# SANDIA REPORT

# Fast Parallel Algorithms for Short-Range Molecular Dynamics

Steve Plimpton

SF2900Q(8-81)

# Fast Parallel Algorithms

## for

## Short–Range Molecular Dynamics

Steve Plimpton
Department 1421
Sandia National Laboratories
Albuquerque, NM   87185
(505) 845-7873
sjplimp@cs.sandia.gov

## Abstract

Three parallel algorithms for classical molecular dynamics are presented. The first assigns each processor a subset of atoms; the second assigns each a subset of inter-atomic forces to compute; the third assigns each a fixed spatial region. The algorithms are suitable for molecular dynamics models which can be difficult to parallelize efficiently — those with short–range forces where the neighbors of each atom change rapidly. They can be implemented on any distributed–memory parallel machine which allows for message–passing of data between independently executing processors. The algorithms are tested on a standard Lennard–Jones benchmark problem for system sizes ranging from 500 to 10,000,000 atoms on three parallel supercomputers, the nCUBE 2, Intel iPSC/860, and Intel Delta. Comparing the results to the fastest reported vectorized Cray Y–MP and C90 algorithm shows that the current generation of parallel machines is competitive with conventional vector supercomputers even for small problems. For large problems, the spatial algorithm achieves parallel efficiencies of 90% and the Intel Delta performs about 30 times faster than a single Y–MP processor and 12 times faster than a single C90 processor. Trade–offs between the three algorithms and guidelines for adapting them to more complex molecular dynamics simulations are also discussed.

This is a pre–print of a paper submitted to a journal; please contact the author if you would like the latest up–to–date version of the manuscript. For example, when they become available, Intel Paragon and Cray MPP timings will likely be added to the tables in Section 7.

# 1 Introduction

Classical molecular dynamics (MD) is a commonly used computational tool for simulating the properties of liquids, solids, and molecules [1, 2]. Each of the $N$ atoms or molecules in the simulation is treated as a point mass and Newton's equations are integrated to compute their motion. From the motion of the ensemble of atoms a variety of useful microscopic and macroscopic information can be extracted such as transport coefficients, phase diagrams, and structural or conformational properties. The physics of the model is contained in a potential energy functional for the system from which individual force equations for each atom are derived.

MD simulations are typically not memory intensive since only vectors of atom information are stored. Computationally, the simulations are "large" in two domains --- the number of atoms and number of timesteps. The length scale for atomic coordinates is Angstroms; in three dimensions many thousands or millions of atoms must usually be simulated to approach even the microscopic scale. In liquids and solids the timestep size is constrained by the demand that the vibrational motion of the atoms be accurately tracked. This limits timesteps to the femtosecond scale and so tens or hundreds of thousands of timesteps are necessary to simulate even picoseconds of "real" time. Because of these computational demands, considerable effort has been expended by researchers to optimize MD calculations for vector supercomputers [22, 28, 31, 39, 41] and even to build special–purpose hardware for performing MD simulations [3, 4]. The current state-of-the-art is such that simulating ten– to hundred–thousand atom systems for picoseconds takes hours of CPU time on machines such as the Cray Y–MP.

The fact that MD computations are inherently parallel has been extensively discussed in the literature [9, 19]. There has been considerable effort in the last few years by researchers to exploit this parallelism on various machines. The majority of the work that has included implementations of proposed algorithms has been for single–instruction/multiple–data (SIMD) parallel machines such as the CM–2 [10, 46], or for multiple–instruction/multiple–data (MIMD) parallel machines with a few dozens of processors [24, 32, 40]. More recently there have been efforts at creating scalable algorithms that will work well on hundred– to thousand–processor MIMD machines [7, 12, 34, 45]. We are convinced that the MIMD programming model (or the single–program/multiple–data SPMD model as it is sometimes called) is the only one that provides enough flexibility to implement all the data structure and computational enhancements that are commonly exploited in MD codes on serial and vector machines. Also, we have found that it is only the current generation of massively parallel MIMD machines with hundreds to thousands of processors that have the computational power to be competitive with the fastest vector machines for MD calculations.

In this paper we present three parallel algorithms which are appropriate for a general class of MD problems that has two salient characteristics. The first characteristic is that forces are limited in range, meaning each atom interacts only with other atoms that are geometrically nearby. Solids and liquids are often modeled this way due to electronic screening effects or simply to avoid the computational cost of including long–range Coulombic forces. For short–range MD the computational effort per timestep scales as $N$, the number of

atoms, but care must be taken to write efficient parallel algorithms that take full advantage of the local nature of the forces.

The second characteristic is that the atoms can undergo large displacements over the duration of the simulation. This could be due to diffusion in a solid or liquid or conformational changes in a biological molecule. The important feature from a computational standpoint is that each atom's neighbors change as the simulation progresses. While the algorithms we discuss could also be used for fixed–neighbor simulations (e.g. all atoms remain on lattice sites in a solid), it is a harder task to continually track the neighbors of each atom and maintain efficient $O(N)$ scaling for the overall computation on a parallel machine.

Our first goal in this effort was to develop parallel algorithms that would be competitive with the fastest methods on vector supercomputers such as the Cray. Moreover we wanted the algorithms to work well on problems with small numbers of atoms, not just for large problems where parallelism is often easier to exploit. This is because currently the vast majority of MD simulations are performed on systems of a few hundred to several thousand atoms where $N$ is chosen to be as small as possible while still accurate enough to model the desired physical effects [6, 38, 47]. The computational goal in these calculations is to perform each timestep as quickly as possible. This is particularly true in non–equilibrium MD where macroscopic changes in the system may take significant time to evolve, requiring millions of timesteps to model. Thus, it is often more useful to be able to perform a 100,000 timestep simulation of a 1000 atom system fast rather than 1000 timesteps of a 100,000 atom system, though the $O(N)$ scaling means the computational effort is the same for both cases. To this end, we consider model sizes as small as a few hundred atoms in this paper.

For very large MD problems, our second goal in this work was to develop parallel algorithms that would be scalable to larger and faster parallel machines. While the timings we present for large MD models ($10^5$ to $10^7$ atoms) on the current generation of parallel supercomputers (hundreds to thousands of processors) are quite fast compared to vector supercomputers, they are still too slow to allow long–timescale simulations to be done routinely. However, our large–system algorithm scales optimally with respect to $N$ and $P$ (the number of processors) so that as parallel machines become more powerful in the next few years, algorithms similar to it will enable larger problems to be studied.

Our earlier efforts in this area [33] produced algorithms which were fast for systems up to tens of thousands of atoms but did not scale optimally with $N$ for larger systems. After improving on these efforts to create a scalable large–system algorithm [34] we have recently added an idea of Tamayo and Giles [45] that has improved the algorithm's performance on medium–sized problems by reducing the inter–processor communication requirements. We have also recently developed a new parallel algorithm which we present here in the context of MD simulations for the first time. It offers the advantages of both simplicity and speed for small to medium–sized problems.

In this paper we present the culmination of our efforts: several algorithms we have found, through implementing and testing a variety of ideas on different parallel machines, to be the fastest methods for short–range molecular dynamics across a wide range of problem sizes. By implementing the algorithms on

3

machines with hundreds to thousands of processors, we have been able to understand in practical terms what algorithmic features work best and tailor the algorithms accordingly to optimize their performance as a function of $N$ and $P$. Due to their scalability, we can also predict how these algorithms will perform on the faster, larger parallel machines of the future.

In the next section, the computational aspects of MD are highlighted and efforts to speed the calculations on vector and parallel machines are reviewed. In Sections 3, 4, and 5 we describe our three parallel algorithms in detail. A standard Lennard–Jones benchmark calculation is outlined in Section 6. In Section 7, implementation details and timing results for the parallel algorithms on three massively parallel MIMD machines are given and comparisons made to the best Cray Y–MP and C90 timings for the benchmark calculation. Discussion of the scaling properties of the algorithms is also included. Next, in Section 8, issues relevant to using the parallel algorithms in different kinds of MD simulations are discussed. Finally, in Section 9, we draw conclusions and give several guidelines for deciding which parallel algorithm is likely to be fastest for a particular short–range MD simulation.

## 2  Computational Aspects of Molecular Dynamics

The computational task in a MD simulation is to integrate the set of coupled differential equations (Newton's equations) given by

$$m_i \frac{d\vec{v}_i}{dt} = \sum_j F_2(\vec{r}_i, \vec{r}_j) + \sum_j \sum_k F_3(\vec{r}_i, \vec{r}_j, \vec{r}_k) + \cdots \qquad (1)$$

$$\frac{d\vec{r}_i}{dt} = \vec{v}_i$$

where $m_i$ is the mass of atom $i$, $\vec{r}_i$ and $\vec{v}_i$ are its position and velocity vectors, $F_2$ is a force function describing pairwise interactions between atoms, $F_3$ describes three–body interactions, and many–body interactions can be added. The force terms are derivatives of energy expressions in which the energy of atom $i$ is typically written as a function of only the positions of itself and other atoms. In practice, only one or a few terms in equation (1) are kept and $F_2$, $F_3$, etc. are constructed so as to include many–body and quantum effects. To the extent the approximations are accurate these equations give a full description of the time–evolution of the system. Thus, the great computational advantage of classical MD, as compared to *ab initio* electronic structure calculations, is that the dynamic behavior of the atomic system is described empirically without having to solve Schrodinger's equation at each timestep.

The force terms in equation (1) are typically non–linear functions of the distance $r_{ij}$ between pairs of atoms and may be either long–range or short–range in nature. For long–range forces, such as Coulombic interactions in an ionic solid or biological system, each atom interacts with all others. Directly computing these forces scales as $N^2$ and is too costly for large $N$. Various approximate methods overcome this difficulty. They include particle–mesh algorithms [29] which scale as $f(M)N$ where $M$ is the number of mesh points,

4

hierarchical methods [5] which scale as $N \log(N)$, and fast–multipole methods [21] which scale as $N$. Recent parallel implementations of these algorithms [17, 50] have improved their range of applicability for many–body simulations, but because of their expense, long–range force models are not commonly used in classical MD simulations.

By contrast, short–range force models are used extensively in MD and is what we are concerned with in this paper. They are chosen either because electronic screening effectively limits the range of influence of the interatomic forces being modeled or simply to truncate the long–range interactions and lessen the computational load. In either case, the summations in equation (1) are restricted to atoms within some small region surrounding atom $i$. This is typically implemented using a cutoff distance $r_c$, outside of which all interactions are ignored. The work to compute forces now scales linearly with $N$. However, even with this savings, the vast majority of computation time spent in a short–range force MD simulation is in evaluating the force terms in equation (1). The time integration typically requires only 2-3% of the total time. To evaluate the sums efficiently requires knowing which atoms are within the cutoff distance $r_c$ at every timestep. The key is to minimize the number of neighboring atoms that must be checked for possible interactions since calculations performed on neighbors at a distance $r > r_c$ are wasted computation. There are two basic techniques used to accomplish this on serial and vector machines; we discuss them briefly here since our parallel algorithms incorporate similar ideas.

The first idea, that of neighbor lists, was originally proposed by Verlet [49]. For each atom, a list is maintained of nearby atoms. Typically, when the list is formed, all neighboring atoms within an extended cutoff distance $r_s = r_c + \delta$ are stored. The list is used for a few timesteps to calculate all force interactions. Then it is rebuilt before any atom could have moved from a distance $r > r_s$ to $r < r_c$. Though $\delta$ is always chosen to be small relative to $r_c$, an optimal value depends on the parameters (e.g. temperature, diffusivity, density) of the particular simulation. The advantage of the neighbor list is that once it is built, examining it for possible interactions is much faster than checking all atoms in the system.

The second technique commonly used for speeding up MD calculations is known as the link-cell method [30]. At every timestep, all the atoms are binned into 3-D cells of side length $d$ where $d = r_c$ or slightly larger. This reduces the task of finding neighbors of a given atom to checking in 27 bins — the bin the atom is in and the 26 surrounding ones. Since binning the atoms only requires $O(N)$ work, the extra overhead associated with it is acceptable for the savings of only having to check a local region for neighbors.

The fastest MD algorithms on serial and vector machines use a combination of neighbor lists and link–cell binning. In the combined method, atoms are only binned once every few timesteps for the purpose of forming neighbor lists. In this case atoms are binned into cells of size $d \geq r_s$. At intermediate timesteps the neighbor lists alone are used in the usual way to find neighbors within a distance $r_c$ of each atom. This is a significant savings over a conventional link–cell method since there are far fewer atoms to check in a sphere of volume $4\pi r_s^3/3$ than in a cube of volume $27 r_c^3$. Additional savings can be gained due to Newton's 3rd law by only computing a force once for each pair of atoms (rather than once for each atom in the pair). In the combined

method this is done by only searching half the surrounding bins of each atom to form its neighbor list. This has the effect of storing atom $j$ in atom $i$'s list, but not atom $i$ in atom $j$'s list, thus halving the number of force computations that must be done.

Although these ideas are simply described, optimal performance on a vector machine requires careful attention to data structures and loop constructs to insure complete vectorization. The fastest implementation reported in the literature is that of Grest, et. al. [22]. They use the combined neighbor list/link-cell method described above to create long lists of pairs of neighboring atoms. At each timestep, they prune the lists to keep only those pairs within the cutoff distance $r_c$. Finally, they organize the lists into packets in which no atom appears twice [39]. The force computation for each packet can then be completely vectorized, resulting in performance on the benchmark problem described in Section 6 that is from 2 to 10 times faster than other vectorized algorithms [28, 41] over a wide range of simulation sizes.

Recently there has been considerable interest in devising parallel MD algorithms. The natural parallelism in MD is that the force calculations and velocity/position updates can be done simultaneously for all atoms. To date, two basic ideas have been exploited to achieve this parallelism. The goal in each is to divide the force computations in equation (1) evenly across the processors so as to extract maximum parallelism. To our knowledge, all algorithms that have been proposed or implemented (including ours) have been variations on these two methods. References [18, 23, 43] include good overviews of various techniques.

