

Automatic Dynamic Load Balancing for a Crack Propagation Application

Gengbin Zheng[†], Michael S. Breitenfeld[‡], Hari Govind[†], Philippe Geubelle[‡], Laxmikant V. Kalé^{†*}

[†]Department of Compute Science, University of Illinois at Urbana-Champaign

[‡]Department of Aerospace Engineering, University of Illinois at Urbana-Champaign

*Corresponding Author: kale@uiuc.edu

Abstract— Automatic, adaptive load balancing is essential for handling load imbalance that may occur during parallel finite element simulations involving mesh adaptivity, nonlinear material behavior and other localized effects. This paper demonstrates the successful application of a measurement-based dynamic load balancing concept to the finite element analysis of elasto-plastic wave propagation and dynamic fracture events. The simulations are performed with the aid of a parallel framework for unstructured meshes called ParFUM, which is based on Charm++ and Adaptive MPI (AMPI) and involves migratable user-level threads. The performance was analyzed using Projections, a performance analysis and post factum visualization tool. The bottlenecks to scalability are identified and eliminated using a variety of strategies resulting in performance gains ranging from moderate to highly significant.

I. INTRODUCTION

Researchers in the field of structural mechanics have often turned to parallel finite element modeling to model physical phenomena with more detail, sophistication, and accuracy. While parallel computing can provide large amounts of computational power, developing parallel software requires substantial efforts to leverage parallel computers efficiently.

Among the challenges associated with the parallelization of finite element codes, achieving load balance is the key to scaling a dynamic application to a large number of processors. This is especially true for dynamic structural mechanics codes where simulations involve rapidly evolving geometry and physics, often resulting in a load imbalance between processors. As a result of this load imbalance, the application has to run at the speed of the slowest processor with deteriorated performance. Solving load imbalance has triggered various research activities in load balancing techniques [1], [2], [3], [4]. Dynamic load balancing attempts to solve the load balance problem at run-time according to the most up-to-date load situation.

Dynamic load balancing is a challenging software design issue and generally creates a burden for the application developers. For example, a computational analyst working on computational fracture mechanics must include the mechanism to inform the decision-making module concerning load balance the estimated CPU load and the communication structure. In addition, once load imbalance is detected and data migration is requested, a developer has to write complicated code for moving data across processors. The ideal load balancing framework should hide the details of load balancing so that the

application developer can concentrate on modeling the physics of the problem.

In this paper, we present an automatic load balancing method and its application in wave propagation and dynamic crack propagation applications. The parallelization model used in this application is the processor virtualization supported by migratable MPI threads. The application runs on a large number of MPI threads (that exceeds the actual physical number of processors), allowing to perform run-time load balancing by migrating MPI threads. The MPI run-time system automatically collects load information from the execution of the application. Based on this instrumented load data, the run-time module makes decisions on migrating MPI threads from heavily loaded processors to underloaded ones. This approach thus requires minimal efforts from the application developer.

The remainder of the paper is organized as follows: Section II presents the finite element formulation used in this paper, with emphasis on the viscoplastic model and cohesive finite element scheme adopted here to model the dynamic propagation of a crack in a ductile medium. Section III describes the parallelization method and programming environment used to implement the structural mechanics application, while Section IV describes ParFUM, the high-level domain specific library introduced to help developer with the parallelization aspect of the application. ParFUM is based on the CHARM++ load balancing framework summarized in Section V. Sections VI and VII respectively describe the performance of parallel adaptive finite element simulations of an elasto-plastic wave propagation problem and of a dynamic fracture event. Section VIII discusses related work in load balancing research. Finally, Section IX concludes with some future plans.

II. COHESIVE FINITE ELEMENT MODEL OF FRACTURE

To simulate the spontaneous initiation and propagation of a crack in a discretized domain, we use an explicit cohesive-volumetric finite element (CVFE) scheme [5], [6], [7]. As its name indicates, the scheme relies on a combination of volumetric elements used to capture the constitutive response of the continuum medium, and of cohesive interfacial elements used to model the failure process taking place in the vicinity of the advancing crack front. The CVFE concept is illustrated in Figure 1, which presents two 4-node tetrahedral volumetric elements tied together by a 6-node cohesive element shown in

its deformed configuration, as the adjacent nodes are initially superposed and the cohesive element has no volume.

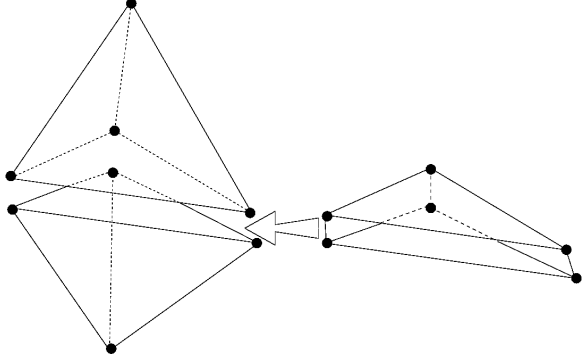


Fig. 1. Two 4-node tetrahedral volumetric elements linked by a 6-node cohesive element.

In the present study, the mechanical response of the cohesive elements is described by the bilinear traction-separation law illustrated in Figure 2 for the case of tensile (Mode I) failure. After an initial stiffening (rising) phase, the cohesive traction T_n reaches a maximum corresponding to the failure strength σ_{max} of the material, followed by a downward phase that represents the progressive failure of the material. Once the critical value Δ_{nc} of the displacement jump is reached, no more traction is exerted across the cohesive interface and a traction-free surface (i.e., a crack) is created in the discretized domain. The emphasis of the dynamic fracture study summarized hereafter is on the simulation of purely mode I failure, although cohesive models have also been proposed for the simulation of mixed-mode fracture events. Also illustrated in Figure 2 is an unloading and reloading path followed by the cohesive traction during an unloading event taking place while the material fails.

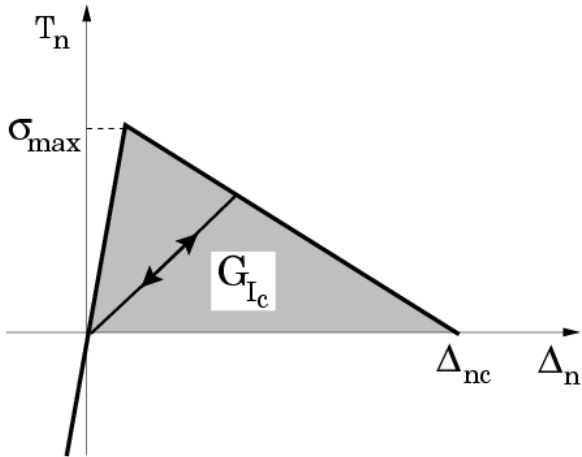


Fig. 2. Bilinear traction-separation law for mode I failure modeling.

The finite element formulation of the CVFE scheme is derived from the following form of the principle of virtual work:

$$\int_V (\rho \ddot{u}_i \delta u_i + S_{ij} \delta E_{ij}) dV = \int_{S_T} T_i^{ex} \delta u_i dS_T + \int_{S_c} T_i \delta \Delta_i dS_c, \quad (1)$$

where the left-hand-side corresponds to the virtual work done by the inertial forces ($\rho \ddot{u}_i$) and the internal stresses (S_{ij}), and the right-hand side denotes the virtual work associated with the externally applied traction (T_i^{ex}) and cohesive traction (T_i) acting along their respective surfaces of application S_T and S_c . In Equation (1), ρ denotes the material density, u_i and E_{ij} are the displacement and strain fields, respectively, and Δ_i denotes the displacement jump across the cohesive surfaces. The implementation relies on an explicit time stepping scheme based on the central difference formulation [5]. A nonlinear kinematics description is used to capture the large deformation and rotation associated with the propagation of the crack. The strain measure used here is the Lagrangian strain tensor \mathbf{E} .

