

## A Cluster Based Load-Balancing Method for Multi-Computer Based Solution-Adaptive Finite Element Graphs

<sup>1</sup>K.M. Sakib, <sup>2</sup>M.H. Kamal and <sup>2</sup>U. Kabir

<sup>1</sup>Institute of Information Technology, <sup>2</sup>Department of Computer Science and Engineering,  
University of Dhaka, Dhaka 1000, Bangladesh

**Abstract:** Load imbalance can degrade the performance of a solution-adaptive finite element application program on a distributed memory multicomputer. To solve the load imbalance problem load of a refined finite element graph can be redistributed based on the current load of each processor. In this study, a load-balancing algorithm is applied to balance the computational load of each processor. A distributed method for load balancing is proposed which is based on the global load balancing information and current load distribution of the system. A simulation model is developed to compare the performance of the proposed method with the existing methods such as MCSTLB, BTLB and CBTLB methods. The execution time and the number of process migration required by different load balancing methods are used for performance evaluation. The experimental result shows that the proposed method is more efficient than that of the existing methods.

**Key words:** Distributed memory multicomputers, load balancing, solution-adaptive finite element graphs

### INTRODUCTION

The finite element method is widely used for the structural modeling of physical systems. In the finite element model, an object can be viewed as a finite element graph, which is a connected and undirected graph that consists of a number of finite elements. Each finite element is composed of a number of nodes. Due to the properties of computation-intensiveness and computation locality, it is very attractive to implement the finite element method on distributed memory multicomputers (Angus *et al.*, 1990; Fox *et al.*, 1988; Simon, 1991; Williams, 1990, 1991). In the context of parallelizing a finite element application program that uses iterative techniques to solve system of equations (Aykanat *et al.*, 1987), a parallel program may be viewed as a collection of tasks represented by nodes of a finite element graph. Each node represents a particular amount of computation and can be executed independently. To efficiently execute a finite element application program on a distributed memory multicomputer, we need to map nodes of the corresponding finite element graph to processors of a distributed memory multicomputer such that each processor has approximately the same amount of computational load and the communication among processors is minimized. Since this mapping problem is known to be NP-complete (Garey and Johnson, 1979), many heuristic method were proposed to find satisfactory

suboptimal solutions (Barnard and Simon, 1994, 1995; Ercal *et al.*, 1990; Fiduccia and Mattheyses, 1982; Gilbert and Zmijewski, 1987; Gilbert *et al.*, 1995; Hendrickson and Leland, 1995a, b; Karypis and Kumar, 1995 a, b; Kernigham and Lin, 1970; Simon, 1991; Williams, 1991).

For a solution-adaptive finite element application program, the number of nodes increases discretely due to the refinement of some finite elements during the execution. This may result in load imbalance of processors. A node remapping or a load-balancing algorithm has to be performed many times in order to balance the computational load of processors while keeping the communication cost among processors as low as possible. For the load balancing approach, some load-balancing algorithms can be used to perform the load balancing process according to the current load of processors. Load-balancing algorithms are performed at run-time; their execution should be fast and efficient.

In this study, a cluster based load-balancing method has been proposed to efficiently deal with the load imbalance problems of solution-adaptive finite element application programs on distributed memory multicomputers. When nodes of a solution-adaptive finite element graph were evenly distributed to processors by some mapping algorithms, according to the communication property of the finite element graph, we can get a processor graph from the partition. For example, Fig. 1 shows a partition of a 21-node finite element graph

