

Hierarchical Partitioning and Dynamic Load Balancing for Scientific Computation

James D. Teresco¹, Jamal Faik², and Joseph E. Flaherty²

¹ Department of Computer Science, Williams College
Williamstown, MA 01267 USA
terescoj@cs.williams.edu

² Department of Computer Science, Rensselaer Polytechnic Institute
Troy, NY 12180, USA
{faikj,flaherje}@cs.rpi.edu

Abstract. Cluster and grid computing has made hierarchical and heterogeneous computing systems increasingly common as target environments for large-scale scientific computation. A cluster may consist of a network of multiprocessors. A grid computation may involve communication across slow interfaces. Modern supercomputers are often large clusters with hierarchical network structures. For maximum efficiency, software must adapt to the computing environment. We focus on partitioning and dynamic load balancing, in particular on hierarchical procedures implemented within the Zoltan Toolkit, guided by DRUM, the Dynamic Resource Utilization Model. Here, different balancing procedures are used in different parts of the domain. Preliminary results show that hierarchical partitionings are competitive with the best traditional methods on a small hierarchical cluster.

Modern three-dimensional scientific computations must execute in parallel to achieve acceptable performance. Target parallel environments range from clusters of workstations to the largest tightly-coupled supercomputers. Hierarchical and heterogeneous systems are increasingly common as symmetric multiprocessing (SMP) nodes are combined to form the relatively small clusters found in many institutions as well as many of today's most powerful supercomputers. Network hierarchies arise as grid technologies make Internet execution more likely and modern supercomputers are built using hierarchical interconnection networks. MPI implementations may exhibit very different performance characteristics depending on the underlying network and message passing implementation (*e.g.*, [32]). Software efficiency may be improved using optimizations based on system characteristics and domain knowledge. Some have accounted for clusters of SMPs by using a hybrid programming model, with message passing for inter-node communication and multithreading for intra-node communication (*e.g.*, [1, 27]), with varying degrees of success, but always with an increased burden on programmers to program both levels of parallelization.

Our focus has been on resource-aware partitioning and dynamic load balancing, achieved by adjusting target partition sizes or the choice of a dynamic

load-balancing procedure or its parameters, or by using a combination of load-balancing procedures. We retain the flat message passing programming model. For hierarchical and heterogeneous systems, different choices of load balancing procedures may be appropriate in different parts of the parallel environment. There are tradeoffs in execution time and partition quality (*e.g.*, surface indices, interprocess connectivity, strictness of load balance) [35] and some may be more important than others in some circumstances. For example, consider a cluster of SMP nodes connected by Ethernet. A more costly graph partitioning can be done to partition among the nodes, to minimize communication across the slow network interface, possibly at the expense of some computational imbalance. Then, a fast geometric algorithm can be used to partition independently within each node.

1 Partitioning and Dynamic Load Balancing

An effective partitioning or dynamic load balancing procedure maximizes efficiency by minimizing processor idle time and interprocessor communication. While some applications can use a static partitioning throughout a computation, others, such as adaptive finite element methods, have dynamic workloads that necessitate dynamic load balancing during the computation. Partitioning and dynamic load balancing can be performed using recursive bisection methods [2, 29, 31, 38], space-filling curve (SFC) partitioning [7, 23–25, 37] and graph partitioning (including spectral [26, 29], multilevel [6, 18, 20, 36], and diffusive methods [8, 19, 22]). Each algorithm has characteristics and requirements that make it appropriate for certain applications; see [4, 35] for examples and [33] for an overview of available methods.

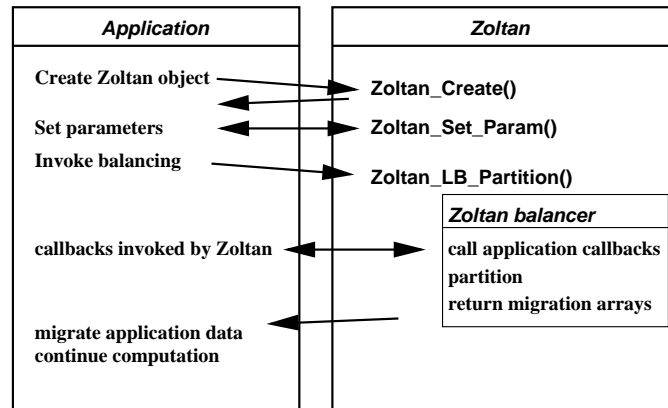


Fig. 1. Interaction between Zoltan and applications.

