



Contents lists available at ScienceDirect

## Future Generation Computer Systems

journal homepage: [www.elsevier.com/locate/fgcs](http://www.elsevier.com/locate/fgcs)

## Exploring dynamic load imbalance solutions with the CoMD proxy application

Olga Pearce<sup>a,\*</sup>, Hadia Ahmed<sup>b</sup>, Rasmus W. Larsen<sup>c</sup>, Peter Pirkelbauer<sup>b</sup>, David F. Richards<sup>a</sup><sup>a</sup> Lawrence Livermore National Laboratory, 7000 East Avenue, Livermore CA 94550, USA<sup>b</sup> Department of Computer and Information Sciences, University of Alabama at Birmingham, USA<sup>c</sup> Department of Computer Science, University of Copenhagen, Denmark

## HIGHLIGHTS

- A proxy application that allows us to study load imbalance and potential solutions.
- Ability to control initial and dynamic load imbalance, flexibility in work migration.
- An evaluation of a load balance algorithm performance on a wide range of scenarios.

## ARTICLE INFO

## Article history:

Received 3 March 2017

Received in revised form 5 October 2017

Accepted 5 December 2017

Available online xxxx

## ABSTRACT

Proxy applications are developed to simplify studying parallel performance of scientific simulations and to test potential solutions for performance problems. However, proxy applications are typically too simple to allow work migration or to represent the load imbalance of their parent applications. To study the ability of load balancing solutions to balance work effectively, we enable work migration in one of the Exascale Co-design Center for Materials in Extreme Environments (ExMatEx) [1] applications, CoMD. We design a methodology to parameterize three key aspects necessary for studying load imbalance correction: (1) the granularity with which work can be migrated; (2) the initial load imbalance; (3) the dynamic load imbalance (how quickly the load changes over time). We present a study of the impact of flexibility in work migration in CoMD on load balance and the associated rebalancing costs for a wide range of initial and dynamic load imbalance scenarios.

© 2018 Elsevier B.V. All rights reserved.

## 1. Introduction

Modern scientific simulations rely on parallel computers to solve state of the art problems requiring vast computational resources. The largest supercomputers have millions of independent processors, and concurrency levels are rapidly increasing. For ideal efficiency, developers of the simulations that run on these machines must ensure that computational work is evenly balanced among processors. Dynamic load balancing is a way to correct the imbalances that arise throughout the application execution, and is increasingly important for overall application performance as a simulation with more processes will waste more resources than a smaller-scale simulation when waiting on a single slow process. To enable dynamic load balancing, applications must implement

a mechanism for work migration. We present a methodology for studying how flexibility in work migration impacts the trade off between the ability to balance the work effectively and the associated costs.

Proxy applications can make it easier to study performance of large simulations and try new solutions. However, as the result of simplification, many proxy applications are not suitable for studying load imbalance solutions because they do not implement work migration. In this work, we extend one of the ExMatEx proxy applications, CoMD, to enable work migration, and to represent dynamic load imbalance found in large scale molecular dynamics simulations. To mimic real simulations, we develop a methodology to set up initial load imbalance, and to dynamically control the imbalance as the simulation progresses. We enable flexibility in the granularity of work migration in CoMD, which allows us to study the impact of the work migration granularity on the ability to lower the imbalance in the simulation and the associated rebalancing costs.

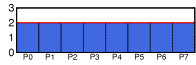
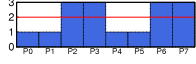
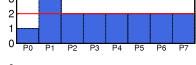
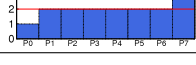
\* Corresponding author.

E-mail addresses: [olga@llnl.gov](mailto:olga@llnl.gov) (O. Pearce), [hadia@cis.uab.edu](mailto:hadia@cis.uab.edu) (H. Ahmed), [rasmuswl@di.ku.dk](mailto:rasmuswl@di.ku.dk) (R.W. Larsen), [pirkelbauer@uab.edu](mailto:pirkelbauer@uab.edu) (P. Pirkelbauer), [richards12@llnl.gov](mailto:richards12@llnl.gov) (D.F. Richards).

<https://doi.org/10.1016/j.future.2017.12.010>

0167-739X/© 2018 Elsevier B.V. All rights reserved.

**Table 1**  
Example load distributions and their imbalance.

	Load on each process	$L_{max}$	$L_{ave}$	Imbalance
(a)		2	2	0%
(b)		3	2	50%
(c)		3	2	50%
(d)		3	2	50%

In this paper, we describe the control mechanisms for (1) initial load imbalance, (2) dynamic imbalance, and (3) work migration granularity we introduced in CoMD. Together, these three aspects enable us to study the effectiveness of load balancing solutions and their costs when the application presents different imbalance behaviors (high/low initial imbalance, fast/slow rate of change in imbalance). We show that our version of CoMD is useful in evaluating work migration granularity and load balance algorithm performance for applications with different imbalance behaviors.

Our contributions in this paper are:

- A proxy application that allows us to study work migration granularity, load imbalance and potential solutions;
- Ability to control initial load imbalance, dynamic load imbalance, and flexibility in work migration;
- An evaluation of a load balance algorithm performance on a wide range of initial and dynamic load imbalance scenarios generated using our new proxy application.

Section 2 describes related work. Section 3 describes the original CoMD proxy application. Section 4 describes how we create initial load imbalance in CoMD, and how we make load imbalance dynamic in CoMD. Section 5 outlines the changes we made to the implementation of CoMD to enable work migration. Section 6 describes our interface for load balance algorithms. Section 7 shows our results.

## 2. Definitions and related work

For this paper, we define load *imbalance* as the scaled maximum load on any process minus the average:

$$Imbalance = \left( \frac{L_{max} - L_{ave}}{L_{ave}} \right) \times 100\%. \quad (1)$$

The definition in Eq. (1) represents opportunity cost of load balancing, as shown in Table 1 [2]. Example (b) shows process loads with 50% imbalance, and because the maximum process load is 3 units, the execution time is 3 units, which is 50% longer than the execution of the balanced example (a) (2 units).

Many simulations implement their own load balance algorithms that are tightly coupled with application data structures (i.e., ParaDiS [3,4]). Others rely on stand-alone libraries like graph partitioners (i.e., ParMetis [5,6], Jostle [7,8], and Zoltan [9,10]). Testing whether their own or stand-alone load balance solutions can correct load imbalance is non-trivial for a production application as enabling work migration can involve significant changes to the application data structures. Our work aims at providing a testbed to enable evaluation of load balance solutions prior to performing the work necessary to enable work migration, and informing decisions about granularity of work migration necessary for the load balance solutions to be successful.

Some applications have the option of using the Adaptive Message Passing Interface (AMPI) developed by the Charm++ group [11, 12] to overdecompose the simulation domain [13] and then balance load by moving virtual processors from overloaded physical processors to the underloaded ones. This approach is application agnostic and based solely on runtime information, but can impose extra communication overhead for tightly coupled applications due to the increase in the surface to volume ratio of the smaller domains. Our work enables testing of load balance algorithms that can rebalance the application based on the input from the application.

LeanMD [14] is a parallel molecular dynamics simulation framework written in Charm++ that exhibits load imbalance. While it is useful for exploring the built-in Charm++ load balance algorithms, it does not provide the control mechanisms to set up different load balance scenarios, which are the main contribution of this work. Similarly, the AMR mini-app implemented in Charm++ and miniAMR from the Mantevo suite [15] do not have a mechanism to control the imbalance.

The Particle-in-cell (PIC) Parallel Research Kernel (PRK) [16] is a paper-and-pencil specification of the computational task to be computed. PIC PRK was developed to help measure the efficiency and effectiveness of dynamic load balance techniques. Similarly to our work, they introduce different imbalance scenarios to test load balance algorithms. Particle distribution in PRK can be exponential, sinusoidal, or linear, resulting in different initial load imbalance. To simulate dynamic imbalance, particles can be uniformly injected or removed throughout execution. Besides representing a different class of simulations, our work goes a step further by explicitly parameterizing the initial and dynamic load imbalance. We also parameterize work assignment granularity to allow studying how flexibility in work assignment impacts the ability to rebalance the application.

## 3. CoMD: ExMatEx proxy app

Molecular dynamics is an important class of simulations that evaluates the forces that the atoms in the system exert on each other over time. CoMD [17] is an ExMatEx [1] classical molecular dynamics proxy application designed to represent this class of simulations. CoMD supports two potential energy models, Lennard-Jones (LJ) [18] and the Embedded Atom Method (EAM) [19,20].

In the simulation, each atom can interact only with other atoms located within a cutoff region defined by the potential function. Because atoms that are far apart have little effect on each other, classical MD simulations frequently use atom interaction models that define the force between two atoms to be zero when their separation distance exceeds a cutoff radius. This reduces the complexity of the force calculation to  $O(n)$ . Fig. 1 illustrates some molecular dynamics definitions in 2D. Fig. 1(a) shows a selected atom and a circle with the cutoff radius which defines the range of interaction. To simplify finding which of the atoms are within the cutoff radius, CoMD divides the simulation space into cells that are no smaller than the cutoff radius, as shown in Fig. 1(b). This cell definition guarantees that the atoms within the cutoff radius from an atom are either in the same cell as the atom, or in the immediately surrounding cells.

