

Load balancing fictions, falsehoods and fallacies [☆]

Bruce Hendrickson

Parallel Computing Sciences Department, Sandia National Labs, Albuquerque, NM 87185-1110, USA

Received 1 June 2000; accepted 1 August 2000

Abstract

Effective use of a parallel computer requires that a calculation be carefully divided among the processors. This *load balancing* problem appears in many guises and has been a fervent area of research for the past decade or more. Although great progress has been made, and useful software tools developed, a number of challenges remain. It is the conviction of the author that these challenges will be easier to address if we first come to terms with some significant shortcomings in our current perspectives. This paper tries to identify several areas in which the prevailing point of view is either mistaken or insufficient. The goal is to motivate new ideas and directions for this important field. © 2000 Elsevier Science Inc. All rights reserved.

Keywords: Load balancing; Graph partitioning; Parallel computer

1. Introduction

Among the factors contributing to the success of parallel computing in the 90s, algorithms and software for load balancing have played a prominent role. Efficient use of a parallel computer requires that the workload be evenly distributed among processors while the amount of inter-processor communication is kept small. This basic challenge presents itself differently in different applications, but the underlying issues are the same. It is quite common to recast this load balancing problem in terms of graphs. Vertices of the graph represent indivisible computations, while edges encode data dependencies between computations. Load balancing is then commonly phrased as a graph partitioning problem – divide the vertices into sets of equal work while cutting as few dependencies (edges) as possible.

Graph partitioning software targeted at parallel computing began to appear in the early 90s and has since become a cottage industry. Among the currently available alternatives are Chaco [1], METIS [2], Jostle [3], PARTY [4] and SCOTCH [5]. Around the same time, complex scientific computing applications began to exhibit excellent performance on parallel machines. A clear conclusion is that the partitioning tools enabled the successful parallelizations.

Although there is some truth in this assertion, there is also a significant degree of falsehood. As we will detail below in Section 2, the standard graph partitioning approach has serious

[☆] This work was supported by the Applied Mathematical Sciences program, US DOE, Office of Energy Research, and was performed at Sandia National Labs, operated for the US DOE under contract No. DE-AC04-94AL85000.

E-mail address: bahendr@cs.sandia.gov (B. Hendrickson).

other side may need it, edge cuts tend to over-estimate the true volume of communication. A simple corollary is that reducing edge cuts might not reduce the true communication volume. This behavior was observed in [7]. But unfortunately, the graph partitioning software in wide use today fails to minimize a correct metric of communication volume.

Two alternative approaches are possible. One is to keep the same graph description, but to explicitly try to minimize the true communication volume. For an unweighted graph, this would merely be the number of vertices which have neighbors in another partition. This approach has been adopted in recent versions of METIS [8]. Unfortunately, this approach is not well suited to problems in which the amount of data associated with the different edges varies. In this case, a more accurate and elegant approach is the hypergraph model of Çatalyürek and Aykanat [9,10]. In this model, the set of vertices which generate and consume a particular value are all connected by a single *hyperedge*. If that hyperedge is split between two processors then a single communication cost is incurred. Each hyperedge can be assigned a different weight if the communication costs are non-uniform. In this way the total number (or weight) of cut hyperedges exactly equals the volume of communication.

Not only have the graph partitioning tools been minimizing an inaccurate measure of communication volume, the true cost of a communication operation is not well modeled by volume. For a distributed memory, message passing computer, the cost of a single message is usually well modeled by two terms: a startup or latency cost which is independent of message size and a bandwidth term which is proportional to the size of the message. Although graph partitioning models try to minimize the latter term, they do nothing for the former. If messages are of only modest size, the latency cost can dominate.

A second effect which is not well captured by standard graph partitioning approaches is the effect of communication congestion. Typically in scientific computations, many messages are simultaneously competing for network resources. Messages which use independent paths through the network can proceed at the same time, but those that compete for paths must take turns. Thus, the pattern of collective communication can have a significant impact on the communication cost. Current graph partitioners fail to accurately deal with this issue.