In the first class of methods a pre-determined set of force computations is assigned to each processor. The assignment remains fixed for the duration of the simulation. The simplest way of doing this is to give a subgroup of atoms to each processor. We call this method an *atom-decomposition* of the workload, since the processor computes forces on its atoms no matter where they move in the simulation domain. More generally, a subset of the force loops inherent in equation (1) can be assigned to each processor. We term this a *force-decomposition* and describe a new algorithm of this type later in the paper. Both of these decompositions are analogous to Lagrangian gridding in a fluids simulations where the grid cells (computational elements) move with the fluid (atoms in MD). By contrast, in the second general class of methods, which we call a *spatial-decomposition* of the workload, each processor is assigned a portion of the physical simulation domain. Each processor computes only the forces on atoms in its sub-domain. As the simulation progresses processors exchange atoms as they move from one sub-domain to another. This is analogous to an Eulerian gridding for a fluids simulation where the grid remains fixed in space as fluid moves through it.

Within the two classes of methods for parallelization of MD, a variety of algorithms have been proposed and implemented by various researchers. The details of the algorithms vary widely from one parallel machine to another since there are numerous problem-dependent and machine-dependent trade-offs to consider, such as the relative speeds of computation and communication. A brief review of some notable efforts follows.

Atom-decomposition methods, also called replicated-data methods [43] because vectors of atom information are replicated across all processors, are often used in MD simulations of molecular systems. This is because the duplication of information makes for straight-forward computation of additional three- and four-

body force terms. Parallel implementations of state–of–the–art biological MD programs such as CHARMM and GROMOS using this technique are discussed in [11, 15]. Force–decomposition methods which systolically cycle atom data around a ring or through a grid of processors have been used on MIMD [24, 43] and SIMD machines [14, 51]. Other force–decomposition methods that use the force–matrix formalism we discuss in Sections 3 and 4 have been presented in [10] and [13]. Boyer and Pawley [10] block–decompose the force matrix in a manner similar to that explained in Section 4, while the method of Brunet, et. al. [13] partitions the matrix element by element. Both of the methods are designed for long–range force systems requiring all–pairs calculations (no neighbor lists) on SIMD machines. Thus the scaling of these algorithms is different from what is presented in Section 4 as is the way they distribute the atom data among processors and perform inter–processor communication.

Spatial–decomposition methods, also called geometric methods [18, 23], are more common in the literature because they are well–suited to very large MD simulations. Recent parallel MIMD implementations for the Intel iPSC/2 hypercube [32, 40, 43], CM–5 [7, 45], and Fujitsu AP1000 [12] have some features in common with the spatial–decomposition algorithm we present in Section 5. The fastest published algorithms for SIMD machines also employ spatial–decomposition techniques [46]. However, the SIMD programming model, which requires processors executing each statement to operate simultaneously on a global data structure, introduces inefficiencies in short–range MD algorithms, particularly when coding the construction and access of variable–length neighbor lists via indirect addressing. Thus the timings in [46] for the benchmark problem discussed in Section 6 on a 32K–processor CM–2 [46] are slower than the single-processor Cray Y–MP timings presented in Section 7. By contrast, the timings for the MIMD parallel algorithms in this paper and references [7, 12, 45] are considerably faster, indicating the advantage a MIMD capability offers for exploiting parallelism in short–range MD simulations.

# 3   Atom–Decomposition Algorithm

In our first parallel algorithm each of the $P$ processors is assigned a group of $N/P$ atoms at the beginning of the simulation. Atoms in a group need not have any special spatial relationship to each other. For ease of exposition, we assume $N$ is a multiple of $P$, though it is simple to relax this constraint. A processor will compute forces on only its $N/P$ atoms and will update their positions and velocities for the duration of the simulation no matter where they move in the physical domain. As discussed in the previous section, this is an *atom–decomposition* of the computational workload.

A useful construct for representing the computational work involved in the algorithm is the $N \times N$ force matrix $F$. The $(ij)$ element of $F$ represents the force on atom $i$ due to atom $j$. Note that $F$ is sparse due to short–range forces and skew–symmetric, i.e. $F_{ij} = -F_{ji}$, due to Newton's 3rd law. We also define $x$ and $f$ as vectors of length $N$ which store the position and total force on each atom. For a 3–D simulation, $x_i$ would store the three coordinates of atom $i$. With these definitions, the atom–decomposition algorithm assigns each processor a sub–block of $F$ which consists of $N/P$ rows. This is shown in Figure 1 where we let

the $z$ subscript denote the processor number from 0 to $P - 1$. Thus, processor $P_z$ computes matrix elements in the $F_z$ block of rows. It also is assigned the corresponding sub–vectors of length $N/P$ denoted by $x_z$ and $f_z$.



Figure 1: The division of the force matrix among processors in the atom–decomposition algorithm. Processor $z$ is assigned a group of $N/P$ rows of the matrix and corresponding pieces of the position and force vectors, $x$ and $f$.

Assume the computation of matrix element $F_{ij}$ requires only the two atom positions $x_i$ and $x_j$. (We relax this assumption in Section 8.) To compute all the elements in $F_z$, processor $P_z$ will need the positions of many atoms owned by other processors. In the atom–decomposition algorithm, this is accomplished by having each processor send its updated atom positions to all the other processors after each timestep, an operation called all–to–all communication. Various algorithms have been developed for performing this operation efficiently on different parallel machines and architectures [19, 48]. We use an idea due to Fox, et. al. [19] that is simple, portable, and works well on a variety of machines. We describe it here because it is the chief communication component of both the atom–decomposition algorithms of this section and the force–decomposition algorithms presented in the next section.

Following Fox's nomenclature, we term the all–to–all communication procedure an *expand* operation. Each processor allocates memory of length $N$ to store the entire $x$ vector. At the beginning of the expand, processor $P_z$ has $x_z$, an updated piece of $x$ of length $N/P$. Each processor needs to acquire all the other processor's pieces, storing them in the correct places in its copy of $x$. Figure 2 illustrates the steps that accomplish this for an 8 processor example. The processors are mapped consecutively to the sub–pieces of the vector. In the first communication step, each processor exchanges its piece with an adjacent processor

in the vector. Processor 2 exchanges with processor 3 in the figure. Now, every processor has a contiguous piece of $x$ that is of length $2N/P$. In the second step, each processor exchanges this piece with a processor two positions away (2 exchanges with 0). Each processor now has a $4N/P$-length piece of $x$. In the last step, each processor exchanges an $N/2$-length piece of $x$ with a processor $P/2$ positions away (2 exchanges with 6); the entire vector now resides on each processor.



Figure 2: An *expand* operation among 8 processors. Processor 2 exchanges successively longer sub–vectors with processors 3, 0, and 6.

A pseudo–code version of the expand operation is given in Figure 3. For simplicity we again assume a power–of–two number of processors; relaxing this assumption is straightforward. The expand proceeds in $\log_2(P)$ steps. At each step $P_z$ performs a data exchange with a partner processor $P'$. The new processor number $P'$ is obtained by flipping one bit in $z$, which itself is a string of $\log_2(P)$ bits. The sub–vector $y$ is sent to $P'$ and the received sub–vector $w$ is concatenated with $y$ (the "|" operation) in the proper order. Thus $y$ doubles in length at every step; at the end of the expand $y$ has become the full $N$–length vector $x$. Costs for a communication algorithm are typically quantified by the number of messages and the total volume of data sent and received. On both these accounts the expand is optimal; each processor performs $\log_2(P)$ sends and receives and exchanges $N - N/P$ data values. This is the reason the expand operation works well on many machines. A drawback is that it requires $O(N)$ storage on every processor. Alternative methods for performing all–to–all communication require less storage at the cost of more sends and receives. This is usually not a good trade–off for MD simulations because, as we shall see, quite large problems can be run with an atom–decomposition algorithm in the many Mbytes of local memory available on current–generation processors.

```
y := x_z
FOR k = 0, ..., log_2(P) - 1
    P' := P_z with k^{th} bit of z flipped
    SEND y to processor P'
    RECEIVE w from processor P'
    IF bit k of z is 0 THEN
        y := y|w
    ELSE
        y := w|y
x := y
```

Figure 3: The *expand* operation for processor $P_z$.

A communication operation that is essentially the inverse of the expand will also prove useful in the atom– and force–decomposition algorithms. Assume each processor has stored new force values throughout its copy of the force vector $f$. Processor $P_z$ needs to know the $N/P$ values in $f_z$, where each of the values is summed across all $P$ processors. A procedure for doing this is known as a *fold* operation [19] and is outlined in Figure 4. Again the operation proceeds in $\log_2(P)$ steps. At each step, $y$ represents a portion of the force vector $f$, and is split into two pieces, $u$ and $v$. One of the pieces is sent to a partner processor $P'$. The received sub–vector $w$ is summed element by element with the retained piece. This summed sub–vector becomes $y$ in the next step, so that $y$ is halving in length at each iteration of the loop. When the fold is finished, $y$ has become $f_z$, with values summed across all $P$ processors. Like the expand, the fold operation requires $\log_2(P)$ sends and receives and $N - N/P$ data to be exchanged by each processor. Additionally it requires $N - N/P$ floating point operations to do the summations, typically a small extra cost.

Having defined the expand and fold operations, we now present two versions of the atom–decomposition algorithm. The first is simpler and does not take advantage of Newton's 3rd law. We call this algorithm A1; it is outlined in Figure 5 with the dominating term(s) in the computation or communication cost of each step listed on the right. We assume at the beginning of the timestep that each processor knows the current positions of all $N$ atoms, i.e. each has a copy of the entire $x$ vector. Step (1) of the algorithm is to construct neighbor lists for all the pairwise interactions that must be computed in block $F_z$. Typically this will only be done once every few timesteps. If the ratio of the physical domain diameter $D$ to the extended force cutoff length $r_s$ is relatively small, it is quicker for $P_z$ to construct the lists by checking all $N^2/P$ pairs in its $F_z$ block. When the simulation is large enough that 4 or more bins can be created in each dimension, it is quicker for each processor to bin all $N$ atoms, then check the 27 surrounding bins of each of its $N/P$ atoms to form the lists. This checking scales as $N/P$ but has a large coefficient, so the overall scaling of the binned neighbor list construction is recorded as $N/P + N$.

10

```
y := f
FOR k = log₂(P) − 1, . . . , 0
    u := top half of y vector
    v := bottom half of y vector
    P' := P_z with k^{th} bit of z flipped
    IF bit k of z is 0 THEN
        SEND v to processor P'
        RECEIVE w from processor P'
        y := u + w
    ELSE
        SEND u to processor P'
        RECEIVE w from processor P'
        y := v + w
f_z := y
```

Figure 4: The *fold* operation for processor $P_z$.

In step (2) of the algorithm, the neighbor lists are used to compute the non-zero matrix elements in $F_z$. As each pairwise force interaction is computed, the force components are summed into $f_z$, so that $F_z$ is never actually stored as a matrix. At the completion of the step, each processor knows the total force $f_z$ on each of its $N/P$ atoms. This is used to update their positions and velocities in step (4). (A step (3) will be added to other algorithms in this and the following sections.) Finally, in step (5) the updated atom positions in $x_z$ are shared among all $P$ processors in preparation for the next timestep via the expand operation of Figure 3. As discussed above, this operation scales as $N$, the volume of data in the position vector $x$.

| | | |
|---|---|---|
| (1) | **Construct** neighbor lists of non-zero interactions in $F_z$ | |
| | ($D < 4r_s$) All pairs | $\frac{N^2}{P}$ |
| | ($D \geq 4r_s$) Binning | $\frac{N}{P} + N$ |
| (2) | **Compute** elements of $F_z$, summing results into $f_z$ | $\frac{N}{P}$ |
| (4) | **Update** atom positions in $x_z$ using $f_z$ | $\frac{N}{P}$ |
| (5) | **Expand** $x_z$ among all processors, result is $x$ | $N$ |

Figure 5: Single timestep of atom-decomposition algorithm **A1** for processor $P_z$.

11

As mentioned above, algorithm **A1** ignores Newton's 3rd law. If different processors own atoms $i$ and $j$ as is usually the case, both processors compute the $(ij)$ interaction and store the resulting force on their atom. This can be avoided at the cost of more communication by using a modified force matrix $G$ which references each pairwise interaction only once. There are several ways to do this by striping the force matrix [42]; we choose instead to form $G$ as follows. Let $G_{ij} = F_{ij}$, except that $G_{ij} = 0$ when $i > j$ and $i + j$ is even, and likewise $G_{ij} = 0$ when $i < j$ and $i + j$ is odd. Conceptually, $G$ is colored like a checkerboard with red squares above the diagonal set to zero and black squares below the diagonal also set to zero. A modified atom-decomposition algorithm **A2** that uses $G$ to take advantage of Newton's 3rd law is outlined in Figure 6.

| | | |
|---|---|---|
| (1) | **Construct** neighbor lists of non-zero interactions in $G_z$ | |
| | $(D < 4r_s)$ All pairs | $\frac{N^2}{2P}$ |
| | $(D \geq 4r_s)$ Binning | $\frac{N}{2P} + N$ |
| (2) | **Compute** elements of $G_z$, | |
| | doubly summing results into local copy of $f$ | $\frac{N}{2P}$ |
| (3) | **Fold** $f$ among all processors, result is $f_z$ | $N$ |
| (4) | **Update** atom positions in $x_z$ using $f_z$ | $\frac{N}{P}$ |
| (5) | **Expand** $x_z$ among all processors, result is $x$ | $N$ |

Figure 6: Single timestep of atom-decomposition algorithm **A2** for processor $P_z$, which takes advantage of Newton's 3rd law.