CoMD uses a Cartesian spatial decomposition of atoms across processes, with each process responsible for computing forces and evolving velocities and positions of all atoms within its domain boundaries. As atoms move across domain boundaries, they are handed off from one process to the logically neighboring process. To compute forces between atoms on different domains, CoMD uses ghost cells, defined as local copies of the immediately neighboring cells residing on the logically neighboring processes. Similar to other molecular dynamic simulations, CoMD uses periodic

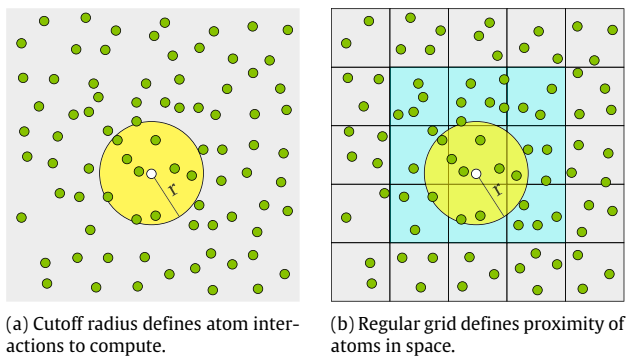


Fig. 1. Molecular dynamics definitions.

boundaries, allowing the atoms near the boundaries to interact with the atoms on the opposite side of the simulated space.

CoMD has been implemented in many different programming models, including MPI+X, CUDA, OpenCL, OpenACC, and X10. For this work we focus on load balancing across nodes in the system, therefore we use the MPI version of CoMD.

#### 4. Introducing load imbalance in CoMD

In this section, we describe how we introduce initial load imbalance in CoMD, and how we ensure the load imbalance changes throughout the simulation.

##### 4.1. Initial load imbalance

The reference implementation of CoMD evenly divides the simulated space between processes. Because the atoms are uniformly

spaced at setup, the atoms and the number of atom interactions that need to be computed are also evenly divided between processes, making the simulation inherently load balanced. To introduce load imbalance, we introduce voids in the simulated structure by removing some atoms from the simulation. Although CoMD simulations are 3D, for simplicity, Fig. 2 shows 2D examples.

Fig. 2(a) illustrates a four-process problem with four cells assigned to each process; the atoms are shown in green. CoMD places atoms on a lattice at the beginning of the simulation, resulting in a load balanced decomposition. We next place spheres with some diameter at random coordinates in the simulated space, as shown in blue in Fig. 2(b). Next, we remove the atoms whose coordinates are inside the spheres, as shown in Fig. 2(c). While the spheres are not part of the simulation, we use the concept of the spherical voids in this work to reason about the uniformity of atom distribution in the simulated space. After we remove the atoms and create the spherical voids, the number of atoms assigned to each process is different, and so is the number of atom interactions the processes will compute, as shown in Fig. 2(d). We expose the control mechanism for the initial load imbalance in CoMD by adding the following user-specified runtime parameters: (1) the spherical void size, (2) the sphere count, and (3) a random seed for generating the coordinates for the sphere center.

##### 4.2. Dynamic load imbalance

To create dynamic load imbalance in CoMD, we introduce the ability to set a non-zero center of mass velocity for the simulation. This causes the atoms to gradually shift in any of the three spatial dimensions (or combination of them) by a fixed distance at every time step.

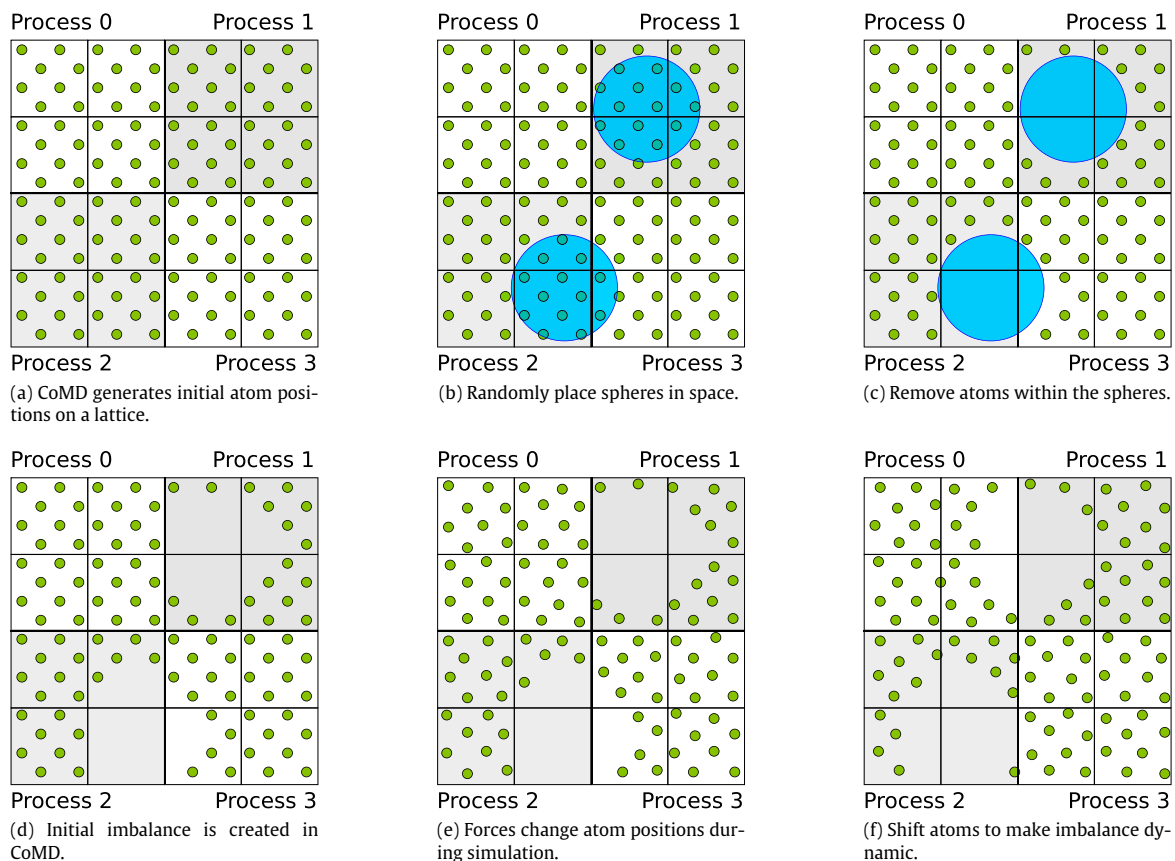


Fig. 2. Introducing load imbalance in CoMD. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

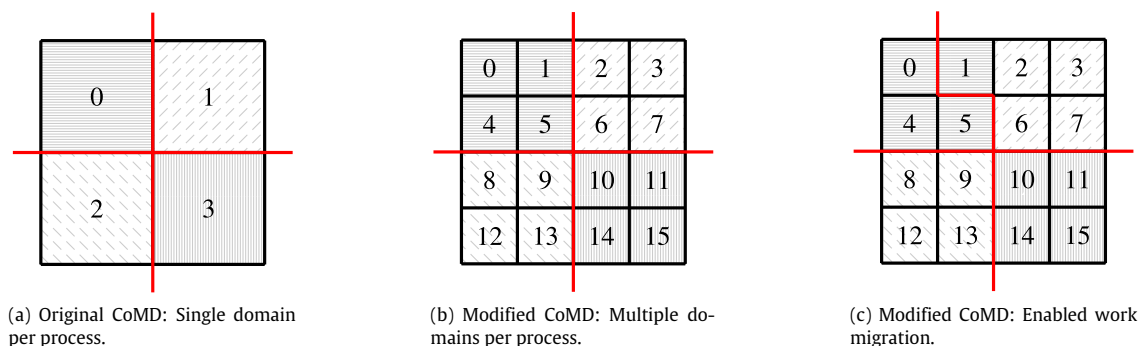


Fig. 3. Using domain overdecomposition to enable load balancing.

- 1: **for** timesteps **do**
- 2:   Execute application iteration on the local domain
- 3:   Communicate with neighbors

Fig. 4. Algorithm: Standard proxy application.

- 1: **for** timesteps **do**
- 2:   **for all** local domains **do**
- 3:     Execute application iteration on the local domain
- 4:     Communicate with neighbors (local and remote)

Fig. 5. Algorithm: Overdecomposed proxy application.

We start with the atom distribution we described in the previous section and demonstrated in Fig. 2(d), where process 1 and process 2 have fewer atoms than process 0 and 3. Fig. 2(e) demonstrates the simulation a timestep later; the atoms have changed positions due to the forces applied to them. Next, in this example, the non-zero center of mass velocity shifts the atoms horizontally, as shown in Fig. 2(f). Now process 0 has fewer atoms, while process 1 has more.