To complete the CVFE scheme, we need to model the constitutive response of the material, i.e., to describe the response of the volumetric elements. In the present study, we use an explicit elasto-visco-plastic update scheme, which is compatible with the nonlinear kinematic description and relies on the multiplicative decomposition of the deformation gradient \mathbf{F} into elastic and plastic parts as

$$\mathbf{F} = \mathbf{F}^e \mathbf{F}^p. \quad (2)$$

The update of the plastic component \mathbf{F}^p of the deformation gradient at the $(n+1)^{th}$ time step is obtained by

$$\mathbf{F}_{n+1}^p = exp \left[\sum_A \frac{\Delta \gamma}{\sqrt{2} \tilde{\sigma}} \left(\sigma^A - \frac{I_1^\sigma}{3} \right) \mathbf{N}^A \otimes \mathbf{N}^A \right] \bullet \mathbf{F}_n^p, \quad (3)$$

where \mathbf{N}^A ($A=1, 2, 3$) denote the Lagrangian axes defined in the initial configuration, $\Delta \gamma$ is the discretized plastic strain increment, I_1^σ is the first Cauchy stress invariant, and $\tilde{\sigma} = \sqrt{(\sigma' : \sigma')/2}$ is the effective stress, with σ' denoting the Cauchy stress deviator whose spectral decomposition is

$$\sigma' = \sum_A \left(\sigma^A - \frac{I_1^\sigma}{3} \right) \mathbf{N}^A \otimes \mathbf{N}^A. \quad (4)$$

The plastic strain increment is given by $\Delta \gamma = \Delta t \dot{\gamma}$, where the plastic strain rate is described in this study by the classical Persyna two-parameter model

$$\dot{\gamma} = \eta \left(\frac{f(\sigma)}{\sigma_Y} \right)^n, \quad (5)$$

in which n and η are material constants, σ_Y is the current yield stress, and $f(\sigma) = (\tilde{\sigma} - \sigma_Y)$ is the overstress. Strain hardening is captured by introducing a tangent modulus E_t relating the increment of the yield stress, $\Delta \sigma_Y$, to the plastic strain increment $\Delta \gamma$. Finally, the linear relation

$$\mathbf{S} = \mathbf{L} \mathbf{E} \quad (6)$$

between the second Piola-Kirchhoff stresses \mathbf{S} and the Lagrangian strains \mathbf{E} is used to describe the elastic response. Assuming material isotropy, the stiffness tensor \mathbf{L} is defined by the Young's modulus E and Poisson's ratio ν .

The main source of load imbalance comes from the very different computational costs associated with the elastic and visco-plastic constitutive updates. As long as the effective stress remains below a given level (chosen in this paper as 80% of the yield stress), only the elastic relation (6) is computed. Once this threshold is reached for the first time, the visco-plastic update is performed, which typically represents a doubling in the computational cost. As the crack propagates through the discretized domain, the load associated with each processor can be substantially heterogeneous, suggesting the need for the robust dynamic load balancing scheme described in Section V.

III. PARALLELIZATION WITH AMPI

The parallel program for simulation of fracture dynamics is written and parallelized using Adaptive MPI.

Adaptive MPI (AMPI) [8], [9] is an MPI implementation and extension based on CHARM++ [10] programming model. CHARM++ is a parallel C++ programming language that embodies the concept of *processor virtualization* [11]. This idea of processor virtualization is that the programmer decomposes the computation, without regard to the physical number of processors available, into a large number of logical work units and data units, which are encapsulated in *virtual processors* (VPs) [11]. The programmer leaves the assignment of VPs to physical processors to the run-time system, which incorporates intelligent optimization strategies and automatic runtime adaptation. These virtual processors themselves can be programmed using any programming paradigm: e.g. they can be organized as indexed collections of C++ objects that interact via asynchronous method invocations, as in CHARM++ [12]. Alternatively, they can be MPI virtual processors implemented as user-level, extremely lightweight threads (NOT to be confused with system level threads or Pthreads), that interact with each other via messages, as in AMPI (illustrated in Figure 3).

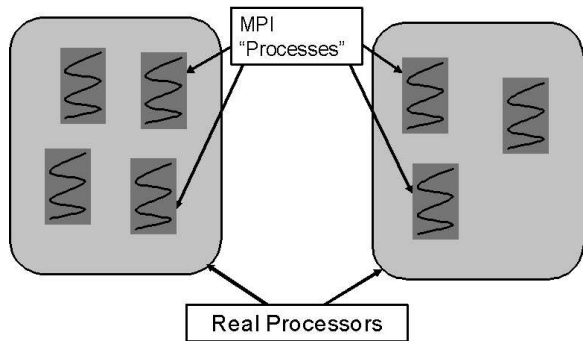


Fig. 3. Implementation of AMPI virtual processors

This idea of processor virtualization brings significant benefits to both parallel programming productivity and parallel

performance [13]. It empowers the run-time system to incorporate intelligent optimization strategies and automatic runtime adaptation. The following is a list of the benefits we have demonstrated in many projects.

Automatic load balancing: AMPI threads (the virtual processors) are decoupled from real processors, thus they are location independent and can migrate from processors to processors. Thread migration provides basic mechanism for load balancing — if some of the physical processors become overloaded, the run-time system can migrate a few of their AMPI threads to underloaded physical processors. The AMPI run-time system provides transparent support of message forwarding after thread migration.

Adaptive overlapping of communication and computation: If one of the AMPI threads is blocked on a receive, another AMPI thread on the same physical processor can run. This largely eliminates the need for the programmer to manually specify some static computation/communication overlapping, as is often required in MPI.

Optimized communication library support: Besides the communication optimization inherited from CHARM++, AMPI supports asynchronous, or non-blocking interfaces to collective communication operations. This allows the overlapping between time-consuming collective operations and other useful computation.

Better cache performance: A virtual processor handles a smaller set of data than a physical processor, so a virtual processor will have better memory locality. This blocking effect is the same method many *serial* cache optimizations employ, and AMPI programs get this benefit automatically.

Flexibility to run on an arbitrary number of processors: Since more than one VPs can be executed on one physical processor, AMPI is capable of running MPI programs on any arbitrary number of processors. This feature proves to be useful in application development and debugging phases.

In many applications, we have demonstrated that the processor virtualization does not incur much cost in parallel performance [13], due to low scheduling overheads of user-level threads. In fact, it often improves cache performance significantly because of its blocking effect.

CHARM++ and AMPI have been used as mature parallelization tools and run-time systems for a variety of real world applications for scalability [14], [15], [16], [17]. To further enhance programmer productivity, we have developed domain-specific frameworks on top of CHARM++ and AMPI to automate the parallelization process, which produces reusable libraries for parallel algorithms. A parallel framework for unstructured meshing called ParFUM that this work is based on is described in the next section.

