

# Unstructured Tree Search on SIMD Parallel Computers: Experimental Results \*

George Karypis and Vipin Kumar  
Department of Computer Science  
University of Minnesota  
Minneapolis, MN 55455

## Abstract

In this paper, we present new methods for load balancing of unstructured tree computations on large-scale SIMD machines, and analyze the scalability of these and other existing schemes. An efficient formulation of tree search on a SIMD machine comprises of two major components: (i) a triggering mechanism, which determines when the search space redistribution must occur to balance search space over processors; and (ii) a scheme to redistribute the search space. We have devised a new redistribution mechanism and a new triggering mechanism. Either of these can be used in conjunction with triggering and redistribution mechanisms developed by other researchers. We analyze the scalability of these mechanisms, and verify the results experimentally. The analysis and experiments show that our new load balancing methods are highly scalable on SIMD architectures. Their scalability is shown to be no worse than that of the best load balancing schemes on MIMD architectures. We verify our theoretical results by implementing the 15-puzzle problem on a CM-2<sup>1</sup> SIMD parallel computer.

## 1 Introduction

Tree search is central to solving a variety of problems in artificial intelligence [11, 25], combinatorial optimization [10, 20], operations research [24] and Monte-Carlo evaluations of functional integrals [30]. The trees that need to be searched for most practical prob-

lems happen to be quite large, and for many tree search algorithms, different parts can be searched relatively independently. These trees tend to be highly irregular in nature and hence, a naive scheme for partitioning the search space can result in highly uneven distribution of work among processors and lead to poor overall performance. The job of partitioning irregular search spaces is particularly difficult for SIMD parallel computers such as the CM-2, in which all processors work in lock-step to execute the same program. The reason is that in SIMD machines, work distribution needs to be done on a global scale (*i.e.* if a processor becomes idle, then it has to wait until the entire machine enters a work distribution phase). In contrast, on MIMD machines, an idle processor can request work from another busy processor without any other processor being involved. Many efficient load balancing schemes have already been developed for dynamically partitioning large irregular trees for MIMD parallel computers [1, 3, 5, 22, 28, 31, 32], whereas until recently, it was common wisdom that such irregular problems cannot be solved on large-scale SIMD parallel computers [20].

Recent research has shown that data parallel SIMD architectures can also be used to implement parallel tree search algorithms effectively. Powley, Korf and Ferguson [26, 27] and Mahanti and Daniels [2, 21] present parallel formulations of a tree search algorithm IDA\*, for solving the 15 puzzle problem on CM-2. Frye and Myczkowski [4] presents an implementation of a depth-first tree search algorithm on the CM-2 for a block puzzle.

The load balancing mechanisms used in the implementations of Frye, Powley, and Mahanti are different

---

\*This work was supported by IST/SDIO through the Army Research Office grant #28408-MA-SDI and by the Army High Performance Computing Research Center at the University of Minnesota.

<sup>1</sup>CM-2 is a registered trademark of Thinking Machines Corporation.

from each other. From the experimental results presented, it is difficult to ascertain the relative merits of these different mechanisms. This is because the performance of different schemes may be impacted quite differently by changes in hardware characteristics (such as interconnection network, CPU speed, speed of communication channels etc.), number of processors, and the size of the problem instance being solved [16]. Hence any conclusions drawn on a set of experimental results may become invalid by changes in any one of the above parameters. Scalability analysis of a parallel algorithm and architecture combination is very useful in extrapolating these conclusions [8, 16, 18]. The iso-efficiency metric has been found to be quite useful in characterizing scalability of a number of algorithms [7, 19]. In particular, it has helped determine optimal load balancing schemes for tree search for a variety of MIMD architectures [18, 6, 15].

In this paper, we present new methods for load balancing of unstructured tree computations on large-scale SIMD machines. The experiments show that our new load balancing methods are highly scalable on SIMD architectures. In particular, from the analysis presented in [12], their scalability is no worse than that of the best load balancing schemes on MIMD architectures.

Section 2 provides a description of existing load balancing schemes and the new schemes we have developed. Section 3 and 4 present the experimental evaluation of static and dynamic triggering. Section 5 comments on other related work in this area.