Even if latency and congestion considerations are included, graph partitioning models invariably try to minimize some version of the total communication cost. But the determining factor in parallel performance is the slowest processor. Graph partitioning generally produces sets with widely varying communication loads. A better objective would be to minimize the maximum communication load, or perhaps to reduce the communication on the most computationally burdened processors. Current partitioning formulations and implementations are not well suited for these objectives.

In summary, current partitioning tools try to optimize an objective that is far removed from the true cost incurred by a parallel computation. An obvious question is why they have proved so successful at enabling large calculations. The reason is that the vast majority of the applications using graph partitioners come from differential equations and so the underlying graph is a mesh. This has several implications. First, the geometric properties of meshes ensure that good partitions exist. If a mesh in d dimensions has n vertices it will have separators of size $n^{(d-1)/d}$ [11]. This ensures that for sufficiently large problems, the ratio of communication to computation will be small, so any reasonable partition should lead to good parallel performance. Second, the vertices associated with computational meshes typically have a bounded number of neighbors. This limits the harm caused by the edge cut approximation. Finally, meshes are generally fairly homogeneous, so the various sets are similar. This limits the communication inhomogeneity possible in these applications.

Unfortunately, other applications are less forgiving. Matrices arising from interior point methods, the singular value decomposition for latent semantic indexing, and other applications

are generally much less structured than mesh matrices. In particular, the variation in the number of edges associated with a vertex can be quite large. Çatalyürek and Aykanat find that for such matrices the hypergraph model can reduce communication by 30–40% [10]. The corresponding percentage improvement for meshes is only in the single digits [12]. But without the guarantee of good partitions, these highly unstructured matrices are liable to require more communication than meshes, and so the partition quality will more directly translate into increased runtime. So although our flawed models have been sufficient for parallelizing the solution of differential equations, future applications will require much more careful models and approaches.

3. Fiction II: simple graphs are sufficient

As described above, the standard manner in which graph partitioning is applied to load balancing involves an undirected graph with weighted vertices and edges. In addition to the problems discussed in Section 2 with accurately representing communication cost in this model, these simple graphs also suffer from a lack of expressibility. Many important classes of computations cannot be accurately described by this model. Examples include the following.

Multi-stage calculations. Many applications consist of a sequence of different, time-consuming calculations. One important example is multi-physics simulations which consist of interleaved computations of, for example, fluids and structures. More mundane examples include the application of a matrix and a preconditioner during an iterative solver, or the different grid calculations in a multi-grid scheme. Although each individual stage of such calculations might be describable by a simple graph, the union of stages often cannot be. This is particularly true if there is a barrier between the different stages, so that overall load balance is only achieved if each individual stage is balanced. When a single stage dominates the runtime, a standard graph may be sufficient, but the need to balance several stages is an objective which cannot be described this way.

Complex constraints. When calculations consist of multiple stages, there are often complicated couplings between them. For instance, the different grids comprising a multigrid solver are related. In some (but not all) cases, the coarse-grid vertices are a subset of the fine-grid vertices. Mechanisms to express these relationships to the partitioner are needed. Another important example arises in finite element calculations where some computations occur on the nodes and others on the elements. A partition which balances each of these computations would be desirable, and the intimate ties between nodes and elements needs to be encoded. For mesh-based applications, the DRAMA interface provides a uniquely rich model for this problem [13]. Unfortunately, current graph partitioning tools are not yet up to this complexity.

Unsymmetric dependencies. In the standard graph model, if a vertex i needs data from a vertex j , then j also needs data from i . Not all dependencies exhibit this symmetry. For instance, matrix–vector multiplication leads to symmetric dependencies only if the matrix is structurally symmetric. Although many methodologies for solving differential equations lead to structurally symmetric matrices, other applications do not. Important examples include latent semantic indexing, least squares problems for data analysis and interior point methods for optimization. The standard graph model does not describe these kinds of problems well.