IV. PARALLEL FRAMEWORK FOR UNSTRUCTURED MESHES

A wide variety of applications involve explicit computations on unstructured grids. As described earlier, one of the key objectives of this work is to create a flexible framework to

perform this type of simulations on parallel computing platforms. Parallel programming introduces several complications:

- Simply expressing a computation in parallel requires the use of either a specialized language such as HPF [18] or an additional library such as MPI.
- Parallel execution makes race conditions and nondeterministic execution possible. Some languages, such as HPF, have a simple lockstep control structure and are thus relatively immune to this problem; while in others, such as pthreads, they are more common.
- Computation and communication must be overlapped to achieve optimal performance. However, few languages provide good support for this overlap, and even simple static schemes can be painfully difficult to implement.
- Load imbalance can severely restrict performance, especially for dynamic applications. Automatic or application-independent load-balancing capabilities are rare (Section VIII).

Our approach to managing the complexity of parallel programming is based on a simple division of labor. In this approach, parallel programming specialists in computer science provide a simple but efficient *parallel framework* for the computation; while application specialists provide the numerics and physics. The parallel framework described hereafter abstracts away the details of its parallel implementation.

Since the parallel framework is application independent, it can be reused across multiple projects. This reuse amortizes the effort and time spent developing the framework and makes it feasible to invest in sophisticated capabilities such as adaptive computation and communication overlap and automatic measurement-based load balancing. Overall, this approach has proven quite effective, leveraging skills in both computer science and engineering to solve problems neither could solve independently.

A. ParFUM Framework

This section describes our parallel framework (called ParFUM) for performing explicit computations on unstructured grids. The framework has been used for finite-element computations, solving partial differential equations, computational fluid dynamics, and other problems.

The basic abstraction provided is very simple — the computational domain consists of an irregular mesh of nodes and elements. The elements are divided into partitions or chunks, normally using the graph partitioning library Metis [19], or ParMetis [20]. These chunks reside in AMPI migratable virtual processors, thereby taking advantage of run-time optimizations including dynamic load balancing. The chunks of meshes and AMPI virtual processors are then distributed across the processors of the parallel machine. There is normally at least one chunk per processor; and often even more. Nodes can be either private, adjacent to the elements of a single partition; or shared, adjacent to the elements of different partitions.

ParFUM application has two main subroutines: the *init* and the *driver*. The *init* subroutine executes only on processor 0 and is used to read the input mesh and physical data and

register it with the framework. The framework then partitions the mesh into as many regions as requested, each partition being a virtual processor. It then executes the *driver* routine on each virtual processor. This routine computes the solution over the local partition of the mesh.

The solution loop for most applications involves a calculation in which each node or element requires data from its neighboring entities. Thus entities on the boundary of a partition need data from entities on other partitions. ParFUM provides a flexible and scalable approach to meet an application's communication requirements. ParFUM adds local read-only copies of remote entities to the partition boundary. These read-only copies are referred to as *ghosts*. A single collective call to ParFUM allows the user to update all ghost entities with data from the original copies on neighboring partitions. This lets application code have effortless access to data from neighboring entities on other partitions. Since the definition of “neighboring” can vary from one application to another, ParFUM provides a flexible mechanism for generating ghost layers. For example, an application might consider two tetrahedra that share a face as neighbors. In another application, tetrahedra that share edges might be considered neighbors. ParFUM users can specify the type of ghost layer required by defining the “neighboring” relationship in the *init* routine and adding multiple layers of ghosts according to the neighboring relationship for applications that require them. In addition, the definition of “neighboring” can vary for different layers. User-specified ghost layers are automatically added after partitioning the input mesh provided during the *init* routine. ParFUM also updates the connectivity and adjacency information of a partition's entities to reflect the additional layers of ghosts. Thus ParFUM satisfies the communication needs of a wide range of applications by allowing the user to add arbitrary ghost layers. After the communication for ghost layers, each local partition is nearly self contained; a serial numerics routine can be run on the partition with only a minor modification to the boundary conditions.

With the above design, ParFUM framework enables straightforward conversion of serial codes into parallel applications. For example, in an explicit structural dynamics computation, each iteration of the time loop has the following structure:

- 1) Compute element strains based on nodal displacements.
- 2) Compute element stresses based on element strains.
- 3) Compute nodal forces based on element stresses.
- 4) Apply external boundary conditions.
- 5) Compute new nodal displacements based on Newtonian physics.

In a serial code, these operations apply over the entire mesh. However, since each operation is local, depending only on a node or element's immediate neighbors, we can partition the mesh and run the same code on each partition.

The only problem is ensuring that the boundary conditions of the different partitions match. The solution we choose is to duplicate the nodes along the boundary and then sum up the nodal forces during step 3, which amounts to this simple change:

- 1) Compute element strains based on nodal displacements.
- 2) Compute element stresses based on element strains.
- 3) Compute nodal forces based on element stresses.
- 4) Apply external *and internal* boundary conditions.
- 5) Compute new nodal displacements based on Newtonian physics.

For existing codes that have already parallelized with MPI, the conversion to ParFUM is even faster, thereby taking advantage of features including dynamic load balancing.

V. LOAD BALANCING FRAMEWORK

Many ParFUM applications involve simulations with dynamic geometry, and use adaptive techniques to solve highly irregular problems. In these applications, load balancing is required to achieve the desired high performance on large parallel machines. It is especially essential for applications where the amount of computation on a mesh partition can increase significantly as the number of elements comprising the partition increases with refinement and/or cohesive element insertion. It is also useful in applications where the computational load for subsets of elements varies over the duration of the simulation.

ParFUM directly utilizes the load balancing framework in CHARM++ and AMPI load balancing framework [3], [21]. The load balancing involves four distinct steps: (1) load evaluation; (2) load balancing initiation which determines when to start a new load balancing; (3) load balancing decision making and (4) task and data migration. These steps are automatic and require minimal effort from the developers.

CHARM++ load balancing framework adopts a unique measurement-based strategy for load evaluation. This scheme is based on the run-time instrumentation, which is feasible due to the *principle of persistence* that can be found in most physical simulations: the communication patterns between objects as well as the computational load of each of them tend to persist over time, even in the case of dynamic applications. This implies that the recent past behavior of a system can be used as a good predictor of the near future. The load instrumentation is fully automatic at runtime. During the execution of a ParFUM application, the run-time times the computation load for each object and records communication pattern into a load “database” on each processor.

The runtime then assesses the load database periodically and determines if load imbalance occurs. The load imbalance can be computed as:

$$\sigma = \frac{L_{max}}{L_{avg}} - 1, \quad (7)$$

where L_{max} is the maximum load across all processors, and L_{avg} is the average load of all the processors. Note that even when load imbalance occurs ($\sigma > 0$), it may not be profitable to start a new load balancing step due to the overhead of load balancing itself. In practice, a load imbalance threshold can be chosen based on a heuristic that the gain of the load balancing ($L_{max} - L_{avg}$) after the load balancing is at least greater than

the estimated cost of load balancing (C_{lb}). That is:

$$\sigma > \frac{C_{lb}}{L_{avg}}. \quad (8)$$

When load balancing is triggered, the load balancing decision module uses the load database to compute a new assignment of virtual processors to physical processors and informs the run-time to execute the migration decision.