## 2 Dynamic Load Balancing Algorithms for Parallel Search

Specification of a tree search problem includes a description of the root node of the tree and a successor-generation-function that can be used to generate successors of any given node. Given these two, the entire tree can be generated and searched for goal nodes. Often strong heuristics are available to prune the tree at various nodes. The tree can be generated using different methods. Depth-first method is used in many important tree search algorithms such as Depth-First Branch and Bound [14], IDA\* [13], Backtracking [10]. In this paper we consider parallel depth-first-search on SIMD machines.

A common method used for parallel depth-first-search of dynamically generated trees on a SIMD machine [26, 21, 4] is as follows. At any time, all the

processors are either in a *search* phase or in a *load balancing* phase. In the search phase, each processor searches a disjoint part of the search space in a depth-first-search (DFS) fashion by performing node expansion cycles in lock-step. When a processor has finished searching its part of the search space, it stays idle until it gets additional work during the next load balancing phase. All processors switch from the search phase to the load balancing phase when a triggering condition is satisfied. In the load balancing phase, busy processors split their work and share it with idle processors. When a goal node is found, all processors quit. If the search space is finite and has no solutions, then eventually all the processors would run out of work, and parallel search will terminate.

Since each processor searches the space in a depth-first manner, the (part of) state space to be searched is efficiently represented by a stack. The depth of the stack is the depth of the node being currently explored; and each level of the stack keeps track of untried alternatives. Each processor maintains its own local stack on which it executes depth-first-search. The current unsearched tree space, assigned to any processor, can be partitioned into two parts by simply partitioning untried alternatives (on the current stack) into two parts. A processor is considered to be busy if it can split its work into two non empty parts, one for itself and one to give away. In the rest of this paper, a processor is considered to be busy if it has at least two nodes on its stack. We denote the number of idle processors by  $I$ , the number of busy processors by  $A$  and the total number of processors by  $P$ . Also, the terms busy and active processors will be used interchangeably.

### 2.1 Previous Schemes for Load Balancing

The first scheme we study is similar to one of the schemes proposed in [21]. In this algorithm, the triggering condition is computed after each node expansion cycle in the search phase. If this condition is satisfied, then a load balancing phase is initiated. In the load balancing phase, idle processors are matched one-to-one with busy processors. This is done by enumerating both the idle and the busy processors; then each busy processor is matched with the idle processor that received the same value during this enumeration. The busy processors split their work into two

parts and transfer one part to their corresponding idle processors<sup>2</sup>. If  $I > A$  then only the first  $A$  idle processors are matched to busy ones and the remaining  $I - A$  processors receive no work. After each load balancing phase, at least one node expansion cycle is completed before the triggering condition is tested again.

A very simple and intuitive scheme [26, 4] is to trigger a load balancing phase when the ratio of active to the total number of processors falls below a fixed threshold. Formally, let  $x$  be a number such that  $0 \leq x \leq 1$ , then the triggering condition for this scheme is:

$$A \leq xP \quad (1)$$

For the rest of this paper we will refer to this triggering scheme as the **static triggering** scheme with threshold  $x$  (in short the  $S^x$ -triggering scheme).

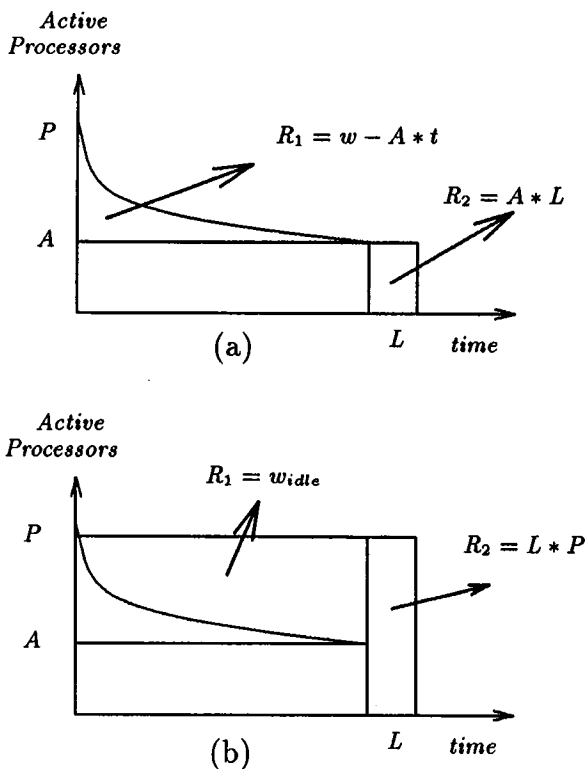


Figure 1: A graphical representation of the triggering conditions for the  $D^P$ -triggering and for the  $D^K$ -triggering schemes.