Step (1) is the same as in algorithm **A1** except only half as many neighbor list entries are made by each processor since $G_z$ has only half the non-zero entries of $F_z$. This is reflected in the factors-of-two included in the scaling entries. For neighbor lists formed by binning, each processor must still bin all $N$ atoms, but only need check half the surrounding bins of each of its $N/P$ atoms. In step (2) the neighbor lists are used to compute elements of $G_z$. For an interaction between atoms $i$ and $j$, the resulting forces on atoms $i$ and $j$ are summed into both the $i$ and $j$ locations of force vector $f$. This means each processor must store a copy of the entire force vector, as opposed to just storing $f_z$ as in algorithm **A1**. When all the matrix elements have been computed, $f$ is folded across all $P$ processors using the algorithm in Figure 4. Each processor ends up with $f_z$, the total forces on its atoms. Steps (4) and (5) then proceed the same as in **A1**.

Note that implementing Newton's 3rd law essentially halved the computational cost in steps (1) and (2), at the expense of doubling the communication cost. There are now two communication steps (3) and (5), each of which scale as $N$. This will only be a net gain if the communication cost in **A1** is less than a third of the overall run time. As we shall see, this will usually not be the case on large numbers of processors, so in practice we almost always choose **A1** instead of **A2** for an atom-decomposition algorithm. However, for

small $P$ or expensive force models, **A2** can be faster.

Finally, we discuss the issue of load–balance. Each processor will has an equal a amount of work if each $F_z$ or $G_z$ block has roughly the same number of non–zero elements. This will be the case if the atom density is uniform across the simulation domain. However non–uniform densities can arise if, for example, there are free surfaces so that some atoms border on vacuum, or phase changes are occurring within a liquid or solid. This is only a problem for load–balancing of the atom–decomposition computation across processors if the $N$ atoms are ordered in a geometric sense as is typically the case. Then a group of $N/P$ atoms near a surface, for example, will have fewer neighbors than groups in the interior. This can be overcome by randomly permuting the atom ordering at the beginning of the simulation, which is equivalent to permuting rows and columns of $F$ or $G$. This insures that every $F_z$ or $G_z$ will have roughly the same number of non–zeros even if the atom density is non–uniform. A random permutation also has the advantage that the load–balance will likely persist as atoms move about during the simulation. Note that this permutation need only be done once, as a pre–processing step before beginning the dynamics.

In summary, the atom–decomposition algorithms divide the MD force computation and integration evenly across the processors (ignoring the $O(N)$ component of binned neighbor list construction which is usually not significant). However, the algorithms require global communication, as each processor must acquire information held by all the other processors. This communication scales as $N$, independent of $P$, so it limits the number of processors that can be used effectively. The chief advantage of the algorithms is that of simplicity. Steps (1), (2), and (4) can be implemented by simply modifying the loops and data structures in a serial or vector code to treat $N/P$ atoms instead of $N$. The expand and fold communication operations can be treated as black–box routines and inserted at the proper locations in steps (3) and (5). Few other changes are typically necessary to parallelize an existing code.

# 4 Force–Decomposition Algorithm

Our next parallel MD algorithm is based on a block–decomposition of the force matrix $F$ rather than a row–wise decomposition as used in the previous section. We term this a *force–decomposition* of the workload. As we shall see, this improves the $O(N)$ scaling of the communication cost to $O(N/\sqrt{P})$. Block–decompositions of matrices are common in linear algebra algorithms [8, 26] for parallel machines which sparked our interest in the idea, but to our knowledge we are the first to apply this idea to short–range MD simulations [27, 36, 37]. The assignment of sub–blocks of $F$ to processors is depicted in Figure 7. We assume for ease of exposition that $P$ is an even power of 2 and that $N$ is a multiple of $P$, although again it is straightforward to relax these constraints. The block owned by each processor is thus of size $(N/\sqrt{P}) \times (N/\sqrt{P})$. We use the Greek subscripts $\alpha$ and $\beta$ to index the row and column blocks of $F$ running from 0 to $\sqrt{P} - 1$. A sub–block of $F$ is denoted as $F_{\alpha\beta}$, and the processor owning it is $P_{\alpha\beta}$. We note that $\alpha$ and $\beta$ also index sub–vectors of $x$ and $f$ of length $N/\sqrt{P}$. To compute the matrix elements in $F_{\alpha\beta}$, processor $P_{\alpha\beta}$ must know the $x_\alpha$ and $x_\beta$ pieces of $x$. As these elements are computed they will be stored in local copies of the force sub–vectors, namely $f_\alpha$
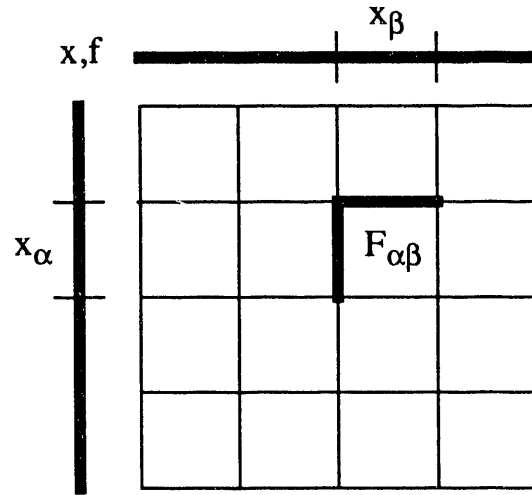
13

and $f_\beta$.



Figure 7: The division of the force matrix among processors in the force–decomposition algorithm. Processor $P_{\alpha\beta}$ is assigned a sub–block $F_{\alpha\beta}$ of size $N/\sqrt{P}$ by $N/\sqrt{P}$. Likewise it stores the corresponding length $N/\sqrt{P}$ pieces of the position and force vectors.

In addition to computing the matrix elements in $F_{\alpha\beta}$, each processor will be responsible for updating the positions and velocities of $N/P$ atoms, as in the atom–decomposition algorithm. These atoms are a sub–vector of $x_\alpha$; that is, the $\sqrt{P}$ processors in row $\alpha$ divide $x_\alpha$ among them, so each is responsible for a contiguous piece of length $N/P$. Numbering these pieces with the column index $\beta$ of the processor, we denote each processor's piece with a superscript, $x_\alpha^\beta$. Similarly, the total force acting on these atoms is the $N/P$–length sub–vector $f_\alpha^\beta$. As in the atom–decomposition case, an element of $f_\alpha^\beta$ is the sum of all the matrix elements across the corresponding row of $F$.

Our first force–decomposition algorithm **F1** is outlined in Figure 8. As before, each processor has updated copies of the needed atom positions at the beginning of the timestep. In this case it is the current sub–vectors $x_\alpha$ and $x_\beta$. In step (1) neighbor lists are constructed. Again, for small problems this is most quickly done be checking all $N^2/P$ possible pairs in $F_{\alpha\beta}$. For large problems, the $N/\sqrt{P}$ atoms in $x_\beta$ are binned, then the 27 surrounding bins of each atom in $x_\alpha$ is checked. The total number of interactions stored in each processor's lists is still $O(N/P)$. The scaling of the binned neighbor list construction is thus $N/P + N/\sqrt{P}$. In step (2) the neighbor lists are used to compute the matrix elements in $F_{\alpha\beta}$. As before the elements are summed into a local copy of $f_\alpha$ as they are computed, so $F_{\alpha\beta}$ need never be stored in matrix form. In step (3) a fold operation is performed within each row of processors so that processor $P_{\alpha\beta}$ obtains the total

14

forces on its $N/P$ atoms, $f_\alpha^\beta$. Although the fold algorithm used is the same as in the preceding section, there is a key difference. In this case the vector $f_\alpha$ being folded is only of length $N/\sqrt{P}$ and only the $\sqrt{P}$ processors in one row are participating in the fold. Thus this operation scales as $N/\sqrt{P}$ instead of $N$ as in the atom-decomposition communication steps.

In step (4), $f_\alpha^\beta$ is used by $P_{\alpha\beta}$ to update the $N/P$ atom positions in $x_\alpha^\beta$. Steps (5a-5d) share these updated positions with all the processors that will need them for the next timestep. These are the processors which share a row or column with $P_{\alpha\beta}$. First, in (5a), the processors in row $\alpha$ perform an expand of their $x_\alpha^\beta$ sub-vectors so that each acquires the entire $x_\alpha$. As with the fold, this operation scales as the $N/\sqrt{P}$ length of $x_\alpha$ instead of as $N$ as it did in algorithms **A1** and **A2**. In step (5b), each processor exchanges its updated atom positions with processor $P_{\beta\alpha}$ which owns the block of $F$ in the transpose position of the matrix. The cost of this operation scales as the $N/P$ length of the data being exchanged. Finally, in step (5c), the processors in each column $\beta$ perform an expand of the received sub-vector $x_\beta^\alpha$. As a result they all acquire $x_\beta$ and are ready to begin the next timestep.

| | | |
|---|---|---|
| (1) | **Construct** neighbor lists of non-zero interactions in $F_{\alpha\beta}$ | |
| | $(D < 4r_s)$ All pairs | $\frac{N^2}{P}$ |
| | $(D \geq 4r_s)$ Binning | $\frac{N}{P} + \frac{N}{\sqrt{P}}$ |
| (2) | **Compute** elements of $F_{\alpha\beta}$, storing results in $f_\alpha$ | $\frac{N}{P}$ |
| (3) | **Fold** $f_\alpha$ within row $\alpha$, result is $f_\alpha^\beta$ | $\frac{N}{\sqrt{P}}$ |
| (4) | **Update** atom positions in $x_\alpha^\beta$ using $f_\alpha^\beta$ | $\frac{N}{P}$ |
| (5a) | **Expand** $x_\alpha^\beta$ within row $\alpha$, result is $x_\alpha$ | $\frac{N}{\sqrt{P}}$ |
| (5b) | **Exchange** atom positions with transpose processor $P_{\beta\alpha}$ | $\frac{N}{P}$ |
| | **Send** $x_\alpha^\beta$ to $P_{\beta\alpha}$ | |
| | **Receive** $x_\beta^\alpha$ from $P_{\beta\alpha}$ | |
| (5c) | **Expand** $x_\beta^\alpha$ within column $\beta$, result is $x_\beta$ | $\frac{N}{\sqrt{P}}$ |

Figure 8: Single timestep of force-decomposition algorithm **F1** for processor $P_{\alpha\beta}$.

As with algorithm **A1**, algorithm **F1** does not take advantage of Newton's 3rd law; each pairwise force interaction is computed twice. Algorithm **F2** avoids this duplicated effort by using the same checkerboarded matrix $G$ that was defined in the preceding section. Note that now the total force on atom $i$ is the sum of all matrix elements in row $i$ minus the sum of all elements in column $i$. The modified force-decomposition algorithm **F2** is outlined in Figure 9. Step (1) is the same as in F1, except that half as many interactions are stored in the neighbor lists. Likewise, step (2) requires only half as many matrix elements be computed. For each $(ij)$ element, the computed force components are now summed into two force vectors instead of one. The force on atom $i$ is summed into $f_\alpha$ in the location corresponding to row $i$. The same force on atom

15

$j$ is also summed into $f_\beta$ in the location corresponding to column $j$. Steps (3a-3d) accumulate these forces so that processor $P_{\alpha\beta}$ ends up with the total force on its $N/P$ atoms. First, in step (3a), the $\sqrt{P}$ processors in column $\beta$ fold their local copies of $f_\beta$. The result is $\overline{f}_\beta^\alpha$. Each element of this $N/P$-length sub-vector is the sum of an entire column of $G$. In step (3b) this sub-vector is exchanged with the transpose-position processor $P_{\beta\alpha}$. The values in the sub-vector each processor receives in this transpose operation are the column contribution to the forces on its $N/P$ atoms. Next, in step (3c), the row contributions to the forces are summed by performing a fold of the $f_\alpha$ vector within each row $\alpha$. The result is $\underline{f}_{\alpha}^\beta$, each element of which is the sum across a row of $G$. Finally, in step (3d) the column and row contributions are subtracted element by element to yield the total forces $f_\alpha^\beta$ on the atoms owned by processor $P_{\alpha\beta}$. The processor can now update the positions and velocities of its atoms; steps 4 and 5 are identical to those of **F1**.

| | | |
|---|---|---|
| (1) **Construct** neighbor lists of non-zero interactions in $G_{\alpha\beta}$ | | |
| $\quad$ $(D < 4r_s)$ All pairs | | $\frac{N^2}{2P}$ |
| $\quad$ $(D \geq 4r_s)$ Binning | | $\frac{N}{2P} + \frac{N}{\sqrt{P}}$ |
| (2) **Compute** elements of $G_{\alpha\beta}$, | | |
| $\quad$ storing results in local copies of $f_\alpha$ and $f_\beta$ | | $\frac{N}{2P}$ |
| (3a) **Fold** $f_\beta$ within column $\beta$, result is $\overline{f}_\beta^\alpha$ | | $\frac{N}{\sqrt{P}}$ |
| (3b) **Exchange** column forces with transpose processor $P_{\beta\alpha}$ | | $\frac{N}{P}$ |
| $\quad$ **Send** $\overline{f}_\beta^\alpha$ to $P_{\beta\alpha}$ | | |
| $\quad$ **Receive** $\overline{f}_\alpha^\beta$ from $P_{\beta\alpha}$ | | |
| (3c) **Fold** $f_\alpha$ within row $\alpha$, result is $\underline{f}_{\alpha}^\beta$ | | $\frac{N}{\sqrt{P}}$ |
| (3d) **Subtract** received $\overline{f}_\alpha^\beta$ from folded $\underline{f}_{\alpha}^\beta$, result is total $f_\alpha^\beta$ | | $\frac{N}{P}$ |
| (4) **Update** atom positions in $x_\alpha^\beta$ using $f_\alpha^\beta$ | | $\frac{N}{P}$ |
| (5a) **Expand** $x_\alpha^\beta$ within row $\alpha$, result is $x_\alpha$ | | $\frac{N}{\sqrt{P}}$ |
| (5b) **Exchange** atom positions with transpose processor $P_{\beta\alpha}$ | | $\frac{N}{P}$ |
| $\quad$ **Send** $x_\alpha^\beta$ to $P_{\beta\alpha}$ | | |
| $\quad$ **Receive** $x_\beta^\alpha$ from $P_{\beta\alpha}$ | | |
| (5c) **Expand** $x_\beta^\alpha$ within column $\beta$, result is $x_\beta$ | | $\frac{N}{\sqrt{P}}$ |

Figure 9: Single timestep of force–decomposition algorithm **F2** for processor $P_{\alpha\beta}$, which takes advantage of Newton's 3rd law.