A. Run-time Support for Thread Migration

In ParFUM applications, load balancing is achieved by migrating AMPI threads that host mesh partitions from overloaded processors to underloaded ones. When an AMPI thread migrates between processors, it must move all the associated data, including its stack and heap-allocated data. The CHARM++ runtime supports both fully automated thread migration and flexible user-controlled migration of data by additional helper functions.

In fully automatic mode, the AMPI run-time system automatically transfers a thread’s stack and heap data which are allocated by special memory allocator called *isomalloc* [9] in a manner similar to that of *PM2* [2]. It is portable on most platforms except for those where the *mmap* system call is unavailable. *Isomalloc* allocates data with a globally unique virtual address, reserving the same virtual space on all processors. With this mechanism, *isomalloc* data can be moved to a new processor without changing the address. This provides a clean way to move a thread’s stack and heap data to a new machine automatically. In this case, migration is transparent to the user code.

Alternatively, users can write their own helper functions to pack and unpack heap data for migrating an AMPI thread. This is useful when application developers wish to have more control in reducing the data volume by using application specific knowledge and/or by packing only variables that are live at the time of migration. The PUP (Pack/UnPack) library [22] was written to simplify this process and reduce the amount of code the developers have to write. The developers only need to write a single PUP routine to traverse the data structure and this routine is used for both packing and unpacking.

B. Load Balancing Strategies

In the step that makes the load balancing decision, the CHARM++ run-time assigns AMPI threads on physical processors, so as to minimize the maximum load (makespan) on the processors. This is known as the Makespan minimization problem, which has been shown as an *NP*-hard optimization problem [23]. However, many combinatorial algorithms have been developed that find a reasonably good approximate solution. CHARM++ load balancing framework provides a spectrum of simple to sophisticated heuristic-based load balancing algorithms, some of which are described in more details below:

- Greedy Strategy: This simple strategy organizes all the objects in decreasing order of their computation times. The algorithm repeatedly selects the heaviest un-assigned

object, and assigns it to the least loaded processor. This algorithm may lead to a large number of migrations. However, it works effectively in most cases.

- **Refinement Strategy:** The refinement strategy is an algorithm which improves the load balance by incrementally adjusting the existing object distribution, especially on highly loaded processors. The computational cost of this algorithm is low because only a subset of processors are examined. Further, this algorithm results in only a few objects being migrated, which makes it suitable for fine-tuning the load balance.
- **Metis-based Strategy:** This strategy uses the METIS graph partitioning library [24] to partition the object-communication graph. The objective of this strategy is to find a reasonable load balance, while minimizing the communication among processors.

CHARM++ load balancing framework also allows a developer to implement his own load balancing strategies based on heuristics specific to the target application (such as in NAMD [14] molecular simulation code).

Load balancing can be done in either centralized or distributed approach depending on how the load balancing decisions are made. In the centralized approach, one central processor makes the decisions globally. The load databases of all processors are collected to the central processor, which may incur high communication overhead and memory usage for the central processor. In the distributed approach, load balance decisions are made in distributed fashion. The load databases are only exchanged among neighboring processors. Due to the lack of the global information and aging of the load data, distributed load balancing inherently is difficult to achieve good load balance as quickly as the centralized approach. In this paper, we use a centralized load balancing strategy.

C. Agile Load Balancing

Applications with fast changing load requires frequent load balancing, which demands fast load balancing with minimal overhead. Normal load balancing strategies in CHARM++ occur in *synchronous* mode, as shown in Figure 4. At load balancing time, the application on each processor stops after it finishes its designated iterations and hands over the control to the load balancing framework to make load balancing decisions. The application can only resume when the load balancing step finishes and all AMPI threads migrate to the destination processors. In practice, this “stop and go” load balancing scheme is simple to implement, and has one important advantage — the AMPI thread migration happens under user control, so that a user can choose a convenient time for the thread migration to possibly minimize the migration data size. However, this scheme is not efficient due to the effect of the global barrier. It suffers from high overhead due to the fact that load balancing process on the central processor has to wait for the slowest processor to join load balancing, and thus wasting CPU cycles on other processors. This motivated us to develop an agile load balancing strategy that performs

asynchronous load balancing which allows overlapping of load balancing time and normal computation.

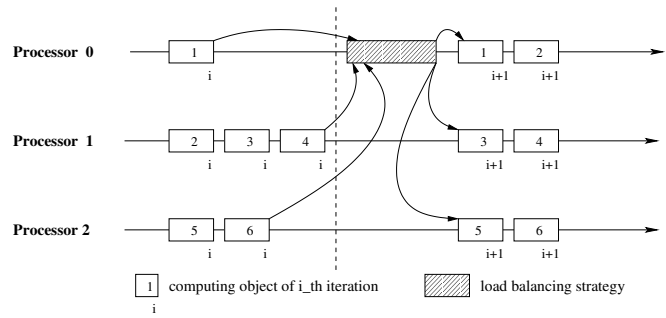


Fig. 4. Traditional synchronous load balancing

The asynchronous load balancing scheme fully takes advantage of CHARM++’s intelligent run-time support for concurrent compositionality [13] that allows dynamic overlapping of the execution of different composed modules in time and space. In the asynchronous scheme, load balancing process occurs concurrently or in the background of normal computation. When it is time for load balancing, each processor sends its load database to the central processor and continues its normal computation without waiting for load balancing to start. When a migration decision is calculated at the background on the central processor, the AMPI threads are instructed to migrate to their new processors in the middle of their computation.

There are a few advantages of asynchronous load balancing over the synchronous scheme. First, eliminating the global barrier helps reducing the idle time on faster processors which otherwise would have to wait for the slower processors to join the load balancing step. Second, it allows the overlapping of load balancing decision making time and computation in an application, which potentially could help improve the overall performance. Finally, each thread can have more flexible control on when to migrate to the designated processor. For example, a thread can choose to migrate when it is about to be idle, which potentially allows overlapping of the thread migration and computation of other threads.

Asynchronous load balancing however imposes a significant challenge to the thread migration in the AMPI run-time system. AMPI threads may migrate *at any time*, whenever they receive the migration notification. In practice, it is not trivial for an AMPI thread to migrate at any time due to the complex run-time state involved, for example when a thread is suspended in the middle of pending receives. In order to support any-time migration of AMPI threads, we extended the AMPI run-time to be able to transfer a complete runtime state associated with the AMPI threads including the pending receive requests and buffered messages for future receives. With the help of isomalloc stack and heap, AMPI threads can be migrated to a new processor transparently at any time without worrying about the scenario that the stack becomes invalid when the threads are resumed on a different processor. For AMPI threads with pending receives, incoming messages are redirected automatically to the destination processors by

the run-time system.

In the next two sections, we present two case studies of simulations to demonstrate the effectiveness of our load balancing strategies.

VI. CASE STUDY 1: ELASTO-PLASTIC WAVE PROPAGATION

The first application is the quasi-one-dimensional elasto-plastic wave propagation problem depicted in Figure 5. It consists of a rectangular bar of length $L = 10\text{ m}$ and cross-section $A = 1\text{ m}^2$. The bar is initially at rest and stress free. It is fixed at one end and subjected at the other end to an applied velocity V ramped linearly from 0 to 20 m/s over $.16\text{ ms}$ and held constant thereafter. The material properties are chosen as follows: yield stress $\sigma_Y = 480\text{ MPa}$, stiffness $E = 73\text{ GPa}$ and $E_t = 7.3\text{ GPa}$, exponent $n = 0.5$, fluidity $\eta = 10^{-6}/\text{s}$, Poisson's ratio $\nu = .33$ and density $\rho = 2800\text{ kg/m}^3$.

