

Benefits of hybrid parallelization in reducing load imbalance in molecular dynamics simulations

Julian Morillo^{a,*}, Maxime Vassaux^b, Peter V. Coveney^b, Marta Garcia-Gasulla^a

^aBarcelona Supercomputing Center, c/Jordi Girona 29, 08034 Barcelona (Spain)

^bCentre for Computational Sciences, University College London, 20 Gordon Street, London, WC1H 0AJ (United Kingdom)

Abstract

The most widely used technique to allow for parallel simulations in molecular dynamics is spatial decomposition, where the physical geometry is divided in boxes, one per processor. This technique can inherently produce computational load imbalance when either the spatial distribution of particles or the computational cost per particle is not uniform. This paper shows the benefits of using a hybrid MPI+OpenMP model to deal with this load imbalance. We consider the LAMMPS, a molecular dynamics simulator that provides its own balancing mechanism and an OpenMP implementation for many of its modules, allowing for a hybrid setup. In this work we extend the current OpenMP implementation of LAMMPS and optimize it and evaluate three different setups: MPI-only, MPI with the LAMMPS balance mechanism, and hybrid setup using our improved OpenMP version. This comparison is made using the five standard benchmarks included with LAMMPS distribution plus two additional test cases. Results show that the hybrid approach can deal with load balance problems better and straightforwardly than the LAMMPS balance mechanism and improve simulations with issues other than load imbalance.

Keywords: load balance, parallel computing, molecular dynamics, MPI, OpenMP, hybrid programming model

1. Introduction and Related Work

LAMMPS (Large-scale Atomic/Molecular Massively Parallel Simulator, [1, 2]) is a classical molecular dynamics code with a focus on materials modeling. It has potentials for solid-state materials (metals, semiconductors) and soft matter (biomolecules, polymers), and coarse-grained or mesoscopic systems. It can be used to model atoms or, more generically, as a parallel particle simulator at the atomic, meso, or continuum scale.

LAMMPS can be run in parallel using MPI and a spatial-decomposition of the simulation domain. The basic idea of a spatial decomposition method is to divide the physical geometry into small boxes, one per processor. Each processor will compute primarily on atoms within its box. This may induce load imbalance in problems with non uniform atom densities. The problem of load imbalance in MPI programs is well known [3] and in particular in molecular dynamics simulations is widely recognized [4, 5].

The challenge of the MPI load imbalance problem comes from the nature of MPI programming, where each process has its own data that can only be shared by explicit message passing. Simultaneously, the nature of load imbalance is dynamic and affected by many factors, therefore difficult to predict. Traditionally, solutions to load imbalance can be divided into two groups; The ones that are applied before execution; in this group, we can consider different mesh partitioners [6, 7]. These solutions are static and cannot address

*Corresponding author

Email addresses: julian.morillo@bsc.es (Julian Morillo), m.vassaux@ucl.ac.uk (Maxime Vassaux), p.v.coveney@ucl.ac.uk (Peter V. Coveney), marta.garcia@bsc.es (Marta Garcia-Gasulla)

32 load changes during the execution. Moreover, they need to be tuned for new architectures, algorithms, or
33 simulations.

34 The approaches applied during the execution can be classified as solutions that "move" data and solutions
35 that change computational resources. The methods that redistribute data [8, 9], usually execute a load
36 balancing algorithm with a given frequency. This algorithm determines if there is a load imbalance problem,
37 when necessary, computes a new partition, and finally redistributes the data as needed. These approaches
38 are not able to deal with very dynamic load imbalance. They also need to be able to measure load and
39 decide how frequently the load balancing algorithm is executed because the cost of redistributing the data
40 is not negligible. Usually, these solutions are implemented within each application; LAMMPS provides its
41 own balancing mechanism [10].