An alternative to static triggering is to use a trigger value that changes dynamically in order to adapt itself to the characteristics of the problem. We call this kind

<sup>2</sup>This is done using the rendezvous allocation scheme described in [9].

of triggering scheme **dynamic triggering**  $D$ . A dynamic triggering scheme was presented and analyzed by Powley, Korf and Ferguson in [26]. For the rest of this paper we will refer to it as the  $D^P$ -triggering scheme. The  $D^P$ -triggering works as follows: Let  $w$  be the sum of the time spent by processors, let  $t$  be the elapsed time since the beginning of the current search phase and let  $L$  be the time required to perform the next load balancing phase. After every node expansion cycle, the ratio  $\frac{w}{t+L}$  is compared against the number of active processors  $A$ , and a load balance is initiated as soon as that ratio is greater or equal to  $A$ . In other words the condition that triggers a load balance is:

$$\frac{w}{t+L} \geq A \quad (2)$$

Because the value of  $L$  cannot be known (it requires knowledge of the future), it is approximated by the cost of the previous load balancing phase.  $D^P$  is a locally greedy approach that tries to maximize the average rate of work over a search and load balancing phase. Polwey *et. al.* also describe variants of  $D^P$ -triggering in [26].

Another way of stating the triggering condition for  $D^P$  is to rewrite eqn (2) as:

$$w - A * t \geq A * L \quad (3)$$

From this equation and Fig. 1(a) we see that the  $D^P$ -triggering scheme will trigger a load balancing phase as soon as the area  $R_1$  is greater or equal to area  $R_2$ .

## 2.2 Our New Schemes for Load Balancing

We have derived a new matching scheme for mapping idle to busy processors in the load balancing phase. This method can be used with either the static or the dynamic triggering schemes. We have also derived a new dynamic triggering scheme.

The new mapping algorithm is similar to the one described earlier but with the following modification. We now keep a pointer that points to the last processor that gave work during the last load balancing phase. Every time we need to load balance, we start matching busy processors to idle processors, starting from the first busy processor after the one pointed by this pointer. When the pointer reaches the last processor, it starts again from the beginning. For the rest

of this paper we will call this pointer **global pointer** and this mapping scheme *GP*. Also, due to the absence of the global pointer we will name the mapping scheme of Section 2.1, *nGP*.

Figure 2 illustrates the *GP* and the *nGP* matching schemes with an example. Assume that at the time when a load balancing phase is triggered, processors 6 and 7 are idle and the others are busy. Also, assume that the global pointer points to processor 5. Now, *nGP* will match processors 6 and 7 to processors 1 and 2 respectively, whereas *GP* will match them to processors 8 and 1 respectively and it will advance the global pointer to processor 1. If after the next search phase, processors 6 and 7 are idle again and the others remain busy, then *nGP* will match them exactly as before where *GP* will match them to processors 2 and 3. The above example also provides the motiva-

Processors	1 2 3 4 5 6 7 8
<b>example 1</b>	
state	BBBBBIIB
global pointer	↑
<i>nGP</i> enumeration of busy processors	1 2 3 4 5 6
<i>GP</i> enumeration of busy processors	2 3 4 5 6 1
enumeration of idle processors	12
<b>example 2</b>	
state	BBBBBIIB
global pointer	↑
<i>nGP</i> enumeration of busy processors	1 2 3 4 5 6
<i>GP</i> enumeration of busy processors	6 1 2 3 4 5
enumeration of idle processors	12

Figure 2: Illustration of the *GP* and *nGP* matching schemes. *B* is used to denote busy processors while *I* is used to denote idle ones.

tion behind *GP*, which is to try to evenly distribute the burden of sharing work among the processors. As shown in [12] the upper bound on the number of load balancing phases required for *GP* is much smaller than that for *nGP*. When  $x \leq 0.5$  both schemes are similar.