Because MD is translationally invariant, this overall shifting the atoms preserves the force interactions between the atoms but changes which process has to compute the interactions. We implemented the center of mass velocity as a runtime parameter; the user indicates the dimension(s) along which the atoms should be shifted, and the shift distance in Angstroms per time step for the selected dimension(s). Atoms can be shifted by a different distance in each dimension. This simulates a production application where the amount of work per process changes over time throughout the application runtime.

## 5. Enabling work migration in CoMD

The reference implementation of CoMD does not allow the correction of load imbalance because of several design decisions. For one, CoMD decomposes the simulated space into the same number of domains as processes, and assigns one domain per process. Also, the domains are rectangular prisms of the same size, and the user cannot control the size or the shape of the domains. Fig. 3(a) shows the simulation space in CoMD decomposed into four equal-sized rectangular domains; one domain is assigned to each process. Additionally, once the domains are assigned to processes, the assignment cannot be changed, so dynamic changes in work cannot be reflected in the work assignment.

We relaxed these design decisions by overdecomposing the application into more domains than available processes. Fig. 3(b) shows the same simulation space decomposed into 16 domains;

the domains are then assigned to four processes. However, if one process has more work than others, the assignment might be more balanced if we move one of the domains to another process, as shown in Fig. 3(c).

In our implementation, we introduce a data structure to store the assignment of domains to processes. Our *domain graph* consists of vertices which represent domains, and edges which represent boundaries with logically neighboring domains. The domain graph is distributed between processes, and each process stores local domains and edges to local and remote logically neighboring domains.

Fig. 4 shows pseudocode for a typical proxy application; each process performs computation on a single domain, and processes communicate boundary results with the logically neighboring processes. Fig. 5 shows pseudocode for an overdecomposed implementation. Each process is now responsible for computing the forces on one or more domains, and communicating the velocity and position updates accordingly.

We compare performance of our overdecomposed version of CoMD with the original implementation to ensure that we have not substantially changed the performance of the proxy application. For these experiments, we use load balanced CoMD, without introducing load imbalance yet. CoMD implements two common potentials, Lennard–Jones (LJ) and the Embedded Atom Method (EAM). Within each timestep, CoMD computes the forces between the particles (using one of the potentials), exchanges ghost cells, and updates atoms. In Figs. 6–9, we show the runtimes for these main computational phases for the original and overdecomposed versions of CoMD. We show both weak scaling (with 256 K atoms per process) and strong scaling (256 M atoms overall).

To measure the penalty of overdecomposition, for these experiments we set the domains to be the smallest possible size, which is equal to a single box in CoMD. In the weak scaling case, this means we have 13,824 domains per process. In the strong scaling case, we have 13 M domains overall.

Fig. 6 shows the runtime of the force computation for both potentials, for the original implementation of CoMD and our overdecomposed version. Both figures demonstrate that the runtime is slightly shorter for our overdecomposed implementation. This may be because in the original implementation, there are many link cells per domain, and in our overdecomposed approach there are fewer link cells per domain. Because we are computing atom–atom interactions for atoms in different link cells, our overdecomposed approach improves cache reuse similar to blocking in matrix–matrix multiplication.

Fig. 7 shows the runtime of the ghost exchange for Lennard–Jones and EAM potentials for the original implementation of CoMD and our overdecomposed version. In the ghost exchange phase, each process gathers information about the atoms that interact with atoms on the neighboring processes, and sends it to the

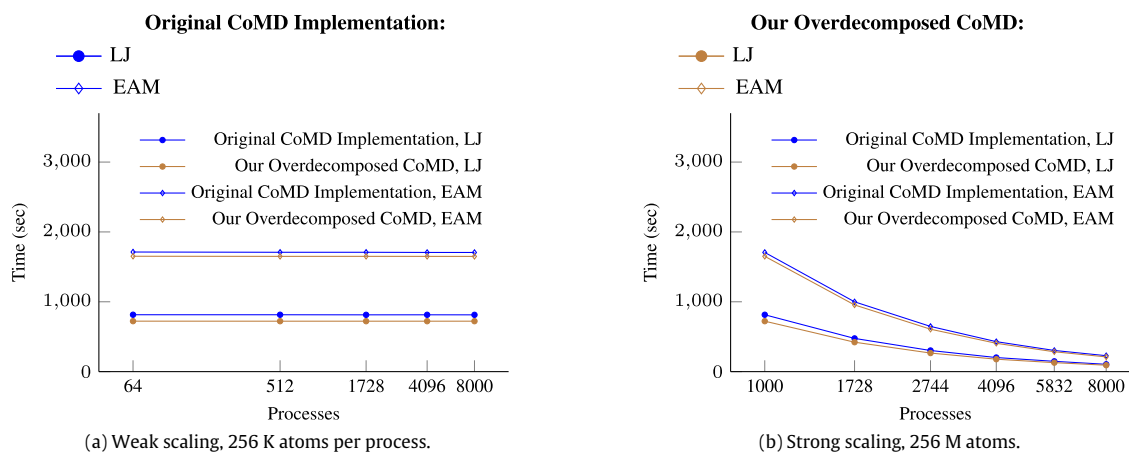


Fig. 6. Force computation (1000 timesteps) for original CoMD vs. overdecomposed CoMD.

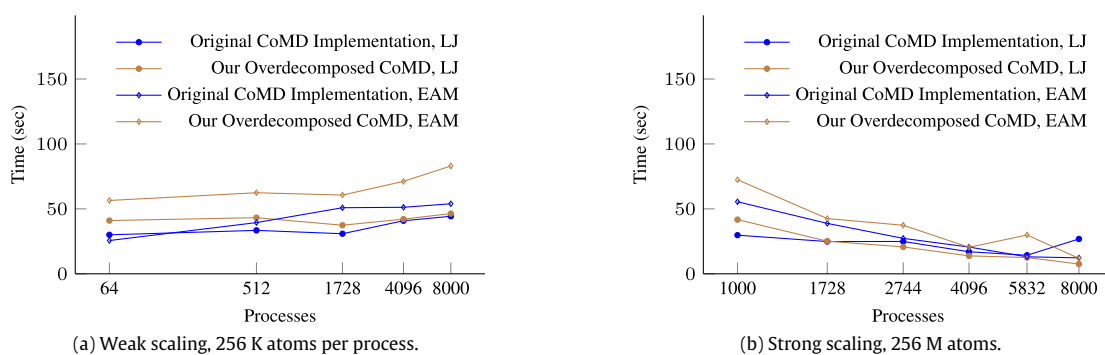


Fig. 7. Ghost exchange (1000 timesteps) for original CoMD vs. overdecomposed CoMD.

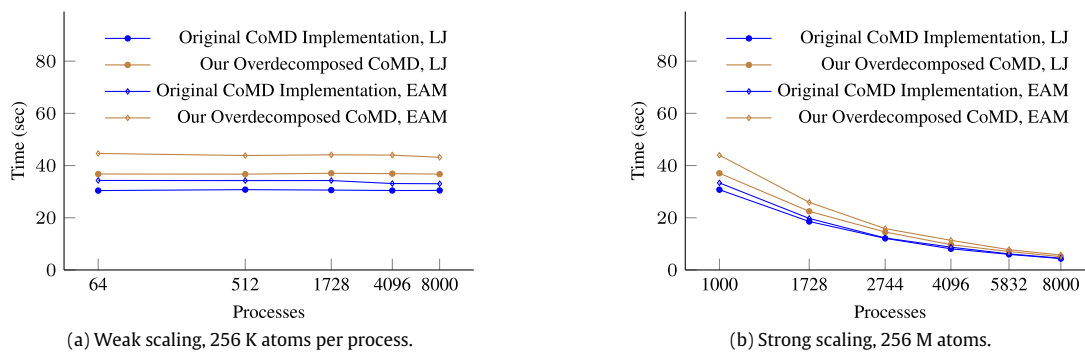


Fig. 8. Atom update (1000 timesteps) for original CoMD vs. overdecomposed CoMD.

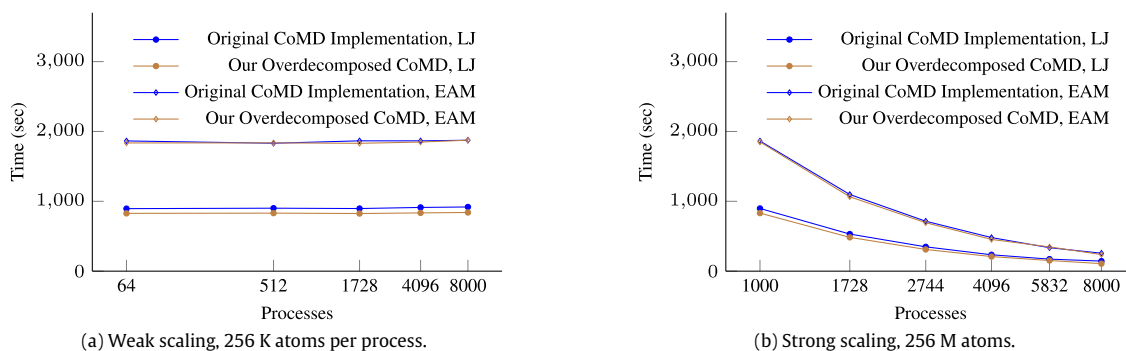


Fig. 9. Total runtime (1000 timesteps) for original CoMD vs. overdecomposed CoMD.



logically neighboring processes. In our implementation, because there are several domains assigned to a process, the process has to go through more link cells to gather the appropriate information, which results in a longer runtime of this phase. However, the cost is still comparable to the original implementation of CoMD, and scales similarly for both strong and weak scaling.