42 In the category of solutions applied at runtime that change the computational resources, we find different
43 approaches. Adaptive MPI [11], for example, rely on virtualized processes, and the runtime is in charge
44 of scheduling them to achieve a good load balance. They run on top of CHARM++ [12] which implies a
45 change in the programming language and model. Etinski et al. [13] propose to use the Dynamic Voltage
46 and Frequency Scaling (DVFS) reducing the frequency of less loaded processes to save power and use the
47 turbo for more loaded ones. Also, in this category, we find the Dynamic Load Balancing library [14, 15]; this
48 library changes the computational resources assigned to the different MPI processes to help balance their
49 load.

50 We propose to use the hybrid programming model MPI+OpenMP [16, 17] to alleviate the load balance
51 problem. This approach offers the advantages of a hybrid code: improves the load balance, and at the
52 same time reduces the pressure on the communication between MPI processes. In contrast with other
53 approaches that need to be programmed adhoc for each input or architecture, an OpenMP parallelization
54 can be exploited in many situations without needing to tune the code specifically.

55 Deng et al. [4], for example, describes an adaptive method for achieving load balance in parallel compu-
56 tations that is tested on standard short-ranged parallel molecular dynamics calculations. Our proposal, in
57 contrast, is to use a hybrid (MPI+OpenMP) approach. We argue that the use of OpenMP can help alleviate
58 MPI scaling issues, especially the ones related to load balance, and that this can be done straightforwardly by
59 leveraging on the OpenMP characteristics. Moreover, our evaluation is not limited to short-ranged molec-
60 ular dynamics calculations: mid-range and long-range simulations are also considered, including all the
61 benchmarks provided by the LAMMPS distribution, together with two extra testcases with quite different
62 characteristics regarding load balance.

63 Many LAMMPS modules have OpenMP versions for shared-memory parallelism, allowing for hybrid
64 setups in which MPI+OpenMP configurations can be run. Although there are OpenMP versions of many
65 LAMMPS modules, many of them lack an OpenMP implementation. Other ones present a parallelization
66 pattern that is not optimum for performance or programmability. This leaves MPI as the only parallel
67 option for these parts of the code. Nonetheless, the code is designed to be easily modified or extended with
68 new functionality. In this paper, we use such a feature to parallelize with OpenMP some code regions that
69 lack this parallelism and improve the original OpenMP implementation in other sections of code.

70 The main contributions of this paper are: i) we present some additions/improvements to the LAMMPS
71 OpenMP implementation; ii) we provide an extensive evaluation of this improved LAMMPS hybrid version
72 against the MPI-only version as a baseline case but also against the LAMMPS balance mechanism mentioned
73 previously; iii) we show how the hybrid approach can deal with imbalance issues better than the balance
74 mechanism and, furthermore, it can improve performance in cases where load imbalance is not the main
75 problem.

76 The document is organized as follows. Section 2 presents a description of LAMMPS and the bench-
77 marks/testcases used for the evaluation together with the performance analysis tools and efficiency metrics
78 employed. In Section 3 we explain why MPI load imbalance is a problem, the difficulties to address it, and
79 why it is very common in molecular dynamics simulations, together with the two compared approaches to
80 solve it: the LAMMPS balancing mechanism and the use of a hybrid model. Section 4 includes a description
81 of our proposed additions to the LAMMPS OpenMP implementation. Section 5 contains the environment
82 employed for the evaluation together with a characterization of the benchmarks/testcases used. Finally, a
83 complete performance comparison of the three evaluated scenarios is done for all the considered benchmark-

84 s/testcases. Section 6 concludes our study with comments and remarks.