In the force–decomposition algorithms, exploiting Newton's 3rd law again halves the computation required in steps (1) and (2). However, the communication cost in steps (3) and (5) does not double. Rather there are 4 expands and folds required in **F2** versus 3 in **F1**. There are also two transpose operations instead of one. Thus, in practice, it is usually faster to use algorithm **F2** with its reduced computational cost and slightly

increased communication cost rather than **F1**. The key point is that all the expand and fold operations in **F1** and **F2** scale as $N/\sqrt{P}$ rather than as $N$ as was the case in algorithms **A1** and **A2**. As we shall see, when run on large numbers of processors this significantly reduces the time the force–decomposition algorithms spend on communication as compared to the atom–decomposition algorithms.

Finally, the issue of load-balance is a more serious concern for the force–decomposition algorithms. Processors will have equal work to do only if all the matrix blocks $F_{\alpha\beta}$ or $G_{\alpha\beta}$ are uniformly sparse. If the atoms are ordered geometrically this will not be the case even for problems with uniform density. This is because such an ordering creates a force matrix with diagonal bands of non–zero elements. As in the atom–decomposition case, a random permutation of the atom ordering produces the desired effect. Only now the permutation should be done as a pre–processing step for all problems, even those with uniform atom densities.

In summary, algorithms **F1** and **F2** divide the MD computations evenly across processors as did the atom-decomposition algorithms. But the block–decomposition of the force matrix means each processor only needs $O(N/\sqrt{P})$ information to perform its computations. Thus the communication and memory costs are reduced by a factor of $\sqrt{P}$ versus algorithms **A1** and **A2**. The force–decomposition strategy retains the simplicity of the atom–decomposition technique; **F1** and **F2** can be implemented using the same "black-box" communication routines as **A1** and **A2**. The force–decomposition algorithms also need no geometric information about the physical problem being modeled to perform optimally. In fact, for load–balancing purposes the algorithms intentionally ignore such information by using a random atom ordering.

# 5 Spatial–Decomposition Algorithm

In our final parallel algorithm the physical simulation domain is subdivided into small 3–D boxes, one for each processor. We call this a *spatial-decomposition* of the workload. Each processor computes forces on and updates the positions and velocities of all atoms within its box at each timestep. Atoms are reassigned to new processors as they move through the physical domain. In order to compute forces on its atoms, a processor need only know positions of atoms in nearby boxes. The communication required in the spatial-decomposition algorithm is thus local in nature as compared to global in the atom- and force-decomposition cases.

The size and shape of the box assigned to each processor will depend on $N$, $P$, and the shape of the physical domain, which we assume to be a 3–D rectangular parallelepiped. Within these constraints the number of processors in each dimension is chosen so as to make each processor's box as "cubic" as possible. This is to minimize communication since in the large $N$ limit the communication cost of the spatial-decomposition algorithm will turn out to be proportional to the surface area of the boxes. An important point to note is that in contrast to the link–cell method for conventional MD described in Section 2, the box lengths may now be smaller or larger than the force cutoff lengths $r_c$ and $r_s$.

Each processor in our spatial–decomposition algorithm maintains two data structures, one for the $N/P$

atoms in its box and one for atoms in nearby boxes. In the first data structure, each processor stores complete information — positions, velocities, neighbor lists, etc. This data is stored in a linked list to allow insertions and deletions as atoms move to new boxes. In the second data structure only atom positions are stored. Interprocessor communication at each timestep keeps this information current.

The communication scheme we use to acquire this information from processors owning the nearby boxes is shown in Figure 10. The first step (a) is for each processor to pair up with an adjacent processor in the east/west dimension, 2 pairs with 1 for example. Processor 2 fills a message buffer with atom positions it owns that are within a force cutoff length $r_s$ of processor 1's box. (The reason for using $r_s$ instead of $r_c$ will be made clear below.) If $d < r_s$, where $d$ is the box length in the east/west direction, this will be all of processor 2's atoms; otherwise it will be those nearest to box 1. Now processors 2 and 1 exchange messages. Processor 2 puts the information it receives into its second data structure. Now the processors pair up in the opposite east–west direction, 2 with 3 in this case, and perform the same operation. If $d > r_s$, all needed atom positions in the east–west dimension have now been acquired by each processor. If $d < r_s$, this procedure is repeated with each processor sending more needed atom positions to its adjacent processors. For example, processor 2 sends processor 1 atom positions from box 3 (which processor 2 now has in its second data structure). This can be repeated until each processor knows all atom positions within a distance $r_s$ of its box, as indicated by the dotted boxes in the figure. The same process is now repeated in the north/south dimension; see step (b) of the figure. The only difference is that messages sent to the adjacent processor now contain not only atoms the processor owns (in its first data structure), but also any atom positions in its second data structure that are needed by the adjacent processor. For $d = r_s$ this has the effect of sending 3 boxes worth of atom positions in one message as shown in (b). Finally, in step (c) the process is repeated in the up/down dimension. Now atom positions from an entire plane of boxes (9 in the figure) are effectively being exchanged in each message.

There are several key advantages to this scheme, all of which reduce the overall cost of communication in our algorithm. First, for $d \geq r_s$, needed atom positions from all 26 surrounding boxes are obtained in just 6 data exchanges. Moreover, as will be discussed in the results section, if the parallel machine is a hypercube, the processors can be mapped to the boxes in such a way that all 6 of these processors will be directly connected to the center processor. Thus message passing will be fast and contention-free. Even if $d < r_s$ so that atom information is needed from more distant boxes, this occurs with only a few extra data exchanges, all of which are still with the 6 immediate neighbor processors. Second, the amount of data communicated is minimized. Each processor acquires only the atom positions that are within a distance $r_s$ of its box. Third, all of the received atom positions can be placed as contiguous data directly into the processor's second data structure. No time is spent rearranging data, except to create the buffered messages that need to be sent. Finally, as will be discussed in more detail below, this message creation can be done very quickly. A full scan of the two data structures is only done once every few timesteps, when the neighbor lists are created, to decide which atom positions to send in each message. The scan procedure creates a list of atoms that make
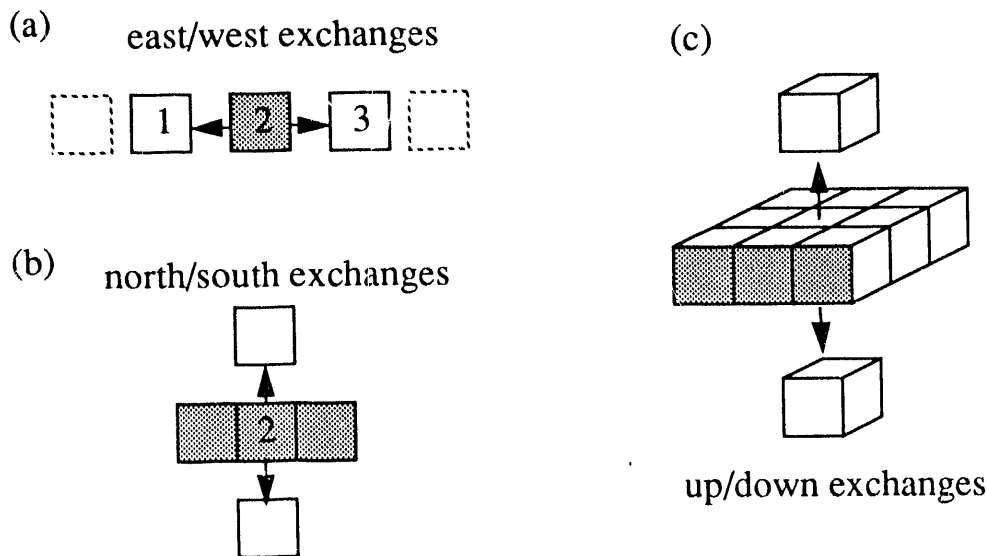
**(a)** east/west exchanges

**(b)** north/south exchanges

**(c)** up/down exchanges

Figure 10: Method by which a processor acquires nearby atom positions in the spatial–decomposition algorithm. In 6 data exchanges all atom positions in adjacent boxes in the (a) east/west, (b) north/south, and (c) up/down directions can be communicated.

up each message. During all the other timesteps, the lists can be used, in lieu of scanning the full atom list, to directly index the referenced atoms and buffer up the messages quickly. This is the equivalent of a gather operation.

We now outline our spatial-decomposition algorithm **S1** in Figure 11. Box $z$ is assigned to processor $P_z$, where $z$ runs from 0 to $P - 1$ as before. Processor $P_z$ stores the atom positions of its $N/P$ atoms in $x_z$ and the forces on those atoms in $f_z$. Steps (1a-1c) are the neighbor list construction, performed once every few timesteps. This is somewhat more complex than in the other algorithms because, as discussed above, it includes the making of lists of atoms that will be communicated at every timestep. First, in step (1a) the positions, velocities, and any other identifying information of atoms that are no longer inside box $z$ are deleted from $x_z$ and stored in a message buffer. These atoms are exchanged with the 6 adjacent processors via the communication pattern of Figure 10. As the information routes through each dimension, processor $P_z$ checks for new atoms that are now inside its box boundaries, adding them to $x_z$. Next, in step (1b), all atom positions within a distance $r_s$ of box $z$ are acquired by the communication scheme described above. As the different messages are buffered by scanning through the two data structures, lists of included atoms are made. The lists will be used in step (5). The scaling factor $\Delta$ for steps (1a) and (1b) will be explained below.

When steps (1a) and (1b) are complete, both of the processor's data structures are current. Neighbor

19

| | |
|---|---|
| **(1a) Move** necessary atoms to new boxes | $\Delta$ |
| **(1b) Make** lists of all atoms that will need to be exchanged | $\Delta$ |
| **(1c) Construct** neighbor lists of interaction pairs in box $z$ | |
| $\quad (d < 2r_s)$ All pairs | $\frac{N}{P}(\frac{N}{2P} + \Delta)$ |
| $\quad (d \geq 2r_s)$ Binning | $\frac{N}{2P} + \Delta$ |
| **(2) Compute** forces on atoms in box $z$, doubly storing results in $f_z$ | $\frac{N}{2P} + \Delta$ |
| **(4) Update** atom positions $x_z$ in box $z$ using $f_z$ | $\frac{N}{P}$ |
| **(5) Exchange** atom positions across box boundaries | |
| $\quad$ with neighboring processors | $\frac{N}{P}(1 + 2r_s/d)^3$ |
| $\quad (d < r_s)$ Send $N/P$ positions to many neighbors | $r_s{}^3$ |
| $\quad (d \approx r_s)$ Send $N/P$ positions to nearest neighbors | $\frac{N}{P}$ |
| $\quad (d > r_s)$ Send positions near box surface to nearest neighbors | $(\frac{N}{P})^{2/3}$ |

Figure 11: Single timestep of spatial-decomposition algorithm **S1** for processor $P_z$.

lists for its $N/P$ atoms can now be constructed in step (1c). If atoms $i$ and $j$ are both in box $z$ (an inner-box interaction), the $(ij)$ pair is only stored once in the neighbor list. If $i$ and $j$ are in different boxes (a two-box interaction), both processors store the interaction in their respective neighbor lists. If this were not done, processors would compute forces on atoms they do not own and communication of the forces back to the processors owning the atoms would be required. A modified algorithm which performs this communication to avoid the duplicated force computation of two-box interactions is discussed below. When $d$, the length of box $z$, is less than two cutoff distances, it is quicker to find neighbor interactions by checking each atom inside box $z$ against all the atoms in both of the processor's data structures. This scales as the square of $N/P$. If $d > 2r_s$, then with the shell of atoms around box $z$, there are 4 or more bins in each dimension. In this case, as with the other algorithms, it is quicker to perform the neighbor list construction by binning. All the atoms in both data structures are hashed into bins of size $r_s$. The surrounding bins of each atom in box $z$ are then checked for possible neighbors.

Processor $P_z$ can now compute all the forces on its atoms in step (2) using the neighbor lists. When the interaction is between two atoms inside box $z$, the resulting force is stored twice in $f_z$, once for atom $i$ and once for atom $j$. For two-box interactions, only the force on the processor's own atom is stored. After computing $f_z$, the atom positions are updated in step (4). Finally, these updated positions must be communicated to the surrounding processors in preparation for the next timestep. This occurs in step (5) using the previously made lists to create each message and the communication pattern of Figure 10. The amount of data exchanged in this operation is a function of the relative values of the force cutoff distance and box length and is discussed in the next paragraph. Also, we note that on the timesteps that neighbor lists are constructed, step (5) does not have to be performed since step (1b) has the same effect.

20

The communication operations in algorithm **S1** occur in steps (1a), (1b), and (5). The communication in the latter two steps is identical. The cost of these steps scales as the volume of data exchanged. For step (5), if we assume uniform atom density, this is proportional to the physical volume of the shell of thickness $r_s$ around box $z$, namely $(d + 2r_s)^3 - d^3$. Note there are roughly $N/P$ atoms in a volume of $d^3$, since $d^3$ is the size of box $z$. There are 3 cases to consider. First, if $d < r_s$ data from many neighboring boxes must be exchanged and the operation scales as $8r_s$ . Second, if $d \approx r_s$, the data in all 26 surrounding boxes is exchanged and the operation scales as $27N/P$. Finally, if $d$ is much larger than $r_s$, only atom positions near the 6 faces of box $z$ will be exchanged. The communication then scales as the surface area of box $z$, namely $6r_s(N/P)^{2/3}$. These 3 cases are explicitly listed in the scaling of step (5). Elsewhere in Figure 11, we use the term $\Delta$ to represent whichever of the three is applicable for a given $N$, $P$, and $r_s$. We note that step (1a) involves less communication since not all the atoms within a cutoff distance of a box face will move out of the box. But this operation still scales as the surface area of box $z$, so we list its scaling as $\Delta$.