The Zoltan Parallel Data Services Toolkit [9, 11] provides dynamic load balancing and related capabilities to a wide range of dynamic, unstructured and/or adaptive applications. Using Zoltan, application developers can switch partitioners simply by changing a run-time parameter, facilitating comparisons of the partitioners’ effect on the applications. Zoltan has a simple interface, and its design is “data-structure neutral.” That is, Zoltan does not require applications to construct or use specific data structures. It operates on generic “objects” that are specified by calls to application-provided callback functions. These callbacks are simple functions that return to Zoltan information such as the lists of objects to be partitioned, coordinates of objects, and topological connectivity of objects. Figure 1 illustrates the interaction between Zoltan and an application.

We focus here on the hierarchical balancing procedures we have implemented within Zoltan, where different procedures are used in different parts of the computing environment.

2 Hierarchical Balancing

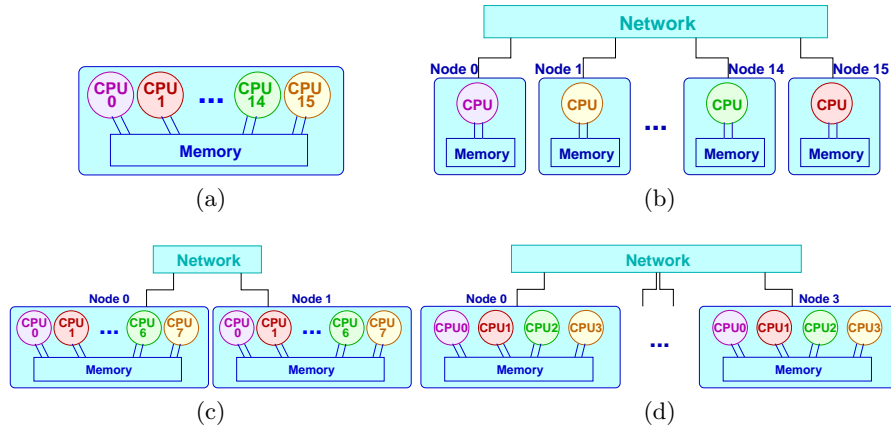


Fig. 2. Examples of parallel computing environments with 16 processors: (a) a 16-way SMP workstation; (b) a 16-node computer with all uniprocessor nodes, connected by a network; (c) two 8-way SMP workstations connected by a network; and (d) four 4-way SMP workstations connected by a network.

Consider four different 16-way partitionings of a 1,103,018-element mesh used in a simulation of blood flow in a human aorta [30]. Here, the geometric method recursive inertial bisection (RIB) [28] and/or the multilevel graph partitioner in ParMetis [20] are used for partitioning. Only two partition quality factors are considered: computational balance, and two *surface index* measures, although other factors [5] should be considered. The *global surface index* (*GSI*) measures

the overall percentage of mesh faces inter-partition boundaries, and the *maximum local surface index* (*MLSI*) measures the maximum percentage of any one partition's faces that are on a partition boundary. The surface indices are essentially normalized edge cuts when considering mesh connectivity in terms of a graph. Partitioning using RIB achieves an excellent computational balance, with partitions differing by no more than one element, and medium-quality surface indices with $MLSI = 1.77$ and $GSI = 0.61$. This would be useful when the primary concern is a good computational balance, as in a shared-memory environment (Figure 2b). Using only ParMetis achieves excellent surface index values, $MLSI = 0.40$ and $GSI = 0.20$, but at the expense of a large computational imbalance, where partition sizes range from 49,389 to 89,302 regions. For a computation running on a network of workstations (NOW) (Figure 2b), it may be worth accepting the significant load imbalance to achieve the smaller communication volume.