Our new dynamic triggering scheme, called  $D^K$ -triggering, takes a different approach than the  $D^P$ -triggering scheme. Formally, let  $w_{idle}$  be the sum of the idle time of all the processors since the beginning of the current search phase and let  $L * P$  be the cost of the next load balancing phase, then the condition that will trigger a load balance is:

$$w_{idle} \geq L * P \quad (4)$$

Fig. 1(b) illustrates this condition,  $R_1$  is  $w_{idle}$  and  $R_2$  is  $L * P$ . This scheme will trigger a load balancing phase as soon as  $R_1 \geq R_2$ . Note that if triggering takes

place earlier than this point, then the load balancing overhead will be higher than the overhead due to idling and vice versa. Thus, our triggering scheme balances the idle time of the processors during the search phase and the cost of the next load balancing phase.

### 3 Static Triggering: Experimental Results

We solved various instances of the 15-puzzle problem [23] taken from [13], on a CM-2 massively parallel SIMD computer. 15-puzzle is a  $4 \times 4$  square tray containing 15 square tiles. The remaining sixteenth square is uncovered. Each tile has a number on it. A tile that is adjacent to the blank space can be slid into that space. An instance of the problem consists of an initial position and a specified goal position. The goal is to transform the initial position into the goal position by sliding the tiles around. The 15-puzzle problem is particularly suited for testing the effectiveness of dynamic load balancing schemes, as it is possible to create search spaces of different sizes ( $W$ ) by choosing appropriate initial positions. IDA\* is the best known sequential depth-first-search algorithm to find optimal solution paths for the 15-puzzle problem [13], and generates highly irregular search trees. We have parallelized IDA\* to test the effectiveness of the various load balancing algorithms. The same algorithm was also used in [26, 21]. Our parallel implementations of IDA\* find all the solutions of the puzzle up to a given tree depth. This ensures that the number of nodes expanded by the serial and the parallel search is the same, and thus we avoid having to consider superlinear speedup effects [29, 26, 21].

We obtained experimental results using both the *nGP* and the *GP* matching schemes for different values of static threshold  $x$ . In our implementation, each node expansion cycle takes about 10ms while each load balancing cycle takes about 34ms. Every time work is split we transfer the node at the bottom of the stack. For the 15-puzzle problem, this appears to provide a reasonable alpha-splitting mechanism. In calculating efficiencies, we used the average node expansion cycle time of parallel IDA\* as an approximation of the sequential node expansion cost. Because of higher node expansion cost associated with SIMD parallel computers [26], the actual efficiencies are lower by a constant ratio than those presented here. However, this does not affect the relative comparison of any of these schemes.

Static Trigger		0.50		0.60		0.70		0.80		0.90	
W	Metric	nGP	GP	nGP	GP	nGP	GP	nGP	GP	nGP	GP
2488958	$N_{expand}$	547	547	479	483	438	438	400	406	384	379
	$N_{lb}$	62	62	105	69	179	78	309	102	376	154
	$E$	0.41	0.41	0.37	0.43	0.29	0.43	0.21	0.41	0.19	0.35
9076121	$N_{expand}$	1957	1957	1708	1730	1520	1563	1364	1429	1320	1325
	$N_{lb}$	72	72	245	84	560	102	1095	134	1317	226
	$E$	0.51	0.51	0.43	0.55	0.32	0.58	0.22	0.59	0.19	0.54
21540929	$N_{expand}$	4629	4629	4078	4091	3588	3687	3234	3376	3104	3113
	$N_{lb}$	73	73	407	90	1251	107	2339	151	3044	250
	$E$	0.54	0.54	0.48	0.60	0.33	0.65	0.24	0.68	0.20	0.67
45584793	$N_{expand}$	9510	9510	8457	8450	7565	7637	6881	7009	6475	6494
	$N_{lb}$	73	73	515	84	1880	106	3636	150	5911	260
	$E$	0.57	0.57	0.54	0.64	0.39	0.70	0.29	0.74	0.21	0.76

Table 1: Experimental results obtained using 8192 CM-2 processors.  $N_{expand}$  is the number of node expansion cycles,  $N_{lb}$  is the number of load balancing phases and  $E$  is the efficiency.

Some of these results are shown in Table 1. All the timings in this table have been taken on 8k processors. From the results shown in this table, we clearly see how  $GP$  and  $nGP$  relate to each other. When  $x = 0.50$  both algorithms perform similarly, where, as predicted by our theoretical analysis presented in [12], the difference between the performance of  $GP$  and  $nGP$  increases as  $x$  increases. From the results shown in Table 1, we also see that the relative performance of  $GP$  versus  $nGP$  increases as  $W$  increases. The reason for that is explained in [12].