85 2. Background

86 2.1. LAMMPS

87 The Large-scale Atomic/Molecular Massively Parallel Simulator (LAMMPS, [1, 2]) is a highly parallelised
88 code for the simulation of classical molecular dynamics. LAMMPS is mainly and widely used by the
89 materials science community. In classical molecular dynamics, atoms or molecules are described as particles
90 which dynamics are controlled by Newton's equations of motion. Particles interactions are determined by
91 potentials describing pairwise, multi-body and long-range interactions. Last, the thermodynamic properties
92 of the ensemble of particles modelled during the course of a molecular dynamics simulation are controlled
93 by the integration of global constraints on the particles position and velocities. LAMMPS, as its name best
94 describes, is intended for parallelism and runs as well on single processors as in parallel using message-passing
95 interface (MPI) and a spatial-decomposition of the simulation domain. Many modules provided by LAMMPS
96 to integrate global constraints or interatomic potentials have versions that provide accelerated performance
97 on CPUs, GPUs, and Intel Xeon Phis. LAMMPS is distributed by Sandia National Laboratories, a US
98 Department of Energy laboratory as an open source code under the terms of the GPL, and its design
99 is meant to be easily extended with new functionalities, which makes it easy for external developers to
100 contribute to the improvement of the code.

101 LAMMPS is used to simulate the dynamics and the properties of a wide range of systems including
102 amorphous and crystallised materials, proteins and much more. The need of computational chemists to
103 simulate larger systems for longer periods of time has continuously pushed the improvement of LAMMPS
104 scalability. Besides, LAMMPS is no also frequently coupled with other tools such as machine-learning or
105 continuum model simulators for scale-bridging purposes. As a result, large ensembles (up to thousands)
106 of molecular dynamics simulations can be simulated simultaneously. LAMMPS has already been used to
107 simulate the dynamics of tens of billions of atoms. On what is referred to as the "Lenard-Jones" benchmark,
108 the highest throughput recorded was 4.34 TFlops in 2005 on a 40×10^9 atoms simulation. In a more recent
109 attempt LAMMPS was shown to reach 2.35×10^{-8} s/atom/timestep.

110 2.2. Benchmarks

111 The evaluation reported in Section 5 is performed executing a combination of LAMMPS standard bench-
112 marking scenarios and a couple of additional scenarios triggering more specifically load-balancing issues.
113 LAMMPS features a set of five standard benchmark representative of the diversity of systems that can
114 be simulated. We assume that parallel efficiency is highly impacted by the range of interatomic potential
115 interactions. We therefore can classify the five scenarios into one of the following three classes of problems:
116 (i) short-range, (ii) mid-range and (iii) long-range interactions. Each scenario simulates the dynamics of
117 32,000 atoms. (i) The so-called "Granular chute" scenario simulates the convective flow of falling particles
118 interacting via a frictional history potential. The "Polymer chain" melt scenario simulates the thermal fluc-
119 tuations of hundred monomers long chains. The two scenarios are representative of short-range interacting
120 systems, each particle interacting on average with respectively 7 and 5 neighbours. (ii) The "EAM" sce-
121 nario simulates the thermodynamic fluctuations of a metallic copper bulk solid which atoms interact via the
122 embedded atom method (EAM) potential. The "Lennard-Jonnes" scenario simulates the thermodynamics
123 of an atomic fluid. The two scenarios are representative of a mid-range interacting system, each particle
124 interacting on average with respectively 45 and 55 neighbours. (iii) The "Rhodopsin" scenario simulates
125 the conformation changes of the rhodopsin protein in a solvated lipid bilayer, the CHARMM force-field is
126 used to described atoms pairwise and multi-body interactions. The "Rhodopsin" scenario also integrates
127 long-range Coulomb interactions, resulting in each particle interacting on average with 440 neighbours. The
128 data required for the simulation of these benchmarks is included in the distribution of LAMMPS. Further
129 details on the constraints applied during the simulation of the scenarios can be found on the Benchmark's
130 page on LAMMPS website (<https://lammps.sandia.gov/bench.html>).

