



A hierarchical partition model for adaptive finite element computation

J.D. Teresco^{*}, M.W. Beall, J.E. Flaherty, M.S. Shephard

Scientific Computation Research Center (SCOREC), Department of Computer Science, Rensselaer Polytechnic Institute, Troy, NY 12180, USA

Received 28 August 1998

Abstract

Software tools for the solution of partial differential equations using parallel adaptive finite element methods have been developed. We describe the design and implementation of the parallel mesh structures within an adaptive framework. The most fundamental concept is that of a hierarchical partition model used to distribute finite element meshes and associated data on a parallel computer. The hierarchical model represents heterogeneous processor and network speeds, and may be used to represent processes in any parallel computing environment, including an SMP, a distributed-memory computer, a network of workstations, or some combination of these. Using this model to segment the computation into chunks which can fit into cache memory provides a potential efficiency gain from an increased cache hit rate, even in a single processor environment. The information about different processor speeds, memory sizes, and the corresponding interconnection network can be useful in a dynamic load balancing algorithm which seeks to achieve a good balance with minimal interprocessor communication penalties when a slow interconnection network is involved. © 2000 Elsevier Science S.A. All rights reserved.

Keywords: Parallel computation; Hierarchical and distributed systems

1. Introduction

The finite element method (FEM) has become a standard analysis tool for solving partial differential equations (PDEs). Computationally demanding three-dimensional problems make adaptive methods and parallel computation essential. Adaptive FEMs provide reliability, robustness, and time and space efficiency. In such a method, the computational domain is discretized into a mesh. During the adaptive solution process, portions of the mesh may be refined or coarsened (*h*-refinement) or moved to follow evolving phenomena (*r*-refinement). The method order may also be varied (*p*-refinement). Each adaptive process concentrates the computational effort in areas where the solution resolution would otherwise be inadequate [7]. Conventional array-based data representations, which work well for fixed-mesh solutions, are not well suited to solutions involving mesh adaptivity [1]. Traversal of the data must be efficient in all cases, but when adaptivity is introduced, modification of the mesh structures and corresponding solution data must also be efficient.

Parallel computation introduces complications such as the need to balance processor loading, coordinate interprocessor communication, and manage distributed data. The standard methodology for optimizing parallel FEM programs relies on a static partitioning of the mesh across the cooperating processors. This is insufficient in an adaptive computation, since load imbalance will be introduced by adaptive enrichment,

^{*} Corresponding author.

E-mail address: ferescoj@cs.rpi.edu (J.D. Teresco).

necessitating dynamic partitioning and redistribution of data. Data structures must support this dynamic mesh migration.

Reusable software libraries are essential for the development of specific applications to solve diverse problems without concern for the details of the underlying mesh structures or parallelization. A variety of such tools have been developed to facilitate the design, implementation, and use of parallel adaptive finite element software. Static data structures cannot be used with adaptive computations, which makes automatic detection by a parallelizing compiler difficult. Parallelism is generally explicit, achieved through the use of a message passing library such as the message passing interface (MPI) [18], and requires a partitioning algorithm to distribute data among processing nodes. Parallel adaptive FEM computations can be distributed in a natural way by a domain decomposition of the mesh underlying the computation, but being adaptive, these meshes will change and the system must account for this.

The approaches can be categorized by a low-level view of the data structures. Two basic paradigms have been used as the underlying structure for these computations. An array-based approach, such as the scalable distributed dynamic array [11], uses the distribution of arrays as the fundamental unit of parallelism. The mesh-based approach [1,4,14,21,24–27] uses mesh connectivity and distribution of mesh entities to achieve parallelism.

Our recent efforts have focused on implementing a framework for the reliable automated solution of problems in science and engineering over arbitrary domains using scalable parallel adaptive techniques [2]. This framework provides common interfaces to various pieces of software typically needed for scientific computation, and in particular when writing finite element analysis procedures. Central to this framework is a distributed mesh structure, described in Section 3.

While parallel adaptive methods have been used successfully in the solutions of many problems [3,5,8,9,11,14,20], it is difficult to take full advantage of a particular parallel computing environment. A partitioning and dynamic load balancing strategy which is efficient in a shared-memory environment may not be in a distributed-memory environment. This is especially the case with heterogeneous or hierarchical computers which combine multiprocessing shared memory nodes connected by networks of various speeds. A partitioner must have some knowledge of the parallel environment to achieve a good distribution. In Section 3, we describe a hierarchical partition model which makes such information available.

2. Target parallel environments

The large range of problem sizes to be addressed and user computing environments available requires software that runs efficiently on computers ranging from single workstations to the largest parallel supercomputers. In order for a parallel program to be efficient on a specific computer, it must obtain and use knowledge of the parallel computational environment on which it is running. Information that is useful in this respect includes available memory, available processing power, and availability and properties (bandwidth and latency) of communication resources. MPI [18] does not provide such information (nor does the MPI-2 specification [19]). Some of this information can be computed through environment queries, but other information will need to be provided manually.

Consider five specific parallel computing environments, represented in Fig. 1. The first (Fig. 1(a)) is an eight-node computer with all uniprocessor nodes connected by a fast, nonblocking switch. Next (Fig. 1(b)), is an eight-processor symmetric multiprocessing (SMP) workstation. The third (Fig. 1(c)), is a collection of eight workstations connected by ethernet, in a typical network of workstations (NOW) configuration. Fig. 1(d) shows a hybrid system, a six-node computer which has four single-processor nodes and two four-way SMP nodes, all connected by a high-speed switch. And finally, Fig. 1(e) shows a more complex networking scheme, with two clusters of asynchronous transfer mode-connected (ATM) four-way SMP workstations, where the two clusters are connected by ethernet.

Consider a parallel program that has been developed and optimized for the environment of Fig. 1(a), which is characterized by a distributed-memory environment with communication provided by a high-speed, nonblocking switch. The computational domain is distributed among the processors by a partitioning algorithm such that the workload is approximately equidistributed and the interfaces between the

