

# SCALABLE PARALLEL OCTREE MESHING FOR TERASCALE APPLICATIONS

TIANKAI TU <sup>\*</sup>, DAVID R. O'HALLARON <sup>†</sup>, AND OMAR GHATTAS <sup>‡</sup>

**Abstract.** This paper presents a new methodology for generating and adapting octree meshes for terascale applications. Our approach combines existing methods, such as parallel octree decomposition and space-filling curves, with a set of new methods that address the special needs of parallel octree meshing. We have implemented these techniques in a parallel meshing tool called *Octor*. Performance evaluations on up to 2,000 processors show that *Octor* has good isogranular scalability, fixed-size scalability, and absolute running time. *Octor* also provides a novel data access interface to parallel PDE solvers and parallel visualization pipelines, making it possible to develop tightly coupled end-to-end finite element simulations on terascale systems.

**1. Introduction.** The emergence of terascale computing has created unprecedented new opportunities for scientists and engineers to simulate complex physical phenomena on a larger scale and at a higher resolution than heretofore possible. This trend has introduced new challenges to applications, algorithms, and systems software. In particular, a significant challenge for terascale finite element simulations is how to generate and adapt high-resolution unstructured meshes with billions of nodes and elements, and how to deliver such meshes to the processors of the terascale system.

A typical approach for preparing a finite element mesh for simulation is to first generate a large mesh structure offline on a server [29, 30], and then upload the mesh to the supercomputer, where additional steps such as mesh partitioning [18] and data redistribution are performed. The result is a series of large time-consuming file transfers and disk I/O that consume large amounts of network and storage resources while contributing nothing to the applications that will be using the mesh structure. Further, since the meshes are generated offline on a server, the offline algorithm is unable to adapt the mesh dynamically at runtime. We propose a better approach where the meshes are generated *in situ*, on the same processors where they will later be used by applications such as finite element solvers and visualization pipelines.

In this paper, we describe a new methodology for parallel octree meshing for terascale applications. At a first glance, it may appear that parallel octree meshing is simply a direct application of the well-known parallel octree method. In fact, parallel octree meshing is fundamentally different from other parallel octree applications in that it must manipulate the vertices (corners) of the octants, which correspond to the mesh nodes. This is required to support finite element simulations, which associate unknowns to mesh nodes and then solve the resulting linear system at these nodes. Complicated correlations between octants and vertices and between vertices and vertices, either on the same processor or on different processors, must all be identified and tracked. Thus, parallel octree meshing presents a set of new problems that do not exist in other parallel octree applications.

Our work builds on a foundation of previous work on parallel octrees [32, 27, 3, 13, 34, 9], mesh generation [5, 26, 22, 6, 13, 12, 17, 7], parallel adaptive mesh refinement [1, 8, 19, 33], parallel adaptive finite element methods [16, 23] and space-filling curves [4, 14]. To address the special requirements of parallel octree meshing, we have developed a set of new algorithms and techniques that include: (1) A new algorithm called *parallel prioritized ripple propagation* that balances an octree efficiently; (2) A sophisticated memory management scheme that facilitates data migration between processors; (3) A new method of extracting parallel mesh structures that delivers mesh data directly to solvers; and (4) A novel data access interface that allows solvers or visualizers to easily interact with the mesher.

We have implemented our methodology within a new parallel meshing tool called *Octor*. Performance evaluations show that *Octor* has good isogranular scalability, good fixed-size scalability, and good absolute running time. Roughly speaking, the runtime increases only logarithmically with the problem size (as measured by the number of mesh elements).

To our knowledge, *Octor* is the first parallel meshing tool that is capable of both generating massive unstructured meshes statically (that is, before a solver starts running) on thousands of processors, and sup-

---

<sup>\*</sup>Computer Science Department, Carnegie Mellon University, Pittsburgh, PA 15213, USA (tutk@cs.cmu.edu)

<sup>†</sup>Computer Science Department and Department of Electrical and Computer Engineering, Carnegie Mellon University, PA 15213, USA (droh@cs.cmu.edu)

<sup>‡</sup>Mechanics, Algorithms, and Computing Laboratory, Department of Biomedical Engineering and Civil & Environmental Engineering, Carnegie Mellon University, PA 15213, USA (oghattas@cs.cmu.edu)

porting dynamic mesh adaptation (that is, while a solver is still executing). Furthermore, it is the first parallel meshing tool that supports tight coupling with other components downstream in the simulation process, thus making it possible to develop terascale end-to-end finite element simulations [31] without the overheads of intermediate file transfers and disk I/O.

We have run Octor on Lemieux, the HP AlphaServer system at the Pittsburgh Supercomputing Center, to generate unstructured finite element meshes for simulating the 1994 Northridge earthquake in the Greater Los Angeles Basin. The largest mesh, with 1.22 billion elements and 1.37 billion nodes, resolves seismic wave frequencies up to 2Hz. The running time to generate this mesh on 2,000 processors is 333 seconds. Given that the execution time of a highly efficient seismic wave propagation solver [2] would be more than 60 hours on the same 2,000 PEs for a typical 80 second earthquake simulation, the meshing time is inconsequential.

The rest of the paper describes the design, implementation and evaluation of Octor. Section 2 provides the background of octree meshing. Section 3 briefly discusses related work. Section 4 presents an overview of the design of Octor. Section 5 explains the algorithm and implementation details. Section 6 presents the performance evaluation of Octor.

**2. Background.** Among the different types of finite element meshes, octree-based hexahedral meshes provide a trade-off between modeling capability and simplicity. On the one hand, octree meshes provide multiresolution capability by local adaptation. On the other hand, octree meshes simplify the treatment of geometry and boundaries by spatial approximation. Octree mesh methods are employed in at least three ways. (1) For PDEs posed in simple domains characterized by highly heterogeneous media in which solution length scales are known *a priori* (such as in linear wave propagation), octree meshes that resolve local solution features can be generated up front. (2) For PDEs in simple domains having solution features that are known only upon solution of the PDEs, octree meshes — driven by solution error estimates — can be adapted dynamically to track evolving fronts and sharp features at runtime (for example to capture shocks). (3) For PDEs posed on complex domains, octrees meshes in combination with special numerical techniques (such as fictitious domain, embedded boundary, or extended finite element methods) can be used to control geometry approximation errors by adapting the octree mesh in regions of high geometric variability, either *a priori* for fixed geometries, or at runtime for evolving geometries. Large-scale examples of *a priori* adapted octree mesh generation can be found in seismic wave propagation modeling [20], while octree mesh methods for compressible flow around complex aircraft configurations provides an excellent example of geometry- and solution-driven dynamic adaptivity [35].