131 In addition to LAMMPS standard benchmark scenarios, we introduce the simulation of an epoxy resin
132 and the simulation of a graphene-based nanocomposite. The simulation of the epoxy resin is referred to as
133 the "Epoxy" scenario, it simulates the non-equilibrium dynamics of highly-crosslinked epoxy polymer chains
134 under applied stretching. The epoxy resin is constrained with fixed number of atoms and temperature.
135 Meanwhile, the volume is controlled throughout the simulation and varied at fixed strain rate. The second
136 scenario features the simulation of graphene-oxide (GO) sheet embedded in a polymer precursor, it is referred
137 to as the "CG-GO" scenario. The GO sheet is a dense, two-dimensional packing of carbon atoms, while
138 polymer precursors consist in a disordered phase of poly(methyl methacrylate) (PMMA) precursors. The
139 "CG-GO" scenario simulates the dynamics of GO during annealing, the number of atoms and the volume
140 of the system are fixed and the temperature is increasing from 300K to 500K.

141 These two custom systems face specific computational efficiency issues which justified the improvement
142 of the current load balancing methods available in LAMMPS. We will perform efficiency measurements
143 which highlight existing bottlenecks and propose load balancing improvement based on the measurements
144 analysis.

145 2.3. Performance Tools and Efficiency Metrics

146 In this paper we go one step further and we aim at gaining in-sight on the reasons for the perfor-
147 mance achieved in the different situations at study. For this we rely on performance analysis tools and
148 the performance methodology promoted by Center of Excellence (CoE) for Performance Optimization and
149 Productivity (POP) ¹.

150 The performance analysis tools used in this work are the following:

151 Extrae: To obtain traces of the different executions, it supports MPI and OpenMP among other parallel
152 programming models[18, 19].

153 Paraver: To visualize the traces obtained with Extrae. It allows us to analyze in detail the behaviour of the
154 program and also to compute the performance metrics[20, 21].

155 The POP CoE has defined a methodology for performance analysis. This methodology is independent
156 of the tool being used for the analysis and defines a set of performance metrics. This metrics are well
157 defined, accepted by the community and meaningful, pointing the analysts to the main factors affecting
158 the performance and scalability of the code. In this paper we use some of this metrics as they allow us to
159 compare the different LAMMPS benchmarks using a common ground.

160 The POP metrics are hierarchical and multiplicative, meaning that the parent metric can always be
161 computed as the product of its childs. Each metric can get values between 0 and 100, and the metric
162 indicates how well that indicator is performing. For example, a load balance of 70% indicates that 30%
163 of the cpu time used is lost due to load imbalance, and also that addressing the load imbalance problem we
164 will enable to improve the execution by at most 30%. Specifically we are going to use 3 metrics from the
165 POP methodology: Parallel Efficiency, Load Balance and Communication efficiency.

166 These efficiency metrics are based on the simplification of a process into two states: the state in which
167 it is performing computation, which is called Useful (blue), and the state in which it is not performing
168 computation, e.g., communicating to other processes, which is called Not useful (red).

169 An example of this simplification can be seen in Figure 1 where we can see two processes, named p_1 and
170 p_2 running from left to right, the time being represented in the x axis. We can see how their execution
171 changes between the two states from Useful to Not useful and vice versa during their execution.

172 We call $P = \{p_1, \dots, p_n\}$ the set of MPI processes and n the number of MPI processes. For each MPI
173 process p we define the set $U_p = \{u_1^p, u_2^p, \dots, u_{|U|}^p\}$ of the time intervals where the application is performing
174 useful computation (the set of the blue intervals). We define the sum of the durations of all useful time
175 intervals in a process p as shown in Equation 1 and we call it the useful duration of a process.

¹<https://pop-coe.eu/>

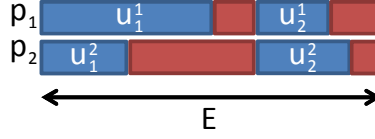


Figure 1: State evolution of two processes

$$D_{U_p} = \sum_{U_p} \blacksquare = \sum_{j=1}^{|U_p|} u_j^p \quad (1)$$

176 Similarly we can define \bar{U}_p and $D_{\bar{U}_p}$ for the red intervals.

177 We also define the elapsed time E as $E = \max_{p=0}^n D_{U_p} + D_{\bar{U}_p}$, the elapsed time is the total duration of the
 178 execution.