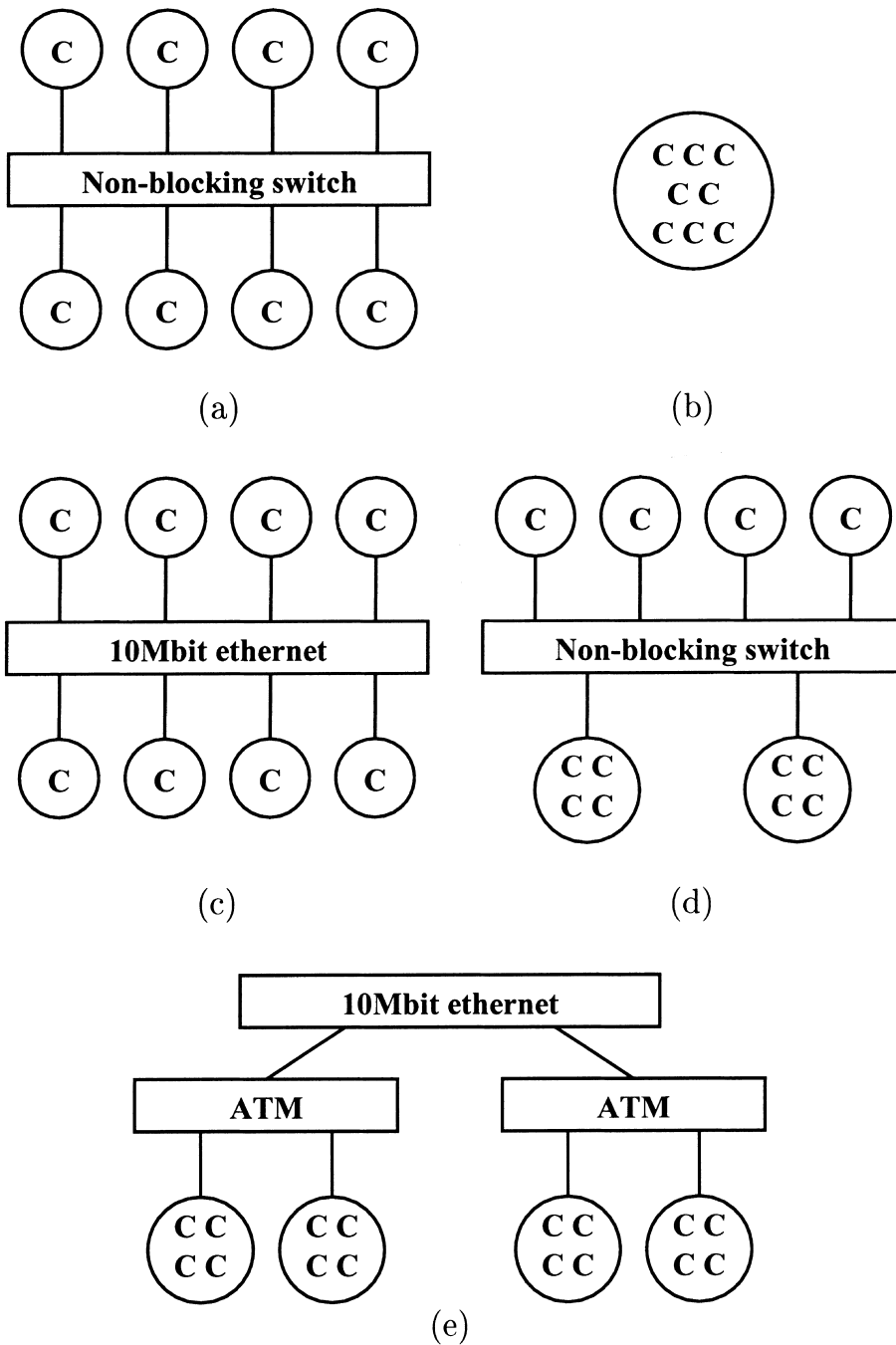


Fig. 1. Examples of parallel computing environments. A C represents a CPU. Multiple Cs in a process node represent an SMP node. (a) An 8-node computer with all uniprocessor nodes, connected by a fast switch. (b) An 8-way SMP workstation. (c) Eight uniprocessor workstations connected by ethernet. (d) An 8-node system with 6-uniprocessor nodes and 2 4-processor SMP nodes connected by a fast switch. (e) Two clusters of ATM connected 4-processor SMP nodes connected by ethernet.

subdomains are small. Perhaps it has been determined that a small load imbalance is tolerable if it will reduce the size of the interfaces; hence, the amount of interprocessor communication across the switch.

Suppose that the program is subsequently targeted for execution on an SMP workstation (Fig. 1(b)), where the processes communicate via shared memory, which is fast relative to a nonblocking switch. The assumption allowing a small computational imbalance to reduce interprocessor communication may no

longer be effective. Since communication is less expensive, the loss of CPU utilization due to the imbalance may significantly degrade performance. A network of eight uniprocessor workstations connected by ethernet (Fig. 1(c)) is at the opposite extreme. In this case, interprocessor communication is very expensive. A more significant computational imbalance may be acceptable if some communication can be avoided. Since the same program can use both environments, the partitioner needs run-time information to determine the existing situation and to decide the optimal relationship between a balanced computational load and minimal communication.

The situation becomes more complex with heterogeneous environments. In the case of ATM clusters (Fig. 1(e)), there are three levels of interconnection network which may be used for one process to communicate with another. Processes on the same SMP node communicate via shared memory. Processes on different SMP nodes communicate via ATM if they are in the same ATM cluster and via ethernet if they are in different clusters. Communication over the shared-memory interface is inexpensive, so there is no need for great expense in partitioning data stored on the same SMP node. Communication over the ATM interface is slower, but still fast enough to use an inexpensive partitioning and to maintain a close balance. When slower communication over ethernet is involved, it is worthwhile to minimize communication over the slow network interface, including use of a more expensive partitioning method and tolerating a larger load imbalance. Heterogeneous processors are also important to consider. Processors within a computation may have different speeds and different amounts of memory available. In cases where computing resources are shared, the full processing power of a given computer may not be available, so a measure of available processing power, depending on CPU demands of other processes, must be maintained.

It is unreasonable to maintain distinct parallel and serial versions of all applications. However, it is also important to eliminate the overhead associated with parallel computation when the program is run on a single processor. Some parallel overhead may be removed at compile time while other overhead is trivial to remove at run time.

3. Parallel mesh infrastructure

We are developing the *Rensselaer partition model* (RPM), which allows a hierarchical partitioning of finite element mesh data and provides additional information about the parallel computational environment in which a program is executing. The model works with the mesh data structures described here; however, many of the ideas may be applied to other systems.

3.1. Basic mesh structures

The *SCOREC mesh database* (MDB) [1] provides an object-oriented, hierarchical representation of a finite element mesh. It includes a set of operators to query and update the mesh data structure. The basic mesh entity hierarchy consists of three-dimensional *regions*, and their bounding *faces*, *edges*, and *vertices*, with bidirectional links between mesh entities of consecutive order (Fig. 2). In three-dimensional meshes, regions are used as finite elements while faces serve as elements in two-dimensions, or as interface elements in three-dimensions. Mesh entities have an explicit *geometric classification* relative to a geometric model of the problem domain to allow for the appropriate representation of the geometry as the mesh is enriched either through *h*- or *p*-refinement. Mesh classification against the geometric domain can be defined more precisely as the unique association of a topological mesh entity of dimension d_i , to a topological entity of the geometric domain of dimension d_j , where $d_i \leq d_j$. Multiple mesh entities can be classified on a geometric entity. A mesh region is classified on the domain region in which it lies. A mesh face is either classified on the domain region or face in which it lies. A mesh edge is classified on the domain region, face, or edge in which it lies. Finally, a mesh vertex is classified on the domain region, face, edge, or vertex in which it lies. Mesh entities are always classified with respect to the lowest order model entity possible. More formal definitions of classification and adjacency can be found in [1].

The full entity hierarchy allows the efficient deletion and creation of mesh entities during *h*-refinement [24] and simplifies attachment of degrees of freedom to the mesh entities and determination of necessary geometric information during *p*-refinement [23]. The database allows for the fast retrieval of adjacency

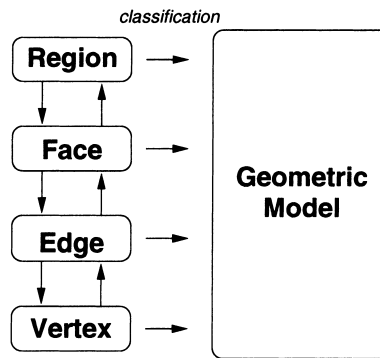


Fig. 2. MDB entity hierarchy with links to a geometric model.

information. Examples of available data include the list of faces bounding an element and the edges sharing a common vertex. All entities can have attached attributes such as solution and error indicator data.