Conceptually, generating or adapting an octree-based hexahedral mesh is straightforward. As shown in Figure 2.1, a problem domain is first refined or coarsened recursively using an octree. We require that two adjacent octants sharing an edge of a face should not differ in edge size by a factor of 2. This constraint is often referred to as the *balance condition* or, more intuitively, *2-to-1 constraint*. Then, octants are mapped to *mesh elements* and vertices are mapped to *mesh nodes*. The nodes hanging at the midpoint of an edge or the center of the face of some element (due to the 2-to-1 constraint) are identified as *dangling nodes*. The remaining nodes are *anchored nodes*. For conforming finite element methods, each dangling node is dependent on the anchored nodes at the endpoints of the edge or the face on which it is hanging through an explicit algebraic constraint. Explicit correlations between dangling nodes and anchored nodes are established.

**3. Related work.** Parallel octree structures have been successfully applied in many areas, especially in the area of N-body simulations [32, 27, 3, 34]. Unlike in octree meshing, these applications do not need to manipulate the vertices of the octants (cells). Nor is it necessary to enforce the balance condition on the octree.

In the mesh generation area, [5, 26, 22] present sequential algorithms that use a quadtree or octree to partition a 2D or 3D point set and then generate triangular or tetrahedral meshes by “warping” the vertices of the leaf octant. [25] presents another sequential algorithm that also uses an octree to partition a domain but converts the leaf octants to hexahedral meshes. [6] introduces a PRAM model of parallelizing quadtree-based triangular mesh generation. [13, 11, 7] present techniques for parallel tetrahedral mesh generation. Most parallel mesh generators, though, run on 32 or 64 processors. To our knowledge, none has run on more than 128 processors. However, it should be noted that generating geometry-conforming 3D meshes (either Delaunay or non-Delaunay) in parallel is a very difficult problem and is being actively researched [10].

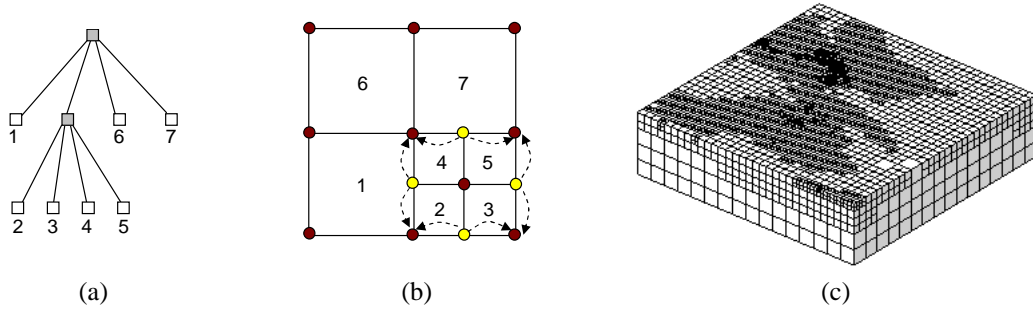


FIG. 2.1. Octree-based hexahedral meshes. (a) Octree domain decomposition. (We use a 2D quadtree to illustrate concepts, but all techniques are designed and implemented for 3D cases.) (b) Octants map to elements and vertices map to mesh nodes. The dark colored dots represent the anchored nodes and the light colored dots represent the dangling nodes. The dashed arrows represent the explicit correlations between dangling nodes and anchored nodes. (c) An example 3D octree mesh.

For parallel mesh adaptation, there are two popular methods: adaptive mesh refinement (AMR) [1, 8, 19, 33] and adaptive finite elements [15, 16, 23]. Both methods start with an initial, coarse mesh, and, as the numerical calculation proceeds, adapt the mesh structure (using different techniques) in those areas where error estimators or indicators dictate a need to do so.

**4. Design overview.** Figure 4.1 shows the process of an end-to-end physical simulation. First, a mesh is generated to model the material property or geometry of the problem. Next, a solver takes the mesh as input and conducts numerical calculations, for example, to solve a system of PDEs. In some applications, the mesh structure needs to be adapted dynamically while the solver is still executing. The numerical results produced by the solver, such as displacements or temperatures, are then correlated to the mesh structure by a visualizer to create 3D images or animations.

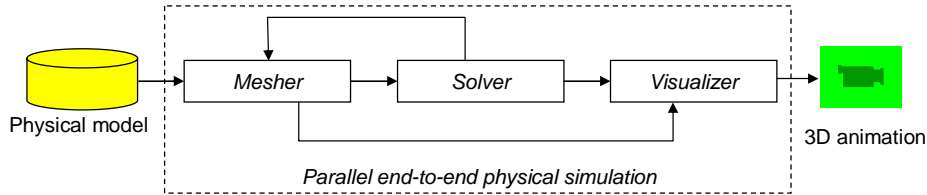


FIG. 4.1. End-to-end physical simulation scenario. The backward arrow from solver to mesh generator represents mesh structure adaptations driven by numerical computations, for example, in adaptive finite element simulations.

The design goal of Octor is to develop a scalable, efficient, and easy-to-use mesher to support such end-to-end simulations for terascale applications. *Scalable* means that the mesher should be able to run on a large number processors without performance degradation. *Efficient* means that mesh elements and nodes should be produced where they will be used instead of on remote processors. *Easy-to-use* means that other simulation components should be able to interact with Octor in a simple way.

To achieve this design goal, we have adopted existing techniques as well as developed new algorithms and methods. We use a distributed, parallel octree to manage the backbone meshing data structure and space-filling curve (SFC) based techniques to partition the meshes. The new methods we have developed are: (1) A new algorithm called *parallel prioritized ripple propagation* that efficiently enforces the 2-to-1 constraint on a parallel octree; (2) A sophisticated memory management scheme that keeps track of application-specific meshing data internally and takes care of data migration between processors (caused by partitioning) automatically; (3) A new method of extracting mesh structures: Instead of extracting a mesh structure first and then partitioning, we partition the backbone parallel octree first and then extract the mesh structure; (4) A novel data access interface that allows a solver or a visualizer to interact easily with the mesher.

Octor provides a small set of API functions (see Appendix) that allow users to write meshing applications. At a high level, Octor meshing consists of the steps shown in Figure 4.2. First, NEWTREE bootstraps

a small octree on each processor. Next, the tree structure is adjusted by REFINETREE and COARSENTREE, either statically or dynamically. While adjusting the tree structure, each processor is responsible only for a small area of the domain. When the adjustment completes, there are many subtrees distributed among the processors. The BALANCETREE step enforces the 2-to-1 constraint on the parallel octree. After a balanced parallel octree is obtained, PARTITIONTREE redistributes the leaf octants among the processors. Finally and most importantly, EXTRACTMESH derives mesh element and node information and determines the various correlations between elements and nodes. The procedure just described provides a general framework for different types of meshing. A particular application may choose not to execute some of the listed steps. For example, for static mesh generation, it might be unnecessary to execute the COARSENTREE step. For dynamic mesh adaptation, it might be unnecessary to run PARTITIONTREE if the numbers of leaf octants on different processors only differ within a small percentage.