Fig. 8 shows the runtime of resorting the particle data and moving the particle information to the link cell that the particles now reside in. Similarly to the ghost exchange phase, in our overdecomposed version for this phase each process has to go through more lists, resulting in longer runtime for this phase than in the original implementation, but similar scaling properties.

Fig. 9 summarizes the impact of overdecomposition on the overall runtime for Lennard–Jones and EAM potentials. Our overdecomposed implementation does not significantly change the overall performance of CoMD.

We implement overdecomposition granularity as a runtime parameter to let us study the trade off between load balance quality and cost. Our hypothesis is that applications with moderate imbalance should be possible to balance even when work assignment granularity is large (i.e., 8 domains per process). Applications with more severe imbalance or more dynamic imbalance may require a finer granularity of work assignment; in this case, the higher cost of rebalancing may be justified by more flexibility in making a more balanced work assignment. Our proxy application allows us to easily study these trade offs.

## 6. Load balance algorithm interface

Our overdecomposed implementation of CoMD maintains a distributed domain graph where the vertices represent domains, and edges represent boundaries with logically neighboring domains. A vertex in the domain graph represents the smallest migratable unit in the application; its optional weight represents the amount of work in that domain. If the weight is not provided, vertices are assumed to be equally weighted. In this paper, we explore two approaches to estimate the work in the domain: (1) as the number of atoms in the domain, and (2) as the number of interactions computed for the domain, which has been shown to be a useful estimate for molecular dynamics simulations [21]. In Section 7.2, we will evaluate these two techniques for load balancing CoMD’s computation of Lennard–Jones and EAM potentials.

Our distributed domain graph can be easily translated to CSR format, and used as an input to generic load balance algorithms like graph partitioners (i.e., ParMetis [5,6], Jostle [7,8], and Zoltan [9, 10]). The output of a load balance algorithm is an assignment of domains to processes. Our implementation sends the domain to its new process and updates the domain graph.

For the results shown in this paper, we used a simple load balance algorithm based on a spatial sort. First, the domains are spatially sorted using a Hilbert curve or Morton curve based on the region of simulated space they represent. Next, the curve is partitioned between the processes, taking domain weights into consideration. This algorithm has the complexity of a sorting algorithm in terms of the number of domains sorted. So far, we used a sequential implementation of the algorithm, necessitating reduction of the relevant information to a single process. In the future, we plan to use a variety of more sophisticated and parallel load balance methods.

Our domain graph is a suitable interface with load balance algorithms since many load balance algorithms use graphs as a representation of the application communication.

## 7. Results

In this section, we evaluate how different work assignment (overdecomposition) granularity, initial load imbalance, and velocity of shifting the atoms in the simulation impact the ability of a load balance algorithm to correct the imbalance. We study the tradeoffs of load balancing accuracy and rebalancing costs.

For our experiments, we use a Linux cluster with nodes consisting of two 2.8 GHz Hex-core Intel Xeon EP X5660 processors, twelve cores per node. All nodes are connected by QDR Infiniband. We use GCC 4.4.7 and MVAPICH v0.99 on top of CHAOS, an HPC variant of RedHat Enterprise Linux (RHEL), running at Linux kernel v2.6.32.

For our measurements, we use Caliper [22], a generic context annotation tool developed at LLNL. Using Caliper, we annotate the phases of execution in the application, and measure how much time the phases take at each time step. The per-timestep measurements allow us to observe the performance changes in the simulation over time.

We define the load of an MPI process in a timestep as the total time minus the time waiting at MPI synchronizations:

$$L_{process} = T_{total} - T_{sync}. \quad (2)$$

We use Caliper’s MPI service to measure the time each process spends in MPI\_Barrier, MPI\_Allreduce, MPI\_Alltoall, and MPI\_Waitall functions at each timestep.  $T_{sync}$  in Eq. (2) is simply the total time the process spends in MPI synchronizations at each timestep. We use the measurements for process loads compute load imbalance using Eq. (1). We use these definitions for describing our experiments and results.

### 7.1. Creating initial load imbalance

The first control mechanism for our study is introducing initial load imbalance before running the simulation. For these experiments, we started by generating 23,328,000 atoms on 64 processes. The simulated space is a cube with a side length of 650.24 Angstroms; each process starts with a cube with a side length of 162.56 Angstroms. We varied the number and size of spherical voids to generate a set of problems with different number of atoms and load imbalance for our experiments.

Fig. 10 shows a set of problems with different configurations for spherical voids. We varied the number of voids from 2 to 10 per process. We varied the radius of the voids from 19.3 to 52.4 Angstroms. The parameterized atom removal resulted in a different number of atoms for each problem. Fig. 10(a) shows the resulting initial imbalance for each problem, as defined by Eq. (1). Generally, large spherical voids lead to less uniform distribution of atoms in the problem and therefore higher imbalance and difficulty in load balancing.

We overdecomposed the generated set of problems and load balanced them. We experimented with three granularities of overdecomposition: 8 domains per process (8x overdecomposition), 64 domains per process (64x overdecomposition), and 512 domains per process (512x overdecomposition). Decomposing the problem into more domains gives us more flexibility in terms of work assignment since it allows us to migrate smaller portions of the problem space. Fig. 10(c) shows the load imbalance for each granularity of overdecomposition after rebalancing. The problems with higher initial imbalance still had higher imbalance even after the rebalancing step, as compared to the problems with lower initial imbalance. This is due to the fact that large spherical voids result in problems with less uniform atom distribution, which are more challenging to load balance. Different granularities of overdecomposition result in roughly the same load imbalance for each problem, with smaller granularities of overdecomposition mostly

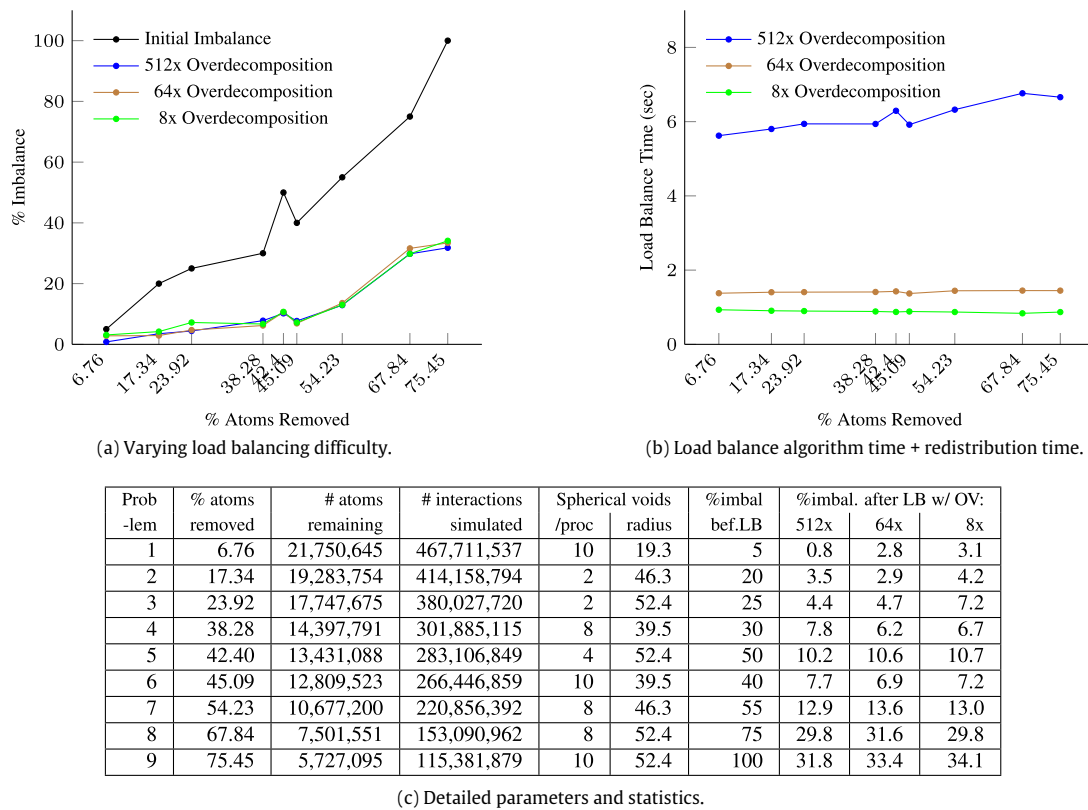


Fig. 10. Initial load imbalance scenarios in CoMD (64 processes, 23 M atoms initially).

achieving slightly lower imbalance due to having more flexibility in work assignment. In some instances, smaller granularity of overdecomposition does not result in lower imbalance, although the difference is negligible and is due to the fact that the work cannot be evenly divided between processes even with smaller granularity.