3.2. Mesh data organization

Mesh entities are listed with the geometric model entities on which they are classified (Fig. 3). When mesh entities are added to or removed from a mesh, or when mesh entities are traversed, the appropriate lists of regions, faces, edges, and vertices stored on the geometric model entities are used. The appropriate geometric entity is readily available from the mesh entity, so no additional searching is needed. Traversal of the entire mesh requires some additional bookkeeping, since lists of geometric entities and the mesh entity lists within them must be traversed. This additional overhead is not significant for performance and the complications are hidden within the iterator implementation. The main advantage of this design is that the inverse classification information (which mesh entities are classified on a given model entity) is readily available. This is useful, for example, when applying a boundary condition on a model face. Rather than traversing all faces in the mesh and querying each to check if it is on the desired boundary, a list of the needed entities is readily available.

3.3. Partition model

The mesh data in a parallel computing environment is distributed. Each entity in the finite element mesh is uniquely assigned to a *partition model entity* or, more simply, to a *partition*. The direct assignment of one

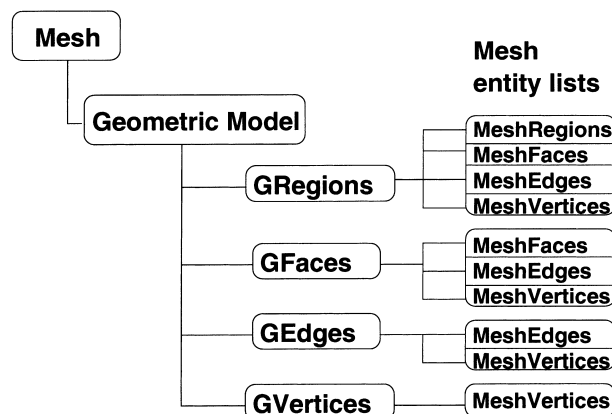


Fig. 3. Data organization with mesh data stored relative to a geometric model. GRegions, GFaces, GEdges, and GVertices refer to lists of geometric model regions, faces, edges, and vertices, respectively.

partition to each process has been a limitation of prior systems [14]. While this is straightforward, there are advantages in allowing multiple partitions on a single process. Our approach to mesh migration (Section 3.9) requires this. Other motivations include the use of the partitions as a coarser structure for a load balancer, as a method to encourage locality of reference to enhance cache performance, and as a means of performing out-of-core computations by swapping partitions in and out of memory. Each of these will be discussed further in Section 6. Each partition is assigned to a specific *process*. Multiple partition model entities may be assigned to a single process. “Process” in this context refers to an address space. The model is hierarchical, with partition model entities assigned to a *process model* and process model entities assigned to a *machine model*.

The machine model represents the physical processing nodes involved in the computation, and information about the communication methods to other involved computers. The process model maps processes to the computer, and maps interprocess communication to inter-computer networks or, perhaps, to a shared-memory interface. The mesh is represented at the partition model level. Partitions must know the mesh entities that they contain, and a mesh entity must know its partition assignment, which we refer to as its *partition model classification*.

3.4. Augmented model

The simplest implementation of the partition model introduces pointers on each mesh entity to store: (i) partition classification and (ii) geometric classification. For space efficiency, these may be combined into a single topological unit called the *augmented model*. Additional topological entities are added to the model which correspond to partitions and partition boundaries, forming a topological structure which contains geometric and partition information. Consider the simple geometry and corresponding mesh of Fig. 4. The geometric model contains 2 faces, 8 edges, and 7 vertices. Mesh entities are classified onto these geometric entities as shown on the right.

In Fig. 5, we divide the mesh into 5 partitions (shown as solid lines). The partitioning is done independently of the geometric model, but the partitioning induces a similar topological structure (Fig. 5, right). When this information is combined with the geometry of Fig. 4, left, the augmented model (Fig. 6) is obtained. The augmented topology has 7 faces, 19 edges, and 15 vertices.

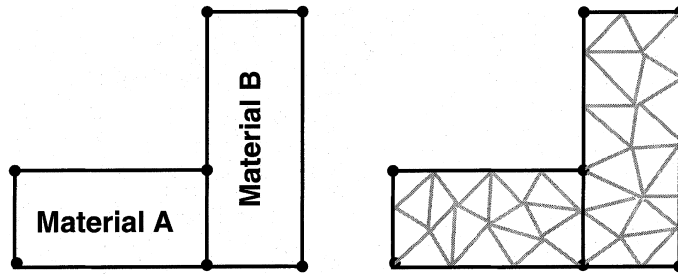


Fig. 4. Example domain geometry (left) and mesh classified on that domain (right). The distinct geometric faces could represent two different materials, for example.

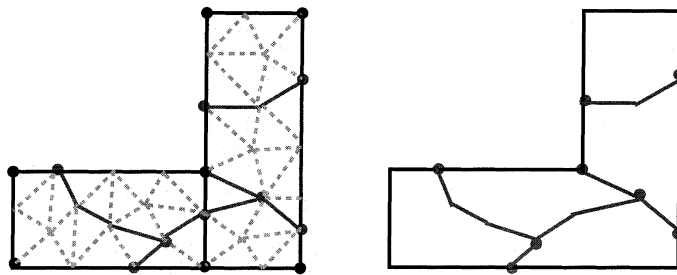


Fig. 5. Mesh partitioning into 5 parts (left) and the associated partition model (right).

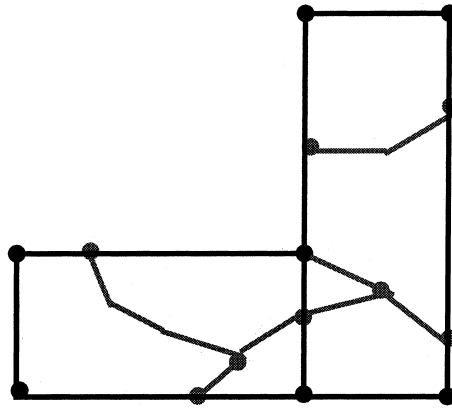


Fig. 6. Topological structure resulting from augmented model combining the geometric and partition model.

The size of the augmented structure will grow with the number of partitions. Such a structure would need to be distributed. The implementation of the augmented model becomes complex, because the CAD packages underlying the geometric model do not support distribution. However, the abstraction is still useful, and the same functionality is achieved by the data organization described in Section 3.5. This can be done because the topological model represents an independent abstraction of the domain, in which the definition of the basic entities can be based on any desired criteria as long as the rules on entity adjacencies are maintained.