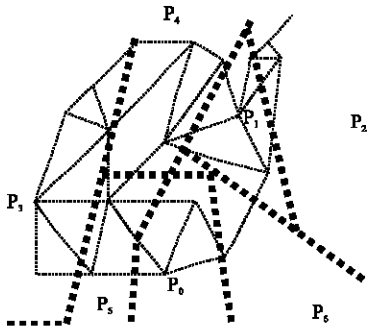


Fig. 1: A partition of 21-node finite element graph on 7 processors

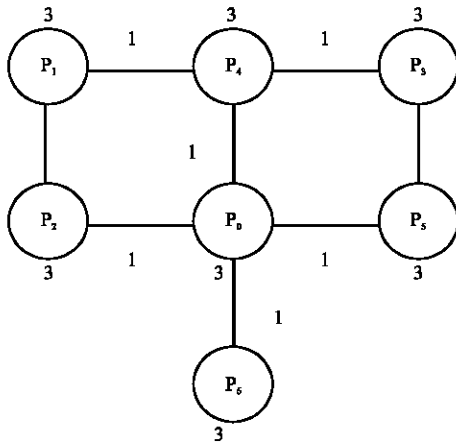


Fig. 2: The corresponding processor graph of Fig. 1

on 7 processors. The corresponding processor graph of Fig. 1 is shown in Fig. 2. In a processor graph, nodes represent the processors and edges represent the communication needed among processors. The weights associated with nodes and edges denote the computation and the communication costs, respectively.

When a finite element graph is refined during runtime, it will result in load imbalance of processors. To balance the computational load of processors, the Cluster method first builds up clusters of processors. Based on clusters, the global load balancing information is calculated by the Tree-Walking Algorithm (TWA) (Wu, 1997). According to the global load balancing information and the current load distribution, a load transfer algorithm is performed to balance the computational load of processors and minimize the communication cost among processors.

To evaluate the performance of the proposed method, this method was implemented along with three other tree-based parallel load-balancing methods, the

MCSTLB method (Chung and Liao, 1999), BTLB method (Chung and Liao, 1999) and CBTLB method (Chung and Liao, 1999). The experimental results show that the execution time and the number of process migration of an application program under a cluster based load-balancing method are always smaller than those of the other methods.

### THE PARALLEL LOAD BALANCING METHODS

**The Maximum Cost Spanning Tree Load-Balancing (MCSTLB) method:** The main idea of the MCSTLB method (Chung and Liao, 1999) is to find a maximum cost spanning tree from the processor graph that is obtained from the initial partitioned finite element graph. The MCSTLB method can be divided into the following 4 phases:

**Phase 1:** Obtain a processor graph  $G$  from the initial partition.

**Phase 2:** Use a similar Kruskal's (1956) algorithm to find a maximum cost spanning tree  $T = (V, E)$  from  $G$ . There are many ways to determine the shape of  $T$ . In this method, the shape of  $T$  is constructed as follows:

- The processor with the largest degree in  $V$  is selected as the root of  $T$ .
- For each nonterminal processor  $v$  in  $T$ , if  $\{u_1, \dots, u_m\}$  are the  $m$  children of  $v$  and  $|u_1| < |u_2| < \dots < |u_m|$ , then  $u_1$  will be the leftmost child of  $v$ ,  $u_2$  will be the second leftmost child of  $v$  and so on, where  $|u_i|$  is the degree of  $u_i$  and  $i = 1, \dots, m$ . If the depth of  $T$  is greater than  $\log M$ , where  $M$  is the number of processors, we will try to adjust the depth of  $T$ . The adjusted method is first to find the longest path (from a terminal processor to another terminal processor) of  $T$ . After the longest path is determined, the middle processor of the path is selected as the root of the tree and the tree is reconstructed according to the above construction process. If the depth of the reconstructed tree is less than that of  $T$ , the reconstructed tree is the desired tree. Other-wise,  $T$  is the desired tree. The purpose of the adjustment is trying to reduce the load balancing steps among processors.

**Phase 3:** Calculate the global load balancing information and schedule the load transfer sequence of processors by using the TWA (Wu, 1997). Assume that there are  $M$  processors in a tree and  $N$  nodes in a refined finite element graph. We define  $N/M$  as the average weight of

a processor. In the TWA method, the quota and the load of each processor in a tree are calculated, where the quota is the sum of the average weights of a processor and its children processors and the load is the sum of the weights of a processor and its children processors. The difference of the quota and the load of a processor is the number of nodes that a processor should send to or receive from its parent. If the difference is negative, a processor should send nodes to its parent. Otherwise, a processor should receive nodes from its parent. According to the global load balancing information, a schedule can be determined.