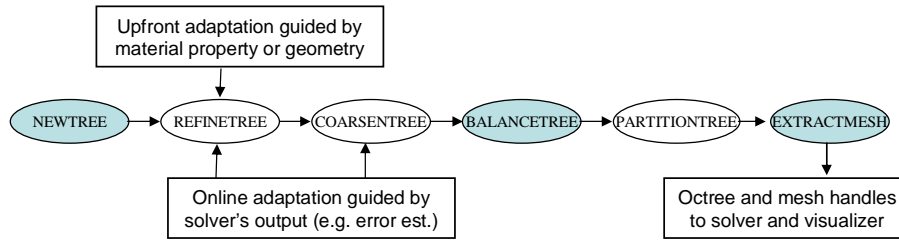


FIG. 4.2. Octree meshing steps. Shaded ovals are required steps. Unshaded ovals are optional steps.

An application such as a finite element solver or a parallel visualization package can interact with Octor in one of two ways. (1) *Indirect interaction*. The application can guide the manipulation of Octor’s internal data structures indirectly using *callback functions*. The function prototypes for REFINETREE, COARSENTREE, and BALANCETREE contain predefined callback function types. The application implements these callback functions and passes the function pointers as arguments when performing any of the three meshing steps. (2) *Direct interaction*. Octor also returns abstract data types (ADT) such as `octree_t`, `octant_t` and `mesh_t` to an application, and provides a set of macros that allow the application to manipulate these ADTs directly. For example, an application can attach to an ADT a pointer that references an arbitrary data structure in memory whose semantics are completely irrelevant to meshing. Although simple, these two data access methods have proven to be sufficiently powerful for integrating a parallel mesher, solver and visualizer [31].

**5. Internals.** This section explains the internal mechanisms of Octor. We highlight only the important algorithm and implementation features, omitting many of the technical details.

**5.1. NEWTREE.** This step operates in three stages: *tree expansion*, *task assignment*, and *tree pruning*.

*Tree expansion.* Each processor starts with a root octant that encloses the entire domain, and then independently expands this tree until a certain level is reached. The level is determined by the number of processors and is usually a small value (3 or 4). During the expansion process, any octant that falls outside of the problem domain is discarded. When the expansion stops, each processor has an identical copy of an initial octree.

*Task assignment.* In this context, a *task* corresponds to the refinement of a leaf octant in the initial octree. Using the preorder traversal as the total ordering, we perform a parallel block data decomposition of the leaf octants. Since all of the processors have complete information about the initial octree, no communication is required during this step. Each processor is assigned a contiguous chunk (in preorder) of leaf octants to refine. These octants are marked as local, and the remaining octants are marked as remote.

*Tree pruning.* On each processor, octants marked as remote are pruned from the initial octree. This is implemented as an octree aggregation operation.

One caveat: because the refinements of the cubic areas (leaf octants) are determined by an application, the workload of each task can be highly uneven, for example, due to the large heterogeneity within the problem domain. There is no way to determine an optimal task assignment scheme. The best we can do is to assign an approximately equal number of tasks to each processor.

**5.2. REFINETREE/COARSENTREE.** In the REFINETREE step, each processor traverses its local leaf octants. On visiting an octant, Octor invokes an application-supplied callback function to determine whether the octant needs to be decomposed. In case a decomposition is triggered, eight new children octants are created and linked into the tree structure automatically. In addition, extra internal memory is allocated by Octor for each new octant to store application-specific meshing data such as material density or velocity. A value initialization callback function is invoked to set the values properly.

The COARSENTREE step follows the reverse procedure. However, there may be cases where a coarsening cannot be carried out. This happens when the eight children of an interior octant are spread over two processors. See Figure 5.1(a) for an example. In these cases, we skip the coarsening of that particular region. Since it occurs only on the boundary of two processors, most internal coarsening can still be performed without problem. Further, since the mesh is likely to be re-partitioned after the REFINETREECOARSENTREE step, those octants spread over a boundary are very likely to be re-distributed to the same processor. Then in the next round of coarsening, they can be successfully merged.

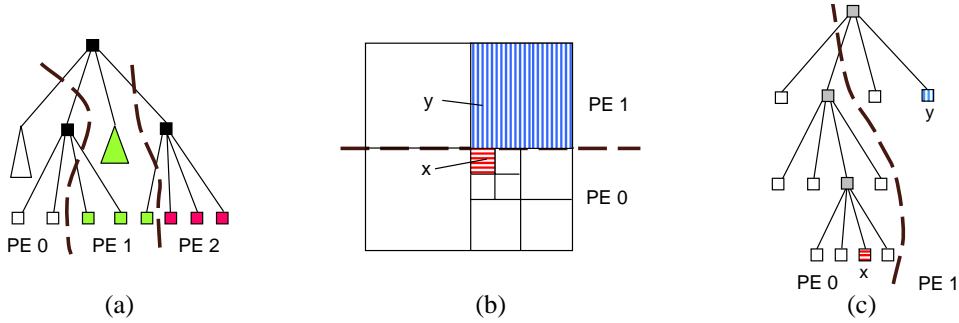


FIG. 5.1. Operations across processor boundaries. The dashed lines represent the boundaries between processors. (a) Coarsening is skipped on processor boundaries. (b) A violation of the balance condition. Octant  $x$  and  $y$  are adjacent to each other spatially but differ in size by a factor of 4. (c) The two spatially adjacent octants ( $x$  and  $y$ ) are distributed on two processors.

**5.3. BALANCETREE.** We have proposed an efficient algorithm called *prioritized ripple propagation* in [29]. The basic idea is to eliminate the effect of small octants on large octants level by level. The algorithm first visits all leaf octants at the lowest level of an octree. For each leaf octant, a series of neighbor-findings are carried out. Any neighboring leaf octants that are more than twice as large as the current octant are decomposed repeatedly until the 2-to-1 constraint is satisfied. When all the leaf octants on the lowest octree level are processed, the algorithm finishes one iteration and moves one level up to start another iteration that processes the leaf octants at a higher level. The algorithm terminates when all the levels are processed. The neighbor-finding operation is implemented using a variant of the standard pointer-based neighbor-finding algorithm [24].

This algorithm works well on a single processor but would run into trouble on a multiprocessor. Figures 5.1(b) and 5.1(c) illustrate the problem. If octant  $x$  and  $y$  are spatially adjacent to each other but are distributed on PE 0 and PE 1, respectively, then when  $x$  searches its neighbor in the north direction, the neighbor-finding algorithm will get stuck.