Differing inputs and outputs. When the set of dependencies is symmetric, then the set of inputs to a calculation is equivalent to the set of outputs. But when the dependencies are non-symmetric the inputs and outputs can differ. For instance, consider the problem of matrix–vector multiplication for non-square matrices. In this case, the set of inputs is of different size than the set of outputs. For such examples, a model is needed which allows the inputs and outputs to be handled separately. The standard graph model is not able to do this.

To address these deficiencies, several alternative partitioning models have been proposed in recent years. One of these is the hypergraph model of Çatalyürek and Aykanat introduced in Section 2. In addition to correctly accounting for communication volume, this model can encode unsymmetric dependencies and, to a certain extent, differing inputs and outputs.

Closely related is the bipartite graph model of Hendrickson and Kolda [14,15]. In this model, inputs are designated as one set of vertices and outputs as a distinct set of vertices. An edge connects an input to an output that depends upon it. This model was devised specifically to handle the problems of unsymmetric dependencies and differing inputs and outputs. It can also encode certain limited classes of multi-stage calculations like two matrices acting on a vector. But it still minimizes the non-optimal metric of edge cuts.

A more general approach than the bipartite model is the multi-constraint, multi-objective partitioning approach of Karypis et al. [16,17]. This model was designed to partition multi-stage calculations, but its generality allows for other applications. It begins with a standard graph model, but augments it in two ways. First, instead of having a single weight for each vertex, vertices are given a vector of weights. A partition is only considered balanced if each of the components of the weight vector is balanced. For instance, value k in the weight vector could reflect the computational cost associated with stage k of the computation. The partitioner would then try to find a single decomposition that partitions each stage evenly. The second augmentation involves edge weights. As with vertices, each edge is allowed to have several weights, reflecting the communication cost associated with different stages.

The generality of this model is appealing, but it leads to very challenging partitioning problems. In many cases, a more focused model may be easier to work with. However, in many ways even the multi-objective, multi-constraint model is not general enough. It cannot, for instance, address the problem of minimizing both the volume of communication and the number of messages. And the model still suffers from all the shortcomings of the edge-cut metric.

Despite the recent activity in alternative partitioning models, much work remains to be done. New approaches which address some of the remaining deficiencies are urgently needed.

4. Fiction III: partition quality is paramount

When advocating a new partitioning algorithm, researchers invariably argue their case by comparing their partition quality and (sometimes) runtime against existing tools. While these two criteria are important, they are far from comprehensive. As discussed in more detail in [18], desirable properties of a partitioner include the following.

Balance the load. This requirement may seem obvious, but it is not always simple. For example, the challenge of balancing multi-stage calculations was discussed in Section 3. More generally, how should the work load be measured? How important is precise balance?

Minimize communication cost. As argued in Section 2, communication cost is a somewhat elusive quantity. The most appropriate metric certainly depends upon the properties of the parallel computer, and probably on those of the application as well.

Run fast in parallel. When static partitioning can be performed as a serial preprocessing step, efficiency is not crucial. But when the partitioning is being performed on an expensive parallel machine instead of a cheap workstation, runtime can be critical. Time spent partitioning is time lost to the application, and can only be justified by a resulting improvement in application performance. This is particularly true for dynamic or adaptive calculations which need to be re-partitioned repeatedly.

Be incremental. The issue of incrementality is critical for dynamic problems, but has no counterpart in static partitioning. When an adaptive or dynamic calculation becomes unbalanced and invokes a repartitioner, the data is already distributed across processors. If the partitioner dictates a new decomposition which is quite different from the current one, then significant quantities of data will need to be transferred between processors. Thus, partitioners which adjust the decomposition in an incremental way are desirable. As an example of the importance of this issue, Touheed et al. [19] report results of several different partitioning algorithms for an adaptive flow calculation on 32 processors of an SGI Origin. In their experiments, the time for data migration was significantly larger than the time spent determining the new partition. Obviously, the relative importance of partition quality, runtime and incrementality depends upon the application. But it may be worthwhile to sacrifice partition quality to improve incrementality.