For hierarchical systems, a hybrid partitioning may be desirable. Consider the two SMP configurations connected by a slow network as shown in Figure 2c,d. In the two 8-way node configuration (Figure 2c), ParMetis is used to divide the computation between the two SMP nodes, resulting in partitions of 532,063 and 570,955 regions, with $MLSI = 0.06$ and $GSI = 0.03$. Within each SMP, the mesh is partitioned eight ways using RIB, producing partitions within each SMP balanced to within one element, and with overall $MLSI = 1.88$ and $GSI = 0.56$. Since communication across the slow network is minimized, this is an appropriate partitioning for this environment. A similar strategy for the four 4-way SMP configuration (Figure 2d) results in a ParMetis partitioning across the four SMP nodes with 265,897, 272,976, 291,207 and 272,938 elements and $MLSI = 0.23$ and $GSI = 0.07$. Within each SMP, partitions are again balanced to within one element, with overall $MLSI = 1.32$ and $GSI = 0.32$. We first presented an example of this type of hybrid partition in [32], although the partitions presented therein were not generated by a fully automatic procedure.

Zoltan's hierarchical balancing automates the creation of such partitions. It can be used directly by an application or be guided by the tree representation of the computational environment created and maintained by the Dynamic Resource Utilization Model (DRUM) [10, 13, 34]. DRUM is a software system that supports automatic resource-aware partitioning and dynamic load balancing for heterogeneous, non-dedicated, and hierarchical computing environments. DRUM dynamically models the computing environment using a tree structure that encapsulates the capabilities and performance of communication and processing resources. The tree is populated with performance data obtained from *a priori* benchmarks and dynamic monitoring agents that run concurrently with the application. It is then used to guide partition-weighted and hierarchical partitioning and dynamic load balancing. Partition-weighted balancing is discussed further in [13]. DRUM's graphical configuration program (Figure 3) may be used to facilitate the specification of hierarchical balancing parameters at each network and multiprocessing node.

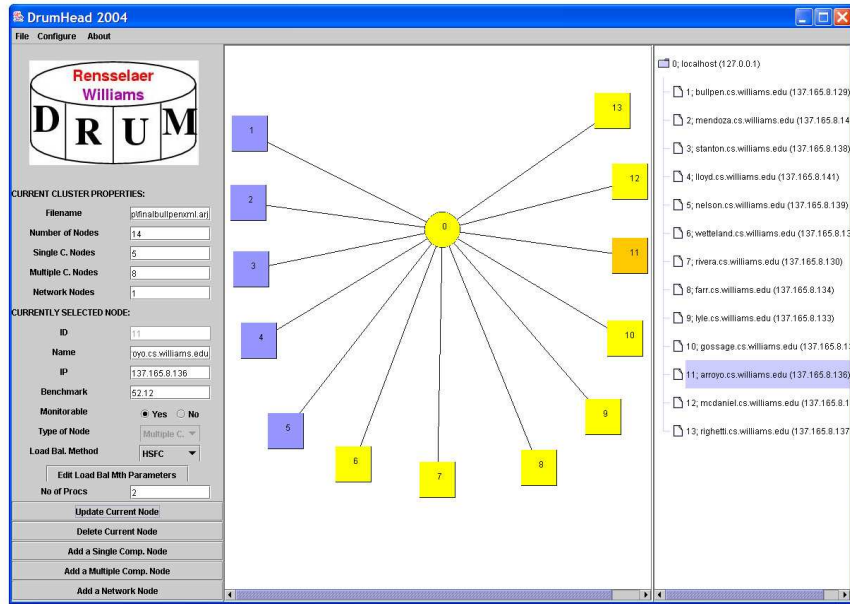


Fig. 3. DRUM's graphical configuration program being used to edit a description of the “Bullpen Cluster” at Williams College. This program may be used to specify hierarchical balancing parameters for Zoltan.