The *parallel prioritized ripple propagation* algorithm overcomes this difficulty in the following way. First, we build a small array (linearly proportional to the number of processors) on each processor that contains the coordinate information of all the *first local* octants of all the processors. When the standard neighbor-finding algorithm gets stuck, we search the small array using the SFC value of the remote neighbor to find which processor hosts the neighbor. However, we do not forward the ongoing neighbor-finding operation to the remote processor immediately. Instead, a request (containing the information of the current octants as well as partial tree-traversal information) is stored in a request buffer that is destined for the host processor. When all processors finish processing all local octants at the current tree level, the request buffers are exchanged between processors. Each processor then services the neighbor-finding requests from other processors and decomposes its large octants as necessary. When all of the processors finish servicing remote requests, the algorithm proceeds to the next higher octree level. The key implementation technique

of the parallel prioritized ripple propagation algorithm is a concise but expressive encoding scheme for the neighbor-finding requests.

**5.4. PARTITIONTREE.** This step creates a new partition for the parallel octree and then migrates octants between processors according to the partition.

The partition is computed in a very simple way. We sort all the leaf octants in Z-order, split them into equal-size chunks and assign the  $i$ th chunk to the  $i$ th processor. Since the preorder traversal of octree leaves is the same as Z-order, there is no need to physically sort the leaf octants. All the leaf octants on processor  $i$  must be sorted in between the leaf octants on processor  $(i - 1)$  and the leaf octants on processor  $(i + 1)$ . See Figure 5.1(a) for an example. Therefore, the only information needed to compute an even partition is the number of leaf octants on each processor. After a few reduction-style communications, each processor receives the information it needs to compute the partition.

The more difficult part of PARTITIONTREE step is actually migrating octants between processors. We have created an internal memory manager for Octor to facilitate this procedure. When octants are created in the refinement process, internal memory is allocated to hold the application-specific data such as velocity or density. While executing the PARTITIONTREE step, the Octor internal memory manager marshals the data that needs to be migrated and sends it to the destination processor in bulk. The unused memory on the sending processor is immediately reclaimed. On the receiving end, the Octor memory manager installs the newly arrived leaf octants in the local octree, and allocates extra memory space to store any application-specific meshing data. There are two advantages of relying on the Octor internal memory manager. First, the function call interface for PARTITIONTREE is very simple to use (see Appendix). Second, Octor is more robust, encapsulating and hiding the complex memory management issues from the application.

A side note: an interesting design feature of Octor is that we have chosen to partition the tree *before* extracting the mesh structure. By migrating octants to their destinations first, each processor can obtain all necessary information for extracting the mesh structures for the simulation or visualization code that is going run on that processor.

**5.5. EXTRACTMESH.** As mentioned earlier, parallel octree meshing is fundamentally different from other parallel octree applications in that there are complicated operations associated with the vertices of the octants, i.e., the mesh nodes.

All of the complexities of handling mesh nodes and other mesh data structures are encapsulated in the EXTRACTMESH step, which performs the following operations: (1) Extracting mesh nodes from the parallel (balanced) octree and computing their coordinates; (2) Discovering dangling nodes; (3) Assigning each element and each node to an owner processor; (4) Establishing sharing relationships between nodes on different processors; (5) Allocating unique *global* element ids and node ids; (6) Assigning per-processor *local* element ids and node ids, (7) Establishing the mapping between global ids with local ids; (8) Correlating local element ids with local node ids (mesh connectivity information); (9) Correlating local dangling node ids with local anchored node ids. These operations have been implemented in an intertwined fashion in Octor in order to reduce communication cost and algorithm complexity. The rest of this section focuses on three interesting issues that will shed some light on the complexity unique to parallel octree meshing.

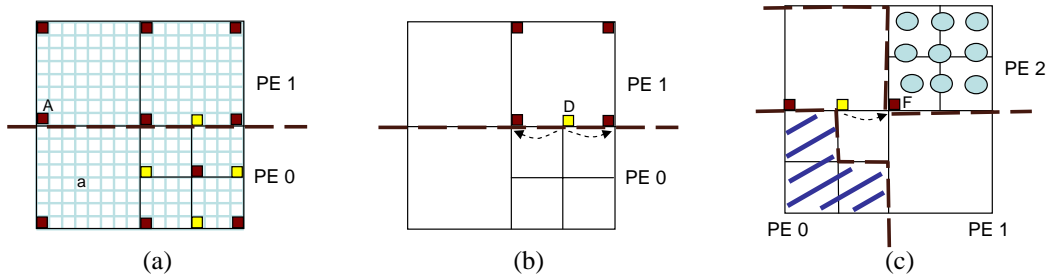


FIG. 5.2. How to extract mesh nodes and identify indirect sharing. (a) Mesh nodes are mapped to the smallest octants in the domain and assigned to processors that cover the respective SFC segments. (b) Subtle sharing case 1: PE 1 is assigned the ownership of a dangling node but does not use it. (c) Subtle sharing case 2: The patterns in the figure illustrate the partitioning of the mesh. PE 0 indirectly uses an anchored node that is on PE 2.

*Determining the ownership of mesh nodes.* Finite elements with Lagrange bases associate unknowns to mesh nodes. If two or more elements sharing a node are distributed on different processors, then each element (and thus its processor) will make some contribution to the values associated with the shared node. The question is which processor will accumulate the value contributions that are distributed among different processors? A commonly used method is to assign each node to an owner processor who will become responsible for all the computation related to a node. A naive implementation would let all the processors sharing that node enter a negotiation process and somehow select a processor as the owner. Our implementation employs a different technique. The basic idea is to treat each mesh node as a smallest octant in the domain. Then each mesh node must fall in the area covered by some processor, as shown in Figure 5.2(a). The processor that encloses the mesh node is assigned as the owner. This procedure can be easily implemented, using well-known space-filling curve techniques, with no interprocessor communications.

*Determining the processors that share a mesh node.* The dual problem of assigning a node to an owner processor is to determine which other processors share a particular node. From the standpoint of an owner processor of a particular node, sharing means that a remote processor needs to access the values associated with this node. A simple case of sharing is shown in Figure 5.2(a). The mesh node marked A is owned by PE 1. PE 0 has element a that will contribute to the values associated with A. Thus, PE 1 and PE 0 should both expect to exchange data with each other. However, there are more subtle cases of sharing. Figure 5.2(b) shows a case where PE 1 does not know the existence of the dangling node D because none of the octants on PE 1 has D as a vertex. However, on PE 0, node D is indeed a vertex of some octant and is produced as a mesh node. Using the node ownership assignment algorithm just described, PE 0 determines PE 1 is the owner of node D and will send data associated with node D to PE 1. Now a discrepancy in communication occurs. Another subtle case is illustrated in Figure 5.2(c). PE 0 has a dangling node that is dependent on node F, which is owned by PE2. So PE 0 needs to make use of the values associated with F. But PE 2 would not be able to determine that PE 0 is sharing F with it. Note that PE 0 and PE 2 are not adjacent to each other in the partition. To resolve such subtle sharing situations, we use a single round of all-to-all communication. Fortunately, since the subtle sharing cases only occur when dangling nodes are involved in the processor boundary area, they are much less common than the simple cases.