The computational portion of algorithm **S1** is in steps (1c), (2), and (4). All of these scale as $N/P$ with additional work in steps (1c) and (2) for atoms that are neighboring box $z$ and stored in the second data structure. The number of these atoms is proportional to $\Delta$ so it is included in the scaling of those steps. The leading term in the scaling of steps (1c) and (2) is listed as $N/2P$ as in algorithms **A2** and **F2**, since inner-box interactions are only stored and computed once for each pair of atoms in algorithm **S1**. Note that as $d$ grows large relative to $r_s$ as it will for very large simulations, the $\Delta$ contribution to the overall computation time decreases and the overall scaling of algorithm **S1** approaches the optimal $N/2P$. In essence, each processor spends nearly all its time working in its own box and only exchanges a relatively small amount of information with neighboring processors to update its boundary conditions.

An important feature of algorithm **S1** is that the lists and structure of the data are only changed once every few timesteps when neighbor lists are constructed. In particular, even if an atom moves outside box $z$'s boundaries it is not reassigned to a new processor until step (1a) is executed [45]. Processor $P_z$ can still compute correct forces for the atom so long as two criteria are met. The first is that an atom cannot move farther than $d$ between two neighbor list constructions since this would cause problems for step (1a). The second is that all nearby atoms within a distance $r_s$, instead of $r_c$, must be updated every timestep. The alternative is to move atoms to their new processors at every timestep [34]. This has the advantage that only atoms within a distance $r_c$ of box $z$ need be exchanged at all timesteps when neighbor lists are not constructed. This reduces the volume of communication since $r_c < r_s$. However, now the neighbor list of a reassigned atom must also be sent. The information in the neighbor list is atom indices referencing local memory locations where the neighbor atoms are stored. If atoms are continuously moving to new processors, these local indices become meaningless. To overcome this, our implementation in [34] assigned a global index (1 to $N$) to each atom which moved with the atom from processor to processor. A mapping of global index to local memory must then be stored in a vector of size $N$ by each processor or the global indices must be sorted and searched to find the correct atoms when they are referenced in a neighbor list. The former solution

limits the size of problems that can be run; the latter solution incurs a considerable cost for the sort and search operations. We found that implementing the Tamayo and Giles idea [45] in our algorithm S1 made the resulting code less complex and reduced the computational and communication overhead. This did not affect the timings for simulations with large $N$, but improved the algorithm's performance for medium-sized problems.

A modified version of **S1** that takes full advantage of Newton's 3rd law can also be devised, call it algorithm **S2**. If processor $P_z$ acquires atoms only from its west, south, and down directions (and sends its own atoms only in the east, north, and up directions), then each pairwise interaction need only be computed once, even when the two atoms reside in different boxes. This requires sending computed force results back in the opposite directions to the processors who own the atoms, as a step (3) in the algorithm. This scheme does not reduce communication costs, since half as much information is communicated twice as often, but does eliminate the duplicated force computations for two-box interactions. An algorithm similar to this is detailed in [12] with excellent results for the Fujitsu AP1000 machine that we note in the next section. Two points are worth noting. First, the overall savings of **S2** over **S1** is small, particularly for large $N$. Only the $\Delta$ term in steps (1c) and (2) is saved. Second, as we mention in our conclusions, the real speed to be gained in spatial-decomposition algorithms for large systems is by improving the single-processor performance of force computation in step (2). As floating point processors in parallel machines become more sophisticated more attention must be paid to data structures and loop orderings in the force and neighbor-list construction routines to achieve high single-processor flop rates. Implementing **S2** requires special-case coding for atoms near box edges and corners to insure all interactions are counted only once [12] and thus affects this optimization process.

Finally, the issue of load-balance is an important concern in any spatial-decomposition algorithm. Algorithm **S1** will be load-balanced only if all boxes have a roughly equal number of atoms (and surrounding atoms). This will not be the case if the physical atom density is non-uniform. Additionally, if the physical domain is not a rectangular parallelepiped, it can be difficult to split into $P$ equal-sized pieces. Sophisticated load-balancing algorithms have been developed [25] to partition an irregular physical domain or non-uniformly dense clusters of atoms, but in general they create sub-domains which are irregular in shape or are connected in an irregular fashion to their neighboring sub-domains. In either case, the task of assigning atoms to sub-domains and communicating with neighbors becomes more costly and complex. If the physical atom density changes over time during the MD simulation, the load-balance problem is compounded. Any dynamic load-balancing scheme requires additional computational overhead and data movement.

In summary, the spatial-decomposition algorithm, like the atom- and force-decomposition algorithms, evenly divides the MD computations across all the processors. Its chief benefit is that it takes full advantage of the local nature of the interatomic forces by performing only local communication. Thus, in the large $N$ limit, it achieves optimal $O(N/P)$ scaling and is clearly the fastest algorithm. However, this is only if good load-balance is also achievable. Since its performance is sensitive to the problem geometry, algorithm **S1**

22

is more restrictive than **A2** and **F2** whose performance is geometry-independent. A second drawback of algorithm **S1** is its complexity; it is more difficult to implement efficiently than the simpler atom- and force-decomposition algorithms. In particular the communication scheme requires extra coding and bookkeeping to create messages and access data received from neighboring boxes. In practice, integrating algorithm **S1** into an existing serial MD code can require a substantial reworking of data structures and code.

# 6 Benchmark Problem

The test case used to benchmark our three parallel algorithms is a MD problem that has been used extensively by various researchers [12, 22, 28, 34, 41, 46, 45]. It models atom interactions with a Lennard-Jones potential energy between pairs of atoms separated by a distance $r$ as

$$\Phi(r) = 4\epsilon \left[ (\frac{\sigma}{r})^{12} - (\frac{\sigma}{r})^6 \right] \tag{2}$$

where $\epsilon$ and $\sigma$ are constants. The derivative of this energy expression with respect to $r$ is the $F_2$ term in equation (1); $F_3$ and higher-order terms are ignored.

The $N$ atoms are simulated in a 3-D parallelepiped with periodic boundary conditions at the Lennard Jones state point defined by the reduced density $\rho^* = 0.8442$ and reduced temperature $T^* = 0.72$. This is a liquid state near the Lennard-Jones triple point. The simulation is begun with the atoms on an $fcc$ lattice with randomized velocities chosen from a Boltzmann distribution. The solid quickly melts as the system evolves to its natural liquid state. A roughly uniform spatial density persists for the duration of the simulation. The simulation is run at constant $N$, volume $V$, and energy $E$, a statistical sampling from the microcanonical ensemble. Force computations using the potential in equation (2) are truncated at a distance $r_c = 2.5\sigma$. The integration timestep is 0.00462 in reduced units. For simplicity we use a leapfrog scheme to integrate equation (1) as in [2]. Other implementations of the benchmark [22] have used predictor-corrector schemes; this only slows their performance by 2-3%.

For timing purposes, the critical features of the benchmark for a given problem size $N$ are $\rho^*$ and $r_c$. These determine how many force interactions must be computed at every timestep. The number of atoms in a sphere of radius $r^* = r/\sigma$ is given by $4\pi\rho^*(r^*)^3/3$. For this benchmark, using $r_c = 2.5\sigma$, each atom has on average 55 neighbors. If neighbor lists are used, the benchmark also defines an extended cutoff length $r_s = 2.8\sigma$ (encompassing about 78 atoms) for forming the neighbor lists and specifies that the lists be created or updated every 20 timesteps. Timings for the benchmark are usually reported in CPU seconds/timestep. If neighbor lists are used then the cost of creating them every 20 steps is amortized over the per timestep timing.

It is worth noting that without running a standard benchmark problem it can be difficult to accurately assess the performance of a parallel algorithm. In particular, it can be misleading to only compare performance of a parallel version of a code to the original vectorized or serial code because, as we have learned from

23

our codes as well as other's results, the vector code performance may well be far from optimal. Even when problem specifications are reported, it can be difficult to compare two algorithm's relative performance when two different benchmark problems are used. This is because of the wide variability in the cost of calculating force equations, the number of neighbors included in cutoff distances, and the frequency of neighbor list building as a function of temperature, atom density, cutoff distances, etc.

# 7  Results

The parallel algorithms of Sections 3, 4, and 5 were tested on three parallel MIMD supercomputers, a nCUBE 2, an Intel iPSC/860, and the Intel Delta. The first two machines are at Sandia; the Delta is at Cal Tech. The nCUBE 2 is a 1024-processor hypercube. Each processor is capable of about 2 Mflops and has 4 Gbytes of memory. Sandia's iPSC/860 has 64 i860 processors connected in a hypercube topology. Its processors have 8 Mbytes of memory and are capable of about 60 Mflops, but in practice 5-10 Mflops is the typical compiled Fortran performance. The Intel Delta has 512 processors configured as a 2-D mesh. The individual processors have 16 Mbytes of memory and are identical to those in the iPSC/860, though the communication network is somewhat faster.

Because the algorithms were implemented in standard Fortran with calls to vendor-supplied message-passing subroutines, only minor changes were required to implement the benchmark codes on the different machines. The algorithms as described do not specify a mapping of processors to the computational elements (force matrix sub-blocks, 3-D boxes, etc.). The mapping could potentially be tailored for a particular machine architecture to minimize message contention (multiple messages using the same communication wire) and the distance messages have to travel between pairs of processors that are not directly connected by a communication wire. We chose mappings that are simple and good choices for hypercubes. For code portability we used the same mappings on the mesh-architecture Delta.

For the atom-decomposition algorithm we simply assign the processors in ascending order to the row-blocks of the force matrix in Figure 1. The expands and folds then take place exactly as in Figure 2. For the force-decomposition algorithm we use a natural calendar ordering of the processors to the force matrix in Figure 7. On a hypercube this means each row and column of the matrix is a sub-cube of processors so that expands and folds within rows and columns can be done optimally. However, the transpose operations in algorithms **F1** and **F2** now require communication between pairs of processors that are architecturally distant. Thus with this mapping there will be some message contention during the transposes as multiple messages route to their distant destinations simultaneously. Since the transpose operations scale as the volume of data exchanged or $N/P$, even with some slow-down due to message congestion, the overall $N/\sqrt{P}$ scaling of the communication portion of the force-decomposition algorithms is not affected. Although we did not implement it for this work, a mapping of processors to the force matrix that produces contention-free transposes for a hypercube is possible and is described in [27]. We have also recently developed a modified force-decomposition algorithm that does not require transpose operations and thus runs slightly faster on

24

mesh–architecture parallel machines [35].

For the spatial–decomposition algorithm, we use a processor mapping that essentially configures a hypercube as a 3–D mesh. Such a mapping is done using a Gray–coded ordering [19] of the processors. This insures each processor's box in Figure 10 has 6 spatial neighbors (boxes in the east, west, north, south, up, down directions) that are assigned to processors which are also nearest neighbors in the hypercube topology. Communication with these neighbors is thus contention–free and as fast as possible. Gray–coding also provides naturally for periodic boundary conditions in the MD simulation since processors at the edge of the 3–D mesh are topological nearest neighbors to those on the opposite edge. The only restriction the Gray–coding imposes is that the number of processors assigned to each dimension of the 3–D mesh be a power–of–two. For the Intel Delta there is no obvious best way to map a 3–D problem to its 2–D mesh of processors. We use the same 3–D Gray–coding assignment scheme for code portability.

Timing results for the benchmark problem on the different parallel machines are shown in Tables I, II, and III for the atom–, force , and spatial–decomposition algorithms. A wide range of problem sizes are considered from $N = 500$ atoms to $N = 10^7$ atoms. The lattice size for each problem is also specified; there are 4 atoms per unit cell for the $fcc$ lattices. Entries with a dashed line are for problems that would not fit in available memory. The last entries in each table are roughly the largest problem sizes for this benchmark that can be run with each algorithm due to memory restrictions on the three parallel machines.

For comparison, we also implemented the vectorized algorithm of Grest, et. al. [22] on single processors of Sandia's Cray Y–MP and a Cray C90 at Cray Research. Our version is slightly different from the original Grest code, using a simpler integrator and allowing for non–cubic physical domains. The timings in reference [22] were for a Cray X–MP. We believe these timings for the faster Y–MP and C90 architectures are the fastest that have been reported for this benchmark problem on a single processor of a convential vector supercomputer. They show a C90 processor to be about 2.5 times faster than a Y–MP processor for this algorithm. The starred Cray timings in the tables are estimates for problems too large to fit in memory on the machines accessible to us. They are extrapolations of the $N = 10^5$ system timing based on the observed linear scaling of the Cray algorithm. It is also worth noting that ideas similar to those used in the parallel algorithms of the previous sections could be used to create efficient parallel Cray codes for multiple processors of a Y–MP or C90. For example, a speed–up of 6.8 on a 8–processor Cray Y–MP has been achieved by Kremer with the Grest, et. al. algorithm [20].