The applied velocity generates a one-dimensional stress wave that propagates through the bar and reflects from the fixed end. At every wave reflection, the stress level in the bar increases as the end of the bar is continuously pulled at a velocity V . During the initial stage of the dynamic event, the material response is elastic as the first stress wave travels through the bar at the dilatational wave speed $C_d = 6215\text{ m/s}$ with an amplitude

$$\sigma = \rho C_d V = 348\text{ MPa} < \sigma_Y. \quad (9)$$

After one reflection of the wave from the fixed end, the stress level in the bar exceeds the yield stress of the material and the material becomes plastic. A snapshot of the location of the elasto-plastic stress wave is shown in Figure 5. The computational overload associated with the plastic update routine (approximately a factor of two increase compared to the elastic case) leads to a significant dynamic load imbalance while the bar transforms from elastic to plastic. In these simulations, the plastic check and update subroutine is called upon when the equivalent stress level exceeds 80% of the yield stress.

The unstructured 400,000-element tetrahedral mesh that spans the bar is initially partitioned into chunks using METIS and these chunks are then mapped to the processors. During the simulation, the processors advance in lockstep with frequent synchronizing communications required by exchanging of boundary conditions, which may lead to bad performance when load imbalance occurs.

The simulation was run on Tungsten Xeon Linux cluster at the National Center for Supercomputing Applications (NCSA). This cluster is based on Dell PowerEdge 1750 servers, each with two Intel Xeon 3.2 GHz processors, running Red Hat Linux and Myrinet interconnect network. The test ran on 32 processors with 160 AMPI virtual processors. Figure 7 shows the results without load balancing in a *Projections* CPU utilization graph over a certain time interval. The figure was generated by *Projections* [25], which is a performance visualization and analysis tool associated with CHARM++ that supplies application-level visual and analytical performance

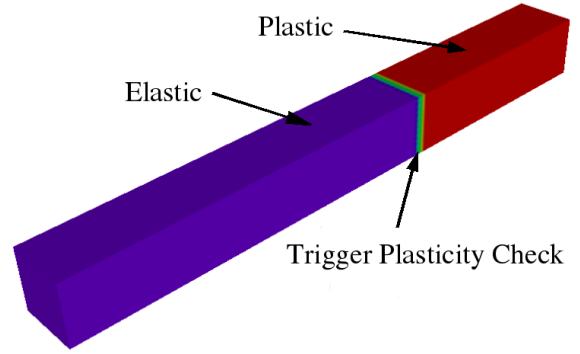


Fig. 5. Location of the traveling elasto-plastic wave at time $C_d t/L = 1.3$.

feedback. This utilization graph shows how the overall utilization changes as the wave propagates through the bar. The total run time was 177 seconds for this run.

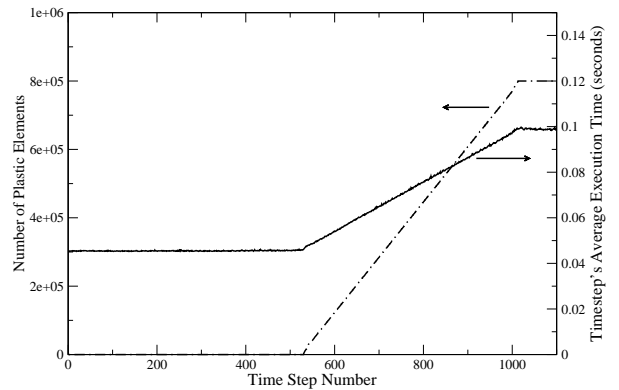


Fig. 6. Evolution of the number of plastic elements.

A separate interest, although not investigated further in this paper, is the period of initial load imbalance (as shown in Figure 7) caused by the quiet generation of subnormal numbers (floating-point numbers that are very close to zero) during the initial propagation of the elastic wave along the initially quiescent bar. This phenomenon is discussed by Lawlor *et al.* [26], who propose an approach to mitigate such performance effects caused by the inherent processor design. This paper is only concerned with the load imbalance associated with the transformation of the bar from elastic to plastic.

As indicated earlier, the load imbalance in this problem is highly transient, as elements at the wave front change from an elastic to a plastic state. In Figure 6, the effects of the plasticity calculations are clearly noticeable in terms of execution time which linearly ramps from the condition of fully elastic to fully plastic resulting in a doubling of the execution time. This leads to a load imbalance between the processor when the simulation is within this linearly ramping region. The load

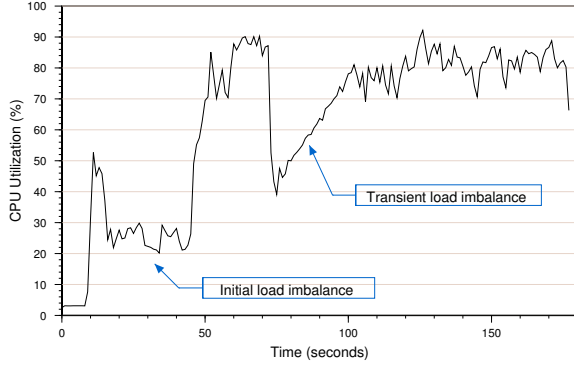


Fig. 7. CPU utilization graph without load balancing (Tungsten Xeon).

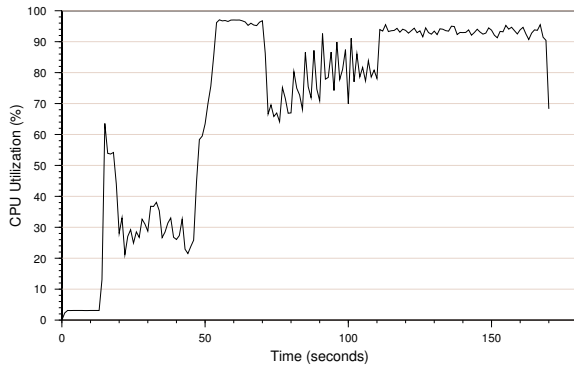


Fig. 8. CPU utilization graph with synchronous load balancing (Tungsten Xeon).

balancer here has to migrate these objects aggressively.

Even though we used a variety of methods and time frames, the problem was not considerably sped up by load balancing. The transition time was too fast for load balancer to significantly speed up the simulation. Also the period of imbalance is a very small portion compared to the total run time. Therefore, a performance improvement here necessitates that the overhead and delays associated with the invocation of the load balancer be minimal. Nevertheless, we managed to speed up the simulation by 7 seconds as shown in Figure 8. The time required for completion reduces to 170 seconds, which yields a 4 percent of overall improvement by the load balancing.

We repeated the same test on 32 processors of the SGI Altix (IA64) at NCSA with the same 160 AMPI virtual processors. Figure 9 shows the result without load balancing in the Projections utilization graph. The total execution time was 207 seconds and a more severe effect of subnormal numbers on this machine was observed in the first hundred seconds of execution time.

In the second run, we ran the same test with the greedy load balancing scheme described in the previous section. The result is shown in Figure 10 in the same utilization graph.