The hierarchical balancing implementation utilizes a lightweight “intermediate hierarchical balancing structure” (IHBS) and a set of callback functions. This permits an automated and efficient hierarchical balancing which can use any of the procedures available within Zoltan without modification and in any combination. Hierarchical balancing is invoked by an application in the same way as other Zoltan procedures. A hierarchical balancing step begins by building an IHBS using these callbacks. The IHBS is an augmented version of the distributed graph structure that Zoltan builds to make use of the ParMetis [21] and Jostle [36] libraries. The hierarchical balancing procedure then provides its own callback functions to allow existing Zoltan procedures to be used to query and update the IHBS at each level of a hierarchical balancing. After all levels of the hierarchical balancing have been completed, Zoltan's usual migration arrays are constructed and returned to the application. Thus, only lightweight objects are migrated internally between levels, not the (larger and more costly) application data. Figure 4 shows the interaction between Zoltan and an application when hierarchical balancing is used.

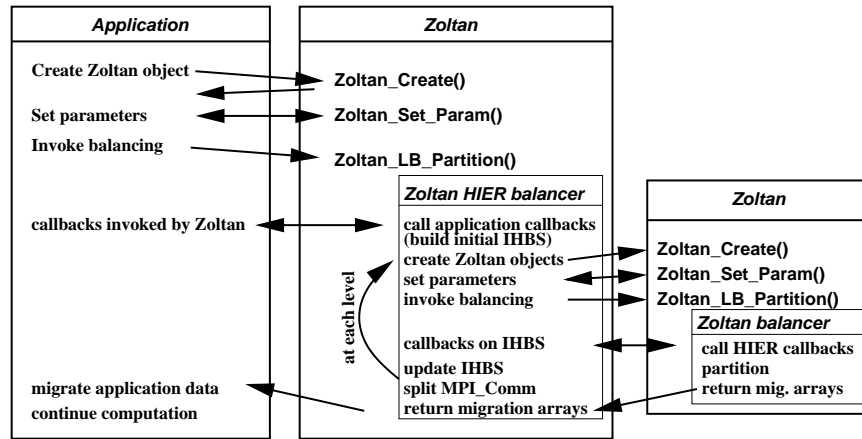


Fig. 4. Interaction between Zoltan and applications when hierarchical balancing is used.

3 Examples

We have tested our procedures using a software package called *LOCO* [16], which implements a parallel adaptive discontinuous Galerkin [3] solution of the compressible Euler equations. We consider the “perforated shock tube” problem, which models the three-dimensional unsteady compressible flow in a cylinder containing a cylindrical vent [14]. This problem was motivated by flow studies in perforated muzzle brakes for large calibre guns [12]. The initial mesh contains 69,572 tetrahedral elements. For these experiments, we stop the computation after 4 adaptive steps, when the mesh contains 254,510 elements.

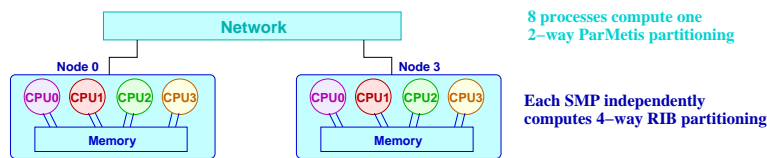


Fig. 5. Hierarchical balancing algorithm selection for two 4-way SMP nodes connected by a network.

Preliminary results are show that hierarchical partitioning can be competitive with the best results from a graph partitioner alone. We use two eight-processor computing environments: one with four Sun Enterprise 220R servers, each with two 450MHz Sparc UltraII processors, the other with two Sun Enterprise 420R servers, each with four 450MHz Sparc UltraII processors. The nodes

are dedicated to computation and do not perform file service. In both cases, inter-node communication is across fast (100 Mbit) Ethernet. A comparison of running times for the perforated shock tube in these computing environments for all combinations of traditional and hierarchical procedures shows that while ParMetis multilevel graph partitioning alone often achieves the fastest computation times, there is some benefit to using hierarchical load balancing where ParMetis is used for inter-node partitioning and inertial recursive bisection is used within each node. For example, in the four-node environment (Figure 5), the computation time following the fourth adaptive step is 571.7 seconds for the hierarchical procedure with ParMetis and RIB, compared with 574.9 seconds for ParMetis alone, 702.7 seconds for Hilbert SFC partitioning alone, 1508.2 seconds for recursive coordinate bisection alone, and 822.9 seconds for RIB alone. It is higher for other hierarchical combinations of methods.