The parallel timings in the three tables are all for single–precision (32–bit) implementations of the benchmark. The Cray timings are, of course, for 64–bit arithmetic since that is the only option. MD simulations do not typically require double precision accuracy since there is a much coarser approximation inherent in the potential model and the integrator. This is particularly true of Lennard–Jones systems since the $\epsilon$ and $\sigma$ coefficients are only specified to a few digits of accuracy as an approximate model of the interatomic energies in a real material. With this said, double precision timings can be easily estimated for the parallel algorithms. The processors in all three of the parallel machines compute about 20–30% slower in double–precision arith-

| Problem Size | | C90 | Y-MP | nCUBE 2 | | Intel iPSC/860 | | Intel Delta |
|---|---|---|---|---|---|---|---|---|
| N | Lattice | P=1 | P=1 | P=512 | P=1024 | P=32 | P=64 | P=256 |
| 500 | 5x5x5 | .00373 | .00930 | .00724 | --- | .0111 | .00880 | .00518 |
| 2048 | 8x8x8 | .0154 | .0369 | .0252 | .0217 | .0446 | .0336 | .0172 |
| 4000 | 10x10x10 | .0232 | .0610 | .0458 | .0394 | .0807 | .0616 | .0314 |
| 6912 | 12x12x12 | .0425 | .106 | .0780 | .0669 | .138 | .103 | .0532 |
| 10976 | 14x14x14 | .0657 | .167 | .124 | .106 | .220 | .164 | .0863 |
| 16384 | 16x16x16 | .103 | .250 | .182 | .155 | .337 | .249 | .130 |
| 32000 | 20x20x20 | .202 | .470 | .351 | .301 | .635 | .474 | .256 |
| 50000 | 20x25x25 | .286 | .733 | .546 | .469 | .993 | .740 | .399 |
| 100,000 | 25x25x40 | .592 | 1.47 | 1.09 | .935 | 1.98 | 1.48 | .820 |

Table I: CPU seconds/timestep for the atom–decomposition algorithm **A1** on several parallel machines for the benchmark simulation. Single processor Cray Y–MP and C90 timings using a fully vectorized algorithm are also given for comparison.

metic than single, so the time spent computing would be increased by that amount. Communication costs in each of the algorithms would essentially double, since the volume of information being exchanged in messages would increase by a factor of two. Thus depending on the fraction of time being spent in communication for a particular $N$ and $P$ (see the scaling discussion below), the overall timings typically increase by 20–50% for double–precision runs.

The tables show the parallel machines to be competitive with the Cray Y–MP and C90 machines across the entire range of problem sizes for all three parallel algorithms. The force–decomposition algorithm is fastest for the smallest problem sizes; spatial–decomposition is fastest for large $N$. The Intel Delta is the fastest of the three parallel machines, up to 30 times faster than a single Y–MP processor on the largest problem sizes using the spatial–decomposition algorithm and 12 times faster than a C90 processor. The nCUBE 2 and Intel Delta can perform million atom simulations of the benchmark problem at 1.17 and .498 seconds/timestep respectively. A surprising result is that the parallel machines are competitive with a single processor of the Cray machines even for the smallest problem sizes. One typically does not think of there being enough parallelism to exploit when there are only a few atoms per processor. The hardware

| Problem Size | | C90 | Y-MP | nCUBE 2 | | Intel iPSC/860 | | Intel Delta | |
|---|---|---|---|---|---|---|---|---|---|
| N | Lattice | P=1 | P=1 | P=512 | P=1024 | P=32 | P=64 | P=256 | P=512 |
| 500 | 5x5x5 | .00373 | .00930 | .00592 | --- | .00980 | .00695 | .00480 | .00455 |
| 2048 | 8x8x8 | .0154 | .0369 | .0110 | .00864 | .0359 | .0250 | .00894 | .00677 |
| 6912 | 12x12x12 | .0425 | .106 | .0245 | .0179 | .112 | .0759 | .0250 | .0160 |
| 10976 | 14x14x14 | .0657 | .167 | .0394 | .0277 | .180 | .122 | .0399 | .0244 |
| 32000 | 20x20x20 | .202 | .470 | .0890 | .0603 | .521 | .349 | .115 | .0667 |
| 50000 | 20x25x25 | .286 | .733 | .162 | .112 | .828 | .544 | .179 | .103 |
| 100,000 | 25x25x40 | .592 | 1.47 | .251 | .171 | 1.75 | 1.10 | .369 | .210 |
| 500,000 | 50x50x50 | 2.86* | 7.33* | 2.47 | 1.66 | --- | 6.04 | 1.96 | 1.17 |
| 1,000,000 | 50x50x100 | 5.92* | 14.7* | --- | 3.29 | --- | --- | 4.04 | 2.41 |

Table II: CPU seconds/timestep for the force–decomposition algorithm **F2** on several parallel machines and the Cray Y–MP and C90.

performance monitor on the Cray also provides the useful metric that the vectorized algorithm runs at about 260 Mflops for large $N$ on a C90 processor. Though the floating point operation count is not identical in the parallel codes, this implies the 512–processor Intel Delta runs the spatial–decomposition algorithm at about 3 Gflops for the largest problem sizes, or about 6 Mflops/processor which is typical of i860 chip performance on compiled Fortran code. However, since basic linear algebra subroutines (BLAS) run at 20-30 Mflops on the i860 (out of a peak speed of 60 Mflops), this indicates much higher performance could still be achieved in parallel MD algorithms by writing code optimized for a particular parallel machine's floating point processor.

Timings for this benchmark on other parallel machines are also discussed in [7, 12, 46, 45], all for spatial–decomposition algorithms. Tamayo, et. al. [46] implemented several SIMD algorithms on a 32K–processor CM–2 (1024 floating point processors). Their fastest algorithm ran at 0.57 sec/timestep for a $N = 18000$ atom system, about a factor of two slower than the single processor Y–MP timing in the tables here. More recently, Tamayo and Giles [45] achieved a time of 0.4 sec/timestep on a $N = 51,200$ atom simulation on 256 processors of a CM–5. This is roughly a factor of two faster than a Y–MP processor and was for a CM–5 without vector units programmed in MIMD mode with explicit message passing; the timings should improve dramatically with the vector units. In other work on the CM–5, Beazley and Lomdahl [7] report a time of 1.00 sec/timestep for a $N = 10^6$ atom simulation on 1024 processors, also without vector units. This however, was for a shorter cutoff distance with only an average of 21 neighbors/atom instead of 55 as in this benchmark. Finally, Brown, et. al. [12] detail an algorithm similar to the **S2** algorithm discussed in Section 5. For a $N = 729000$ atom system (at a slightly smaller reduced density of $\rho^* = 0.8$) run on 512 processors

| Problem Size | | C90 | Y-MP | nCUBE 2 | | Intel iPSC/860 | | Intel Delta | |
| N | Lattice | P=1 | P=1 | P=512 | P=1024 | P=32 | P=64 | P=256 | P=512 |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| 500 | 5x5x5 | .00373 | .00930 | .0130 | .0119 | .0129 | .0106 | .00706 | .00592 |
| 2048 | 8x8x8 | .0154 | .0369 | .0173 | .0148 | .0321 | .0189 | .00837 | .00650 |
| 6912 | 12x12x12 | .0425 | .106 | .0374 | .0250 | .0768 | .0436 | .0159 | .0111 |
| 16384 | 16x16x16 | .103 | .250 | .0650 | .0407 | .161 | .0874 | .0275 | .0167 |
| 50000 | 20x25x25 | .286 | .733 | .160 | .0967 | .420 | .224 | .0664 | .0380 |
| 100,000 | 25x25x40 | .592 | 1.47 | .298 | .165 | .798 | .418 | .119 | .0678 |
| 500,000 | 50x50x50 | 2.86* | 7.33* | 1.17 | .650 | 3.66 | 1.88 | .501 | .261 |
| 1,000,000 | 50x50x100 | 5.92* | 14.7* | 2.23 | 1.17 | --- | 3.68 | .951 | .498 |
| 5,000,000 | 100x100x125 | 28.6* | 73.3* | 10.2 | 5.28 | --- | --- | 4.45 | 2.31 |
| 10,000,000 | 100x125x200 | 59.2* | 147.* | --- | 10.2 | --- | --- | --- | 4.60 |

Table III: CPU seconds/timestep for the spatial-decomposition algorithm S1.

of the Fujitsu AP1000 they report a time of 0.927 ser /timestep.

The timings in Table I show that communication c. ¬.s have begun to dominate in the atom–decomposition algorithm by the time hundreds of processors are used. There is little speed up gained by doubling the number of processors used. By contrast timings in Table II show the force–decomposition algorithm is speeding up by roughly 30% when the number of processors is doubled. The timings for the largest problem sizes in Table III evidence excellent scaling properties. Doubling $P$ nearly halves the run times for a given $N$. Similarly, as $N$ increases for fixed $P$, the run times per atom actually become faster as the surface–to–volume ratio of each processor's box is reduced. We note, however, that this scaling depends on uniform atom density within a simple domain such as the rectangular parallelepiped of the benchmark problem.

A comparison of the different algorithm's performance using data from all 3 tables can be better displayed in graphical form. Figure 12 shows the nCUBE 2's performance on the benchmark simulation on 1024 processors as a function of problem size. Single processor Y–MP timings are also included. The linear scaling of all the algorithms when $N$ is large is evident. Note that force–decomposition is faster than atom–decomposition across all problem sizes due to its reduced communication costs. On this many processors, the spatial–decomposition algorithm has significant overhead costs for small $N$. This is because the $d/r_s$ ratio is so small that each processor has to communicate with a large number of neighboring processors to acquire all its needed information. As $N$ increases, this overhead is reduced relative to the computation performed

28

inside the processor's box, and the algorithm's performance asymptotically approaches its optimal $O(N/P)$ performance. Thus there is a cross–over size $N$ at which the spatial–decomposition algorithm becomes faster than force–decomposition. We return to this point in the conclusions section.
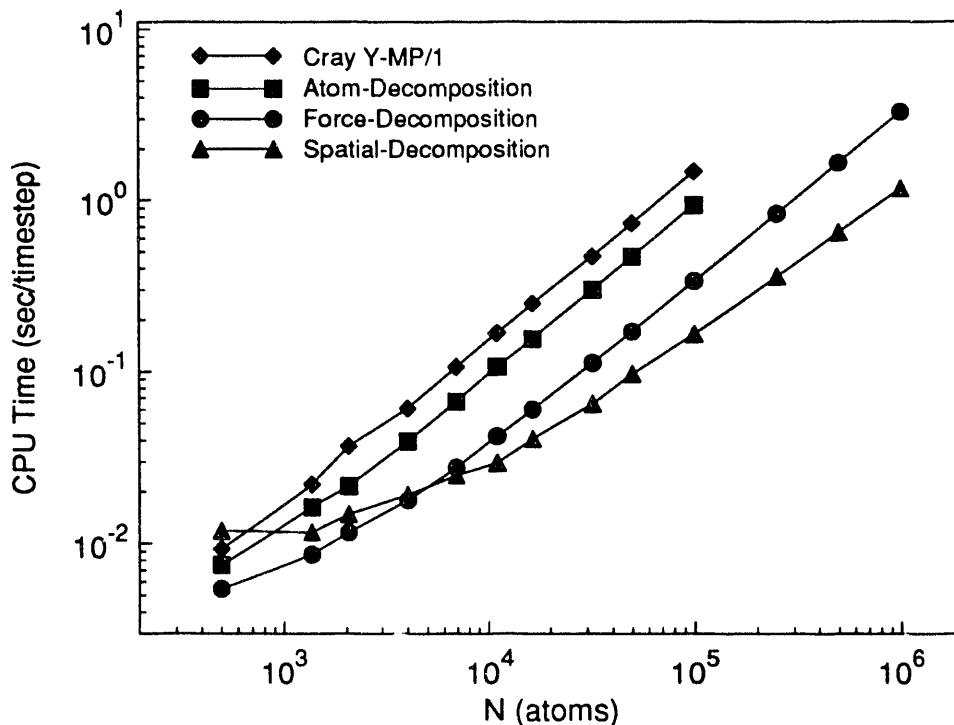


Figure 12: CPU timings (seconds/timestep) for the three parallel algorithms on 1024 processors of the nCUBE 2 for different problem sizes. Single-processor Cray Y–MP timings are also given for comparison.

In Figure 13 we plot the Intel Delta's performance on the $N = 10976$ atom benchmark as a function of number of processors. The single–processor i860 and Y–MP timings are also shown; the Y–MP is about 13.3 times faster than a i860 processor on this problem. The dotted line is the maximum achievable speed of the Delta if any of the algorithms were 100% efficient. Parallel efficiency is defined as the run time on 1 processor divided by the quantity $(P \times$ run time on $P$ processors). Thus if the 512–processor timing is 256 times as fast as the 1–processor timing, the algorithm is 50% efficient. On small numbers of processors communication is not a significant factor and all the algorithms perform similarly. But as $P$ increases, the algorithms become less efficient. The atom–decomposition falls off most rapidly due to the $O(N)$ scaling of its communication. On the Delta's large 2–D mesh the all–to–all communication this algorithm requires is particularly inefficient (because of message contention), causing a slow–down when going from 256 to

512 processors. Force–decomposition is next most efficient due to its $O(N/\sqrt{P})$ communication scaling. It remains competitive with the spatial–decomposition algorithm across a wide range of numbers of processors. When hundreds or thousands of processors are used, even the spatial–decomposition algorithm becomes less efficient since now the box size is small relative to the force cutoff distance for this $N$. It is worth noting that the trends in the plots of Figures 12 and 13 are the same for the other machines and problem sizes tested in this study. Though the absolute data values are functions of $N$, $P$, and the benchmark attributes, the relative trade–offs between the various algorithms are consistently the same.

Using one–node timings on the nCUBE and Intel machines as reference points, parallel efficiencies can be computed for all the algorithms. The nCUBE 2 one–processor timing is $9.15 \times 10^{-4}$ seconds/timestep/atom. Both Intel machines give a one–processor timing of $2.03 \times 10^{-4}$ seconds/timestep/atom. Because the algorithms and codes scale so linearly, these values can be used to predict optimal timings for problems larger than will fit on a single processor. For the million–atom simulation, the timings in Table III show the spatial–decomposition algorithm **S1** thus has a parallel efficiency of 76% on 1024 processors of the nCUBE and 80% on 512 processors of the Intel Delta. The larger simulations approach roughly a 90% parallel efficiency. To put these numbers in context, consider that on the nCUBE, the million–atom simulation means each processor has about 1000 atoms in its box. But the range of the cutoff distance in the benchmark is such that about 2600 atoms from surrounding boxes are still needed at every timestep to compute forces. Thus the spatial–decomposition algorithm **S1** is 76% efficient even though two–and–a–half times as many atom positions are communicated as are updated locally by each processor.