*Returning a mesh structure that can be directly used by a solver.* Since we are already extracting mesh structures on the same processors where solvers will be running, we have gone one step further to extract a mesh structure that can be directly used by finite element solvers. By *directly*, we mean that when a handle to the local mesh is returned to an application, a solver can immediately build its communication schedule without incurring any communication. Operations (6)–(9), listed earlier in this section, accomplish this task on a solver’s behalf. The advantage of carrying out all these mapping operations within Octor is that we can make efficient use of the various internal data structures built during the EXTRACTMESH step. For example, we use the same node hash table repeatedly for various operations such as discovering dangling nodes and correlating mesh elements to mesh nodes.

**6. Performance evaluation.** In this section, we present a performance evaluation of Octor for statically generating large unstructured finite element meshes.<sup>1</sup> The meshes are generated for modeling earthquake simulations, and are directly used by an explicit seismic wave-propagation solver. The target region of these simulations is the Greater Los Angeles Basin, which comprises a 3D volume of 100 km x 100 km x 37.5 km. The material model we used to drive the mesh generation process is the Southern California Earthquake Center (SCEC) 3D velocity model [21] (Version 3, 2002). A sequence of meshes is generated to satisfy different simulation frequency requirements. Roughly speaking, the higher the frequency, the finer (larger) the mesh. Performance data presented in this section were obtained from production and experimentation runs of earthquake simulations on Lemieux, the HP AlphaServer system at Pittsburgh Supercomputing Center, on a number of processors ranging from 1 to 2,000.

Our performance evaluation focuses on two issues: (1) The isogranular scalability [28] of Octor. That is, how scalable is Octor when we generate larger meshes on more processors while keeping the average number of elements on each processor more or less the same? (2) The fixed-size scalability of Octor. That is, how

<sup>1</sup>We have not yet evaluated the performance of Octor for real-world dynamic meshing applications. Thus, the performance of COARSENTREE is not presented.

scalable is Octor when we fix the problem size and increase the number of processors?

**6.1. Isogranular scalability study.** Figure 6.1 summarizes the results of generating meshes of various sizes on different number of processors. Since the earth is highly heterogeneous, the largest elements of the meshes are 64 times as large in edge size as the smallest ones (the difference between “max leaf level” and “min leaf level”). Also, because of the multi-resolution in the mesh, there are a large number of dangling nodes in the mesh (between 11% and 20%).

The larger meshes (1Hz, 1.5Hz, and 2Hz) are for terascale finite element simulations that run on thousands of processors [2]. However, for the purpose of this study, we have chosen the numbers of processors to ensure that the memory utilization per processor across the different experiments is roughly the same. Given the unstructured nature of the meshes, it is impossible to guarantee the per-processor element (node) number to be exactly the same over different runs. Nevertheless, we have contained the difference to within 10% (as shown on the “Element/PE” row). The smallest mesh (0.23Hz) has 661K elements and is generated on one processor, while the largest one (2Hz) has 1.22B elements and is generated on 2,000 processors.

PEs	1	16	52	184	748	2000
Mesh name	0.23Hz	0.5Hz	0.75Hz	1Hz	1.5Hz	2Hz
Elements	6.61E+5	9.92E+6	3.13E+7	1.14E+8	4.62E+8	1.22E+9
Nodes	8.11E+5	1.13E+7	3.57E+7	1.34E+8	5.34E+8	1.37E+9
Anchored	6.48E+5	9.87E+6	3.12E+7	1.14E+8	4.61E+8	1.22E+9
Dangling	1.63E+5	1.44E+6	4.57E+6	2.03E+7	7.32E+7	1.48E+8
Max leaf level	11	13	13	14	14	15
Min leaf level	6	7	8	8	9	9
Elements/PE	6.61E+5	6.20E+5	6.02E+5	6.20E+5	6.18E+5	6.12E+5
Octor time(sec)	19.98	74.79	127.95	150.26	305.35	332.85

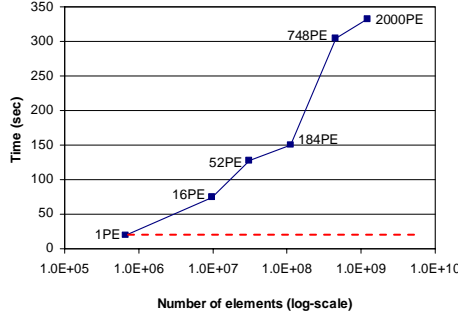
FIG. 6.1. Summary of the meshes of different resolutions for earthquake wave-propagation simulations and the time to generate them on different number of processors.

Figure 6.2(a) shows how the Octor meshing time increases as the problem size increases (along with the number of processors used). Note that the horizontal axis (problem size) is in log-scale, while the vertical axis (time) is in linear scale. It can be seen that, roughly, the meshing time increases as just a logarithm of problem size. As shown in Figure 6.2(b), the increased running time is mostly due to the PARTITIONTREE step. Although Octor does not achieve the theoretical optimal isogranular performance (the dashed line), the linear meshing time trend suggests that Octor has good isogranular scalability. We can speculate that if Octor were to be used to mesh a problem 10 times as large as the 2Hz mesh on 20,000 processors, then the meshing time would only increase moderately. Here is another way to appreciate the efficiency of Octor: Assume that a 2Hz mesh (80 GB in size) already exists on a lab server and that we have to move it across the network to a supercomputer. Transferring the mesh at peak rates over a gigabit ethernet connection would require more than 10 minutes. In comparison, Octor takes less than 6 minutes to generate the mesh from scratch *in situ*.

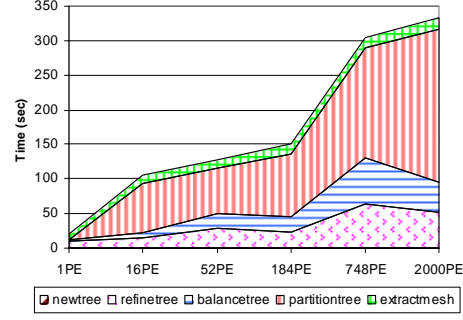
We now examine the running time of different meshing steps in more detail for the isogranular runs. Figure 6.3(a) shows the execution time breakdown for generating the different meshes (annotated by the number of processors used). Each bar represents the contribution of the five meshing steps as a percentage of the total execution time. From the fact that percentage breakdowns (except for single-PE run that does not invoke PARTITIONTREE) do not change significantly between different runs, we can deduce that none of the meshing steps is a scalability bottleneck. Otherwise, we would have seen disproportionately large fluctuations in the percentage breakdowns. We can also see that the dominant cost of meshing is associated with PARTITIONTREE, which is probably the most bandwidth (memory and network) intensive operation. However, further research is needed to quantitatively determine the reason.

