implemented on the tree machine that we studied in Section 2. We will use the tree machine to implement the NonFragment set; specifically, every record in the NonFragment set resides in a different square processor of the tree machine. We will assume that there is a small conventional computer "driving" the tree machine; that computer will also store the TreeEdges set. We will now describe how each part of Algorithm 1 can be implemented in this scheme.

Setting TreeEdges to be empty is trivial, and initializing NonFragment to be empty can be accomplished by broadcasting a "clear" command to all of the square processors in the tree machine. The for statement is implemented by performing the calculations in the "driver" computer, and then inserting each record into the tree machine. Since insertions can be pipelined, the total time required for this step is linear in N .

We now come to implementing the code in the while loop. We can find the record of NonFragment with the minimum NNDist field in logarithmic time by letting the triangles compute the minimum of those fields (and also returning the name of the record achieving that minimum). We then fetch the Node and NNinFrag fields of that record, insert the edge into TreeEdges, and delete the record from NonFragment (this amounts to turning off a processor). We must now perform the foreach statement, and it is here that we employ the massive parallelism of the tree machine. To every record currently stored in the tree, we (simultaneously) ship

down a description of `b.Node` (sufficient to compute the distance function) and then compute `ThisDist`. If `ThisDist` is greater than the `NNDist` field of the particular processor, we (temporarily) turn it off. We now broadcast the assign commands to alter the `NNinFrag` and `NNDist` fields, which are executed only by the listening processors. We then turn back on all the processors currently representing elements on `NonFragment`, and repeat the while loop until `NonFragment` is empty.

The method that we just described is easy to program on the VLSI tree machine of Bentley and Kung [1979], assuming that the distance function between pairs of nodes can be computed by the square processors in constant time. The initialization steps (assigning sets to be empty and initializing `NonFragment`) require linear time on the tree machine. Each iteration of the while loop requires $O(\lg N)$ time. (Data must travel from the input node to the output node three times: once to find the minimum `NNDist`, once to find the associated values and delete the old record, and once to perform the updates of the foreach statement.) Since this step is executed $N-1$ times, the total time used by the step is $O(N \lg N)$. That cost dominates the cost of the initialization, and the total running time of the Prim-Dijkstra MST algorithm implemented on the tree machine is therefore proportional to $N \lg N$.

To illustrate the operation of this parallel algorithm on a more concrete example, let us consider the problem of constructing the MST of a set of N points in k -space (using the Euclidean metric). Initializing the set `NonFragment` involves storing k reals in each of N processors; this can be accomplished in Nk time. In the while statement, finding the minimum can be accomplished in $\lg N$ time, and the name of the record realizing the minimum can also be retrieved (and deleted) that quickly. We must now execute the foreach statement. This involves entering the k coordinates of the point most recently added to the fragment to all of the N processors (and computing the sum of the squares of the coordinate-wise distances); this can be accomplished in time proportional to k . Updating the appropriate fields requires constant time, so the total time per step of this algorithm is proportional to $k + \lg N$. The total running time of the algorithm is therefore

proportional to $N \lg N + Nk$.¹

The MST algorithm of this section can be used to yield an $N \lg N$ time algorithm on a tree machine with $N/(\lg N)$ square processors, each of which has $\lg N$ storage. To accomplish this, we store $\lg N$ vertices at each square node, and implement the foreach statement by executing its code $\lg N$ times. Each iteration of the while loop still costs $\lg N$ time, and the total cost of the algorithm is still proportional to $N \lg N$. The processor-time product of this modified algorithm is N^2 , rather than the $N^2 \lg N$ of the previous version. Note that this scheme can be used to give an $N \cdot G(N)$ time MST algorithm on an $N/G(N)$ -processor machine for any function G that grows more quickly than $\lg N$ (the trick is to store $G(N)$ vertices in every square processor).

5. Conclusions

Sollin's [1977] algorithm is another method that has been proposed for the parallel construction of MSTs in $N \lg N$ time on an N -processor machine. The method that we saw in Section 4 is a fundamentally different approach to the problem than Sollin's: our method performs N iterations of logarithmic cost, whereas Sollin's performs $\lg N$ iterations of linear cost. A particular advantage of our method over Sollin's is that while his performs $N^2 \lg N$ distance calculations in the worst case, ours performs only $\binom{N}{2}$, which is optimal. A related advantage of ours is that it performs only N^2 operations altogether (unlike Sollin's $N^2 \lg N$), and this approach might therefore lead to an N -processor algorithm with linear running time. We also saw that for any function G growing as fast as $\lg N$, our algorithm can be implemented in $N \cdot G(N)$ time on an $N/G(N)$ -processor machine, giving an optimal processor-time product of $O(N^2)$; this is superior to Sollin's processor-time product of $O(N^2 \lg N)$. Yet another parallel MST algorithm has been given by Savage [1977]. Her algorithm runs in $O(\lg^2 N)$ time on a machine with $O(N^2/\lg N)$ processors, and is

¹It is possible to use the tree machine to achieve an asymptotically faster MST algorithm for point sets in k -space by implementing Yao's [1977] "directional nearest neighbor" MST algorithm. The directional nearest neighbors are found by pipelining the tree machine, and the resulting MST is then computed on a uniprocessor. The running time of the resulting algorithm is $O[(N \lg \lg N + Nk) \cdot f(k)]$, where $f(k)$ is the number of directional nearest neighbor searches Yao's algorithm requires in k dimensions. The parallel implementation of the Prim-Dijkstra algorithm would probably be faster than this method in almost every practical application, however.

therefore applicable in an entirely different sphere than the algorithm of Section 4. (Note that the processor-time product of her algorithm is proportional to $N^2 \lg N$.)