Finally, we discuss the scalability of the different parallel algorithms in the large $N$ limit. Table IV shows the overall scaling of the computation and communication portions of the 5 algorithms. This is constructed from the scaling entries for the various steps of the algorithms in Figures 5, 6, 8, 9, and 11, using large $N$ values when there is an option. Some coefficients are included to show contrasts between the various algorithms. The amount of memory required per processor to store atom position and force vectors is also listed in the table.

Computation in the atom–decomposition algorithm **A1** scales as $N/P + N$ where the second term is for binned neighbor list construction. The coefficient on this term is small so it is usually not a significant factor. The communication scales as $N$, as does the memory to store all atom positions. By contrast, atom–decomposition algorithm **A2** implements Newton's 3rd law so its leading computational term is cut in half. Now the communication cost is doubled and the entire force vector must be stored on each processor as well.

Force–decomposition algorithms **F1** and **F2** have the same computational complexity as **A1** and **A2** respectively except the binning for neighbor list construction now scales as $N/\sqrt{P}$, again not typically a significant factor. In **F1** there are 3 expands/folds and one transpose operation for a total communication cost of $3N/\sqrt{P} + N/P$. Similarly **F2** requires 4 expands/folds and 2 transposes. Implementing **F1** requires storing two atom position sub–vectors and one force sub–vector, all of length $N/\sqrt{P}$. **F2** requires an extra
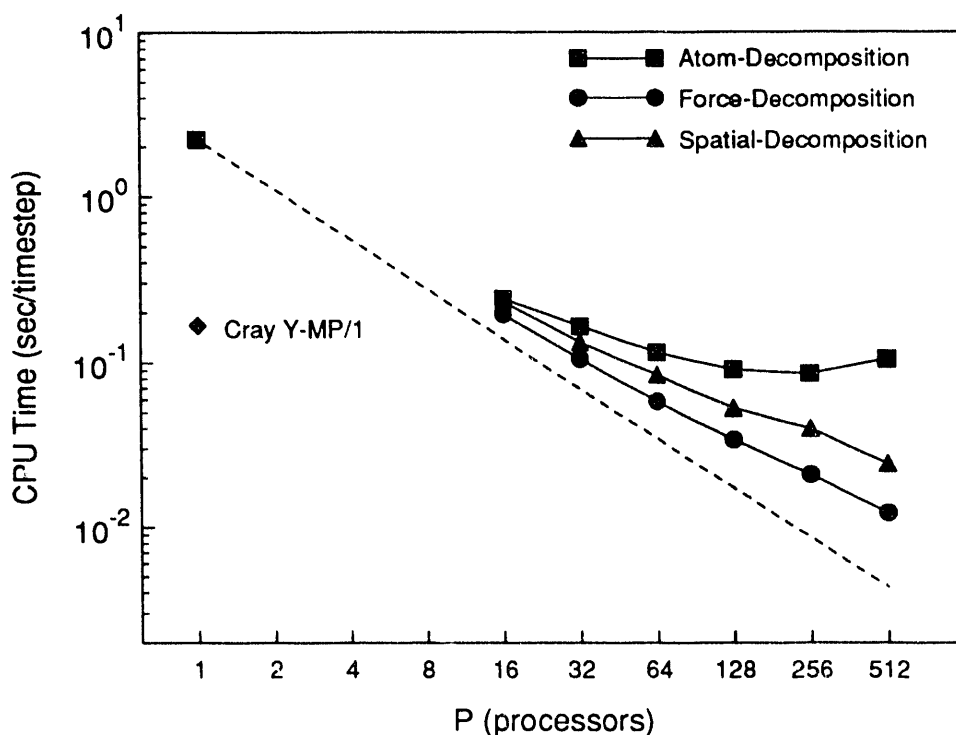
Figure 13: CPU timings (seconds/timestep) for the three parallel algorithms on the Intel Delta for different numbers of processors on a benchmark simulation with $N = 10976$ atoms. Single–processor i860 and Cray Y–MP timings are shown for comparison.

force sub–vector.

Computation in the spatial–decomposition algorithm S1 scales as $N/2P$ since it implements Newton's 3rd law for interactions between atom pairs inside a processor's box. For large N problems there is an extra factor for computations performed on nearby atoms within a distance $r_s$ of the box faces. The number of these atoms is proportional to the surface area of the box face $(N/P^{2/3})$ times $r_s$ for each of the 6 faces. The communication in algorithm S1 scales as the same factor as do the memory requirements for storing the nearby atoms. Additionally, $O(N/P)$ memory must be allocated for storing the atoms in a processor's box.

As mentioned above, the memory entries in Table IV are for the atom position and force vectors used in each algorithm. However, in all of the algorithms processors must store additional $O(N/P)$ information such as neighbor lists and velocities for the $N/P$ atoms they own. In practice, for the force– and spatial– decomposition algorithms, storage of neighbor lists is actually the dominant factor in limiting the size of

| Algorithm | Computation | Communication | Memory |
|-----------|-------------|---------------|--------|
| A1 | $\dfrac{N}{P}+N$ | $N$ | $N$ |
| A2 | $\dfrac{N}{2P}+N$ | $2N$ | $2N$ |
| F1 | $\dfrac{N}{P}+\dfrac{N}{\sqrt{P}}$ | $3\dfrac{N}{\sqrt{P}}+\dfrac{N}{P}$ | $3\dfrac{N}{\sqrt{P}}$ |
| F2 | $\dfrac{N}{2P}+\dfrac{N}{\sqrt{P}}$ | $4\dfrac{N}{\sqrt{P}}+2\dfrac{N}{P}$ | $4\dfrac{N}{\sqrt{P}}$ |
| S1 | $\dfrac{N}{2P}+6r_s\left(\dfrac{N}{P}\right)^{2/3}$ | $6r_s\left(\dfrac{N}{P}\right)^{2/3}$ | $\dfrac{N}{P}+6r_s\left(\dfrac{N}{P}\right)^{2/3}$ |

Table IV: Scaling properties of all 5 parallel algorithms as a function of problem size $N$ and number of processors $P$. Run time scaling for the communication and computation portions of the algorithms as well as their per-processor memory requirements are listed.

problem that can be run. For example, the 4–6 Gbytes of user memory available on the nCUBE 2 and Intel Delta machines used in this study allows a $10^7$ atom simulation to be run with the $r_c = 2.5\sigma$ cutoff of the benchmark problem. Using a shorter cutoff on the 1024-processor CM–5 (which has more user memory) allowed Beazley and Lomdahl to run a 70-million atom 3–D simulation [7]. These figures are indicative of the size of problems that can be run on current-generation parallel machines.

# 8  Application of the Algorithms

While the benchmark problem discussed in Sections 6 and 7 is relatively simple, the parallel algorithms described in this paper can be used in a variety of more complex MD simulations with little modification. We discuss the parallel implications of some common MD issues in the next several paragraphs.

(A) Force models more computationally expensive than Lennard–Jones potentials are often used in MD simulations of various materials. Pairwise forces, even if they are very expensive, can often be pre–computed once and then stored in table form or as a set of interpolating coefficients. Then they turn out to be little

more expensive to compute with than 6-12 potentials. Modern parallel computers have ample memory for storing quite large tables of force values and/or coefficients in duplicate on every processor.

(B) Force models that are functions of atom velocities, or other quantities besides just atom positions, are sometimes used. An example is the embedded atom method (EAM) potentials [16] commonly used in modeling metals and metal alloys where an atom's energy is a function of electron density contributions from neighboring atoms as well as conventional pair-potential interactions. A more general N-body simulation example is vortex methods in fluid dynamics where "particles" of fluid interact via their vorticities. All of the parallel algorithms described here can be augmented in steps (3) and (5) to communicate additional atom-based quantities as needed [36] without affecting their overall parallel scaling.

(C) More sophisticated many-body force models are often used in MD simulations of covalently bonded materials. Examples include angular (three-body) forces for silicon and torsional (four-body) forces for organic polymers or proteins. These forces can be most easily computed in parallel if a single processor knows the positions of all the atoms in a particular bond group. This is simple to implement in the atom-decomposition algorithm since each processor knows the positions of every atom. Likewise, the spatial-decomposition algorithm can also be modified to insure each processor acquires enough information from surrounding boxes to compute all the many-body terms its atoms are a party to, since the bond groups are short-range in nature. However, the force-decomposition algorithm requires special care in this respect. This is because a processor only knows the positions of $2N/\sqrt{P}$ atoms that have no special spatial relationship to each other. One solution is to perform a pre-processing step to reorder the atoms in such a way that one or more processors will know the positions of all the atoms in each bond group. We have developed methods for doing this in organic MD simulations where the connectivity of the bond groups is static [37, 35]. However, we know of no simple way to use the force-decomposition idea for the more general case of dynamically changing connectivities, such as for silicon three-body potentials.

(D) Though force calculation is the key computational kernel in MD simulations, the quantities of interest are often global parameters like pressure, structure factors, and diffusion coefficients. These thermodynamic and transport properties are often calculated once every 50 or 100 timesteps and add little to the overall computational cost of a serial program. The same is true for the parallel case. In short-range MD each processor can compute its partial contribution to one of these quantities from the atom information it already knows. Then the local values can be accumulated quickly as a global sum across all the processors.

(E) In many MD codes, neighbor list construction is triggered by atom movement. For example, lists will only be recreated when some atom has moved half the distance $r_s - r_c$. Again, this can easily be implemented in the parallel algorithms by having each processor check if any of its $N/P$ atoms have met the criterion, then exchanging a global flag to decide if the neighbor list routines should be called. If the list of interacting neighbors is static in a particular MD simulation (e.g. atoms on a lattice), then step (1) in all of the parallel algorithms becomes unnecessary. The remaining steps of the algorithms are still a fast way to parallelize the necessary computation and communication for this special case.

(F) The benchmark problem implements a constant $N$, volume $V$, and energy $E$ microcanonical ensemble. Another common choice is to hold $N$, pressure $P$, and temperature $T$ constant, sampling from the canonical ensemble. This involves rescaling the simulation domain dimensions and velocities at each timestep (or every few timesteps) to hold the pressure and temperature constant. In parallel this requires a small amount of additional communication, a global summation or exchange of the rescaling parameters, similar to the effort involved in (D) and (E) above.

(G) A simple leapfrog integrator was used in our implementation of the benchmark problem. More complex ODE integrators such as Runge–Kutta or predictor–corrector methods can easily be used in the context of any of the parallel algorithms in step (4). These methods are perfectly parallel since they only require information about the $N/P$ atoms already owned by each processor. They may also require extra $O(N/P)$ memory to store old timestep values or work vectors.

(H) Multiple–timescale MD methods have been proposed [44], where work is done at staggered times on different length scales to allow longer timesteps to be taken on average. Only very short–range information is used to compute forces in the smallest timesteps. These schemes are an effort to include longer–range effects while avoiding true long–range force computation. They are typically implemented by a hierarchy of neighbor lists which store information for the different length scales. Since they are still inherently short–range force models, they can be implemented within the general framework of any of the parallel algorithms we have presented. In the limit that the force computation becomes truly long–range in nature, pairwise forces are usually not the computational method of choice as discussed in Section 2. However, if long–range pairwise forces are used, they can still be computed directly in the force–matrix formalism of the atom– and force–decomposition algorithms [27]. By contrast the spatial–decomposition algorithm would now require long–range communication and become an inefficient solution.

# 9  Conclusion

We have detailed the construction and implementation of three kinds of parallel algorithms for MD simulations with short-range forces. Each of them has advantages and disadvantages. The atom–decomposition algorithm is simplest to implement and load–balances automatically, but because it performs all–to–all communication, its communication costs begin to dominate its run time on large numbers of processors. The force–decomposition algorithm is also relatively simple, though it often requires some pre–processing to assure load–balance. It also works well independent of the physical problem's geometry. Its $O(N/\sqrt{P})$ scaling is better than that of the atom–decomposition algorithm, but is not optimal for large simulations. The spatial–decomposition algorithm does exhibit optimal $O(N/P)$ scaling for large problems. However it suffers more easily from load–imbalance and is more difficult to implement efficiently.

In practical terms, how does one choose the "best" parallel algorithm for a particular MD simulation? Assuming one knows the ranges of $N$ and $P$ the simulation will be run with, we find the following four guidelines helpful.

(A) Choose an atom–decomposition algorithm only if the communication cost is expected to be negligible. In this case simplicity outweighs the inefficient communications. Typically this will only be true for small $P$ (say $P \leq 16$ processors) or very expensive forces where computation time dominates communication time.

(B) A force–decomposition approach will be faster than atom–decomposition in all other cases. Both the atom– and force–decomposition algorithms scale linearly with $N$ for fixed $P$. This means for a given $P$, the parallel efficiency of either algorithm is independent of $N$. Moreover, as $P$ doubles, the efficiency of the communication portion of the atom–decomposition algorithm goes down by a factor of 2, while the force–decomposition algorithm's efficiency decreases by a factor of only $\sqrt{2}$. Thus, once $P$ is large enough that force–decomposition is significantly faster than atom–decomposition, it will remain faster as $P$ increases, independent of $N$. For the benchmark problem this was the case for $P \geq 16$ processors.