We used our set of problems with different initial load imbalance and parameterized overdecomposition granularities to look at the costs of rebalancing and redistributing work. Fig. 10(b) shows that larger granularities of overdecomposition result in higher costs of rebalancing, due to the finer granularity of work assignment. Additionally, the cost of rebalancing was slightly higher in problems with more severe imbalance as more domains were migrated during rebalancing.

While Fig. 10(a) suggests slightly better load balance with smaller granularities of overdecomposition, Fig. 10(b) shows the higher costs of overdecomposing into smaller domains. We will look at these costs in more detail in Section 7.3 as we examine weak scaling of our approach and overall runtimes of the problems.

Fig. 10(c) details the runtime parameters (number and size of spherical voids) for each problem shown in Figs. 10(a) and 10(b). It also shows the number of remaining atoms and initial load imbalance as well as the number of interactions each problem computes in the initial timestep. We also list the resulting load imbalance for each granularity of overdecomposition.

Our ability to vary the initial load imbalance allows us to evaluate how well a load balance method can correct the imbalance for a given overdecomposition granularity.

## 7.2. Input to load balance algorithm

We explored two different variants of a load balancing algorithm. One variant uses atom counts, and the other uses interaction

counts to estimate the amount of work in a domain. We evaluated the two variants using the load imbalance scenarios introduced in Fig. 10. Figs. 11 and 12 show our results for CoMD using Lennard-Jones and EAM potentials respectively. We measured the actual load imbalance (measured imbalance), the imbalance in atoms per domain, and the imbalance in interaction count per domain.

Fig. 11(a) depicts the initial load imbalance without overdecomposition for Lennard-Jones potential. We show three metrics: (1) the actual load imbalance (measured imbalance), (2) the imbalance in atoms per domain, and (3) the imbalance in interaction count per domain. All three metrics indicate similar load imbalance, suggesting that it would be reasonable to load balance either atoms or particle interactions. Fig. 11(b) shows the imbalance metrics for Lennard-Jones potential when we use an 8x overdecomposition and apply the two variants of the load balancing algorithm (atom-based and interaction-based). The actual load imbalance did not differ significantly between the two load balancing techniques. However, both atom imbalance and imbalance in number of interactions metrics slightly exaggerate the load imbalance. Figs. 11(c) and 11(d) show the imbalance metrics for 64x and 512x overdecomposition. In all scenarios, both load balancing techniques perform similarly. These figures also indicate that the additional flexibility gained through higher overdecomposition does not result in better load balance.

Fig. 12(a) shows the initial load imbalance for computing the EAM potential. All the three metrics indicate similar load imbalance. Figs. 12(b)–12(d) show the results using 8x, 64x, and 512x overdecomposition. The actual load imbalance did not differ significantly between the two load balancing techniques. As with Lennard-Jones potential, higher overdecomposition does not improve our ability to correct load imbalance. At this scale, creating 8 domains per process is sufficient to load balance the problem in CoMD.

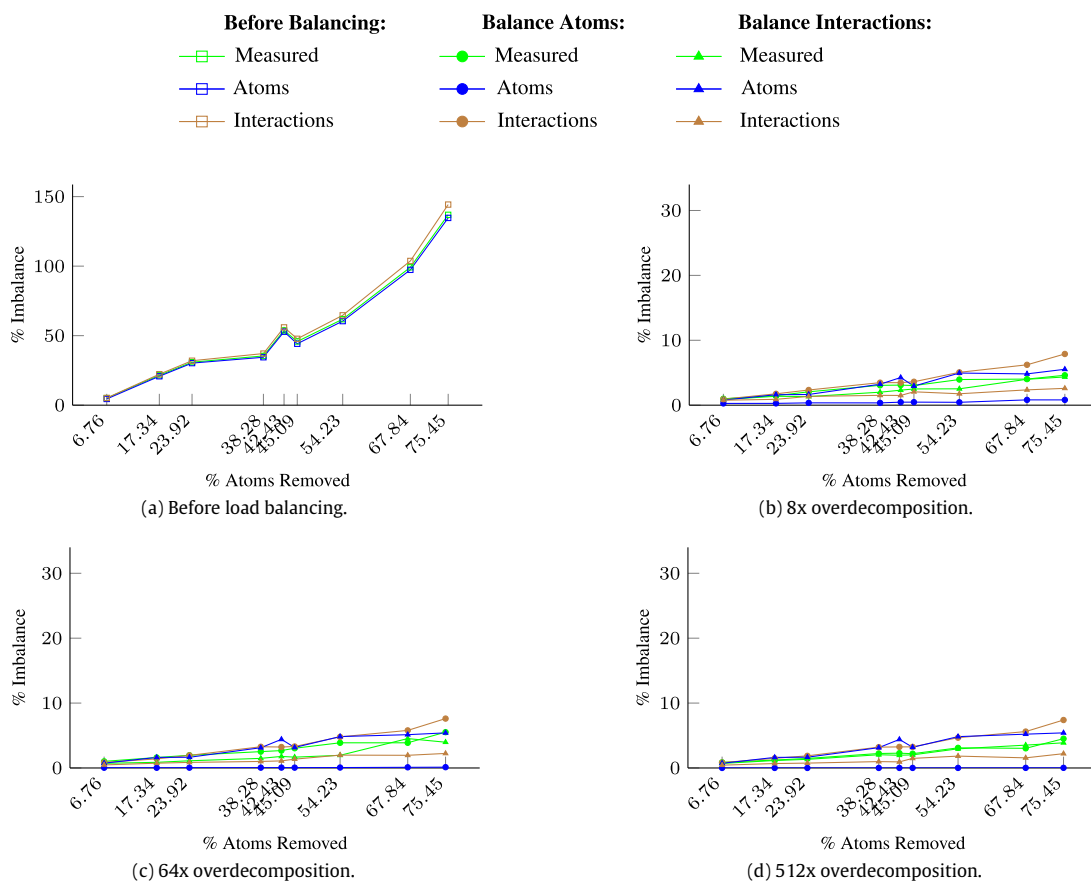


Fig. 11. Load imbalance for the LJ potential (64 processes, 23M atoms initially).

Overall, our work verifies that for the CoMD proxy application, the atom count and the interaction count are both reasonable approximations of the actual load imbalance in the application. We use these accurate load imbalance estimates as input to our load imbalance algorithms, and see effective load imbalance correction as a result.

### 7.3. Weak scaling

To evaluate the impact of scale on the performance and accuracy of a load balance algorithm, we performed a weak scaling study. For these experiments, we use Problem 5 in Fig. 10(c). This problem has 4 spherical voids with a radius of 52.4 Angstroms per process, and an initial imbalance of 50%. We weak scale Problem 5 from 64 to 4,096 processes by proportionally increasing the simulated space and keeping the size and the number of spherical voids per processor the same.

Fig. 13 demonstrates the impact of load imbalance and load balancing costs on the total runtime of the simulation when weak scaling the problem. We ran each problem for 160 timesteps, and rebalanced in the third timestep (once per run). Figs. 14 and 15 provide detailed parameters and statistics for Fig. 13. Fig. 14 lists the number of atoms, interactions, and initial imbalance at each scale. Load imbalance goes up slightly at higher scales because differences in process loads become more pronounced at scale.

Fig. 13(a) illustrates the ability of our load balance algorithm to load balance the problem; we show the load imbalance right after rebalancing when using 8x, 64x, and 512x overdecomposition. Fig. 15 details the load imbalance after the rebalancing. For all three granularities of overdecomposition, the load balance algorithm is successful at reassigning work in a more balanced manner. Because

the difficulty of load balancing increases with scale even in the weak scaling case, we see that the post-rebalancing imbalance is higher at higher scale. While the results are similar for all three levels of overdecomposition, 512x overdecomposition is slightly better due to more flexibility in work assignment. However since 512x overdecomposition deals with fine-grained data, it uses more memory for bookkeeping and runs out of memory on 4,096 processes when using our sequential load balance algorithm.

Fig. 13(b) shows the runtime of our load balance algorithm. Because we only have a sequential algorithm in place at the moment and need to gather/scatter its input/output, its execution time grows as the problem is weak scaled. The execution time grows much faster with scale for finer granularity of overdecomposition because the number of domains in the problem is the input size to the load balance algorithm.

Fig. 13(c) shows the time for redistributing the domains in the simulation after the load balance algorithm determines the new domain assignments. The higher the overdecomposition, the more domains there are in the problem, and therefore the longer it takes to redistribute the simulation.

Fig. 13(d) shows the total runtime of the simulation over 160 timesteps. Fig. 15 details the total runtime without load balancing, and with load balancing for each granularity. Problems run without load balancing take the longest to finish. Problems run with 8x overdecomposition plus load balancing are the fastest especially at higher scale.

Because we only have a sequential load balance algorithm that requires reduction of all necessary data to a single process, 512x overdecomposition used too much memory and resulted in a slow execution time of the load balance algorithm, outweighing the benefits of achieving lower imbalance.