4 Discussion

We have demonstrated the ability to use the hierarchical balancing implemented within Zoltan as both an initial partitioning and a dynamic load balancing procedure for a realistic adaptive computation. While a traditional multilevel graph partitioning is often the most effective for this application and this computing environment, the results to date demonstrate the potential benefits of hierarchical procedures. In cases where the top-level multilevel graph partitioner achieves a decomposition without introducing load imbalance, it provides the fastest time-to-solution in our studies to this point (though only slightly faster than the multilevel graph partitioner alone is used). When imbalance is introduced by the multilevel graph partitioners, using hierarchical balancing to achieve strict balance within an SMP is beneficial.

Studies are underway that utilize hierarchical balancing on larger clusters, on other architectures, and with a wider variety of applications. We expect that hierarchical balancing will be most beneficial when the extreme hierarchies found in grid environments are considered. Significant benefits will depend on an MPI implementation that can provide a very efficient intra-node communication. This is not the case in the MPICH [17] implementation used on the cluster in our experiments. Here, all communication needs to go through a network layer (MPICH's p4 device), even if two processes are executing on the same SMP node. We are eager to run experiments of clusters built from other types of SMP nodes and on computational grids, and specifically those with MPI implementations that do provide appropriate intra-node communication optimizations.

Enhancements to the hierarchical balancing procedures will focus on usability and efficiency. Further enhancements to DRUM's machine model and graphical configuration tool will facilitate and automate the selection of hierarchical procedures. Efficiency may be improved by avoiding unnecessary updates to the intermediate structure, particularly at the lowest level partitioning step. Maintaining the intermediate structure across subsequent rebalancing steps would reduce startup costs, but is complex for adaptive problems. It would also be

beneficial to avoid building redundant structures, such as when ParMetis is used at the highest level of a hierarchical balancing, however this would require some modification of individual Zoltan methods, which we have been careful to avoid thus far.

The IHBS itself has potential benefits outside of hierarchical balancing. It could be used to allow incremental enhancements and post-processing “smoothing” [15] on a decomposition before Zoltan returns its migration arrays to the application. The IHBS could also be used to compute multiple “candidate” decompositions with various algorithms and parameters, allowing Zoltan or the application to compute statistics about each and only accept and use the one deemed best.

Partitioning is only one factor that may be considered for an effective resource-aware computation. Ordering of computation and of communication, data replication to avoid communication across slow interfaces, and use of multithreading are other resource-aware enhancements that may be used. DRUM’s machine model currently includes some information that may be useful for these other types of optimizations, and it will be augmented to include information to support others.

Acknowledgments

The authors were supported in part by Sandia contract PO15162 and the Computer Science Research Institute at Sandia National Laboratories. Sandia is a multiprogram laboratory operated by Sandia Corporation, a Lockheed Martin Company, for the United States Department of Energy under Contract DE-AC04-94AL85000.

The authors would like to thank Erik Boman, Karen Devine, and Bruce Hendrickson (Sandia National Laboratories, Albuquerque, USA) for their valuable input on the design of Zoltan’s hierarchical balancing procedures.

Students Jin Chang (Rensselaer), Laura Effinger-Dean (Williams College), Luis Gervasio (Rensselaer), and Arjun Sharma (Williams College) have contributed to the development of DRUM.

We would like to thank the reviewers of this article for their helpful comments and suggestions.

References

1. S. B. Baden and S. J. Fink. A programming methodology for dual-tier multicomputers. *IEEE Transactions on Software Engineering*, 26(3):212–216, 2000.
2. M. J. Berger and S. H. Bokhari. A partitioning strategy for nonuniform problems on multiprocessors. *IEEE Trans. Computers*, 36:570–580, 1987.
3. R. Biswas, K. D. Devine, and J. E. Flaherty. Parallel, adaptive finite element methods for conservation laws. *Appl. Numer. Math.*, 14:255–283, 1994.
4. E. Boman, K. Devine, R. Heaphy, B. Hendrickson, M. Heroux, and R. Preis. LDRD report: Parallel repartitioning for optimal solver performance. Technical Report SAND2004-0365, Sandia National Laboratories, Albuquerque, NM, February 2004.