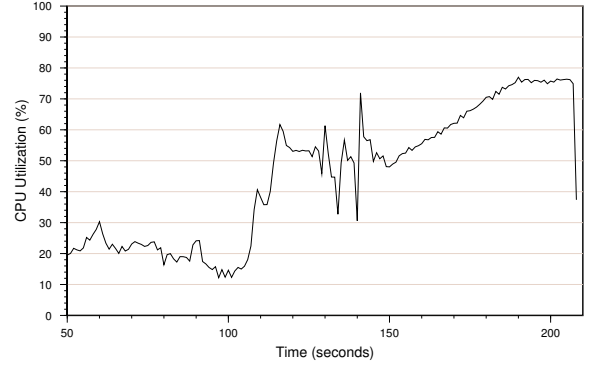


Fig. 9. CPU utilization graph without load balancing (SGI Altix).

The load balancing is invoked around time interval 130 in the figure. After the load balancing, the CPU utilization is slightly improved and the total execution time is now around 198 seconds.

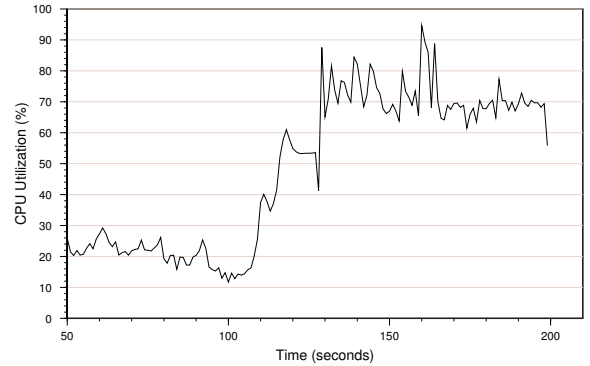


Fig. 10. CPU utilization graph with synchronous load balancing (SGI Altix).

Finally, we ran the same test with the same greedy algorithm in an asynchronous load balancing scheme (Section V-C). Asynchronous load balancing scheme avoids the stall of an application for load balancing and overlaps the computation with the load balancing and migration. The result is shown in Figure 11 in a utilization graph. It can be seen that, after load balancing, the overall CPU utilization was further improved and the total execution time is 187 seconds, which is a 20 second improvement.

VII. CASE STUDY 2: DYNAMIC FRACTURE

The second application involves a single edge notched fracture specimen of width $W = 5m$, height $H = 5m$, thickness $T = 1m$ and initial crack length $a_0 = 1m$, having a weakened plane starting at the crack tip and extending along the crack plane to the opposite edge of the specimen. The material properties used in this simulation are $\sigma_Y = 900 MPa$, $E = 210 GPa$, $E_t = 2.4 GPa$, $n = 0.5$, $\eta = 10^{-6}/s$, $\nu = .3$, and $\rho = 7850 kg/m^3$. A linearly ramped velocity of 0 to 1 m/s over 2.0 ms and held constant thereafter is

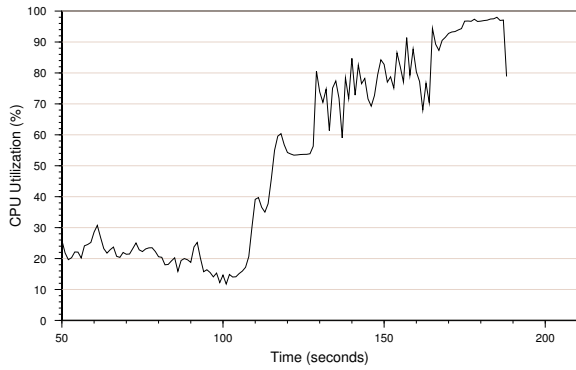


Fig. 11. CPU utilization graph with asynchronous load balancing (SGI Altix).

applied along the top and bottom surfaces of the specimen. A single layer of six-node cohesive elements are placed along the weakened interface, with the failure properties described by a critical crack opening displacement value $\Delta_{nc} = .8\text{ mm}$ and a cohesive failure strength $\sigma_{max} = 95\text{ MPa}$. The mesh consists of 91,292 cohesive elements along the interface plane and 4,198,134 linear strain tetrahedral elements. As the stress wave emanating from the top and bottom edges of the specimen reach the fracture plane, a region of high stress concentration is created around the initial crack tip. In that region, the equivalent stress exceeds the yield stress of the material leading to the creation of a plastic zone. As the stress level continues to build up in the vicinity of the crack front, the cohesive tractions along the fracture plane start to exceed the cohesive failure strength of the weakened plane and a crack starts to propagate rapidly along the fracture plane, surrounded by a plastic zone and leaving behind a plastic wake, as illustrated in Figure 12.

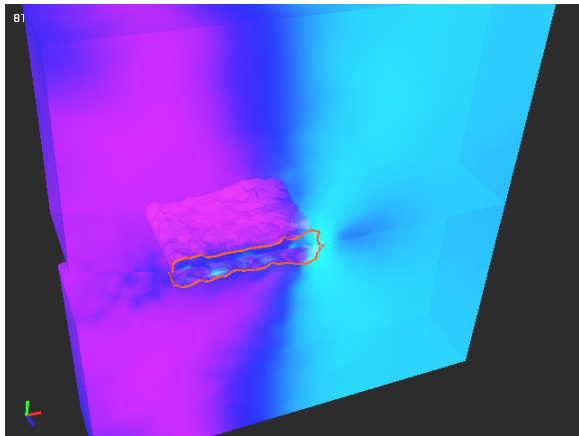


Fig. 12. Snapshot of the plastic zone surrounding the propagating planar crack at time $C_d t/a_0 = ???$. The iso-surfaces denote the extent of the region where the elements have exceeded the yield stress of the material.

This simulation was run on the Turing cluster at the University of Illinois of Urbana-Champaign. The cluster consists of 640 dual Apple G5 nodes connected with Myrinet network.

The simulation without load balancing took about 24 hours on 100 processors. The average processor utilization is shown in the bottom curve of Figure 13. The processor utilization is on the Y axis and the time on the X axis. It can be seen that around time interval 120, the CPU utilization dropped from around 85% to only about 42%. This is due to the advent of the elastic elements transitioning into plastic elements around the crack tip, leading to the beginning of load imbalance. In Figure 14 the number of plastic elements starts to increase dramatically as the crack starts to propagate along the interface. As more elastic elements turn plastic, the CPU utilization slowly increases. The load imbalance can also be easily seen in the CPU utilization graph over processors in Figure 15. While some of the processors have the CPU utilization as high as about 90%, some processors only have about 50% of the CPU utilization during the whole execution.

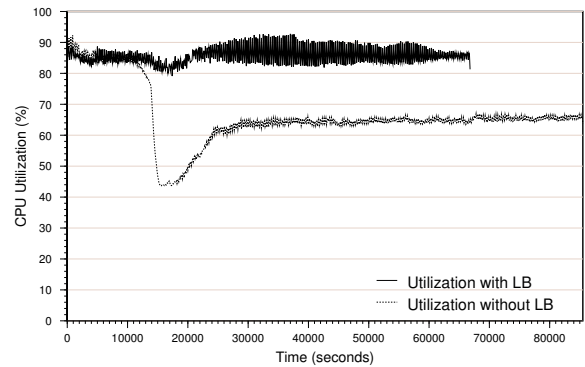


Fig. 13. CPU utilization graph over time without vs. with load balancing for the fracture problem shown in Figure 12 (Turing Apple Cluster).