Be frugal with memory. As with runtime, the partitioner is competing with the application for finite space on the parallel machine. A partitioner which requires only modest memory will allow larger calculations to be performed. Again, the importance of this issue varies enormously between applications.

Support determination of communication pattern. Once a new decomposition has been determined and the data redistributed, the application still needs to figure out how to work with the new partition. Specifically, each processor needs to know what data to exchange with which other processors. In its most general form, this determination can be expensive. Some partitioners, generally those based upon geometric properties of the data, can greatly simplify this task.

Be easy to use. Although difficult to quantify, ease of use is a critical aspect of a good tool. Many, if not most, application developers would gladly trade some performance for simplicity. For static partitioners with their file interfaces, this is not a big issue. But for dynamic partitioners which are invoked via parallel subroutine calls, the design of the interface is very important.

The importance of these different issues varies between applications. But it is insufficient to focus our attention solely on only one or two criteria.

5. Fiction IV: existing tools solve the problem

A number of good partitioning and load balancing tools have been developed in recent years, and new packages continue to appear. Developers of such software (including the author) like to believe that their tools meet the needs of a wide range of parallel applications. Clearly, the usage of these tools supports this belief. But application developers have no choice but to use the available tools. In fact, the existing tools are far from ideal. Two shortcomings in particular are worth mentioning.

The first shortcoming is the inadequacy of the graph partitioning model which was discussed in Section 2. Although it has been a useful abstraction, load balancing is not identical to graph partitioning.

The second shortcoming is in the nature of the software engineering. Here, it is worth differentiating between static and dynamic tools. Most load balancing packages address the static partitioning problem, and run on a serial computer. They are designed to work as part of a preprocessing step in which a large computation is prepared for a parallel machine. The tools are generally invoked as stand-alone codes with file interfaces. This approach is fairly mature, and sufficient for many applications.

However, a growing number of applications are ill-suited to this approach and would benefit from more sophisticated software engineering. Parallel load balancing tools are clearly needed for adaptive calculations, but they are necessary in several other settings as well. For instance, if the

computational problem is generated on a parallel machine (e.g. by a parallel mesh generator) then the load balancing needs to be done in parallel too. Also, if the target parallel computer is heterogeneous, then the partitioning cannot easily be done in advance. Only at runtime will an application know what resources are available to it, and without this information it cannot properly balance the load.

Unfortunately, parallel or dynamic partitioning software is much less mature than its static counterpart. To a large degree this is merely a consequence of the added complexity associated with developing parallel algorithms and software. But several algorithmic and software issues make dynamic load balancing libraries inherently more difficult. First, as discussed in Section 4, the dynamic partitioning problem has to contend with the issue of incrementality. For the static problem, there are only two principle metrics: quality and runtime. And multilevel methods seem to cover this two-dimensional space quite well, providing high quality solutions in modest time. But the need for incrementality adds a third key metric to the dynamic problem. It seems unlikely that a single algorithm will be able to cover this full three-dimensional space adequately. So a good tool will need to provide a suite of partitioners, some with greater incrementality, some with higher quality, etc. It requires more work to build such a toolkit than to implement a single algorithm.

Another important difference between the static and the dynamic cases is that dynamic partitioners cannot be stand-alone codes with file interfaces. Instead, they will be invoked as subroutine libraries, which raises new challenges in software design. The input arguments to most existing parallel partitioners (e.g. [3,13,20]) include a graph in some format. The burden of constructing the graph in the specified format is placed upon the application developer. The exception to this rule is Zoltan [21], in which functions are passed across the interface instead of data. In our view, this object-oriented design has several advantages. It reduces memory requirements, allows for the partitioner to only extract the information it needs, and the tool can be modified without changing the interface. But most importantly, we believe it is easier to use.

In summary, the load balancing community has yet to embrace the software engineering techniques which have simplified the inter-operability of commercial software. We believe that new tools developed with this mindset will significantly improve the utility of load balancing software. This is particularly true for dynamic load balancing tools.