Figure 6.3(b) shows the running time of the meshing steps in a different way. We group together the absolute running times of each step for different runs and compare the running time change of each individual step over different problem sizes. Since the NEWTREE times are below 0.5 seconds for all the runs, we have not presented them as a group in the figure. It is clear that REFINETREE, BALANCETREE, and EXTRACTMESH perform well in an isogranular setting, while the performance of PARTITIONTREE deteriorates somewhat with larger numbers of processors. It is surprising that EXTRACTMESH, which is the most complex step in terms of algorithms and communications, runs faster and scales better than all of the other steps except for the relatively simple NEWTREE step. This is strong evidence that our techniques for extracting mesh structures





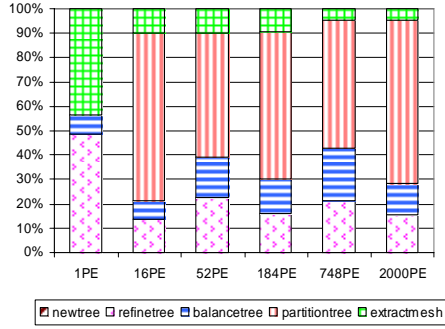
(a) Execution time trend.



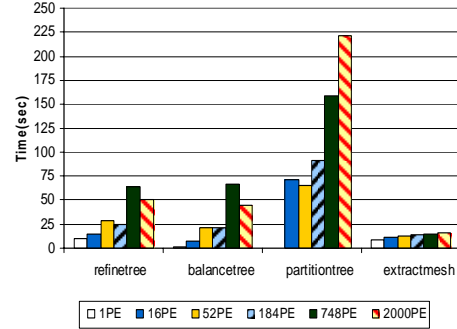
(b) Execution time breakdown for the steps.

FIG. 6.2. Octor meshing execution time and detailed breakdown. (a) The running time of Octor increases roughly as a logarithm of problem size (as determined by the total number of elements). (b) The trend of execution time increases in different steps. The time for each step is plotted with a unique pattern, which is stacked on top of the previous step. From bottom up, the stripes represent time spent in the NEWTREE, REFINERTREE, BALANCETREE, PARTITIONTREE, and EXTRACTMESH step, respectively. Note that since the NEWTREE time is very short (less than 0.5 sec), its stripe is invisible in the figure.

are effective and efficient.



(a) Execution time percentage breakdown.



(b) Comparison of execution times of each step on different number of processors.

FIG. 6.3. Understanding the performance of Octor from two different perspectives. Figure (a) shows where the execution time goes. Figure (b) shows how each step's running time varies on different numbers of processors.

**6.2. Fixed-size scalability study.** For applications that are computation bound, larger numbers of processors might be used to solve a moderately-sized problem just to improve the turn-around time. Memory utilization is not of particular importance to these applications. Therefore, we also need to investigate the performance of Octor when we fix the problem size and increase the number of processors.

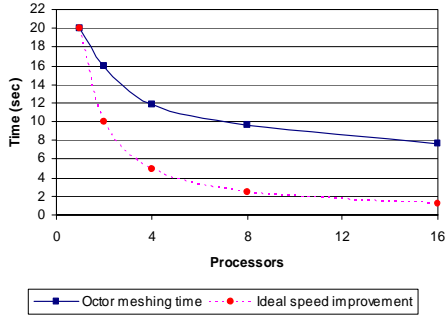
Three sets of fixed-size scalability experiments were conducted, for small size, medium size and large size problems, respectively. The experimental setups are shown in Figure 6.4.

PEs	1	2	4	8	16	32	64	128	256	512	748	1024	2000
Small case (0.23Hz, 661k elements)	x	x	x	x	x								
Medium case (0.5Hz, 9.92m elements)				x	x	x	x	x					
Large case (1Hz, 114m elements)								x	x	x	x	x	x

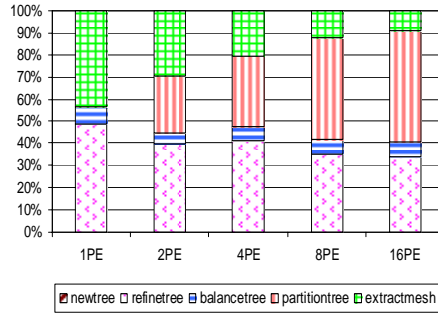
FIG. 6.4. Setup of fixed-size speedup experiments. Entries marked with "x" represent experiment runs.

Figure 6.5 shows the performance of Octor for different fixed-sized problems. The figures on the left side shows the speedup of meshing on more processors. The dashed line represents the ideal speedup. Although not perfect, Octor achieves good speedups while running on a large number of processors. The figures on the right show the execution time percentage breakdowns for different fixed-sized runs on different numbers of

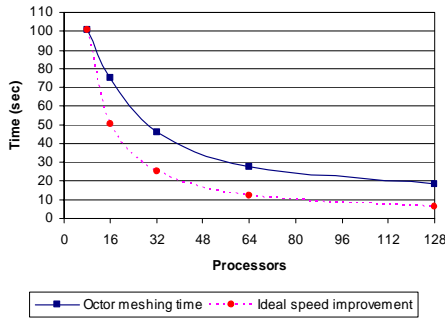
processors. That there is no major fluctuation in the breakdowns (except for the small case problem where 1 PE is involved) is more evidence that there are no scalability bugs among the meshing steps.



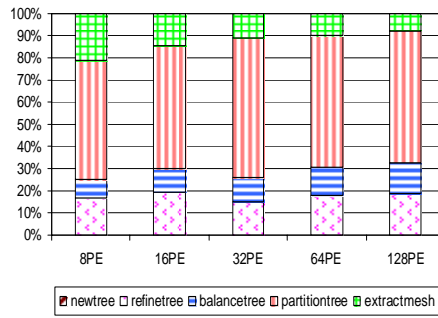
(a) Small case meshing speedup.



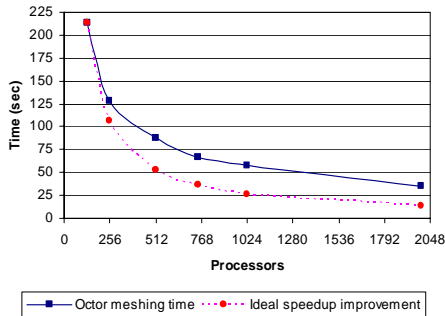
(b) Small case running time percentage breakdown.