179 Parallel Efficiency (PE). The Parallel efficiency indicates the amount of time that is being lost due to the
 180 parallelization of the code. Or, equivalently, the ratio between the time used for useful computation and
 181 the total consumed CPU time. As we said the Parallel efficiency PE can be computed as the product of
 182 its childs, in this case the Load balance LB and Communication efficiency CE and is defined as shown in
 183 Equation 2.

$$PE = \frac{\sum_{U_p} \blacksquare}{E * n}; PE = LB * CE \quad (2)$$

184 Load Balance. Load balance measures the efficiency loss due to different loads (useful computation) for each
 185 process. Its definition can be seen in Equation 3.

$$LB = \frac{\sum_{i=1}^n D_{U_i}}{n * \max_{i=1}^n D_{U_i}} \quad (3)$$

186 Communication efficiency. Finally, the Communication efficiency is the efficiency loss for communicating
 187 data, it can be divided into two child metrics Serialization efficiency and Transfer efficiency. Serialization
 188 corresponds to time lost due to synchronizations between different processes, i.e., when one process needs
 189 to wait for another one. Transfer is the time lost in any kind of MPI overhead, it includes different fac-
 190 tors such as: network bandwidth, communication latency or implementation overheads. The definition of
 191 Communication efficiency can be found in Equation 4.

$$CE = \frac{\max_{i=1}^n D_{U_i}}{E} \quad (4)$$

192 3. Challenges and Proposed Approaches to Load Imbalance

193 3.1. Imbalance in Molecular Dynamics Simulations

194 Molecular dynamics (MD) is a commonly used tool for simulation of the structural, thermodynamic, and
 195 transport properties of biological and polymeric systems on the picosecond to nanosecond timescale. During
 196 a timestep of the MD simulation, forces are computed on each atom due to its interaction with other atoms,
 197 and atoms move by integrating simple Newtonian equations of motion [5].

198 The parallel nature of MD simulations has long been recognized [5, 22]. The overall calculation on
 199 P processors should scale as N/P , being N the total number of atoms in the simulated system. For gen-
 200 eral molecular systems simulated on message-passing machines, most parallel implementations have used
 201 the *replicated - data* technique [23] where a copy of all N atomic positions is stored on each of P proces-
 202 sors. This enables easy computation and load-balancing. However at each timestep, the interprocessor

203 communication needed to globally update a copy of the N -vector of atom positions scales as N , indepen-
204 dent of P . Thus replicated-data methods do not scale to large numbers of processors. An alternative
205 known as *force – decomposition* scales as N/\sqrt{P} but is still sub-optimal [24]. For large N/P ratios, spatial-
206 decomposition methods are clearly the best algorithmic choice. By subdividing the physical volume among
207 processors, most computations become local and communication is minimized so that optimal N/P scaling
208 can be achieved. Such method is the one used by LAMMPS.

209 The basic idea of a spatial decomposition method for MD is to divide the physical geometry into small
210 boxes, one per processor. Each processor will compute primarily on atoms within its box. This may induce
211 load imbalance in problems with non uniform atom densities.

212 3.2. Balancing

213 To alleviate the balancing problem, LAMMPS provides the balance command [10]. This command
214 adjusts the size and shape of processor sub-domains within the simulation box, to attempt to balance
215 the number of atoms or particles and thus indirectly the computational cost (load) more evenly across
216 processors. The load balancing is "static" in the sense that this command performs the balancing once,
217 before or between simulations. The processor sub-domains will then remain static during the subsequent
218 run. To perform "dynamic" balancing, LAMMPS provides the fix balance command, which can adjust
219 processor sub-domain sizes and shapes on-the-fly during a run.