**Phase 4:** Perform load transfer (send/receive) based on the global load balancing information, the schedule and  $T$ . Assume that processor  $P_i$  needs to send  $m$  nodes to processor  $P_j$  and let  $N$  denote the set of nodes in  $P_i$  that are adjacent to those of  $P_j$ . In order to keep the communication cost as low as possible, in the load transfer, nodes in  $N$  are transferred first. If  $|N|$  is less than  $m$ , then nodes adjacent to those in  $N$  are transferred. This process is continued until the number of nodes transferred to  $P_j$  is equal to  $m$ .

**The Binary Tree Load Balancing (BTLB) method:** The BTLB method (Chung and Liao, 1999) is similar to the MCSTLB method (Chung and Liao, 1999). The only difference between these two methods is that the MCSTLB method is based on a maximum cost spanning tree to balance the computational load of processors while the BTLB method is based on a binary tree. The BTLB method can be divided into the following four phases:

**Phase 1:** Obtain a processor graph  $G$  from the initial partition.

**Phase 2:** Use a similar Kruskal's (1956) algorithm to find a binary tree  $T = (V, E)$  from  $G$ , where  $V$  and  $E$  denote the processors and edges of  $T$ , respectively. The method to determine the shape of a binary tree is the same as that of the MCSTLB method.

**Phase 3:** Calculate the global load balancing information and schedule the load transfer sequence of processors by using the TWA.

**Phase 4:** Perform load transfer (send/receive) based on the global load balancing information, the schedule and  $T$ . The load transfer method is the same as that of the MCSTLB method.

**The Condensed Binary Tree Load Balancing (CBTLB) method:** The main idea of the CBTLB method (Chung and Liao, 1999) is to group processors of the processor

graph into metaprocessors. Each metaprocessor is a hypercube. The CBTLB method can be divided into the following 5 phases:

**Phase 1:** Obtain a processor graph  $G$  from the initial partition.

**Phase 2:** Group processors of  $G$  into metaprocessors to obtain a condensed processor graph  $G_c$  incrementally. The metaprocessors in  $G_c$  are constructed as follows: First, a processor  $P_i$  with the smallest degree in  $G$  and a processor  $P_j$  that is a neighbor processor of  $P_i$  and has the smallest degree among those neighbor processors of  $P_i$  are grouped into a metaprocessor. Then, the same construction is applied to other ungrouped processors until there are no processors can be grouped into a hypercube. Repeat the grouping process to each metaprocessor until there are no metaprocessors can be grouped into a higher order hypercube.

**Phase 3:** Find a binary tree  $T = (V, E)$  from  $G_c$ , where  $V$  and  $E$  denote the metaprocessors and edges of  $T$ , respectively. The method of constructing a binary tree is the same as that of the BTLB method.

**Phase 4:** Based on  $T$ , calculate the global load balancing information and schedule the load transfer sequence by using a similar TWA method for metaprocessors. To obtain the global load balancing information, the quota and the load of each processor in a tree are calculated. The quota is defined as the sum of the average weights of processors in a metaprocessor  $C_i$  and processors in children processors of  $C_i$ . The load is defined as the sum of the weights of processors in a metaprocessor  $C_i$  and processors in children metaprocessors of  $C_i$ . The difference of the quota and the load of a metaprocessor is the number of nodes that a metaprocessor should send to or receive from its parent metaprocessor. After calculating the global load balancing information, the schedule is determined as follows. Assume that  $m$  is the number of nodes that a metaprocessor  $C_i$  needs to send to another metaprocessor  $C_j$ . We have the following two cases:

**Case 1:** If the weight of  $C_i$  is less than  $m$ , the schedule of these two metaprocessors is postponed until the weight of  $C_i$  is greater than or equal to  $m$ .