(C) For a given $P$, the scaling of the spatial–decomposition algorithm is not linear with $N$. For small $N$ communication and overhead costs are significant and the efficiency is poor; for large $N$ the efficiency is asymptotically optimal (100%). Thus when compared to a force–decomposition approach, there will be some cross–over point as $N$ increases for a given $P$ where a spatial–decomposition algorithm becomes faster. In the benchmark the cross–over size was several thousands of atoms on hundreds of processors, as in Figure 12. In general, the cross–over size is a function of the complexity of the force model, force cutoff distances, and the computational and communication capabilities of a particular parallel machine. It will also be a function of $P$. For example, if the force cutoff distance is reduced (say to $r_c = 2^{1/6}\sigma$ as is common for fluids modeled with purely repulsive forces), a spatial–decomposition algorithm will require less exchange of information between processors assigned to neighboring boxes and the cross–over size will be reduced as well. A rough estimate of a spatial–decomposition algorithm's efficiency for a given $N$ and $P$ can be made by noting each processor's box has volume $d^3 = N/P$, but it computes and communicates information in an extended volume of $(d + 2r_s)^3$. Comparing the extended volume to the box volume gives a rough measure of the extra (inefficient) work the algorithm is performing.

(D) The preceding paragraph assumes the computation in the spatial–decomposition algorithm is perfectly load–balanced. Load–imbalance imposes an upper bound on the efficiency a spatial–decomposition algorithm can achieve. For example, biological simulations of proteins solvated by water often are performed in a vacuum so that the atoms in the simulation fill a roughly spherical volume. If this domain is treated as a cube and split into $P$ pieces then the sphere fills only a $\pi/6$ fraction of the cube and a 50% parallel inefficiency results. The net effect of load–imbalance is to increase the cross–over size at which a spatial–decomposition algorithm becomes faster than a force–decomposition approach. In practice, we have found the force decomposition algorithm can be faster or at least quite competitive with spatial–decomposition algorithms for simulations of up to many tens of thousands of atoms [37].

In Section 7 we discussed the performance of the parallel algorithms on three different parallel computers, the nCUBE 2, Intel iPSC/860, and Delta. The results show that current–generation parallel machines are competitive with multi–processor Cray–class vector supercomputers for short–range MD simulations. More

generally, these algorithms can be implemented on any parallel computer that allows its processors to execute code independently of each other and exchanges data between processors by standard message–passing techniques. This is the definition of a multiple instruction/multiple data (MIMD) parallel architecture. Most of the current– and next–generation parallel supercomputers support this mode of programming, including the Intel Paragon, CM–5, and Cray MPP machines. Several features of the algorithms take advantage of the flexibility of the MIMD paradigm, including the code to build and access variable–length neighbor lists via indirect addressing, to select/pack/unpack data for messages, and to efficiently exchange variable–length data structures between sub–groups of processors as in Figures 2 and 10.

Finally, we are confident these algorithms or versions based on similar ideas will continue to be good choices for MD simulations on parallel machines of the future. Optimizing their performance for next–generation machines will require improving their single–processor computational performance. As the individual processors used in parallel machines become faster and more complex, high computational rates can only be achieved by writing pipelined or vectorized code. Thus, many of the data reorganization and other optimization techniques that have been developed for MD on vector machines [22] will become important for parallel implementations as well.

# 10 Acknowledgments

# References

[1] F. F. Abraham. Computational statistical mechanics: methodology, applications and supercomputing. *Advances in Physics*, 35:1–111, 1986.

[2] M. P. Allen and D. J. Tildesley. *Computer Simulation of Liquids*. Clarendon Press, Oxford, 1987.

[3] D. J. Auerbach, W. Paul, A. F. Bakker, C. Lutz, W. E. Rudge, and F. F. Abraham. A special purpose parallel computer for molecular dynamics: Motivation, design, implementation, and application. *J. Phys. Chem.*, 91:4881–4890, 1987.

[4] A. F. Bakker, G. H. Gilmer, M. H. Grabow, and K. Thompson. A special purpose computer for molecular dynamics calculations. *J. Comp. Phys.*, 90:313–335, 1990.

[5] J. Barnes and P. Hut. A hierarchical $O(N \log N)$ force-calculation algorithm. *Nature*, 324:446–449, 1986.

[6] M. Baskes, M. Daw, B. Dodson, and S. Foiles. Atomic-scale simulation in materials science. *Materials Research Society Bulletin*, pages 28–34, Feb 1988.

[7] D. M. Beazley and P. S. Lomdahl. Message-passing multi-cell molecular dynamics on the Connection Machine 5. Technical Report LA–UR–92–3158, Los Alamos National Laboratories, Los Alamos, NM, 1993.

[8] R. H. Bisseling and J. G. G. van de Vorst. Parallel LU decomposition on a transputer network. In G. A. van Zee and J. G. G. van de Vorst, editors, *Lecture Notes in Computer Science, Number 384*, pages 61–77. Springer-Verlag, 1989.

[9] B. M. Boghosian. Computational physics on the Connection Machine. *Comp. in Phys.*, Jan/Feb, 1990.

[10] L. L. Boyer and G. S. Pawley. Molecular dynamics of clusters of particles interacting with pairwise forces using a massively parallel computer. *J. Comp. Phys.*, 78:405–423, 1988.

[11] B. R. Brooks and M. Hodošček. Parallelization of CHARMM for MIMD machines. *Chemical Design Automation News*, 7:16–22, 1992.

[12] D. Brown, J. H. R. Clarke, M. Okuda, and T. Yamazaki. A domain decomposition parallelization strategy for molecular dynamics simulations on distributed memory machines. *Comp. Phys. Comm.*, 74:67–80, 1993.

[13] J. P. Brunet, A. Edelman, and J. P. Mesirov. Hypercube algorithms for direct N-body solvers for different granularities. *SIAM J. Sci. Dist. Comp.*, 1993. to appear.

[14] J. P. Brunet, J. P. Mesirov, and A. Edelman. An optimal hypercube direct N-body solver on the Connection Machine. In *Proc. Supercomputing '90*, pages 748–752. IEEE Computer Society Press, 1990.

[15] T. W. Clark, J. A. McCammon, and L. R. Scott. Parallel molecular dynamics. In *Proc. 5th SIAM Conference on Parallel Processing for Scientific Computing*, pages 338–344. SIAM, 1992.

[16] M. S. Daw and M. I. Baskes. Embedded-atom method: Derivation and application to impurities, surfaces, and other defects in metals. *Phys. Rev. B*, 29:6443–6453, 1984.

[17] H. Q. Ding, N. Karasawa, and W. A. Goddard III. Atomic level simulations on a million particles: The cell multipole method for Coulomb and London interactions. *J. Chem. Phys.*, 97:4309, 1992.

[18] D. Fincham. Parallel computers and molecular simulation. *Molec. Sim.*, 1:1–45, 1987.

[19] G. C. Fox, M. A. Johnson, G. A. Lyzenga, S. W. Otto, J. K. Salmon, and D. W. Walker. *Solving Problems on Concurrent Processors: Volume 1*. Prentice Hall, Englewood Cliffs, NJ, 1988.

[20] G. S. Grest at Exxon Research, personal communication, 1993.

[21] L. Greengard and V. Rokhlin. A fast algorithm for particle simulations. *J. Comp. Phys.*, 73:325–348, 1987.

[22] G. S. Grest, B. Dünweg, and K. Kremer. Vectorized link cell Fortran code for molecular dynamics simulations for a large number of particles. *Comp. Phys. Comm.*, 55:269–285, 1989.

[23] S. Gupta. Computing aspects of molecular dynamics simulations. *Comp. Phys. Comm.*, 70:243–270, 1992.

[24] H. Heller, H. Grubmuller, and K. Schulten. Molecular dynamics simulation on a parallel computer. *Molec. Sim.*, 5:133–165, 1990.

[25] B. Hendrickson and R. Leland. An improved spectral graph partitioning algorithm for mapping parallel computations. Technical Report SAND90-1460, Sandia National Laboratories, Albuquerque, NM, 1992.

[26] B. Hendrickson and D. Womble. The torus–wrap mapping for dense matrix calculations on massively parallel computers. Technical Report SAND92-0792, Sandia National Laboratories, Albuquerque, NM, 1992.

[27] B. A. Hendrickson and S. J. Plimpton. Parallel many–body simulations without all–to–all communication. Technical Report SAND92-2766, Sandia National Laboratories, Albuquerque, NM, 1993.

[28] D. M. Heyes and W. Smith. *Inf. Q. Computer Simulation Condensed Phases (Daresbury Laboratory)*, 28:63, 1988.

[29] R. W. Hockney and J. W. Eastwood. *Computer Simulation Using Particles*. Adam Hilger, New York, NY, 1988.

[30] R. W. Hockney, S. P. Goel, and J. W. Eastwood. Quiet high–resolution computer models of a plasma. *J. Comp. Phys.*, 14:148–158, 1974.

[31] J. J. Morales and M. J. Nuevo. Comparison of link–cell and neighbourhood tables on a range of computers. *Comp. Phys. Comm.*, 69:223–228, 1992.

[32] M. R. S Pinches, D. J. Tildesley, and W. Smith. Large–scale molecular dynamics on parallel computers using the link–cell algorithm. *Molec. Sim.*, 6:51–87, 1991.

[33] S. J. Plimpton. Molecular dynamics simulations of short–range force systems on 1024–node hypercubes. In *Proc. 5th Distributed Memory Computing Conference*, pages 478–483. IEEE Computer Society Press, 1990.

[34] S. J. Plimpton. Scalable parallel molecular dynamics on MIMD supercomputers. In *Proc. Scalable High Performance Computing Conference–92*, pages 246–251. IEEE Computer Society Press, 1992.

[35] S. J. Plimpton and B. A. Hendrickson. A new strategy for parallelizing molecular dynamics simulations of organic systems. In preparation.

[36] S. J. Plimpton and B. A. Hendrickson. Parallel molecular dynamics with the embedded atom method. In *Materials Theory and Modeling*, volume 291, pages 37–42. Materials Research Society Symposium Proc., Fall 1992.

[37] S. J. Plimpton, B. A. Hendrickson, and G. S. Heffelfinger. A new decomposition strategy for parallel bonded molecular dynamics. In *Proc. 6th SIAM Conference on Parallel Processing for Scientific Computing*, pages 178–182. SIAM, 1993.

[38] S. J. Plimpton and E. D. Wolf. Effect of interatomic potential on simulated grain–boundary and bulk diffusion: A molecular dynamics study. *Phys. Rev. B*, 41:2712–2721, 1990.

[39] D. C. Rapaport. *Comput. Phys. Rep.*, 9:1, 1988.

[40] D. C. Rapaport. Multi–million particle molecular dynamics. II. Design considerations for distributed processing. *Comp. Phys. Comm.*, 62:217–228, 1991.

[41] M. Schöen. Structure of a simple molecular dynamics Fortran program optimized for Cray vector processing computers. *Comp. Phys. Comm.*, 52:175–185, 1989.

[42] H. Schreiber, O. Steinhauser, and P. Schuster. Parallel molecular dynamics of biomolecules *Parallel Computing*, 18:557–573, 1992.

[43] W. Smith. Molecular dynamics on hypercube parallel computers. *Comp. Phys. Comm.*, 62:229–248, 1991.

[44] W. B. Street, D. J. Tildesley, and G. Saville. Multiple timestep methods in molecular dynamics. *Mol. Phys.*, 35:639–48, 1978.

[45] P. Tamayo and R. Giles. A parallel scalable approach to short-range molecular dynamics on the CM-5. In *Proc. Scalable High Performance Computing Conference-92*, pages 240–245. IEEE Computer Society Press, 1992.

[46] P. Tamayo, J. P. Mesirov, and B. M. Boghosian. Parallel approaches to short-range molecular dynamics simulations In *Proc. Supercomputing '91*, pages 462–470. IEEE Computer Society Press, 1991.

[47] P. A. Taylor, J. S. Nelson, and B. W. Dodson. Adhesion between atomically flat metallic surfaces. *Phys. Rev. B*, 44:5834–5841, 1991.

[48] R. van de Geijn. Efficient global combine operations. In *Proc. 6th Distributed Memory Computing Conference*, pages 291–294. IEEE Computer Society Press, 1991.

[49] L. Verlet. Computer experiments on classical fluids. I. Thermodynamical properties of Lennard-Jones molecules. *Phys. Rev.*, 159:98–103, 1967.

[50] M. S. Warren and J. K. Salmon. A parallel treecode for gravitational N-body simulations with up to 20 million particles Bulletin of the American Astronomical Society, 23:1345, 1991.

[51] A. Windemuth and K. Schulten. Molecular dynamics simulation on the Connection Machine. *Molec. Sim.*, 5:353–361, 1991.

## Internal Distribution

| | |
|---|---|
| Paul Fleury | 1000 |
| Brian Dodson | 1104 |
| Jeff Nelson | 1112 |
| Ed Barsis | 1400 |
| Sudip Dosanjh | 1402 |
| Bill Camp | 1421 |
| Grant Heffelfinger | 1421 |
| Martin Lewitt | 1421 |
| Steve Plimpton (15) | 1421 |
| Mark Sears | 1421 |
| Dick Allen | 1422 |
| Bruce Hendrickson | 1422 |
| Kevin McCurley | 1423 |
| Art Hale | 1424 |
| Rob Leland | 1424 |
| Mike Proicou | 1424 |
| Steve Attaway | 1425 |
| Paul Taylor | 1432 |
| John Curro | 1702 |
| Mike Kent | 1815 |
| Elizabeth Holm | 1831 |
| Dona Crawford | 1900 |
| Ray Cline | 1952 |
| Randy Cygan | 6118 |
| Gary Carlson | 6211 |
| John Shelnutt | 6211 |
| Mike Colvin | 8117 |
| Richard Judson | 8117 |
| Charles Tong | 8117 |
| Joe Schoeniger | 8117 |
| Murray Daw | 8341 |
| Stephen Foiles | 8341 |
| Carl Melius | 8353 |
| Jim Plimpton | 9301 |
| | |
| Technical Library (5) | 7141 |
| Technical Publications | 7151 |
| Document Processing for | |
| DOE/OSTI (10) | 7613-2 |
| Central Technical File | 8523-2 |

# END

# DATE
# FILMED
1/20/94