3.5. Multiple mesh pointers

Partitions are represented by assigning separate “mesh” pointers to each, as shown in Fig. 7. This is equivalent to the augmented model in functionality, but does not require the overhead of the more complex

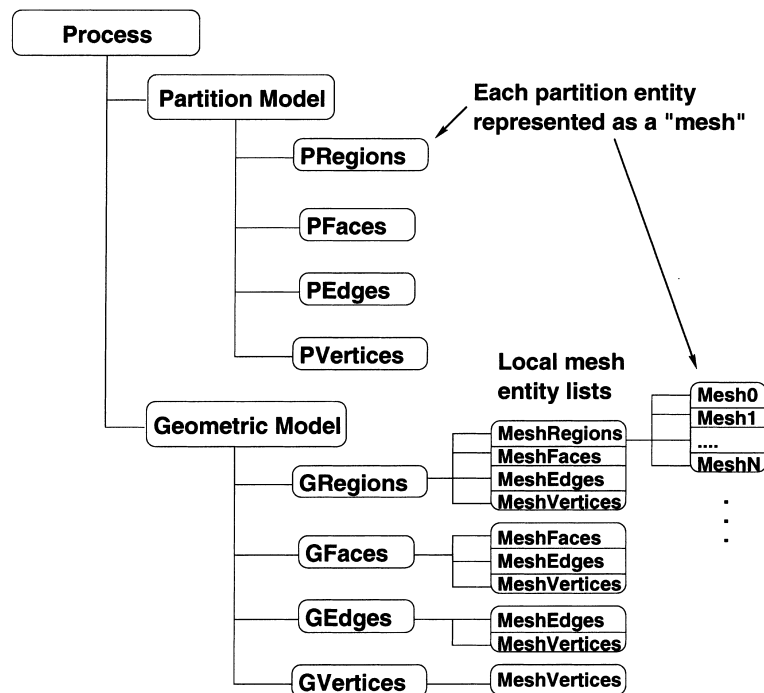


Fig. 7. Data organization with both geometric and partition models.

distributed augmented model structure. In this case, the partition model is only a topological construct. It contains no geometric information and is distinct from the geometric model. The mesh entity lists within each geometric model entity are extended to contain one list for each partition within the local process, but avoid the second per-mesh entity pointer for partition classification. There is a small bit field on each entity which is used to retrieve the mesh pointer to which a mesh entity belongs, and this is used to determine partition classification, since each partition is stored as a separate mesh.

3.6. Process and machine models

Partition model entities are assigned to a process model (Fig. 8). Each process, with a unique address space in a parallel computation, is assigned a process model entity. The process model represents the “live” state of the job and is concerned with its computational and communication performance. Each process tracks the work assigned to it and, at any given time, the amount of this work that has been completed. The actual communication performance, between this process and other processes to which it communicates, is tracked within the *process interconnect* structures associated with the active connection between process objects.

In contrast, the machine model represents the theoretical performance that the hardware environment is capable of delivering while the process model represents the actual performance that is being achieved. Process model entities are assigned to a machine model (Fig. 9).

The machine model describes the computational environment in which the process is running. This information includes available memory, available processing power, available communication resources and the properties of these resources. In addition, the machine model provides capabilities to manage and modify the resources used by the computation. The model is a graph that describes the various hierarchies in the computational environment. A machine model entity represents a processing node, which may have multiple processing units. Multiple processes can be assigned to a single computer when either a computer includes SMP nodes or multiple processes are run on a single CPU. The latter is often desirable when debugging, such as when debugging a test case with a large number of processes on a smaller number of nodes. Multiple processes on a computer are not specifically managed by RPM. They are under the control of the operating system scheduler. *MachineInterconnect* entities can be used to model different network speeds.

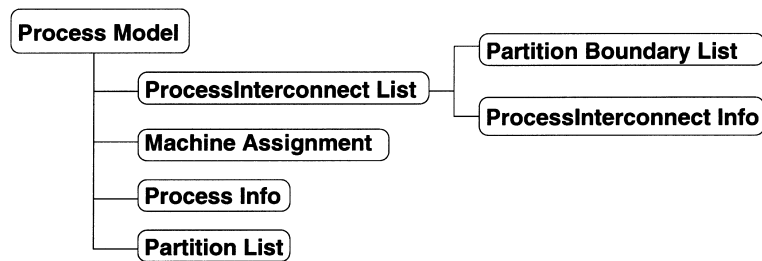


Fig. 8. Process model data organization.

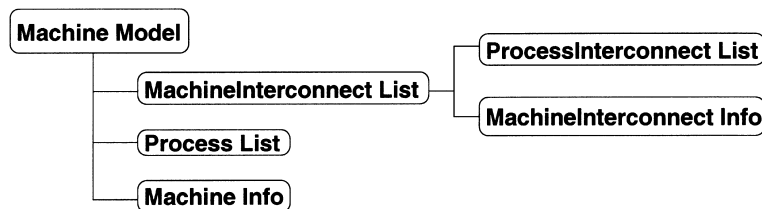


Fig. 9. Machine model data organization.

In Fig. 10, we see a sample two-dimensional mesh (left) and its classification onto the partition model (right). Mesh entities are classified on a partition model entity, which, in this case, is a partition face or boundary. Partition model entities are assigned to process model entities (Fig. 11, left), allowing multiple partitions to be assigned to a single process. And the process model entities are associated with a machine model (Fig. 11, right), allowing multiple processes to be assigned to a single machine “entity”.

3.7. Boundary structures

Mesh entities need only be replicated when on a partition boundary that is also a process boundary. Fig. 12 shows two-dimensional examples of mesh faces which share a common edge across a partition boundary. The shared mesh edge is classified on the partition edge in each case. On the left, the partition edge is local to the process, so the mesh entity need not be replicated and is stored only on the “mesh” associated with the partition edge. On the right, the partition edge is also on the process boundary. In this case, the partition boundary is replicated in each process. The mesh edge classified on the partition boundary must also be replicated. It is still classified on the partition edge.

Consider a mesh adjacency operator that retrieves the list of faces which bound a given region. One or more of the faces returned could be on a partition boundary and may be replicated. The operator must return a pointer to a local entity. There is a unique instance of the duplicated face which is the actual entity (the one located on the *owner* of the entity). The copies returned by nonowners contain appropriate

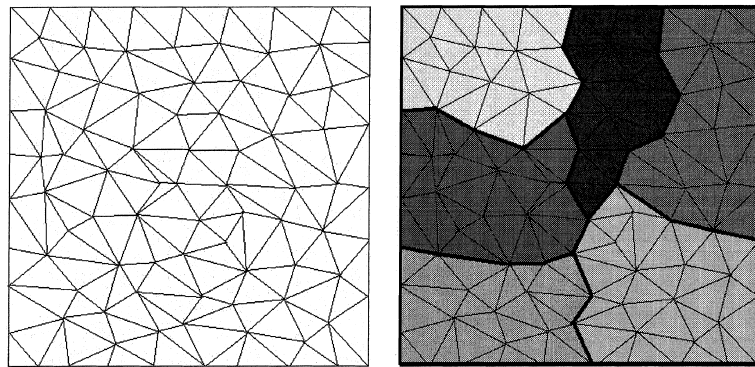


Fig. 10. A sample two-dimensional mesh (left) and its partitioning (right).