**Case 2:** If the weight of  $C_i$  is greater than or equal to  $m$ , a schedule can be made between processors of  $C_i$  and  $C_j$ . Assume that  $ADJ$  denotes the set of processors in  $C_i$  that are adjacent to those in  $C_j$ . If the sum of the weights of processors in  $ADJ$  is less than  $m$ , a schedule is made to transfer nodes of processors in  $C_i$  to processors in  $ADJ$

such that the weights of processors in ADJ is greater than or equal to  $m$ . If the sum of the weights of processors in ADJ is greater than or equal to  $m$ , a schedule is made to send  $m$  nodes from processors in ADJ to those in C.

**Phase 5:** Perform load transfer (send/receive) among metaprocessors based on the global load balancing information, the schedule and T. The load transfer method is similar to that of the BTLB method. After performing load transfer process among metaprocessors, a Dimension Exchange Method (DEM) is performed to balance the computational load of processors in metaprocessors.

### THE PROPOSED CLUSTER BASED LOAD BALANCING METHOD

The main idea of cluster based method is to construct an arrangement of processors, where the processors are combined into groups. After the construction of processor group, the load information for each processor is collected and the load balancing algorithm is performed such that the processor can balance their load by transferring minimum number of processes and the overall load balancing time is also improved.

**Phase 1:** Group construction.

**Step 1:** Divide  $N$  number of processors into  $N/3$  number of groups. In a group there might be one or two or three processors. In each case the group might be constructed as following:

**Case 1:** If a group has three nodes, then one of them is called the parent node and the other two are left and right child, respectively.

**Case 2:** If a group has two nodes, then one of them is called the parent node and the other is the left child.

**Case 3:** If a group has only one node, then it is the parent node.

In each group the children nodes send their state information to the parent node when they try to balance the load.

If there is only one group, then go to Phase 3.

**Step 2:** Group three local groups to form a large group. In this large group, one node acts as parent and other two as left and right child, respectively.

This process of constructing large group is continued until there is only one large group.

**Phase 2:** Load estimation

Each processor in the system has varying number of processes and each process has varying amount of load. To find the average weight or Quota of a processor we first have to calculate the sum of loads of all processors and then we divide the total sum by the number of processors of the system. Thus we obtain the Quota for each processor and from the Quota we calculate the high threshold and low threshold value for each processor, where

High threshold = Quota +  $x$  (where  $x = 5\%$  of quota)

Low threshold = Quota -  $x$  (where  $x = 5\%$  of Quota)

Now a processor's state is defined as follows:

**Case 1:** The processor is in normal state if its load is greater than low threshold and less than high threshold.

**Case 2:** The processor is in underloaded state if its load is below the low threshold.

**Case 3:** The processor is in overloaded state if its load is above the high threshold.

**Phase 3:** Load distribution.

**Step 1:** In this level, for each group the group load and the group quota is calculated. The group load is defined as the sum of loads of each processor in a group, which is not in normal state and the group quota is the sum of the quota for each processor in the group, which is not in normal state. From the group Quota, the high threshold and low threshold is also calculated for the group. Now depending on the group load and threshold values of the group, the following two cases may occur.

**Case 1:** If the group load is greater than low threshold and less than high threshold, then it is possible to balance the load of the group internally. For each member node of the group, the difference of quota and load is the number of processes that a node should send or receive from other nodes. If the difference is negative, a node should transfer load, otherwise it should receive loads.

**Case 2:** If the group load is greater than the high threshold value or less than the low threshold value, then load balancing is not possible within the group. In this case, the parent will contain the group load information.

If the load of all groups, in this level is balanced, then the load distribution process is terminated. Otherwise step 1 is repeated until a higher level large group exists.

**Step 2:** When the largest group is reached, the group load is distributed among the members of the group, which is not in normal state. For each member node of the group, the difference of quota and load is the number of processes that a node should send or receive from other nodes. If the difference is negative, a node should transfer load, otherwise it should receive loads. Then, each group of the next lower level distributes the load among the processors of that group in the same way.