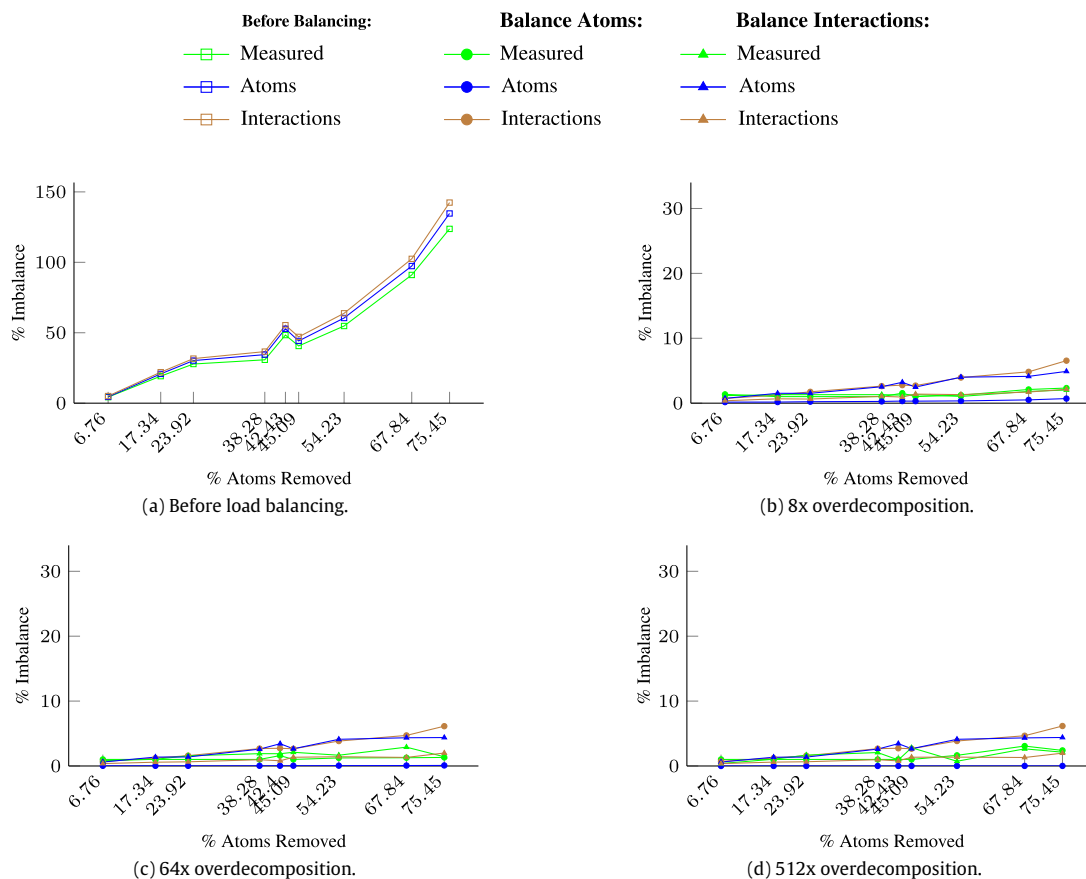


Fig. 12. Load imbalance for the EAM potential (64 processes, 23M atoms initially).

Our proxy application allowed us to study performance, accuracy, and associated cost of a load balance algorithm for an increasing number of processes.

#### 7.4. Dynamic imbalance

To allow dynamic work changes and to study their effect on the ability to load balance the simulation, we experiment with our control mechanism of adjusting a center of mass velocity for the atoms. For these experiments we compare Problems 5 and 9 in Fig. 10(c) on 64 processes. Problem 5 has 4 spherical voids with the radius of 52.4 Angstroms per process; it has 13,431,088 atoms and computes 283,106,849 atom-pair force interactions, and 50% initial imbalance. Problem 9 has 10 spherical voids with the radius of 52.4 Angstroms per process; it has 5,727,095 atoms and computes 115,381,879 atom-pair force interactions, and 100% initial imbalance. In addition to being more imbalanced, Problem 9 performs significantly less computation than Problem 5. For both problems, we use three granularities of overdecomposition, and run the load balance algorithm at even intervals with different frequencies.

Fig. 16 shows the effect of dynamic imbalance. We ran all problems for 160 timesteps. We show three center of mass velocities: 0.1%, 0.5%, and 1%, where 1% velocity means each timestep the atoms were shifted by 1% of the length of the simulation space in one dimension of the problem. The change in the imbalance is small when the atoms are shifted slowly, and a faster shift of atoms results in a larger variation in imbalance. The imbalance can both decrease and increase because the atoms in the problem will shift from one process to the next, changing the amount of computation each process performs. For Problem 5, the variation in imbalance as

the atoms shift is not large, because each process does a significant amount of work even as the atoms shift, as shown in Fig. 16(a). For Problem 9, the variation in imbalance is larger, because there are fewer atoms overall and the impact of them shifting to other processes is greater, as shown in Fig. 16(b).

While it is possible to trigger rebalancing based on a heuristic or a model which determines when load balancing would be beneficial, in this work we only explore load balancing with a certain static frequency. We define frequency of rebalancing as the number of times the problem was rebalanced during the execution, so a frequency of 2x means the problem was rebalanced twice, namely at timestep 1 and timestep 81. Figs. 17 and 18 show the total runtimes for Problems 5 and 9, and how imbalance in Problems 5 and 9 evolves when they are rebalanced with different frequency. We ran the problems for 160 timesteps. We show rebalancing at equal intervals with different frequencies that result in the best runtime (as demonstrated in Fig. 17(a)) for 0.1%, 0.5%, and 1.0% velocities. We show the runtimes for different center of mass velocities, and different frequencies of rebalancing. Minimizing the total runtime requires appropriately evaluating the tradeoffs between the costs and benefit of rebalancing with different frequency; we intend to use our version of CoMD for developing and evaluating models which handle the trade off for different imbalance scenarios and rebalancing costs.

Fig. 17(a) shows the total runtimes for Problem 5. Fig. 17(b) shows that load imbalance increases slowly when center of mass velocity is 0.1%, and Fig. 17(a) confirms that load balancing twice results in the lowest execution time (64x overdecomposition). The imbalance increases faster when the shift velocity is 0.5%, as shown in Fig. 17(c), so load balancing 8 times results in the lowest execution time (8x overdecomposition). Load imbalance increases

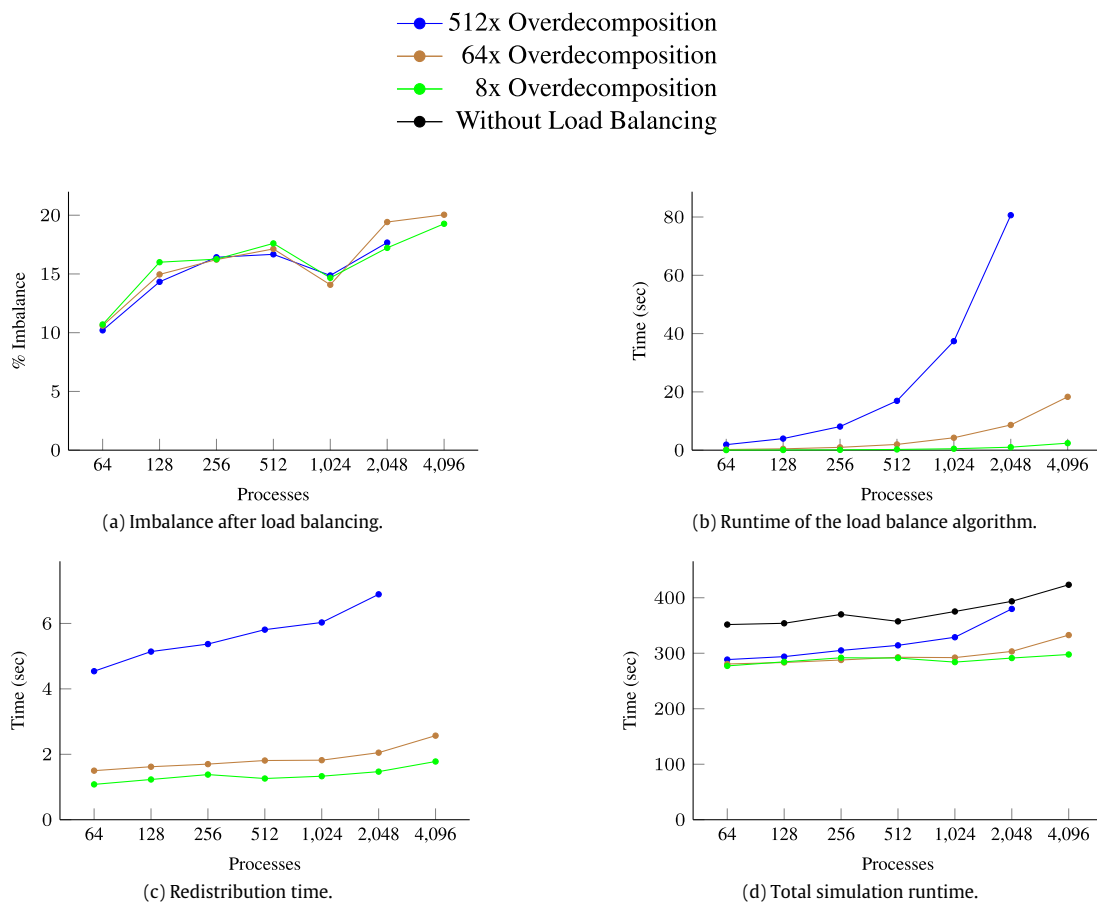


Fig. 13. Impact of load imbalance and rebalancing costs on total runtime (weak scaling)

most dramatically when the atoms are shifted with velocity of 1.0%, as shown in Fig. 17(d), necessitating rebalancing 14 times (8x overdecomposition).