5. C. L. Bottasso, J. E. Flaherty, C. Özturan, M. S. Shephard, B. K. Szymanski, J. D. Teresco, and L. H. Ziantz. The quality of partitions produced by an iterative load balancer. In B. K. Szymanski and B. Sinharoy, editors, *Proc. Third Workshop on Languages, Compilers, and Runtime Systems*, pages 265–277, Troy, 1996.
6. T. Bui and C. Jones. A heuristic for reducing fill in sparse matrix factorization”. In *Proc. 6th SIAM Conf. Parallel Processing for Scientific Computing*, pages 445–452. SIAM, 1993.
7. P. M. Campbell, K. D. Devine, J. E. Flaherty, L. G. Gervasio, and J. D. Teresco. Dynamic octree load balancing using space-filling curves. Technical Report CS-03-01, Williams College Department of Computer Science, 2003.
8. G. Cybenko. Dynamic load balancing for distributed memory multiprocessors. *J. Parallel Distrib. Comput.*, 7:279–301, 1989.
9. K. Devine, E. Boman, R. Heaphy, B. Hendrickson, and C. Vaughan. Zoltan data management services for parallel dynamic applications. *Computing in Science and Engineering*, 4(2):90–97, 2002.
10. K. D. Devine, E. G. Boman, R. T. Heaphy, B. A. Hendrickson, J. D. Teresco, J. Faik, J. E. Flaherty, and L. G. Gervasio. New challenges in dynamic load balancing. *Appl. Numer. Math.*, 52(2–3):133–152, 2005.
11. K. D. Devine, B. A. Hendrickson, E. Boman, M. St. John, and C. Vaughan. *Zoltan: A Dynamic Load Balancing Library for Parallel Applications; User’s Guide*. Sandia National Laboratories, Albuquerque, NM, 1999. Tech. Report SAND99-1377. Open-source software distributed at <http://www.cs.sandia.gov/Zoltan>.
12. R. E. Dillon Jr. A parametric study of perforated muzzle brakes. ARDC Tech. Report ARLCB-TR-84015, Benét Weapons Laboratory, Watervliet, 1984.
13. J. Faik, J. D. Teresco, K. D. Devine, J. E. Flaherty, and L. G. Gervasio. A model for resource-aware load balancing on heterogeneous clusters. Technical Report CS-05-01, Williams College Department of Computer Science, 2005. Submitted to *Transactions on Parallel and Distributed Systems*.
14. J. E. Flaherty, R. M. Loy, M. S. Shephard, M. L. Simone, B. K. Szymanski, J. D. Teresco, and L. H. Ziantz. Distributed octree data structures and local refinement method for the parallel solution of three-dimensional conservation laws. In M. Bern, J. Flaherty, and M. Luskin, editors, *Grid Generation and Adaptive Algorithms*, volume 113 of *The IMA Volumes in Mathematics and its Applications*, pages 113–134, Minneapolis, 1999. Institute for Mathematics and its Applications, Springer.
15. J. E. Flaherty, R. M. Loy, M. S. Shephard, B. K. Szymanski, J. D. Teresco, and L. H. Ziantz. Adaptive local refinement with octree load-balancing for the parallel solution of three-dimensional conservation laws. *J. Parallel Distrib. Comput.*, 47:139–152, 1997.
16. J. E. Flaherty, R. M. Loy, M. S. Shephard, and J. D. Teresco. Software for the parallel adaptive solution of conservation laws by discontinuous Galerkin methods. In B. Cockburn, G. Karniadakis, and S.-W. Shu, editors, *Discontinuous Galerkin Methods Theory, Computation and Applications*, volume 11 of *Lecture Notes in Computational Science and Engineering*, pages 113–124, Berlin, 2000. Springer.
17. W. Gropp, E. Lusk, N. Doss, and A. Skjellum. A high-performance, portable implementation of the MPI message passing interface standard. *Parallel Computing*, 22(6):789–828, Sept. 1996.
18. B. Hendrickson and R. Leland. A multilevel algorithm for partitioning graphs. In *Proc. Supercomputing ’95*, 1995.
19. Y. F. Hu and R. J. Blake. An optimal dynamic load balancing algorithm. Preprint DL-P-95-011, Daresbury Laboratory, Warrington, WA4 4AD, UK, 1995.