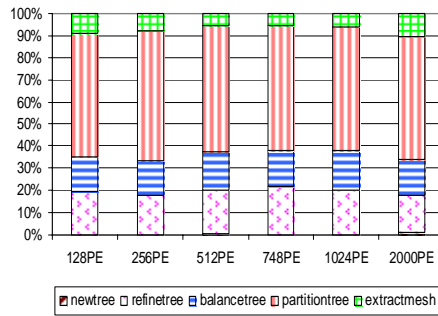
(c) Medium case meshing speedup.



(d) Medium case running time percentage breakdown.



(e) Large case meshing speedup.



(f) Large case running time percentage breakdown.

FIG. 6.5. Execution time of generating a fixed-sized mesh on different numbers of processors and how different steps of meshing contribute to the total running time percentage wise.

An important and closely related question is what happens to a solver when a fixed-sized problem is meshed on larger numbers of processors? Figure 6.6 shows the speedup of an explicit finite element seismic wave propagation solver for earthquake simulation. The solver takes a mesh produced by Octor *in situ* and runs on the same processors. The solver runs 4,000 time steps for the medium case (Figure 6.6(a)) and 2,500 time steps for large case. Each time step requires the same amount work for a fixed-size problem. A complete earthquake simulation would run somewhere between 10,000 to 80,000 time steps.

A surprising finding from these figures is that the solver achieves almost perfect speedup on hundreds and thousands of processors, even though the strategy we used to partition the tree is extremely simple. Also, a closer look at the vertical axes of Figure 6.5(c)(e) and Figure 6.6(a)(b) reveals that the meshing time becomes inconsequential when compared with the running time of the solver. So the bottom line is that as

long as a physical problem can be meshed in a reasonable amount of time, Octor meshing will not constitute a bottleneck for terascale simulations.

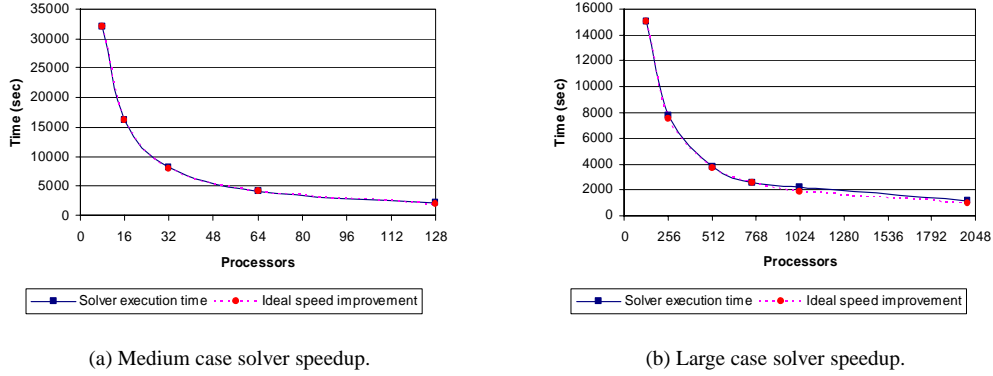


FIG. 6.6. Solver achieves almost perfect speedup using the partitioned meshes generated by Octor. (We have omitted the plotting for the small case, which looks identical to that of the medium case.)

**7. Conclusion.** This paper has demonstrated how to combine existing, well-known techniques, such as parallel octrees and space-filling curves, with a new set of algorithms and methods, such as parallel prioritized ripple propagation, to deliver a new capability for generating and adapting massive octree meshes for high-performance scientific computing.

We have implemented these new ideas in Octor, a parallel octree meshing tool that has been used successfully for terascale finite element earthquake simulations. The good scalability of Octor on up to 2,000 processors indicates that it has the potential to continue to scale on more processors. Future work will focus on applying Octor in the context of terascale applications with adaptive meshing requirements.

**Acknowledgments.** This work is sponsored in part by NSF under grant IIS-0429334, in part by a sub-contract from the Southern California Earthquake Center as part of NSF ITR EAR-0122464, in part by NSF under grant EAR-0326449, and in part by a grant from Intel. Supercomputing time at the Pittsburgh Supercomputing Center is supported under NSF TeraGrid grant MCA04N026P. We would like to thank our Quake project colleagues Jacobo Bielak and Leonardo Ramirez-Guzman, and our SCEC CME partners Tom Jordan and Phil Maechling, for providing such a compelling application context. Thanks also to Hongfeng Yu and Kwan-Liu Ma for helping us understand the data access requirements of parallel visualization packages. Special thanks to John Urbanic, Chad Vizino, and Sergiu Sanielevici at PSC for their outstanding technical support.

## Appendix. Major functions of the Octor API.

```
octree_t * octor_newtree(double x, double y, double z, int rectx, int myid, int groupsize);
int octor_refinetree(octree_t *octree, toexpand_t *toexpand, setrec_t *setrec);
int octor_coarsentree(octree_t *octree, toshrink_t *toshrink, setrec_t *setrec);
int octor_balancetree(octree_t *octree, setrec_t *setrec);
int octor_partitiontree(octree_t *octree);
mesh_t * octor_extractmesh(octree_t *octree);
void octor_deletemesh(mesh_t *mesh);
void octor_deletetree(octree_t * octree);
```

## REFERENCES

- [1] T. ABEL, G. L. BRYAN, AND M. L. NORMAN, *The formation and fragmentation of primordial molecular clouds*, Astron. Astrophys., (2000), pp. 39–44.
- [2] V. AKCELIK, J. BIELAK, G. BIROS, I. IPANOMERITAKIS, ANTONIO FERNANDEZ, O. GHATTAS, E. KIM, J. LOPEZ, D. R. O’HALLARON, T. TU, AND J. URBANIC, *High resolution forward and inverse earthquake modeling on terascale computers*, in SC2003, Phoenix, AZ, Nov. 2003. Gordon Bell Award for Special Achievement.