The contributions of this paper are in two areas. In a "theoretical" setting, we have investigated an implementation of the Prim-Dijkstra MST algorithm on a parallel computer. While this algorithm has the same running time as a previous algorithm, it uses an optimal number of distance calculations, and therefore gives hope for a faster parallel algorithm. A modified version of the algorithm runs on $N/G(N)$ processors, with optimal processor-time product. We also saw a faster parallel MST algorithm for the special case of finding the MST of a point set in Euclidean k -space, answering an open problem in computational geometry (see Shamos [1978, p. 210, Problem 7]). On a more practical level, the algorithm that we have investigated uses a machine with a simple and regular interconnection strategy that can be efficiently implemented in VLSI technology.

Acknowledgements

The helpful comments of Kevin Brown, Dorothea Haken, Jim Saxe, Michael Shamos, and Herb Wilf are gratefully acknowledged.

References

- Bentley, J. L., B. W. Weide and A. C. Yao [1978]. "Optimal expected-time algorithms for closest-point problems," Carnegie-Mellon Computer Science Department Report CMU-CS-79-111 (March 1979).
- Bentley, J. L. and H. T. Kung [1979]. "A tree machine for searching problems," *Proceedings of the 1979 International Conference on Parallel Processing*, (August 1979), pp. 257-266, IEEE.
- Borodin, A. and I. Munro [1975]. *The Computational Complexity of Algebraic and Numeric Problems*, American Elsevier, New York.
- Cheriton, D. and R. E. Tarjan [1976]. "Finding minimum spanning trees," *SIAM J. Computing* 5, 4, (December 1976), pp. 724-742.

- Dijkstra, E. W. [1959]. "A note on two problems in connexion with graphs," *Numerisch Math.* 1, (October 1959), pp. 269-271.
- Kung, H. T. [1979]. "The structure of parallel algorithms," to appear in *Advances in Computers*.
- Prim, R. C. [1957]. "Shortest connection networks and some generalizations," *Bell System Technical Journal* 36, (November 1957), pp. 1389-1401.
- Savage, C. D. [1977]. Parallel algorithms for graph theoretic problems. University of Illinois Coordinated Science Laboratory Report UILU-ENG-77-2231.
- Shamos, M. I. [1978]. Computational Geometry, unpublished Ph.D. Thesis, Yale University, May 1978.
- Sollin [1977]. An algorithm attributed to Sollin in *Introduction to the Design and Analysis of Algorithms* by S. E. Goodman and S. T. Hedetniemi (McGraw-Hill, 1977), Section 5.5. (Also referred to in Yao [1977].)
- Whitney, V. K. M. [1972]. "Algorithm 422, minimal spanning tree [H]," *Comm. of the ACM* 15, 4 (April 1972), pp. 273-274.
- Yao, A. C. [1977]. "On constructing minimum spanning trees in k-dimensional space and related problems," Stanford Computer Science Department Report STAN-CS-77-642 (December 1977).

REPORT DOCUMENTATION PAGE		READ INSTRUCTIONS BEFORE COMPLETING FORM
1. REPORT NUMBER CMU-CS-79-142	2. GOVT ACCESSION NO.	3. RECIPIENT'S CATALOG NUMBER
4. TITLE (and Subtitle) A TREE MACHINE FOR SEARCHING PROBLEMS		5. TYPE OF REPORT & PERIOD COVERED Interim
		6. PERFORMING ORG. REPORT NUMBER
7. AUTHOR(s) JON LOUIS BENTELY		8. CONTRACT OR GRANT NUMBER(s) N00014-76-C-0370
9. PERFORMING ORGANIZATION NAME AND ADDRESS Carnegie-Mellon University Computer Science Department Pittsburgh, PA 15213		10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS
11. CONTROLLING OFFICE NAME AND ADDRESS Office of Naval Research Arlington, VA 22217		12. REPORT DATE August 1979
		13. NUMBER OF PAGES 39
14. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office)		15. SECURITY CLASS. (of this report) UNCLASSIFIED
		15a. DECLASSIFICATION/DOWNGRADING SCHEDULE
16. DISTRIBUTION STATEMENT (of this Report) Approved for public release; distribution unlimited.		
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)		
20. ABSTRACT (Continue on reverse side if necessary and identify by block number)		