220 Load-balancing is typically most useful if the particles in the simulation box have a spatially-varying
221 density distribution or when the computational cost varies significantly between different particles. For
222 example, a model of a vapor/liquid interface, or a solid with an irregular geometry containing void regions.
223 In these cases, LAMMPS default of dividing the simulation box volume into a regular-spaced grid of 3d
224 bricks, with one equal-volume sub-domain per processor, may assign numbers of particles per processor
225 in a way that the computational effort varies significantly. This can lead to poor performance when the
226 simulation is run in parallel.

227 The balancing can be performed with or without per-particle weighting. With no weighting, the balancing
228 attempts to assign an equal number of particles to each processor. With weighting, the balancing attempts
229 to assign an equal aggregate computational weight to each processor, which typically induces a different
230 number of atoms assigned to each processor. The weight assigned to a particle is defined a priori by the user
231 based on his knowledge of the particle, for example the expected number of neighbours and interactions.

232 3.3. Hybridization

233 It is not a trivial task to determine the optimal model (pure MPI vs MPI+OpenMP) to use for some
234 specific application. Although pure MPI can sometimes outperform hybrid, it is not less true that lots of
235 counterexamples do exist and results tend to vary with input data, problem size, etc. even for a given code.
236 In order to get optimal scalability one should in any case try to implement the following strategies:

- 237 • Reduce synchronization overhead
- 238 • Reduce load imbalance
- 239 • Reduce computational overhead and memory consumption
- 240 • Minimize MPI communication

241 Works like [25] pinpoint cases where hybrid programming model (MPI+OpenMP) can indeed be the superior
242 solution because of reduced communication needs and memory consumption, or improved load balance.

243 Hybridizing the code can help alleviate MPI scaling issues, especially the ones related to load balance as
244 the load balance within OpenMP is addressed straightforwardly when using a dynamic schedule with work-
245 sharing or the tasking model (i.e. generating explicit tasks that will be dynamically executed by threads
246 when they become idle).

247 To perform the tests with the hybrid versions of the benchmarks we have made use of OpenMP features
248 already provided by LAMMPS library. Version lammmps-20Nov19 has been used. However, and guided by
249 the Epoxy testcase, some modifications have been added to the OpenMP implementation that are described
250 in the following section.

251 4. Implementation

252 The use of the performance tools described in Section 2.3 to trace and analyze the execution of the Epoxy
253 testcase allowed us to find a source of load imbalance in void NPairHalfBinNewtonTri::build(NeighList *list)
254 function (blue regions in Figure 2). The bottom part of Figure 2 shows a timeline with the execution of two
255 LAMMPS iterations for 48 OpenMP threads. The upper timeline represents the execution with 48 MPI
256 ranks (included for reference, both timelines are at the same timescale). It can be seen how the OpenMP
257 parallelization of void NPairHalfBinNewtonTri::build(NeighList *list) alleviates the load imbalance. Note,
258 however, that the MPI execution is almost two times faster than the OpenMP one due to the sequential
259 parts (parts not OpenMP parallelized) in the later.