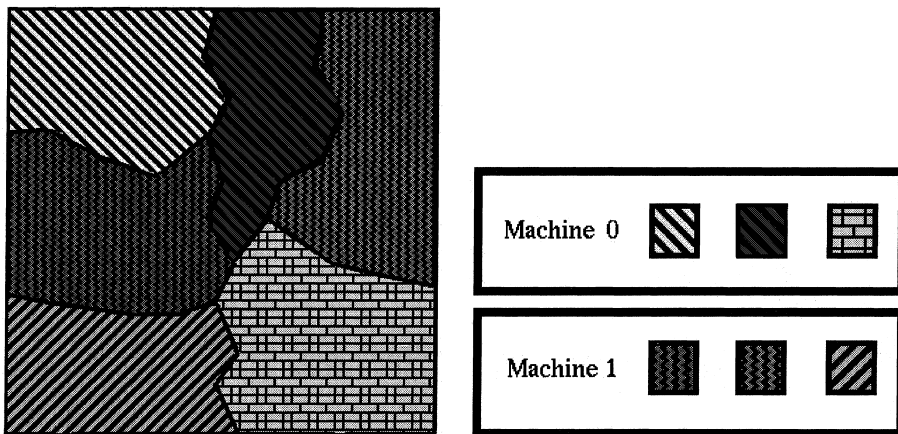


Fig. 11. Assignment of partitions to a process model (left), and partition and process assignments on a machine model (right). Patterns indicate assignments on the process model.

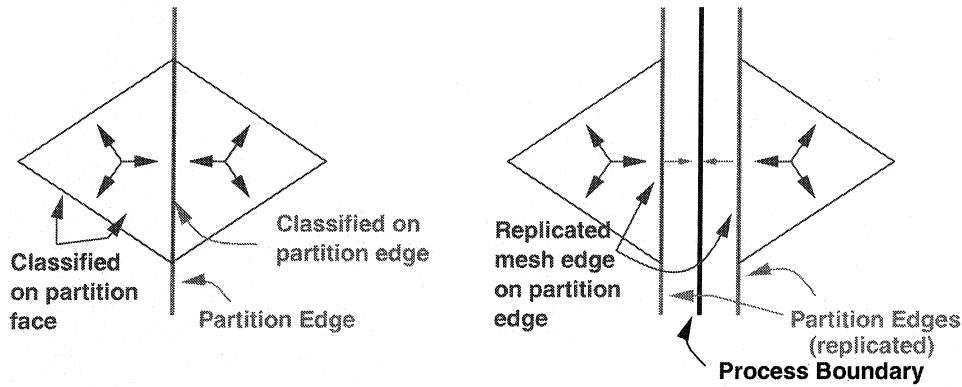


Fig. 12. Partition classification of mesh entities for the case of same-process partition (left) and at the partition boundary which is also a process boundary (right).

information for the process to obtain the “real” mesh entity, which in this case would be located within the address space of another process.

3.8. Mesh entity counts and traversals

A parallel adaptive FEM code must query and traverse a mesh in a number of different ways. A solution procedure will need to traverse all entities where solution data can be stored. This may be restricted to entities associated with a particular portion of the geometry. A linear algebra solver will need to traverse the interiors of partitions to solve local problems and partition boundaries to solve interface problems. A partitioning algorithm will need to query the number of mesh entities in each partition and on partition boundaries.

Given these possibilities, the question “how many mesh faces” can have different meanings. Within RPM, the answer that is returned depends on “who” you ask. The global mesh returns the global count. A partition model entity (a partition or a partition boundary) provides the number classified on that partition model entity. A geometric entity returns the number associated with that geometry. Traversal functions hide the partitioning when possible, and, like entity-count operators, separate operators are provided for each type of traversal.

An important feature is the ability to traverse mesh entities by geometric entity or by partition. This is much more efficient than traversing all entities and checking the classification (geometric or partition) of each to determine whether it is one we wish to use.

Table 1

Parallel environments involved in network performance comparisons. Variations in architecture, interprocess communication vehicles, numbers of processes and computers are examined

Architecture	IPC vehicle	Procs.	Computers	Key
IBM SP2	switch	2–32	2–32	SP2s
IBM SP2	Ethernet	2–32	2–32	SP2e
Sun Ultra 2/2200	shared memory	1–2	1	SUNsh
Sun Ultra 2/2200	local p4	1–2	1	SUNlp4
Sun Ultra 2/2200	Ethernet p4	1–6	1–3	SUNp4(10)
Sun Ultra 2/2200	fast Ethernet p4	1–6	1–3	SUNp4(100)
SGI Onyx II	shared memory	1–8	1	SGIsh
SGI Onyx II	local p4	1–8	1	SGIlp4

3.9. Inter-process migration

The approach to mesh migration is similar to the one described by Flaherty et al. [14], but is implemented at a more abstract level. To migrate a subset of the mesh, the subset is first reclassified on a new partition model entity. This reclassification is very inexpensive: it involves two pointer updates and an entity count adjustment. Most of the partition reclassifications will not involve a change to the topology of the partition model. More extensive operations are only necessary in those cases when the reclassification produces new partition adjacencies. Following reclassification of mesh entities on the new partition, the partition is migrated to the destination process (abstractly, a “reclassification” of the partition on the process model) and migration of the associated mesh follows. Some entities on interpartition boundaries become the interprocessor boundary. Since the new partition boundary structures already exist, the interprocess boundary update is straightforward. This abstraction of migration allows appropriate packing to avoid the inefficiency of using frequent, small messages. This strategy requires the ability to assign multiple partition model entities to a process and dynamic creation and deletion of partition model entities.

4. Network performance comparisons

Specific operations were run in a variety of parallel environments to appraise differences in communication devices at both the hardware and software levels. All cases used MPI for communication. IBM’s implementation of MPI was used on their SP2, while others used the MPICH package [17]. MPICH is freely available and runs on architectures ranging from networks of workstations, where it uses p4 as an underlying communication library, to SMPs where it takes advantage of the availability of shared memory as a communication device, and on high-end parallel machines, where it uses the vendor’s optimized parallel communication libraries.

The test program simulates work done at the mesh level during a parallel adaptive computation. The specific test sequence was (i) load the mesh onto one processor, (ii) partition the mesh onto n processors, (iii) refine a portion of the mesh, and (iv) repartition the mesh to rebalance. Computers involved in the test configurations (Table 1) were a 36-node IBM SP2, with communications over the high-speed switch and over an ethernet network that connects all nodes; an eight-processor SGI Onyx2 with communication via shared memory; and three dual-processor Sun Ultra 2/2200 workstations, using shared memory and both slow (10 Mb/s) and fast (100 Mb/s) ethernet. The initial meshes contained 44 177 elements (coarse case), 79 944 elements (medium case), and 169 733 elements (fine case). In each case, local refinement was applied which increased mesh sizes to 87 170, 158 046, and 453 900 for the coarse, medium, and fine meshes. Tables 2–7 show wall-clock execution times for portions of the code. The “load” time is the time for the initial mesh to be loaded on a single processor. This is a totally serial process and is used as a consistency

Table 2
Wall-clock times for network tests using coarse meshes on the IBM SP2

System	Load	Partition		Refine	Rebalance	
		Total	Migration		Total	Migration
SP2s-2	35.65	14.27	8.37	32.87	15.78	8.79
SP2e-2	36.07	22.83	13.97	33.00	25.44	14.40
SP2s-4	35.69	19.56	12.27	20.60	53.37	4.81
SP2e-4	35.89	44.92	27.79	21.59	72.38	13.87
SP2s-8	36.02	20.99	14.34	11.54	23.01	4.54
SP2e-8	38.00	61.29	33.34	16.32	62.55	20.23
SP2s-16	36.10	21.50	15.39	11.39	12.48	4.58
SP2e-16	36.47	90.42	37.99	59.27	87.09	36.39
SP2s-32	37.21	25.12	16.79	8.99	8.32	3.59
SP2e-32	48.08	131.18	39.22	93.62	269.55	181.43