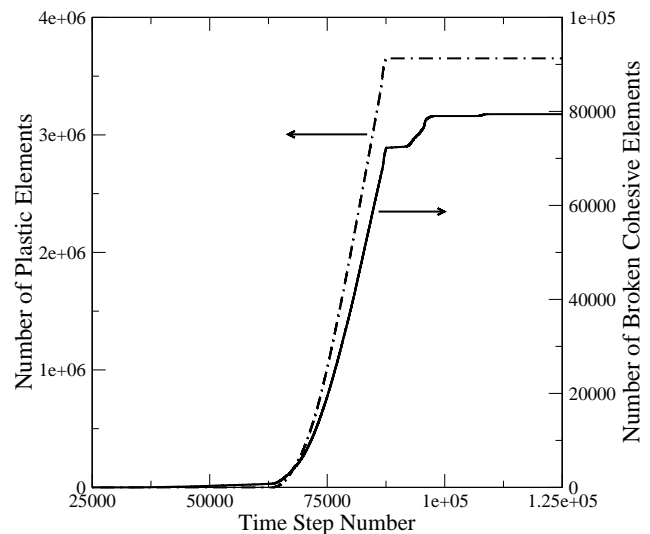


Fig. 14. Evolution of the number of plastic and broken cohesive elements.

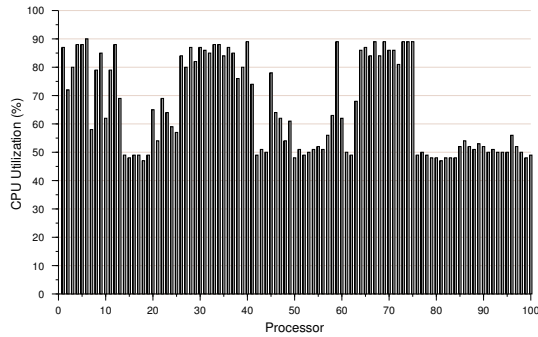


Fig. 15. CPU utilization across processors without load balancing (Turing Apple Cluster).

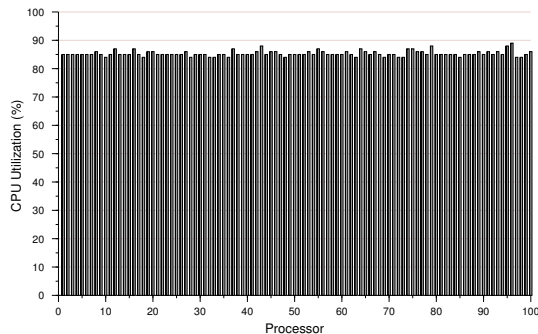


Fig. 16. CPU utilization across processors with load balancing (Turing Apple Cluster).

With the greedy load balancing strategy invoked every 500 time-steps, the simulation finished in 18 hours, a saving of nearly 6 hours or 25% over the same simulation with no load balancing. This increase is caused by the overall increased processor utilization, which can be seen in the upper curve of Figure 13. The peaks correspond to the times when the load balancer is activated during the simulation. There is an immediate improvement in the utilization when the load balancer is invoked. Then the performance slowly deteriorates as more elements become plastic. The next invocation tries to balance the load, all over again. Figure 16 further illustrates that load balance has been improved from Figure 15 in the view of the CPU utilization across processors. It can be seen that a CPU utilization of around 85% is achieved on all processors with negligible load variance.

VIII. RELATED WORK

The goal of our work is a generic load balancing framework that optimizes the load balance of the irregular and highly dynamic applications with an application independent interface, therefore we will focus our discussion in this section to those dynamic load balancing systems for parallel applications. In particular, we wish to distinguish our research using the following criteria:

- Supporting data migration. Migrating data has advantages over migrating “heavy-weight” processes which adds complexity to the runtime system.
- Generality. Load balancing methods are designed to be application independent. They can be used for a wide variety of applications.
- Automatic load estimation. The load balancing framework does not rely on application developer to provide application load information.
- Communication-aware load balancing. The framework takes communication into account explicitly rather than implicitly for example using domain specific knowledge. Communication pattern including multicast and communication volume are directly recorded into load balancing database for load balancing algorithms.
- Adaptive to execution environment. Take background load and non-migratable load into account.

Table I shows the comparison of CHARM++ load balancing framework to several other software systems that support dynamic load balancing. DRAMA [27] is designed specifically to support finite element applications. This specialization enables DRAMA to provide an application “independent” load balancing using its built-in cost functions for the category of applications. Zoltan [28], [29] does not make assumptions about applications’ data, and is designed to be general purpose load balancing library. However, it relies on application developers to provide cost function and communication graph. A recent system PREMA [30], [31] supports very similar idea of migratable objects, however, its load balancing method primarily focuses on task scheduling problem as in non-iterative applications. The Chombo [32] package has been developed by Lawrence Berkeley National Lab. It provides a set of tools including load balancing for implementing finite difference methods for the solution of partial differential equations on block-structured adaptively refined rectangular grids. It requires users to provide input for computational workload in real number for each box (box is a partition of mesh). Charm++ provides the most comprehensive features for load balancing. It is applicable to most scientific and engineering applications that present persistent computation (even dynamic). Charm++ load balancing is also capable of adapting to the change of background load [33].

IX. CONCLUSION

Dynamic and adaptive parallel load balancing is indispensable for handling load imbalance that may arise during a parallel simulation due to mesh adaptation, material nonlinearity, etc. This paper demonstrates the successful application of a measurement-based dynamic load balancing concept to the crack propagation problem, that uses cohesive elements. There are myriad of other problems where the same principle applies.

In the future we plan to enhance the performance of more complex ParFUM applications with the load balancing framework, which were previously considered unsolvable in reasonable amount of time, using ParFUM. Example applications include adaptive insertion and activation of cohesive elements

System Name	Data Migration	Generality	Explicit Comm.	Automatic Cost Est.	Adaptive
DRAMA	Yes	No	No	Yes	No
Zoltan	Yes	Yes	No	No	No
PREMA	Yes	No	No	No	No
Chombo	Yes	No	No	No	No
CHARM++	Yes	Yes	Yes	Yes	Yes

TABLE I
SOFTWARE SYSTEMS THAT SUPPORT DYNAMIC LOAD BALANCING

for dynamic fracture simulations, adaptive mesh adaptation. We also will explore using load balancing framework on very large parallel machines such as 64K processor Blue Gene/L.

ACKNOWLEDGEMENTS

This work was supported by the Center for the Simulation of Advanced Rockets under contract number B341494 by the U.S Department of Energy.