This process of load distribution is repeated until any lower level group exists.

**RESULTS AND DISCUSSION**

This study compares the performance of the load-balancing methods by implementing the algorithm with some simulation programs. The criteria used to evaluate the performance are execution time and the number of processes to be migrated to balance the system load.

**Comparison of execution time of different load balancing methods:** The execution time of different load balancing methods, with 7, 15, 25, 30 and 40 processors are shown in Table 1.

From Table 1, it is notified that among MCSTLB, BTLB and CBTLB method, the execution time of CBTLB method is better than the other two. This is because the CBTLB method can reduce the size of a tree with a large ratio so that the overheads to do the load transfer among the metaprocessors are less than those of the MCSTLB and BTLB method. Thus it can reduce the load transfer time efficiently. We also observe that the execution time for the Cluster method is less than that of CBTLB method. This is because the CBTLB method does not try to balance the load within a metaprocessor after forming the group. As a result a metaprocessor, which can be balanced locally, is grouped into higher-level hypercube.

Table 1: The execution time in seconds of different load balancing method for different load samples with different number of processors

Methods	No. of processes				
	7	15	25	30	40
MCSTLS	1.500549	1.500549	1.500549	1.500549	1.554396
BTLB	1.103846	1.10549	1.100000	1.154396	1.100000
CBTLB	1.500549	1.500549	1.500549	1.500549	1.500549
Cluster	0.659340	0.692308	0.714286	0.714286	0.714286

Table 2: Number of process migration of different load balancing method for different load samples with different number of processors

Methods	No. of processes									
	5	7	10	15	20	25	30	35	40	45
MCSTLB	278	440	719	1270	1841	2370	3113	3550	4135	4562
BTLB	301	509	824	1460	2156	2937	2156	4068	4810	4204
CBTLB	432	782	1389	2286	3509	4426	5556	6617	7355	8126
Cluster	112	166	277	312	464	494	577	601	720	1223

This makes fruitless process transfer and thus takes more time to balance the load. Though, in Cluster methods, grouping is performed in each refinement, it takes less time to balance the system load.

**Comparison of the number of process migration of different methods:** The number of processes to be migrated in different load balancing methods, in a system with 5, 7, 10, 15, 20, 25, 30, 35, 40 and 45 processors are shown in Table 2.

**CONCLUSION**

Different types of load-balancing algorithm for solution-adaptive finite element application program on distributed memory multicomputers were proposed. These are MCSTLB method, the BTLB method, the CBTLB method and the Cluster method. In MCSTLB method, BTLB method and CBTLB method, a logical tree (a maximum cost spanning tree for MCSTLB method, a binary tree for BTLB method and a condensed binary tree for CBTLB method) is constructed from a processor graph. Based on the tree structure and the current load of the system, an existing method tries to balance the system load. But in those methods, the static nature of the logical tree makes a huge number of process migrations, which consume not only time but also the communication network bandwidth.

In this study, a new, improved group-based method is proposed to balance the load among the sites of a distributed memory multicomputer system to overcome the problems associated with the previous methods. In this method, the processors are grouped so that the members of a group can try to balance their load within the group without knowing the states of the other processors belonging to a different group. Otherwise, when balancing the load within the group is not possible, this group tries to balance the load in a large group. Thus, in this method a process is migrated only then, when it finds its suitable destination. So the discussion concludes that the proposed method requires fewer process migration and less execution time than the existing methods.

To evaluate the performance of the existing load balancing methods and the proposed one, the algorithms are implemented with some simulation programs. Two criteria are execution time and the number of process migration of different algorithms for an application program is used for performance evaluation. The experiment result shows that the execution time and number of process to be migrated of the proposed method is better than that of the existing methods.