Fig. 18(a) shows the total runtimes for Problem 9. Fig. 18(b) shows that load imbalance increases slowly when shift velocity is 0.1%, and load balancing twice results in the lowest execution time (8x overdecomposition). The imbalance increases faster when the shift velocity is 0.5%, as shown in Fig. 18(c), so load balancing 4 times results in the lowest execution time (8x overdecomposition). For this instance of the problem (0.5% velocity), the state of the application at step 80 is such that the load balance algorithm does not improve the imbalance; in fact, load imbalance increases; we are still investigating why this happens. Load imbalance increases most dramatically when the atoms are shifted with velocity of 1.0%, as shown in Fig. 18(d), requiring rebalancing 8 times for the best runtime (8x overdecomposition). Similarly to the 0.5% velocity case, for 1.0% velocity case, there are points in the problem when rebalancing increases the load imbalance (i.e., at timestep 40); we are still investigating the causes. Because there is less work per process in Problem 9 as compared to Problem 5, higher imbalance can be tolerated and fewer rebalancing steps can be amortized over the duration of the problem.

Table 2 shows details for the lowest runtimes in each velocity/frequency category. As velocity increases, the benefit of increasing the rate of rebalancing can outweigh the rebalancing cost, resulting in overall improvement in total runtime of the simulation.

Our ability to shift atoms in the simulation over time allows us to study how dynamic work changes effect our ability to load balance the simulation for a given overdecomposition granularity.

Num Procs	Num Atoms Remaining	Number of Interactions	% Imbalance before LB
64	13,431,088	283,106,849	50.70
128	26,505,376	558,243,612	50.30
256	53,190,202	1,120,473,002	55.76
512	106,273,179	2,239,029,303	50.55
1,024	212,741,503	4,482,828,235	58.00
2,048	425,692,895	8,971,434,404	65.00
4,096	852,665,073	17,971,364,315	63.00

Fig. 14. Detailed parameters for Fig. 13.

Table 2  
Details for Best Case Execution for Different Velocities.

Problem	Velocity	Velocity		
		0.1%	0.5%	1.0%
Problem 5 (50%)	Overdecomp.	64x	8x	8x
	LB Frequency	2x	8x	14x
	Time (s)	281.1	286.8	296.9
Problem 9 (100%)	Overdecomp.	8x	8x	8x
	LB Frequency	2x	4x	8x
	Time (s)	169.4	199.16	202.1

## 8. Conclusions

We extended an ExMatEx proxy application, CoMD, to be a suitable testbed for assessing the ability of load balance algorithms

Num Procs	% Imbalance After LB			Total Runtime			
	512x	64x	8x	No LB	512x	64x	8x
64	10.20	10.60	10.70	351.70	288.59	280.66	277.20
128	14.33	14.96	16	353.98	293.91	283.32	284.52
256	16.42	16.22	16.26	370.02	305.10	287.86	291.86
512	16.67	17.13	17.6	357.41	314.23	292.73	291.33
1,024	14.87	14.08	14.67	375.26	328.83	292.16	284.03
2,048	17.67	19.42	17.22	393.54	379.90	303.28	291.36
4,096	NA	20.04	19.27	423.40	NA	332.71	297.77

Fig. 15. Detailed statistics for Fig. 13.

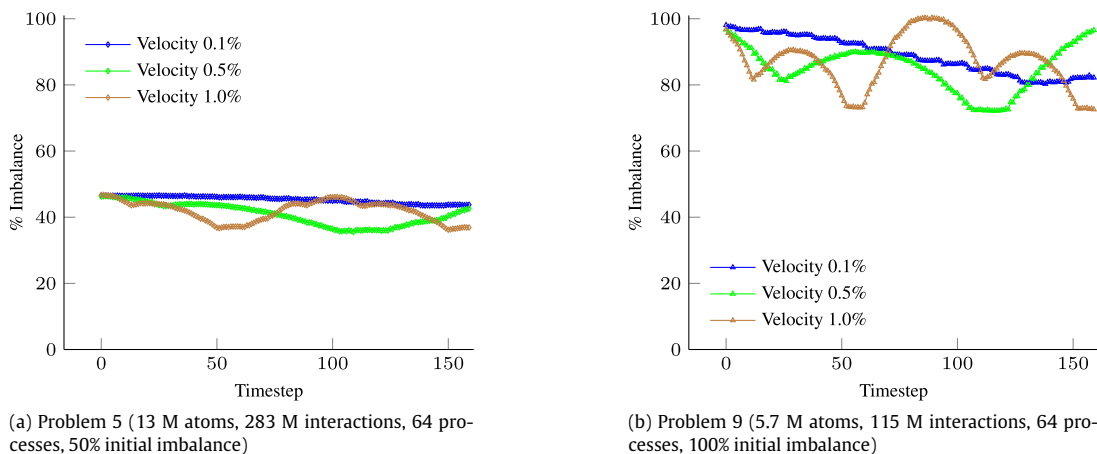


Fig. 16. Effect of velocity with which the atom shift on load imbalance.

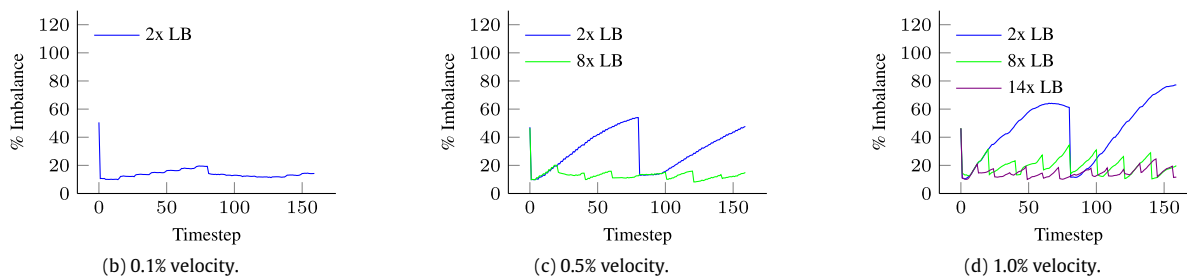
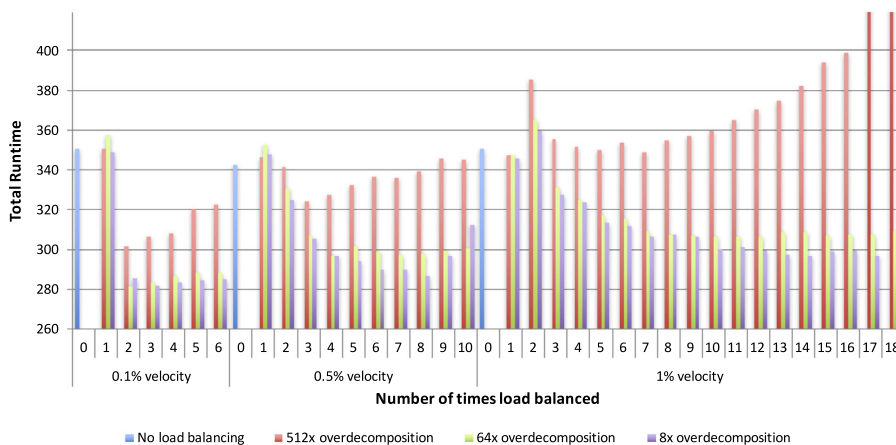
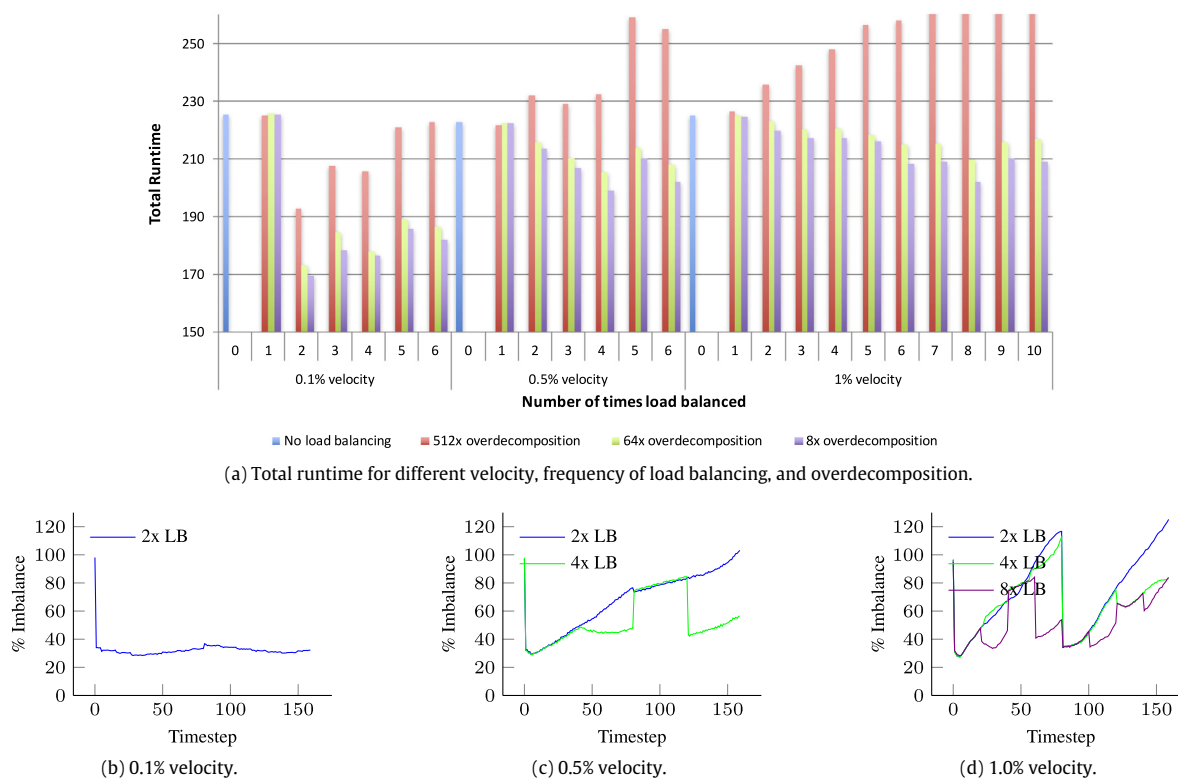