6. Fiction V: load balancing means finding the right partition

A good distribution of work among processors is the key to obtaining high performance on a parallel computer. Each processor should have useful work to do for the duration of the computation, and the overhead of interprocessor communication must be minimized. A natural, but flawed, conclusion is that the best way to parallelize an application is to find a good partitioning.

This conclusion is flawed for two reasons. First, as discussed in Section 4, dynamic and adaptive calculations require the partition to be adjusted on the fly. No single partition is adequate for the duration of the computation. The second, and more interesting, reason arises in the context of multi-stage calculations which were introduced in Section 3. The different stages in these calculations can have conflicting partitioning objectives. An optimal partition for one stage may be far from optimal for another.

As discussed in Section 3, several alternative partitioning models are able to combine two or more stages into a single partitioning problem. However, this may not be the best answer. Perhaps no single decomposition exists which enables good performance from all the stages. An alternative is to support multiple decompositions so that each stage can be performed efficiently,

independent of the partitioning demands of the other stages. The disadvantage of this approach is that data must be moved between the decompositions, incurring a cost in both time and memory. But in some applications, this approach enables complex calculations to run efficiently on many more processors than a single decomposition would allow. One such application is the crash simulation work of Plimpton et al. [22].

Low speed impacts (e.g. car crashes) are usually simulated with Lagrangian techniques. A mesh is constructed of the car in its native geometry. As the simulation proceeds and the car hits, for example, a telephone pole, the mesh distorts to follow the deformation of the bumper. As the bumper deforms, it eventually strikes the radiator. In the simulation, this is revealed when the mesh of the bumper contacts the mesh of the radiator. At this point, new forces need to be included in the simulation.

There are two dominant stages in these calculations. The first is a finite element analysis of the impacting bodies to reveal how they deform under stress. The second is the geometric search for contacts in the mesh. Unfortunately, these two stages have quite different decomposition needs. The finite element analysis is a prototypical example for which graph partitioning works well. The pattern of data dependencies is static and determined by the topological structure of the mesh. The contact detection stage changes every timestep, and data dependencies reflect geometric proximity. In the work by Plimpton et al., graph partitioning is used for the finite element stage and a geometric partitioner is used for contact detection. The resulting crash code was the first to scale well to hundreds and even thousands of processors and has enabled simulations of unprecedented scale and fidelity [23].

Several aspects of crash simulations make them good candidates for the multiple-decomposition approach. These applications are generally not memory intensive, so the duplication of data associated with the multiple decompositions is not a problem. Also, each stage of the computation is expensive, so the cost of the data transfer can be tolerated. And finally, the partitioning needs of the two stages are quite different, so single-decomposition approaches have exhibited only limited scalability. For multi-stage applications with these features, multiple decompositions should be considered.

7. Fiction VI: all the problems are solved

Within the load balancing community, there is a sense of satisfaction with the status quo. The complacency that this engenders makes this the biggest and most damaging fiction of all. As discussed above, there is a compelling need for more accurate and expressive models, new partitioning algorithms to address these models and better designed software tools and interfaces. In addition to these general needs, a number of interesting and important problems have not yet been adequately addressed. Among them are the following.

Partitioning for sparse solve on each subdomain. Some preconditioners based on domain decomposition involve a sparse, direct solver on each subdomain. An important example is the FETI class of preconditioners [24]. The work and memory required by each subdomain for such a calculation is quite difficult to predict in advance. It requires a much more detailed analysis than can be provided by merely summing vertex weights. Current graph partitioning models are unable to balance the work or memory required on each subdomain. New models and partitioning ideas would be very helpful here, and would also be applicable to the parallelization of sparse direct methods.

Partitioning for good aspect ratios. Another property of FETI preconditioners is that they work best when the individual subdomains have small aspect ratios (i.e. not be long and skinny) [25].

Recent work by Diekmann et al. [26,27] has shown an elegant method for encoding this objective into the standard graph model. Further work in this area would be welcome.