- [3] S. ALURU, G. M. PRABHU, AND J. GUSTAFSON, *Truly distribution-independent algorithms for the n-body problem*, in SC1994, Washington, D.C., 1994.
- [4] S. ALURU AND F. E. SEVILGEN, *Parallel domain decomposition and load balancing using space-filling curves*, in Proceedings of the 4th IEEE Conference on High Performance Computing, 1997.
- [5] M. BERN, D. EPPSTEIN, AND J. GILBERT, *Provably good mesh generation*, in Proceedings of 31st Symposium on Foundation of Computer Science, 1990, pp. 231–241.
- [6] M. BERN, D. EPPSTEIN, AND S. TENG, *Parallel construction of quadtrees and quality triangulations*, International Journal of Computational Geometry and Applications, 9 (1999), pp. 517–532.
- [7] D. K. BLANDFORD, G. E. BLELLOCH, AND C. KADOW, *Compact parallel Delaunay tetrahedralization*. Submitted to SC2005.
- [8] A. CALDER, B. C. CURTIS, J. DURSI, B. FRYXELL, G. HENRY, P. MACNEICE, K. OLSON, P. RICKER, R. ROSNER, F. TIMMES, H. TUFO, J. TRURAN, AND M. ZINGALE, *High-performance reactive fluid flow simulations using adaptive mesh refinement on thousands of processors*, in Proceedings of SC2000, 2000. Gordon Bell Award for Special Achievement.
- [9] P. M. CAMPBELL, K. D. DEVINE, J. E. FLAHERTY, L. G. GERVASIO, AND J. D. TERESCO, *Dyanmic octree load balancing using space-filling curves*, Tech. Report CS-03-01, Department of Computer Science, Williams College, jan 2003.
- [10] N. CHRISOCHOIDES, *Parallel mesh generation*, in Numerical Solution of Partial Differential Equations on Parallel Computers, M. Bruaset, P. Bjorstad, and A. Tveit, eds., Springer, 2005. to appear.
- [11] N. CHRISOCHOIDES AND D. NAVE, *Parallel Delaunay mesh generation kernel*, International Journal for Numerical Methods in Engineering, 1 (2001).
- [12] ———, *Parallel Delaunay mesh generation kernel*, Int. J. Num. Methods in Engineering, 58 (2003), pp. 161–176.
- [13] H. L. DE COUGNY AND M. S. SHEPHARD, *Parallel refinement and coarsening of tetrahedral meshes*, International Journal for Numerical Methods in Engineering, 46 (1999), pp. 1101–1125.
- [14] J. M. DENNIS, *Partitioning with space-filling curves on the cubed-sphere*, in Proceedings of Workshop on Massively Parallel Processing at IPDPS’03, Nice, France, 2003.
- [15] H. C. EDWARDS AND J. C. BROWNE, *Scalable dynamic distributed arrays and its application to parallel hp adaptive finite element code*, in Proceedings of POOMA ’96, Santa Fe, New Mexico, 1996.
- [16] J. E. FLAHERTY, R. LOY, C. OZTURAN, M. S. SHEPHARD, B. K. SZYMANSKI, J. D. TERESCO, AND L. H. ZIANTZ, *Parallel structures and dynamic load balancing for adaptive finite element computation*, Applied Numerical Mathematics, (1998).
- [17] CLEMENS KADOW AND NOEL J. WALKINGTON, *Adaptive dynamic projection-based partitioning for parallel Delaunay mesh generation and refinement*, in SIAM Conference on Parallel Processing for Scientific Computing, San Francisco, CA, Feb 2004.
- [18] G. KARYPIS AND V. KUMAR, *A coarse-grain parallel formulation of multi-level k-way graph partitioning algorithm*, in 8th SIAM Conference on Parallel Processing for Scientific Computing, 1997.
- [19] D. J. KERBYSON, H. J. ALME, A. HOISIE, F. PETRINI, H. J. WASSERMAN, AND M. GITTINGS, *Predictive performance and scalability modeling of a large-scale application*, in SC2001, 2001.
- [20] E. KIM, J. BIELAK, AND O. GHATTAS, *Large-scale northridge earthquake simluation using octree-based multiresolution mesh method*, in Proceedings of the 16th ASCE Engineering Mechanics Conference, Seattle, Washington, July 2003.
- [21] H. MAGISTRALE, S. DAY, R. CLAYTON, AND R. GRAVES, *The SCEC Southern California reference three-dimensional seismic velocity model version 2*, Bulletin of the Seismological Society of America, (2000).
- [22] S. A. MITCHELL AND S. A. VAVASIS, *Quality mesh generation in three dimensions*, in Proceedings of the 8th ACM Symposium on Computational Geometry, 1992, pp. 212–221.
- [23] A. K. PATRA, A. LASZLOFFY, AND J. LONG, *Data structures and load balancing for parallel adaptive hp finite-element methods*, International Journal on Computers and Mathematics with Applications, 46 (2003), pp. 105–123.
- [24] HANAN SAMET, *Applications of Spatial Data Structures: Computer Graphics, Image Processing and GIS*, Addison-Wesley Publishing Company, 1990.
- [25] R. SCHNEIDERS, *An algorithm for the generation of hexahedral element meshes based on an octree technique*, in Proceedings of 6th International Meshing Roundtable, 1997.
- [26] M. S. SHEPHARD AND M. K. GEORGES, *Automatic three-dimensional mesh generation by the finite octree technique*, International Journal for Numerical Methods in Engineering, 32 (1991).
- [27] J. P. SINGH, C. HOLT, J. L. HENNESSY, AND ANOOP GUPTA, *A parallel adaptive fast multipole method*, in SC1993, Portland, OR, 1993.
- [28] X. SUN AND L. NI, *Scalable problems and memory-bounded speedup*, Journal of Parallel and Distributed Computing, 19 (1993), pp. 27 – 37.
- [29] T. TU AND D. R. O’HALLARON, *A computational database system for generating unstructured hexahedral meshes with billions of elements*, in SC2004, Pittsburgh, PA, Nov. 2004.
- [30] T. TU, D. R. O’HALLARON, AND J. LOPEZ, *Etree – a database-oriented method for generating large octree meshes*, in Proceedings of the Eleventh International Meshing Roundtable, Ithaca, NY, Sept. 2002, pp. 127– 138. Also in Engineering with Computers (2004) 20:117–128.
- [31] T. TU, H. YU, L. RAMIREZ-GUZMAN, J. BIELAK, O. GHATTAS, K. MA, AND D. R. O’HALLARON, *From physical model to scientific understanding: A scalable online end-to-end supercomputing approach*. Working paper.
- [32] M. S. WARREN AND J. K. SALMON, *A parallel hashed oct-tree n-body algorithm*, in SC1993, Portland, OR, 1993.
- [33] A. WISSINK, R. HORNING, S. KOHN, S. SMITH, AND N. ELLIOTT, *Large scale parallel structured AMR calculations using the SAMRAI framework*, in SC2001, Denver, CO, 2001.
- [34] L. YING, G. BIROS, D. ZORIN, AND H. LANGSTON, *A new parallel kernel-independent fast multipole method*, in SC2003, Phoenix, AZ, 2003.
- [35] D. P. YOUNG, R. G. MELVIN, M. B. BIETERMAN, F. T. JOHNSON, S. S. SAMANT, AND J. E. BUSSOLETTI, *A locally refined rectangular grid finite element: Application to computational fluid dynamics and computational physics*, Journal of

Computational Physics, 92 (1991), pp. 1–66.