Fig. 17. Effect of frequency of rebalancing on load imbalance and simulation runtime. Problem 5 (13 M atoms, 283 M interactions, 64 processes, 50% initial imbalance).



**Fig. 18.** Effect of frequency of rebalancing on load imbalance and simulation runtime. Problem 9 (5.7 M atoms, 115 M interactions, 64 processes, 100% initial imbalance).

to correct dynamic load imbalance. We designed a methodology to parameterize three key aspects necessary for studying load imbalance correction: the granularity with which work can be migrated, the initial load imbalance, and how quickly the imbalance changes throughout the simulation. We presented a study demonstrating our ability to use our modified version of CoMD to evaluate the effectiveness of a load balance algorithm and the associated rebalancing costs for a wide range of initial and dynamic load imbalance scenarios.

## Acknowledgment

This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344 (LLNL-JRNL-725317).

## References

- [1] Exascale Co-design Center for Materials in Extreme Environments (ExMatEx). <http://www.exmatex.org>.
- [2] O. Pearce, T. Gamblin, B.R. de Supinski, M. Schulz, N.M. Amato, Quantifying the effectiveness of load balance algorithms, in: International Conference on Supercomputing (ICS), June 2012.
- [3] V. Bulatov, W. Cai, J. Fier, M. Hiratani, G. Hommes, T. Pierce, M. Tang, M. Rhee, K. Yates, T. Arsenlis, Scalable line dynamics in ParaDiS, in: SC'04, November 2004.
- [4] A. Arsenlis, W. Cai, M. Tang, M. Rhee, T. Opperstrup, G. Hommes, T.G. Pierce, V.V. Bulatov, Enabling strain hardening simulations with dislocation dynamics, *Modelling Simulation Mater. Sci. Eng.* 15 (6) (2007) 553.
- [5] K. Schloegel, G. Karypis, V. Kumar, A unified algorithm for load-balancing adaptive scientific simulations, in: SC'00, August 2000.
- [6] K. Schloegel, G. Karypis, V. Kumar, Parallel multilevel algorithms for multi-constraint graph partitioning, in: International Euro-Par Conference on Parallel Processing, August 2000.
- [7] C. Walshaw, M. Cross, Parallel optimization algorithms for multilevel mesh partitioning, *Parallel Comput.* 26 (12) (2000) 1635–1660.
- [8] C. Walshaw, M. Cross, M.G. Everett, Parallel dynamic graph partitioning for adaptive unstructured meshes, *J. Parallel Distrib. Comput.* 47 (2) (1997) 102–108.
- [9] K. Devine, E. Boman, R. Heaphy, B. Hendrickson, J. Teresco, J. Faik, J. Flaherty, L. Gervasio, New challenges in dynamic load balancing, *Appl. Numer. Math.* 52 (2–3) (2005) 133–152.
- [10] K. Devine, B. Hendrickson, E. Boman, M. St. John, C. Vaughan, Design of dynamic load-balancing tools for parallel applications, in: International Conference on Supercomputing (ICS), May 2000.
- [11] A. Bhatel , L.V. Kal , S. Kumar, Dynamic topology aware load balancing algorithms for MD applications, in: SC'09, November 2009.
- [12] R.K. Brunner, L.V. Kal , Handling application-induced load imbalance using parallel objects, in: International Workshop on Parallel and Distributed Computing for Symbolic and Irregular Applications, 2000.
- [13] E.R. Rodrigues, P.O.A. Navaux, J. Panetta, A. Fazenda, C.L. Mendes, L.V. Kal , A comparative analysis of load balancing algorithms applied to a weather forecast model, in: International Symposium on Computer Architecture and High Performance Computing, (SBAC-PAD), October 2010.
- [14] V. Mehta, *LeanMD: A Charm++ Framework for High Performance Molecular Dynamics Simulation on Large Parallel Machines*, University of Illinois at Urbana-Champaign, 2004.
- [15] M.A. Heroux, D.W. Doerfler, P.S. Crozier, J.M. Willenbring, H.C. Edwards, A. Williams, M. Rajan, E.R. Keiter, H.K. Thornquist, R.W. Numrich, *Improving Performance Via Mini-Applications*, Technical Report SAND2009-5574, Sandia National Laboratories, 2009.
- [16] E. Georganas, R.F.V. der Wijngaart, T. Mattson, Design and Implementation of a Parallel Research Kernel for Assessing Dynamic Load-Balancing Capabilities, in: International Parallel and Distributed Processing Symposium (IPDPS), May 2016.
- [17] Official website for CoMD. <http://www.exmatex.org/comd.html>.
- [18] D. Frenkel, B. Smit, *Understanding Molecular Simulation*, second ed., Academic Press, Inc., Orlando, FL, USA, 2001.
- [19] M.S. Daw, M.I. Baskes, Embedded-atom method: Derivation and application to impurities, surfaces, and other defects in metals, *Phys. Rev. B* 29 (1984) 6443–6453.



- [20] M.S. Daw, S.M. Foiles, M.I. Baskes, *The embedded-atom method: A review of theory and applications*, *Mater. Sci. Rep.* 9 (7) (1993) 251–310.
- [21] O. Pearce, T. Gamblin, B.R. de Supinski, T. Arsenlis, N.M. Amato, Load balancing N-body simulations with highly non-uniform density, in: *International Conference on Supercomputing (ICS)*, June 2014.
- [22] D. Boehme, T. Gamblin, D. Beckingsale, P.-T. Bremer, A. Gimnez, M. LeGendre, O. Pearce, M. Schulz, Caliper: performance introspection for HPC software stacks, in: *The International Conference for High Performance Computing, Networking, Storage, and Analysis, SC'16*, Salt Lake City, Utah, USA, November 13–18, 2016.



**Olga Pearce** is a computer scientist in the Center for Applied Scientific Computing (CASC) at Lawrence Livermore National Laboratory (LLNL). Her research interests include distributed data structures and parallel algorithms; parallel programming models; performance analysis, optimization, and modeling; and application load balancing. Olga joined LLNL in 2009 as a Lawrence Scholar, and joined CASC in 2014. Olga received her Ph.D. in Computer Science in 2014 from Texas A&M University.



**Hadia Ahmed** is a Ph.D. candidate in the Department of Computer and Information Sciences at University of Alabama at Birmingham. Her research interests include: high performance computing, parallel applications, compiler analysis techniques, and bioinformatics.



**Rasmus W. Larsen** is a Master's student at the Department of Computer Science, University of Copenhagen, Denmark. His research interests include functional programming, compilers, GPU programming, and high performance computing. Before starting university, he was an independent game developer.



**Peter Pirkelbauer** is an assistant professor in the department of Computer and Information Sciences at the University of Alabama at Birmingham. His research interests include programming languages, compilers, and parallel systems.



**David Richards** is a computational physicist at Lawrence Livermore National Laboratory. In addition to nearly 20 years experience designing scientific simulation codes, David has also played a leadership role in understanding and guiding the direction of future hardware and software at LLNL. David holds a B.S. in Physics from Harvey Mudd College and a Ph.D. in Physics from the University of Illinois at Urbana-Champaign.