Partitioning for heterogeneous parallel architectures. Most new and emerging parallel architectures have heterogeneous networks. For instance, many machines are clusters of symmetric multi-processors. The communication properties within an SMP are different than those for communication between SMPs. The growth of commodity cluster computing is also leading to parallel machines with heterogeneity among the processors. Existing partitioning and load balancing tools fail to incorporate this information. Several important questions arise including how to model the architecture and how to adapt partitioning algorithms for them. Some interesting early work on these questions has been pursued by Teresco et al. [28], but much remains to be done.

Partitioning to minimize congestion. As discussed in Section 2, the true cost of a collective communication operation can be difficult to predict. In the not-too-distant past, vendors provided specialized high performance networks that efficiently routed complex sets of messages. But for the large number of current machines with commodity networking, network performance is less robust. In particular, if many messages are all competing to use a single wire, then performance will suffer. Methods for generating partitions in which this contention is avoided are needed. Some ideas along this line can be found in the work of Pellegrini [29] or Hendrickson et al. [30], but new insights are needed.

Despite the general feeling that load balancing is a mature area, there is a great need for new ideas and insights. The wide range of open or only imperfectly solved problems provides many opportunities for new research.

Acknowledgements

The ideas in this paper have been influenced by my collaborations with Rob Leland, Tammy Kolda and Karen Devine. I have also benefited from conversations with Alex Pothen, Steve Plimpton and George Karypis. In addition, I am indebted to Guy Lonsdale and the organizers of the Third DRAMA Workshop for an invitation to speak, which motivated the organization of this paper.

References

- [1] B. Hendrickson, R. Leland, The Chaco user's guide, version 2.0. Technical Report SAND95-2344, Sandia National Labs, Albuquerque, NM, 1995.
- [2] G. Karypis, V. Kumar, A fast and high quality multilevel scheme for partitioning irregular graphs, Technical Report 95-035, Department of Computer Science, University of Minnesota, Minneapolis, MN, 1995.
- [3] C. Walshaw, M. Cross, M. Everett, Mesh partitioning and load-balancing for distributed memory parallel systems, in: Proceedings of the Parallel and Distributed Computing for Computational Mechanics: Systems and Tools, Saxe-Coburg Publications, 1998.
- [4] R. Preis, R. Diekmann, The PARTY partitioning-library, user guide – version 1.1, Technical Report tr-rsfb-96-024, Department of Computer Science, University of Paderborn, Paderborn, Germany, 1996.
- [5] F. Pellegrini, SCOTCH 3.1 user's guide, Technical Report 1137-96, LaBRI, Université Bordeaux I, France, August 1996.
- [6] B. Hendrickson, Graph partitioning and parallelsolvers: Has the emperor no clothes? in: Solving Irregularly Structured Problems in Parallel, Irregular '98, Lecture Notes in Computer Science, vol. 1457, Springer, Berlin, 1998, pp. 218–225.
- [7] B. Hendrickson, T. Kolda, Graph partitioning models for parallel computing, *Parallel Comput.* 26 (12) (2000) 1519–1534.