Table 3

Wall-clock times for network tests using coarse meshes on Sun Ultra 2/2200 workstations

System	Load	Partition		Refine	Rebalance	
		Total	Migration		Total	Migration
SUNsh-2	9.27	7.13	4.32	15.79	8.04	4.25
SUNlp4-2	7.48	7.03	4.04	15.81	8.14	4.41
SUNp4(10)-2	7.64	15.62	9.59	20.16	18.81	10.81
SUNp4(10)-4	7.66	19.60	12.34	11.64	21.45	8.46
SUNp4(10)-6	7.77	25.04	15.04	7.54	25.94	10.99
SUNp4(100)-2	8.14	8.38	4.84	16.10	9.58	5.18
SUNp4(100)-4	7.82	10.55	6.76	10.19	11.30	3.23
SUNp4(100)-6	7.67	11.21	7.54	5.98	7.91	2.71

Table 4

Wall-clock times for network tests using coarse meshes on the SGI Onyx

System	Load	Partition		Refine	Rebalance	
		Total	Migration		Total	Migration
SGIsh-2	6.59	7.94	5.01	17.65	9.18	5.37
SGIlp4-2	6.59	8.44	5.31	18.02	9.77	5.62
SGIsh-4	6.69	10.63	7.65	11.35	8.21	3.35
SGIlp4-4	6.64	11.28	7.85	11.59	8.76	3.47
SGIsh-6	6.62	10.80	8.18	6.21	5.44	2.09
SGIlp4-6	6.65	12.12	8.62	6.63	6.78	2.66
SGIsh-8	6.71	11.46	8.77	6.17	5.54	2.84
SGIlp4-8	6.59	13.24	9.05	6.99	7.04	3.42

Table 5

Wall-clock times for network tests using medium meshes on the IBM SP2

System	Load	Partition		Refine	Rebalance	
		Total	Migration		Total	Migration
SP2s-2	64.83	26.07	16.02	58.44	35.41	16.76
SP2e-2	64.83	42.78	26.75	59.11	51.39	26.54
SP2s-4	64.64	41.35	22.58	32.34	162.70	10.02
SP2e-4	65.14	68.11	39.85	32.96	203.01	28.33
SP2s-8	65.03	38.61	25.95	26.66	67.44	11.58
SP2e-8	65.73	110.40	55.05	38.82	179.34	43.08
SP2s-16	65.19	38.56	27.70	24.60	35.86	12.04
SP2e-16	68.26	129.74	63.39	52.15	177.55	102.16
SP2s-32	65.97	41.84	29.82	14.83	16.82	7.20
SP2e-32	77.11	152.45	64.20	99.56	442.96	325.18

check and a gauge of the relative performance of the different processors. “Partition total” time is the total time to partition the initial mesh using parallel sort inertial recursive bisection (PSIRB) [24]. “Partition migration” time is the portion of the “partition total” time that was spent doing message passing for mesh migration during partitioning. “Refine” time is the time spent refining the partitioned mesh. The amount of communication in this step depends on the number of processors that are affected by mesh refinement. “Rebalance total” time is the total time to rebalance the mesh using PSIRB after refinement, and “rebalance migration” is the subset of “rebalance total” time spent in message passing for migration during the rebalancing.

Table 6

Wall-clock timings for network tests using fine meshes on Sun Ultra 2/2200 workstations

System	Load	Partition		Refine	Rebalance	
		Total	Migration		Total	Migration
SUNsh-2	40.56	26.63	15.87	99.33	55.81	30.14
SUNlp4-2	32.91	26.60	15.81	97.11	54.71	29.58
SUNp4(10)-2	33.11	59.28	36.52	121.75	129.95	78.35
SUNp4(10)-4	33.78	102.13	46.55	51.26	208.63	48.03
SUNp4(10)-6	33.45	91.21	58.32	37.29	206.42	60.08
SUNp4(100)-2	33.15	29.64	17.93	95.64	60.79	33.52
SUNp4(100)-4	34.39	70.20	25.91	50.43	156.44	20.65
SUNp4(100)-6	33.27	45.80	30.05	34.44	114.81	16.39

Table 7

Wall-clock times for network tests using fine meshes on the SGI Onyx

System	Load	Partition		Refine	Rebalance	
		Total	Migration		Total	Migration
SGIsh-2	28.04	33.64	20.38	131.38	74.37	38.94
SGIlp4-2	28.93	34.25	21.05	131.61	73.60	40.54
SGIsh-4	29.57	50.63	31.28	66.88	97.59	25.50
SGIlp4-4	28.41	52.39	31.33	61.82	83.65	24.89
SGIsh-6	27.85	45.32	33.21	43.47	55.54	15.68
SGIlp4-6	28.06	47.49	34.76	42.24	59.43	18.40
SGIsh-8	29.34	51.52	36.37	33.88	73.18	25.51
SGIlp4-8	27.54	51.47	36.12	32.28	76.58	28.25

These network comparison tests provide insight into some important issues. Given current MPI implementations, the same communication device must be used for all communications, regardless of other available devices. The shared-memory device should be the fastest within a given computer, and it is interesting to know its performance relative to using or not using it. The results indicate that the shared-memory device is not always fastest and the difference between the two communication devices is not as large as expected. This could be due to a number of factors including system overhead, overhead introduced by MPICH, and the quality of the implementation of the shared-memory device. Further investigation is planned. On the Sun and SGI platforms (Tables 3, 4, 6 and 7), direct comparisons between the shared-memory and p4 network devices were run. The results were unexpected. The parts of the code which require communication typically ran at the same or slower speed relative to the shared-memory device. The serial portion of the code (loading the mesh) was significantly slower when using the shared memory device, especially on the Sun computer. One speculation for the reason for this is that the shared-memory device uses threads, forcing the use of thread-safe libraries which are typically slower due to mutual-exclusion locks [16]. We would like to make use of multithreading in the future, so this somewhat disturbing result will be investigated further. On the SGI platform, no significant difference was observed between the two communication devices.

We are also interested in comparing parallel performance of various interconnection networks. Scalability for clusters of workstations would be to be limited to systems having approximately 10 computers due to bandwidth and latency issues. The procedures tested performed about as well as expected and in some cases a bit better. Comparisons of the runs on the IBM SP2 using the high-speed switch and ethernet (Tables 2 and 5) show that the ethernet performance was fairly good up to about 8 processors. Times on ethernet were higher as expected, but as the number of processors increased, total execution time, especially for refinement, decreased. Above 8 processors, not unexpectedly, the performance of the ethernet falls

rapidly. Similar, but less dramatic effects can be seen in the tests run between the Sun computers on 10 Mbs ethernet and 100 Mbs fast ethernet (Tables 3 and 6).

The portion of partitioning and rebalancing time taken by migration varies from less than 10% to 73%. The PSIRB implementation works in two phases: (i) a partitioning phase which consists of computation with a relatively small amount of message passing to determine element destinations, and (ii) a communication-dominated migration phase. For example, consider the “Rebalance” times on the IBM SP2 (Table 5). With the switch, communication is relatively inexpensive and scales well. The sustained communication of the migration phase performs well relative to the shorter, latency-dominated communication and computation of the partitioning phase. Ethernet is generally slower and performance degrades as processors are added and as multiple processors attempt to send concurrent messages, which is typical in the migration phase.