#### 4 Dynamic Triggering, Experimental Results

We implemented all four combinations of the two dynamic triggering schemes  $D^P$  and  $D^K$ , and the two matching schemes  $nGP$  and  $GP$ , in the parallel IDA\* to solve the 15-puzzle problem on CM-2. In all cases, the root node is given to one of the processors and static triggering with  $x = 0.85$  is used until 85% of the processors became active. Thus in the initial distribution phase, each node expansion cycle was followed by a work distribution cycle until 85% of the processors had work. After the initialization phase, triggering was done using the respective dynamic triggering schemes. The results are summarized in Table 2.

From the results shown in this table we can see that for the  $nGP$  matching scheme, the  $D^K$ -triggering scheme performs slightly worse than the  $D^P$ -triggering scheme for larger problems. For the  $GP$  matching scheme,  $D^K$ -triggering performs consistently better than  $D^P$ -triggering for all problems. Also the  $GP$  matching scheme constantly outperforms  $nGP$  for both dynamic triggering schemes as it does for static triggering. Comparing the two dynamic triggering schemes in Table 2, with the static triggering scheme in Table 1, we see that  $GP$ - $D^K$  performs as good as the  $GP$ - $S^x$  schemes using optimal trigger values for

Dynamic Trigger		$D^P$ -triggering		$D^K$ -triggering	
W	Metric	nGP	GP	nGP	GP
2488958	$N_{expand}$	685	596	552	500
	* $N_{lb}$	137	100	83	68
	$E$	0.29	0.34	0.37	0.42
9076121	$N_{expand}$	2002	1778	1758	1467
	* $N_{lb}$	161	110	186	114
	$E$	0.45	0.52	0.46	0.60
21540929	$N_{expand}$	4436	3894	4002	3200
	* $N_{lb}$	214	116	378	183
	$E$	0.52	0.62	0.49	0.69
45584793	$N_{expand}$	8682	7517	8014	6406
	* $N_{lb}$	314	137	669	279
	$E$	0.57	0.70	0.54	0.77

Table 2: Experimental results obtained using 8192 CM-2 processors using various dynamic triggering schemes.  $N_{expand}$  is the number of node expansion cycles, \* $N_{lb}$  is the number of work transfers and  $E$  is the efficiency. Note that for the  $D^K$ -triggering scheme \* $N_{lb}$  is equal to the number of load balancing phases.

each problem.

#### 5 Related Work

Powley, Korf and Ferguson [26, 27] present load balancing algorithms which have different triggering and matching schemes. Their matching scheme, min- $f$  ordering, is identical to  $nGP$  or  $GP$  when the number of active processors is less than the number of idle ones; but when the active processors are more than the idle ones, work is given out from processors that have nodes with smaller  $f$ -value (*i.e.*, the lower bound on the cost of the node). The primary motivation behind the min- $f$  ordering is that nodes with smaller  $f$ -values are likely to represent more work than those with larger  $f$ -values. Clearly, the min- $f$  ordering should lead to no more load balancing phases than  $nGP$ . It may even lead to fewer number of load balancing phases than  $GP$  depending upon how good a predictor the  $f$ -value is of the overall load. But unlike  $GP$ , it is difficult to put an upper bound on the number

of load balancing phases for the min- $f$  ordering scheme. Also, implementation of min- $f$  requires a sorting step which is more time consuming than simple enumeration required by  $GP$ . This will become important if the cost of sorting is high compared with the rest of the load balancing phase. In our experiments with the 15-puzzle problem, we found min- $f$  and  $GP$  to require about the same number of load balancing phases. But, in our experimental setup, the load balancing cost for min- $f$  is about 2.5 times that of  $GP$ . Hence  $GP$  leads to somewhat higher efficiencies. Powley *et. al.* also report experimental results for a matching scheme which randomly matches idle processors to busy processors, and point out that its performance is similar to the min- $f$  scheme.

Powley *et. al.* present three different triggering schemes: static triggering,  $D^P$ -triggering and a variation of  $D^P$ -triggering with the following modifications: A load balancing cycle is triggered when either the  $D^P$ -triggering condition holds or when the number of active processors is less than  $P/2$  and the time spent searching is at least half of the time spent in load balancing. These modifications to  $D^P$  guarantee that at least one third of the total time is spent searching, and alleviate many of the drawbacks of the original  $D^P$ -triggering scheme. Besides load balancing within iterations of IDA\*, Powley *et. al.* perform load balancing during the initial work distribution and between iterations of IDA\*. These steps may not be applicable, in general, to dynamic distribution of unstructured trees on parallel computers.