REFERENCES

- [1] Robert K. Brunner and Laxmikant V. Kalé. Handling application-induced load imbalance using parallel objects. In *Parallel and Distributed Computing for Symbolic and Irregular Applications*, pages 167–181. World Scientific Publishing, 2000.
- [2] Gabriel Antoniu, Luc Bouge, and Raymond Namyst. An efficient and transparent thread migration scheme in the PM^2 runtime system. In *Proc. 3rd Workshop on Runtime Systems for Parallel Programming (RTSPP) San Juan, Puerto Rico. Lecture Notes in Computer Science 1586*, pages 496–510. Springer-Verlag, April 1999.
- [3] Gengbin Zheng. *Achieving High Performance on Extremely Large Parallel Machines: Performance Prediction and Load Balancing*. PhD thesis, Department of Computer Science, University of Illinois at Urbana-Champaign, 2005.
- [4] Amnon Barak, Shai Geday, and Richard G. Wheeler. The mosaic distributed operating system. In *LNCS 672*. Springer, 1993.
- [5] X.-P. Xu and A. Needleman. Numerical simulation of fast crack growth in brittle solids. *Journal of the Mechanics and Physics of Solids*, 42:1397–1434, 1994.
- [6] G. T. Camacho and M. Ortiz. Computational modeling of impact damage in brittle materials. *Int. J. Solids Struct.*, 33:2899–2938, 1996.
- [7] P. H. Geubelle and J. Baylor. Impact-induced delamination of composites: a 2d simulation. *Composites B*, 29:589–602, 1998.
- [8] Chao Huang, Gengbin Zheng, Sameer Kumar, and Laxmikant V. Kalé. Performance evaluation of adaptive MPI. In *Proceedings of ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming 2006*, March 2006.
- [9] Chao Huang, Orion Lawlor, and L. V. Kalé. Adaptive MPI. In *Proceedings of the 16th International Workshop on Languages and Compilers for Parallel Computing (LCPC 2003)*, LNCS 2958, pages 306–322. College Station, Texas, October 2003.
- [10] L. V. Kale and Sanjeev Krishnan. Charm++: Parallel Programming with Message-Driven Objects. In Gregory V. Wilson and Paul Lu, editors, *Parallel Programming using C++*, pages 175–213. MIT Press, 1996.
- [11] Laxmikant V. Kalé. The virtualization model of parallel programming : Runtime optimizations and the state of art. In *LACSI 2002*, Albuquerque, October 2002.
- [12] Orion Sky Lawlor and L. V. Kalé. Supporting dynamic parallel object arrays. *Concurrency and Computation: Practice and Experience*, 15:371–393, 2003.
- [13] Laxmikant V. Kalé. Performance and productivity in parallel programming via processor virtualization. In *Proc. of the First Intl. Workshop on Productivity and Performance in High-End Computing (at HPCA 10)*, Madrid, Spain, February 2004.
- [14] James C. Phillips, Gengbin Zheng, Sameer Kumar, and Laxmikant V. Kalé. NAMD: Biomolecular simulation on thousands of processors. In *Proceedings of the 2002 ACM/IEEE conference on Supercomputing*, pages 1–18, Baltimore, MD, September 2002.
- [15] Sameer Kumar, Chao Huang, Gheorge Almasi, and Laxmikant V. Kalé. Achieving strong scaling with NAMD on Blue Gene/L. In *Proceedings of IEEE International Parallel and Distributed Processing Symposium 2006*, April 2006.
- [16] Ramkumar V. Vadali, Yan Shi, Sameer Kumar, L. V. Kale, Mark E. Tuckerman, and Glenn J. Martyna. Scalable fine-grained parallelization of plane-wave-based ab initio molecular dynamics for large supercomputers. *Journal of Computational Chemistry*, 25(16):2006–2022, Oct. 2004.
- [17] Filippo Gioachin, Amit Sharma, Sayantan Chackravorty, Celso Mendes, Laxmikant V. Kale, and Thomas R. Quinn. Scalable cosmology simulations on parallel machines. In *7th International Meeting on High Performance Computing for Computational Science (VECPAR)*, July 2006.
- [18] C.H. Koelbel, D.B. Loveman, R.S. Schreiber, G.L. Steele Jr., and M.E. Zosel. *The High Performance Fortran Handbook*. MIT Press, 1994.
- [19] George Karypis and Vipin Kumar. Multilevel k-way partitioning scheme for irregular graphs. *Journal of Parallel and Distributed Computing*, 48:96 – 129, 1998.
- [20] George Karypis and Vipin Kumar. A fast and high quality multilevel scheme for partitioning irregular graphs. *SIAM J. Sci. Comput.*, 20(1):359–392, 1998.
- [21] Milind Bhandarkar, L. V. Kale, Eric de Sturler, and Jay Hoeflinger. Object-Based Adaptive Load Balancing for MPI Programs. In *Proceedings of the International Conference on Computational Science, San Francisco, CA, LNCS 2074*, pages 108–117, May 2001.
- [22] Rashmi Jyothi, Orion Sky Lawlor, and L. V. Kale. Debugging support for Charm++. In *PADTAD Workshop for IPDPS 2004*, page 294. IEEE Press, 2004.
- [23] J. K. Lenstra, D. B. Shmoys, and E. Tardos. Approximation algorithms for scheduling unrelated parallel machines. *Math. Program.*, 46(3):259–271, 1990.
- [24] George Karypis and Vipin Kumar. Parallel multilevel k-way partitioning scheme for irregular graphs. In *Supercomputing '96: Proceedings of the 1996 ACM/IEEE conference on Supercomputing (CDROM)*, page 35, 1996.
- [25] Laxmikant V. Kale, Gengbin Zheng, Chee Wai Lee, and Sameer Kumar. Scaling applications to massively parallel machines using projections performance analysis tool. In *Future Generation Computer Systems Special Issue on: Large-Scale System Performance Modeling and Analysis*, volume 22, pages 347–358, February 2006.
- [26] Orion Lawlor, Hari Govind, Isaac Dooley, Michael Breitenfeld, and Laxmikant Kale. Performance degradation in the presence of subnormal floating-point values. In *Proceedings of the International Workshop on Operating System Interference in High Performance Applications*, September 2005.
- [27] A. Basermann, J. Clinckemaillie, T. Coupeuz, J. Fingberg, H. Dignonnet, R. Ducloux, J.-M. Gratién, U. Hartmann, G. Lonsdale, B. Maerten, D. Roose, and C. Walshaw. Dynamic load balancing of finite element applications with the DRAMA Library. In *Applied Math. Modeling*, volume 25, pages 83–98, 2000.
- [28] K. Devine, B. Hendrickson, E. Boman, M. St. John, and C. Vaughan. Design of Dynamic Load-Balancing Tools for Parallel Applications. In *Proc. Intl. Conf. Supercomputing*, May 2000.
- [29] Karen D. Devine, Erik G. Boman, Robert T. Heaphy, Bruce A. Hendrickson, James D. Teresco, Jamal Faik, Joseph E. Flaherty, and Luis G. Gervasio. New challenges in dynamic load balancing. *Appl. Numer. Math.*, 52(2–3):133–152, 2005.
- [30] Kevin Barker, Andrey Chernikov, Nikos Chrisochoides, and Keshav Pingali. A Load Balancing Framework for Adaptive and Asynchronous

Applications. In *IEEE Transactions on Parallel and Distributed Systems*, volume 15, pages 183–192, 2003.

- [31] Kevin J. Barker and Nikos P. Chrisochoides. An Evaluation of a Framework for the Dynamic Load Balancing of Highly Adaptive and Irregular Parallel Applications. In *Proceedings of SC 2003*, Phoenix, AZ, 2003.
- [32] Chombo Software Package for AMR Applications. <http://seesar.lbl.gov/anag/chombo/>.
- [33] Robert K. Brunner and Laxmikant V. Kalé. Adapting to load on workstation clusters. In *The Seventh Symposium on the Frontiers of Massively Parallel Computation*, pages 106–112. IEEE Computer Society Press, February 1999.