In summary, these tests indicate that the mesh refinement and repartitioning codes scale well on reasonably large workstation clusters connected with, e.g., fast ethernet. The differences in the communication networks on this test case underscores the need for knowledge of these differences when trying to achieve a good partition for an adaptive computation.

5. Load balancing within RPM

Communication costs associated with interprocess boundary entities are useful in a partitioner. The availability of the information in the process and machine models allows a partitioner to avoid large communication volume on interprocess boundaries which communicate across a slow network. Many partitioners attempt to minimize the total size of the interprocessor boundaries, but that is not always equivalent to minimizing the actual communication cost. Like the computational costs of elements, the communication costs associated with boundary entities are not necessarily uniform.

Fig. 13 shows three 4-way partitionings of the initial mesh used in a simulation of a perforated shock tube [13]. The mesh contains 44177 regions. We consider only two factors in partition quality for this example, though other factors [6] must be considered for a more thorough analysis. We consider computational balance, and two *surface index* measures. The *global surface index* (GSI) measures the overall percentage of mesh faces which are on interpartition boundaries, and the *maximum local surface index* (MLSI) measures the maximum percentage of interpartition boundary faces on any one partition. Partitioning using PSIRB (Fig. 13(a)) achieves an excellent computational balance, with three processors assigned 11044 regions, and the fourth 11045 regions, and reasonable surface indices with $MLSI = 4.67$ and $GSI = 1.53$. A partition such as this would be useful when the primary concern is a good computational balance, such as in a shared-memory environment. The partition using ParMetis [22] (Fig. 13(b)) achieves excellent surface index values, $MLSI = 2.72$ and $GSI = 1.00$, but at the expense of a larger computational imbalance. Here, the partitions contain 10731, 11369, 10732, and 11345 regions. For a computation running on a NOW, it may be worth accepting the load imbalance to achieve the smaller communication volume.

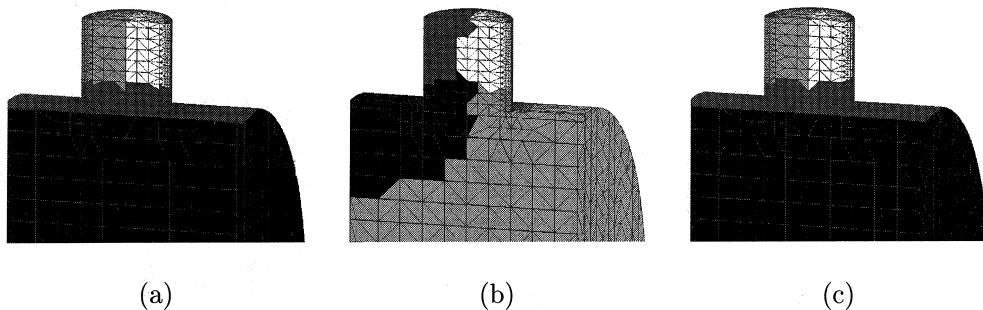


Fig. 13. Examples of partitioning using additional knowledge of parallel environment. Partitioning using (a) PSIRB, (b) ParMetis, (c) multilevel combination of PSIRB and ParMetis. Shading indicates processor assignments.

When executing on a pair of two-way SMP nodes connected by a network, a partition like the one shown in Fig. 13(c) may be desirable. Here, a combination of ParMetis and PSIRB is used. ParMetis is used to divide the computation among the two SMP nodes, resulting in partitions of 21 414 and 22 763 regions, with $MLSI = 0.99$ and $GSI = 0.48$. Within each SMP, the mesh is partitioned with PSIRB, resulting in the partitioning of Fig. 13(c), where the partitions contain 10 707, 10 707, 1, 381, and 11 382 regions, and the overall $MLSI = 5.18$ and $GSI = 1.69$. However, since the extra effort was made to minimize communication across the slow network, this will be a very effective partitioning for this environment.

A common operation when optimizing partitions is an interprocess boundary smoothing [10,12,14,21] which exchanges elements across interprocess boundaries to decrease communication volume. This operation can provide significant improvement in surface index values, especially when applied to partitions using spatial cuts to introduce a partition boundary (PSIRB or an octree-based partitioner [15,20]). However, this improvement is achieved at the cost of the smoothing operation itself and the small load imbalances introduced. Depending on the parallel environment, this may or may not be worthwhile and represents an additional way of taking advantage of information provided by RPM.

The measurement of load imbalance in our existing algorithms is based on a weighted number of elements. Some sources of heterogeneity, such as the extra work on higher-order elements in adaptive p -refinement or the extra time steps taken on smaller elements in a local refinement method [15], can be accounted for using the weighted element scheme. Others, like the run-time imbalance caused by other processes in a nondedicated parallel environment, are not known a priori with respect to a given step and do not fit well into this model. Even when it is possible to estimate appropriate weights, there may be considerable effort involved. It will be possible to measure load imbalance dynamically by instrumenting the communication libraries to measure synchronization delays.

6. Discussion

We present the architecture of RPM: a hierarchical partition model to provide a distributed mesh data organization and to quantify information about the parallel computing environment. Within RPM, mesh data is organized to allow efficient traversals of entities of the entire mesh, a partition, or a geometric model entity on which mesh entities are classified. Each traversal is important during phases of a parallel adaptive computation. Support for multiple partitions within a process allows for a natural way to represent entities on partition boundaries and an elegant mesh migration strategy. The process and machine models provide information about the parallel execution environment. In a preliminary study, we present execution times for a representative subset of our parallel adaptive finite element software on different computers and interconnection networks. There are significant differences in performance with, e.g., network bandwidth and latency. Partitioning algorithms that work well for one configuration may not be efficient for another. An example of a mesh partitioned by three different procedures on three different architectures illustrates this point. Further enhancements to existing partitioning algorithms and development of new algorithms becomes possible with the availability of the information from RPM.

In the future, the idea of a “partition” may be extended to add functionality for overlapping mesh migration with computation (called a *PipePartition*), and for computations on meshes too large to fit in memory (the *FilePartition*). A *StandardPartition* implements the behavior of a partition described previously. For a *StandardPartition* the mesh resides in memory within a single process. The *PipePartition* allows migration to be done incrementally as mesh entities are moved into the partition, while the *FilePartition* allows entities to be migrated to disk or other off-line storage. Each is derived from a generic *Partition*.

When a *PipePartition* is created it is given a destination partition within another process. The entities to be migrated are moved from their existing partition to this *PipePartition*. However, in the background, the entities are actually moved to the destination partition (in appropriately sized chunks) to maximize communication efficiency while the rest of the entities are being moved into the partition. In the final step, the partition is “sent” to the other process; however, at this point most of the data should already be sent. This step just packs and sends the remaining mesh entities. The most effective way to do this would

be to have the `PipePartition` spawn a separate thread to transfer the entities. This would require the transfer of entities to the partition to be thread safe.

The `FilePartition` is a similar concept, but its destination is a file rather than another partition. The primary use of this would be to act as a coarse level paging scheme so that a portion of the mesh of a large problem could be loaded into memory and processed while the remainder is kept on disk. This may also be used as the means of storing the entire partitioned mesh to disk.