Mahanti and Daniels proposed two dynamic load balancing algorithms, FESS and FECS, in [21, 2]. In both these schemes a load balancing phase is initiated as soon as one processor becomes idle and the matching scheme used is similar to  $nGP$ . The difference between FESS and FECS is that during each load balancing phase FESS performs a single work transfer while FECS performs as many work transfers as required so that the total number of nodes is evenly distributed among the processors. As our analysis has shown the FESS scheme has poor scalability and because this scheme usually performs as many load balancing phases as node expansion cycles, its performance depends on the relative costs of communication and node expansion. FECS performs better work distribution thus requires fewer number of load balancing phases. Hence it has better performance than FESS, and its scalability may be close to the

$GP$  matching scheme. The memory requirements of FECS is unbounded, and modifications to handle this problem are discussed in [21].

Frye and Myczkowski proposed two dynamic load balancing algorithms in [4]. The first scheme is similar to  $nGP-S^x$  with the difference that each busy processor gives one piece of work to as many idle processors as many pieces of work it has. Clearly this scheme has a poor splitting mechanism. Extending this algorithm in such a way so that the total number of nodes is evenly distributed among the processors results in a scheme similar to FECS of Mahanti *et. al.* The second algorithm is based on nearest neighbor communication. In this scheme after each node expansion cycle the processors that have work check to see if their neighbors are idle. If this is the case then they transfer work to them. This scheme is similar to the nearest neighbor load balancing schemes for MIMD machines thus as shown in in [17] it is sensitive to the quality of the work splitting mechanism. Hence, this algorithm is sensitive to the quality of the alpha-splitting mechanism.

## References

- [1] S. Arvindam, Vipin Kumar, V. Nageshwara Rao, and Vineet Singh. *Automatic test Pattern Generation on Multi-processors*. *Parallel Computing*, 17, number 12:1323-1342, December 1991.
- [2] M. Evett, James Hendler, Ambujashka Mahanti, and Dana Nau. *PRA\*: A Memory-Limited Heuristic Search Procedure for the Connection Machine*. In *Proceedings of the third symposium on the Frontiers of Massively Parallel Computation*, pages 145-149, 1990.
- [3] Raphael A. Finkel and Udi Manber. *DIB - A Distributed implementation of Backtracking*. *ACM Trans. of Progr. Lang. and Systems*, 9 No. 2:235-256, April 1987.
- [4] Roger Frye and Jacek Myczkowski. *Exhaustive Search of Unstructured Trees on the Connection Machine*. In *Thinking Machines Corporation Technical Report*, 1990.
- [5] M. Furuichi, K. Taki, and N. Ichiyoshi. *A Multi-Level Load Balancing Scheme for OR-Parallel Exhaustive Search Programs on the Multi-PSI*. In *Proceedings of the 2nd ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, 1990. pp.50-59.
- [6] Ananth Grama, Vipin Kumar, and V. Nageshwara Rao. *Experimental Evaluation of Load Balancing Techniques for the Hypercube*. In *Proceedings of the Parallel Computing 91 Conference*, 1991.
- [7] Anshul Gupta and Vipin Kumar. *The scalability of FFT on Parallel Computers*. In *Proceedings of the Frontiers 90*