20. G. Karypis and V. Kumar. A fast and high quality multilevel scheme for partitioning irregular graphs. *SIAM J. Scien. Comput.*, 20(1), 1999.
21. G. Karypis and V. Kumar. Parallel multilevel k -way partitioning scheme for irregular graphs. *SIAM Review*, 41(2):278–300, 1999.
22. E. Leiss and H. Reddy. Distributed load balancing: design and performance analysis. *W. M. Kuck Research Computation Laboratory*, 5:205–270, 1989.
23. W. F. Mitchell. Refinement tree based partitioning for adaptive grids. In *Proc. Seventh SIAM Conf. on Parallel Processing for Scientific Computing*, pages 587–592. SIAM, 1995.
24. A. Patra and J. T. Oden. Problem decomposition for adaptive hp finite element methods. *Comp. Sys. Engng.*, 6(2):97–109, 1995.
25. J. R. Pilkington and S. B. Baden. Dynamic partitioning of non-uniform structured workloads with spacefilling curves. *IEEE Trans. on Parallel and Distributed Systems*, 7(3):288–300, 1996.
26. A. Pothen, H. Simon, and K.-P. Liou. Partitioning sparse matrices with eigenvectors of graphs. *SIAM J. Mat. Anal. Appl.*, 11(3):430–452, 1990.
27. R. Rabenseifner and G. Wellein. Comparison of parallel programming models on clusters of SMP nodes. In H. Bock, E. Kostina, H. Phu, and R. Rannacher, editors, *Proc. Intl. Conf. on High Performance Scientific Computing*, pages 409–426, Hanoi, 2004. Springer.
28. M. S. Shephard, J. E. Flaherty, H. L. de Cougny, C. Özturan, C. L. Bottasso, and M. W. Beall. Parallel automated adaptive procedures for unstructured meshes. In *Parallel Comput. in CFD*, number R-807, pages 6.1–6.49. Agard, Neuilly-Sur-Seine, 1995.
29. H. D. Simon. Partitioning of unstructured problems for parallel processing. *Comp. Sys. Engng.*, 2:135–148, 1991.
30. C. A. Taylor, T. J. R. Hugues, and C. K. Zarins. Finite element modeling of blood flow in arteries. *Comput. Methods Appl. Mech. Engrg.*, 158(1–2):155–196, 1998.
31. V. E. Taylor and B. Nour-Omid. A study of the factorization fill-in for a parallel implementation of the finite element method. *Int. J. Numer. Meth. Engrg.*, 37:3809–3823, 1994.
32. J. D. Teresco, M. W. Beall, J. E. Flaherty, and M. S. Shephard. A hierarchical partition model for adaptive finite element computation. *Comput. Methods Appl. Mech. Engrg.*, 184:269–285, 2000.
33. J. D. Teresco, K. D. Devine, and J. E. Flaherty. *Numerical Solution of Partial Differential Equations on Parallel Computers*, chapter Partitioning and Dynamic Load Balancing for the Numerical Solution of Partial Differential Equations. Springer-Verlag, 2005.
34. J. D. Teresco, J. Faik, and J. E. Flaherty. Resource-aware scientific computation on a heterogeneous cluster. *Computing in Science & Engineering*, 7(2):40–50, 2005.
35. J. D. Teresco and L. P. Ungar. A comparison of Zoltan dynamic load balancers for adaptive computation. Technical Report CS-03-02, Williams College Department of Computer Science, 2003. Presented at COMPLAS '03.
36. C. Walshaw and M. Cross. Parallel Optimisation Algorithms for Multilevel Mesh Partitioning. *Parallel Comput.*, 26(12):1635–1660, 2000.
37. M. S. Warren and J. K. Salmon. A parallel hashed oct-tree n -body algorithm. In *Proc. Supercomputing '93*, pages 12–21. IEEE Computer Society, 1993.
38. R. Williams. Performance of dynamic load balancing algorithms for unstructured mesh calculations. *Concurrency*, 3:457–481, October 1991.