It remains to show how well the algorithms used will scale to computers with tens of thousands of processors. Algorithms which frequently compute global reductions and synchronize globally will break down in such an environment. There is very likely to be some hierarchy with groups of processors connected by fast and slow network interfaces, making the hierarchical partition model a valuable tool. Additionally, the ability to store multiple partitions per process could allow partitioning to be inherently multilevel, allowing more global algorithms to work with the coarse structure of the partitions, while balance among partitions could be done via local neighborhood optimizations.

Acknowledgements

Some authors were supported by the National Science Foundation through grants ASC-9318184 and CCR-9527151, by AFOSR grant F49620-95-1-0407, by ARO grant DAAG55-98-1-0200, and by Simmetrix, Inc. Computer systems used include the IBM SP systems at the Maui High Performance Computing Center, and the IBM SP2 at Rensselaer.

References

- [1] M.W. Beall, M.S. Shephard, A general topology-based mesh data structure, *Int. J. Numer. Meth. Engng.* 40 (9) (1997) 1573–1596.
- [2] M.W. Beall, M.S. Shephard, A geometry-based analysis framework, in: *Advances in Computational Engineering Science*, Tech. Science Press, Forsyth, GA, 1997, pp. 557–562.
- [3] R. Biswas, K.D. Devine, J.E. Flaherty, Parallel, adaptive finite element methods for conservation laws, *Appl. Numer. Math.* 14 (1994) 255–283.
- [4] R. Biswas, L. Oliker, A. Sohn, Global load balancing with parallel mesh adaption on distributed-memory systems, in: *Proceedings of Supercomputing '96*, Pittsburgh, 1996.
- [5] C.L. Bottasso, H.L. de Cougny, M. Dindar, J.E. Flaherty, C. Özturan, Z. Rusak, M.S. Shephard, Compressible aerodynamics using a parallel adaptive time-discontinuous Galerkin least-squares finite element method, in: *Proceedings of the 12th AIAA Appl. Aero. Conf.*, No. 94-1888, Colorado Springs, 1994.
- [6] C.L. Bottasso, J.E. Flaherty, C. Özturan, M.S. Shephard, B.K. Szymanski, J.D. Teresco, L.H. Ziantz, The quality of partitions produced by an iterative load balancer, in: B.K. Szymanski, B. Sinharoy (Eds.), *Proceedings of the third workshop on Languages, Compilers, and Runtime Systems*, Troy, 1996, pp. 265–277.
- [7] K. Clark, J.E. Flaherty, M.S. Shephard, *Adaptive Methods for Partial Differential Equations*, *Appl. Numer. Math.*, special ed., 14, 1994.
- [8] K.D. Devine, J.E. Flaherty, Parallel adaptive *hp*-refinement techniques for conservation laws, *Appl. Numer. Math.* 20 (1996) 367–386.
- [9] M. Dindar, A. Lemnios, M. Shephard, J.E. Flaherty, K.E. Jansen, An adaptive procedure for rotorcraft aerodynamics, *AIAA Paper 98-2417*, 16th Applied Aerodynamics Conference, 1998.
- [10] P. Diniz, S. Plimpton, B. Hendrickson, R. Leland, Parallel algorithms for dynamically partitioning unstructured grids, in: D. Bailey (Ed.), *Proceedings of the Seventh SIAM Conference on Parallel Processing for Scientific Computing*, SIAM, February 1995, pp. 615–620.
- [11] H.C. Edwards, A Parallel Infrastructure for Scalable Adaptive Finite element Methods and its Application to Least Squares C^∞ Collocation, Ph.D. Thesis, The University of Texas at Austin, May 1997.
- [12] C. Farhat, M. Lesoinne, Automatic partitioning of unstructured meshes for the parallel solution of problems in computational mechanics, *Int. J. Numer. Meth. Engng.* 36 (1993) 745–764.
- [13] J.E. Flaherty, M. Dindar, R.M. Loy, M.S. Shephard, B.K. Szymanski, J.D. Teresco, L.H. Ziantz, An adaptive and parallel framework for partial differential equations, in: D.F. Griffiths, D.J. Higham, G.A. Watson (Eds.), *Proceedings of the 17th Dundee Biennial Conference on Numerical Analysis*, 1997; *Pitman Research Notes in Mathematics Series*, No. 380, Addison-Wesley, Reading, MA, 1998, pp. 74–90.
- [14] J.E. Flaherty, R.M. Loy, C. Özturan, M.S. Shephard, B.K. Szymanski, J.D. Teresco, L.H. Ziantz, Parallel structures and dynamic load balancing for adaptive finite element computation, *Appl. Numer. Math.* 26 (1998) 241–263.

- [15] J.E. Flaherty, R.M. Loy, M.S. Shephard, B.K. Szymanski, J.D. Teresco, L.H. Ziantz, Adaptive local refinement with octree load-balancing for the parallel solution of three-dimensional conservation laws, IMA Preprint Series 1483, Institute for Mathematics and its Applications, University of Minnesota, 1997; J. Parallel and Dist. Comput., to appear.
- [16] W. Gropp, Personal communication.
- [17] W. Gropp, E. Lusk, N. Doss, A. Skjellum, A high-performance, portable implementation of the MPI message passing interface standard, *Parallel Comput.* 22, 1996.
- [18] Message Passing Interface Forum, University of Tennessee, Knoxville, Tennessee, MPI: A Message Passing Interface Standard, 1st ed., 1994.
- [19] Message Passing Interface Forum, University of Tennessee, Knoxville, Tennessee, MPI-2: Extensions to the Message-Passing Interface, 1st ed., 1997.
- [20] T. Minyard, Y. Kallinderis, Octree partitioning of hybrid grids for parallel adaptive viscous flow simulations, *Int. J. Numer. Meth. Fluids* 26 (1998) 57–78.
- [21] T. Minyard, Y. Kallinderis, K. Schulz, Parallel load-balancing for dynamic execution environments, in: *Proceedings of the 34th Aerospace Sciences Meeting and Exhibit*, No. 96-0295, Reno, 1996.
- [22] K. Schloegel, G. Karypis, V. Kumar, Parallel multilevel diffusion algorithms for repartitioning of adaptive meshes, Tech. Report 97-014, University of Minnesota, Department of Computer Science and Army HPC Center, Minneapolis, MN, 1997.
- [23] M.S. Shephard, S. Dey, J.E. Flaherty, A straight forward structure to construct shape functions for variable p -order meshes, *Comp. Meth. Appl. Mech. and Engng.* 147 (1997) 209–233.
- [24] M.S. Shephard, J.E. Flaherty, C.L. Bottasso, H.L. de Cougny, C. Özturan, M.L. Simone, Parallel automated adaptive analysis, *Parallel Comput.* 23 (1997) 1327–1347.
- [25] V. Vidwans, Y. Kallinderis, V. Venkatakrishnan, Parallel dynamic load-balancing algorithm for three-dimensional adaptive unstructured grids, *AIAA J.* 32 (3) (1994) 497–505.
- [26] C.H. Walshaw, M. Berzins, Dynamic load balancing for PDE solvers on adaptive unstructured meshes, *Concurrency: Practice and Experience* 7 (1) (1995) 17–28.
- [27] C.H. Walshaw, M. Cross, M. Everett, Mesh partitioning and load-balancing for distributed memory parallel systems, in: *Proceedings of Par. Dist. Comput. for Comput. Mech.* Lochinver, Scotland, 1997.