**REFERENCES**

- Angus, I.G., G.C. Fox, J.S. Kim and D.W. Walker, 1990. Solving Problems on Concurrent Processors. Englewood Cliffs, N.J.: Prentice Hall, Vol. 2.
- Aykanat, C., F. Ozgüner, S. Martin and S.M. Doraivelu, 1987. Parallelization of a Finite Element Application Program on a Hypercube Multiprocessor, Hypercube Multiprocessor, pp: 662-673.
- Barnard, S.T. and H.D. Simon, 1995. A Parallel Implementation of Multilevel Recursive Spectral Bisection for Application to Adaptive Unstructured Meshes, Proc. 7th SIAM Conf. Parallel Processing for Scientific Computing, San Francisco, pp: 627-632.
- Barnard, S.T. and H.D. Simon, 1994. Fast Multilevel Implementation of Recursive Spectral Bisection for Partitioning Unstructured Problems, Concurrency: Practice and Experience, 6: 101-117.
- Chung, Y.C. and C.J. Liao, 1999. Tree-Based Parallel Load Balancing Methods for Solution-Adaptive Finite Element Graphs on Distributed Memory Multicomputer. IEEE. Trans. Parallel and Distrib. Sys., 10: 360-370.
- Ercal, F., J. Ramanujam and P. Sadayappan, 1990. Task Allocation onto a Hypercube by Recursive Mincut Bipartitioning. J. Parallel Distrib. Comput., 10: 35-44.
- Fiduccia, C.M. and R.M. Mattheyes, 1982. A Linear-Time Heuristic for Improving Network Partitions. Proc. 19th IEEE. Design Automation Conf., pp: 175-181.
- Fox, C., M. Johnson, G. Lyzenga, S. Otto, J. Salman and D.W. Walker, 1988. Solving Problems on Concurrent Processors. Englewood Cliffs, N.J.: Prentice Hall, Vol. 1.
- Garey, M.R. and D.S. Johnson, 1979. Computers and Intractability, A Guide to Theory of NP-Completeness. San Francisco: Freeman.
- Gilbert, J.R. and E. Zmijewski, 1987. A Parallel Graph Partitioning Algorithm for a Message-Passing Multiprocessor. Int. J. Parallel Programming, 16: 427-449.
- Gilbert, J.R., G.L. Miller and S.H. Teng, 1995. Geometric Mesh Partitioning: Implementation and Experiments. Proc. 9th Int. Parallel Processing Symp., Santa Barbara, Calif., pp: 418-427.
- Hendrickson, B. and R. Leland, 1995. A Multilevel Algorithm for Partitioning Graphs, Proceeding of Supercomputing.
- Hendrickson, B. and R. Leland, 1995. An Improved Spectral Graph Partitioning Algorithm for Mapping Parallel Computations. SIAM J. Scientific Computing, 16: 452-469.
- Karypis, G. and V. Kumar, 1995. Multilevel k-way Partitioning Scheme for Irregular Graphs, Technical Report 95-064, Department of Computer Science, University of Minnesota, Minneapolis.
- Karypis, G. and V. Kumar, 1995. MeTiS-Unstructured Graph Partitioning and Spares Matrix Ordering System. University of Minnesota.
- Kernigham, B.W. and S. Lin, 1970. An Efficient Heuristic Procedure for Partitioning Graphs. Bell Sys. Technol. J., 49: 292-370.
- Kruskal, J.B., 1956. On the Shortest Spanning Subtree of a Graph and the Traveling Salseman Problem. Proc. AMS., 7: 48-50.
- Simon, H.D., 1991. Partitioning of Unstructured Problems for Parallel Processing. Comput. Sys. Eng., 2: 135-148.
- Williams, R.D., 1991. Performance of Dynamic Load Balancing Algorithms for Unstructured Mesh Calculations, Concurrency: Practice and Experience, 3: 457-481.
- Williams, R.D., 1990. DIME: Distributed Irregular Mesh Environment. California Institute of Technology.
- Wu, M.Y., 1997. On Runtime Parallel Scheduling for Processor Load Balancing. IEEE. Trans. Parallel and Distrib. Sys., 8: 173-186.