- Conference on Massively Parallel Computation*, October 1990. An extended version of the paper will appear in *IEEE Transactions on Parallel and Distributed Systems*, 1993.
- [8] John L. Gustafson, Gary R. Montry, and Robert E. Benner. *Development of Parallel Methods for a 1024-Processor Hypercube*. *SIAM Journal on Scientific and Statistical Computing*, 9 No. 4:609–638, 1988.
- [9] W. Daniel Hillis. *The Connection Machine*. MIT Press, 1991.
- [10] Ellis Horowitz and Sartaj Sahni. *Fundamentals of Computer Algorithms*. Computer Science Press, Rockville, Maryland, 1978.
- [11] Laveen Kanal and Vipin Kumar. *Search in Artificial Intelligence*. Springer-Verlag, New York, 1988.
- [12] George Karypis and Vipin Kumar. *Unstructured Tree Search on SIMD Parallel Computers*. Technical Report 92–21, University of Minnesota, 1992.
- [13] Richard E. Korf. *Depth-First Iterative-Deepening: An Optimal Admissible Tree Search*. *Artificial Intelligence*, 27:97–109, 1985.
- [14] Vipin Kumar. *DEPTH-FIRST SEARCH*. In Stuart C. Shapiro, editor, *Encyclopaedia of Artificial Intelligence: Vol 2*, pages 1004–1005. John Wiley and Sons, Inc., New York, 1987. Revised version appears in the second edition of the encyclopedia to be published in 1992.
- [15] Vipin Kumar, Ananth Grama, and V. Nageshwara Rao. *Scalable Load Balancing Techniques for Parallel Computers*. Technical report, Tech Report 91-55, Computer Science Department, University of Minnesota, 1991.
- [16] Vipin Kumar and Anshul Gupta. *Analyzing Scalability of Parallel Algorithms and Architectures*. Technical report, TR-91-18, Computer Science Department, University of Minnesota, June 1991. A short version of the paper appears in the Proceedings of the 1991 International Conference on Supercomputing, Germany, and as an invited paper in the Proc. of 29th Annual Allerton Conference on Communication, Control and Computing, Urbana, IL, October 1991.
- [17] Vipin Kumar, Dana Nau, and Laveen Kanal. *General Branch-and-bound Formulation for AND/OR Graph and Game Tree Search*. In Laveen Kanal and Vipin Kumar, editors, *Search in Artificial Intelligence*. Springer-Verlag, New York, 1988.
- [18] Vipin Kumar and V. Nageshwara Rao. *Parallel Depth-First Search, Part II: Analysis*. *International Journal of Parallel Programming*, 16 (6):501–519, 1987.
- [19] Vipin Kumar and Vineet Singh. *Scalability of Parallel Algorithms for the All-Pairs Shortest Path Problem: A Summary of Results*. In *Proceedings of the International Conference on Parallel Processing*, 1990. Extended version appears in *Journal of Parallel and Distributed Processing (special issue on massively parallel computation)*, Volume 13, 124-138, 1991.
- [20] Karp R. M. *Challenges in Combinatorial Computing*. To appear January 1991.
- [21] A. Mahanti and C. Daniels. *SIMD Parallel Heuristic Search*. To appear in *Artificial Intelligence*, 1992.
- [22] V. Nageshwara Rao and Vipin Kumar. *Parallel Depth-First Search, Part I: Implementation*. *International Journal of Parallel Programming*, 16 (6):479–499, 1987.
- [23] Nils J. Nilsson. *Principles of Artificial Intelligence*. Tioga Press, 1980.
- [24] Christos H. Papadimitriou and Kenneth Steiglitz. *Combinatorial Optimization, Algorithms and Complexity*. Prentice Hall, 1982.
- [25] Judea Pearl. *Heuristics - Intelligent Search Strategies for Computer Problem Solving*. Addison-Wesley, Reading, MA, 1984.
- [26] C. Powley, R. Korf, and C. Ferguson. *IDA\* on the Connection Machine*. To appear in *Artificial Intelligence*, 1992.
- [27] Curt Powley and Richard E. Korf. *SIMD and MIMD Parallel Search*. In *Proceedings of the AAAI Spring Symposium*, pages 49–53, 1989.
- [28] Abhiram Ranade. *Optimal Speedup for Backtrack Search on a Butterfly Network*. In *Proceedings of the Third ACM Symposium on Parallel Algorithms and Architectures*, 1991.
- [29] V. Nageshwara Rao and Vipin Kumar. *On the Efficiency of Parallel Backtracking*. *IEEE Transactions on Parallel and Distributed Systems*, (to appear), 1992. available as a technical report TR 90-55, Computer Science Department, University of Minnesota.
- [30] Jasec Myczkowski Roger Frye. *Load Balancing Algorithms on the Connection Machine and their Use in Monte-Carlo Methods*. In *Proceedings of the Unstructured Scientific Computation on Multiprocessors Conference*, 1992.
- [31] Wei Shu and L. V. Kale. *A Dynamic Scheduling Strategy for the Chare-Kernel System*. In *Proceedings of Supercomputing 89*, pages 389–398, 1989.
- [32] Benjamin W. Wah and Y. W. Eva Ma. *MANIP - A Multicomputer Architecture for Solving Combinatorial Extremum-Search Problems*. *IEEE Transactions on Computers*, c-33, May 1984.