- [8] G. Karypis, personal communication, 1998.
- [9] Ü.V. Çatalyürek, C. Aykanat, Decomposing irregularly sparse matrices for parallel matrix–vector multiplication, in: *Parallel Algorithms for Irregularly Structured Problems, Irregular '96, Lecture Notes in Computer Science*, vol. 1117, Springer, Berlin, 1996, pp. 75–86.
- [10] Ü. Çatalyürek, C. Aykanat, Hypergraph-partitioning based decomposition for parallel sparse-matrix vector multiplication, *IEEE Trans. Parallel Distrib. Syst.* 10 (7) (1999) 673–693.
- [11] S.-H. Teng, Points, spheres and separators: a unified geometric approach to graph partitioning, Ph.D. Thesis, Department of Computer Science, Carnegie Mellon University, Pittsburgh, PA, 1991.
- [12] C. Aykanat, personal communication, 1998.
- [13] <http://www.ccr1-nece.de/DRAMA>.
- [14] T.G. Kolda, Partitioning sparse rectangular matrices for parallel processing, in: *Solving Irregularly Structured Problems in Parallel, Irregular '98, Lecture Notes in Computer Science*, vol. 1457, Springer, Berlin, 1998, pp. 68–79.
- [15] B. Hendrickson, T. Kolda, Partitioning rectangular and structurally nonsymmetric sparse matrices for parallel processing, *SIAM J. Sci. Comput.* 21 (6) (2000) 2048–2072.
- [16] G. Karypis, V. Kumar, Multilevel algorithms for multi-constraint graph partitioning, Technical Report 98-019, Department of Computer Science, University of Minnesota, Minneapolis, MN, 1998.
- [17] K. Schloegel, G. Karypis, V. Kumar, A new algorithm for multi-objective graph partitioning, in: *Proceedings of the EuroPar '99, Lecture Notes in Computer Science*, Springer, Berlin, 1999.
- [18] B. Hendrickson, K. Devine, Dynamic load balancing in computational mechanics, *Comput. Meth. Appl. Mech. Engrg.* 184 (2–4) (2000) 485–500.
- [19] N. Touheed, P. Selwood, P.K. Jimack, M. Berzins, A comparison of some dynamic load-balancing algorithms for a parallel adaptive flow solver application, *Parallel Comput.* 26 (12) (2000) 1535–1554.
- [20] G. Karypis, V. Kumar, Parallel multilevel graph partitioning, Technical Report 95-036, Department of Computer Science, University of Minnesota, Minneapolis, MN, 1995.
- [21] E. Boman, K. Devine, B. Hendrickson, W. Mitchell, M. St. John, C. Vaughan, <http://www.cs.sandia.gov/Zoltan>.
- [22] S. Plimpton, S. Attaway, B. Hendrickson, J. Swegle, C. Vaughan, D. Gardner, Transient dynamics simulations: parallel algorithms for contact detection and smoothed particle hydrodynamics, *J. Parallel Distrib. Comput.* 50 (1998) 104–122 (previous version appeared in *Proc. Supercomputing '96*).
- [23] K. Brown, S. Attaway, S. Plimpton, B. Hendrickson, Parallel strategies for crash and impact simulations, *Comput. Meth. Appl. Mech. Engrg.* 184 (2–4) (2000) 375–390 (invited paper).
- [24] C. Farhat, F.X. Roux, An unconventional domain decomposition method for an efficient parallel solution of large-scale finite element systems, *SIAM J. Sci. Stat. Comp.* 13 (1992) 379–396.
- [25] C. Farhat, N. Maman, G. Brown, Mesh partitioning for implicit computation via domain decomposition: impact and optimization of the subdomain aspect ratio, *Int. J. Numer. Meth. Engrg.* 38 (1995) 989–1000.
- [26] R. Diekmann, R. Preis, F. Schlimbach, C. Walshaw, Aspect-ratio for mesh partitioning, in: D. Pritchard, J. Reeve (Eds.), *Proceedings of the Euro-Par '98, LNCS 1470*, Springer, Berlin, 1998, pp. 347–351.
- [27] R. Diekmann, R. Preis, F. Schlimbach, C. Walshaw, Shape-optimized mesh partitioning and load balancing for adaptive FEM, *Parallel Comput.* 184 (2000) 269–285.
- [28] J.D. Teresco, M.W. Beall, J.E. Flaherty, M.S. Shephard, A hierarchical partition model for adaptive finite element computation, Technical Report, Rensselaer Polytechnic Institute, Troy, NY, 1998, *Comput. Meth. Appl. Mech. Eng.* 184 (2000) 269–285.
- [29] F. Pellegrini, Static mapping by dual recursive bipartitioning of process and architecture graphs, in: *Proceedings of the Scalable High Performance Computational Conference*, IEEE Press, New York, 1994, pp. 486–493.
- [30] B. Hendrickson, R. Leland, R. Van Driessche, Skewed graph partitioning, in: *Proceedings of the Eighth SIAM Conference on Parallel Processing for Scientific Computing*, SIAM, Philadelphia, PA, 1997.