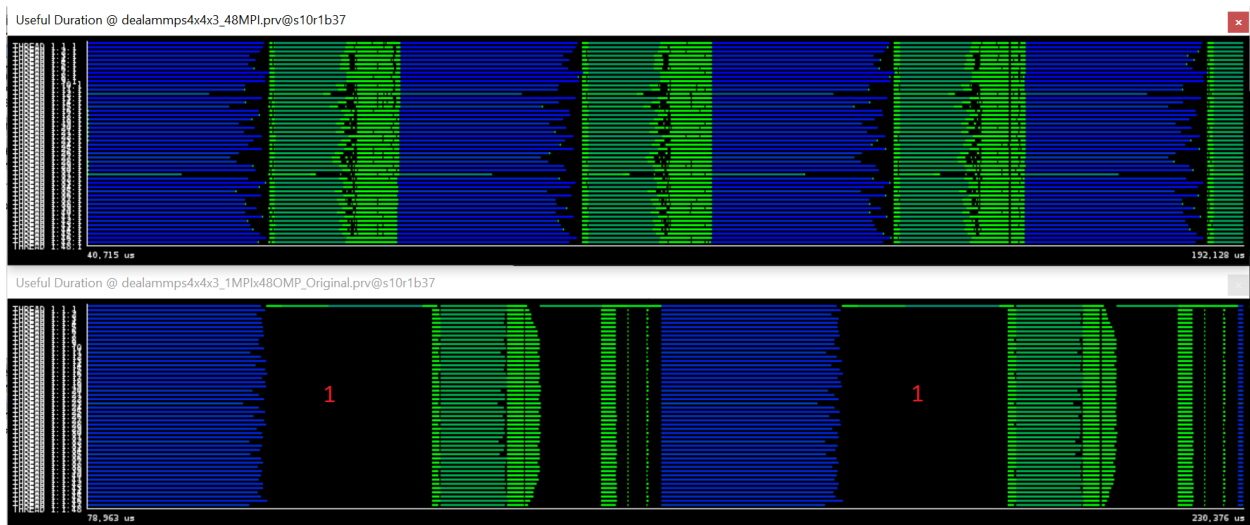


Figure 2: MPI-only (top) vs OpenMP-only (bottom) execution: the load imbalance is alleviated (blue region) but there are significant parts not OpenMP parallelized

260 4.1. OpenMP taskification

261 Given the timeline in Figure 2, our first work was to parallelize with OpenMP the biggest sequential
262 part (marked with a red 1 in the timeline). This part corresponds to Neighbor::build_topology() function.
263 The original LAMMPS code is depicted in Listing 1.

```
264 void Neighbor::build_topology()  
265 {  
266     if (force->bond) {  
267         neigh_bond->build();  
268         ...  
269     }  
270     if (force->angle) {  
271         neigh_angle->build();  
272         ...  
273     }  
274     if (force->dihedral) {  
275         neigh_dihedral->build();  
276         ...  
277     }  
278     if (force->improper) {  
279         neigh_improper->build();  
280         ...  
281     }  
282 }  
283 }
```

Listing 1: void Neighbor::build_topology() LAMMPS original code

285 This function consists in 4 calls to 4 different functions build (each one from a different class). Each of
 286 these 4 functions have a very similar structure that consists in a computation phase that ends up with an
 287 MPI_Allreduce of an int calculated in this computation phase (among other things). The 4 functions work
 288 with different data structures, therefore they are independent of each other and can be run in parallel.

289 The objective was to annotate these 4 calls to build functions with OpenMP tasks in order to allow their
 290 execution in parallel by different threads. Besides, the issue associated with the MPI_Allreduce command
 291 remains at the end of each function. The solution consists in moving these communications out of the
 292 4 functions and put them one level above in the Neighbor::build_topology function. The idea, then, is
 293 to have at the end 4 tasks with the computation of the 4 build functions and one final task with the 4
 294 communications. In this way, we delay as much as possible MPI communications, preventing unnecessary
 295 waiting times if there is imbalance between MPI ranks in some of the 4 computation phases.

296 In order to allow the MPI communications to be executed at the end of the 4 build computations we
 297 need to move them outside of each function. To do that, a change in the signature of the functions is needed:
 298 we need them to return an int (the value shared in the MPI_Allreduce) instead of void. Then, in each of
 299 the build functions, the calculated int in the computation phase is returned by the function instead of being
 300 directly shared with other MPI ranks through the MPI_Allreduce call. These returned values are then used
 301 in the new MPI_Allreduce calls located in void Neighbor::build_topology() function.

302 Code Listing 2 presents a skeleton of the final implementation of Neighbor::build_topology() function.
 303 As it can be seen, 4 new local variables are declared: the variables will be used to store the values returned
 304 by each of the build functions and, in turn, to honor the dependencies between the 4 computation tasks and
 305 the communications task.

```

306 void Neighbor::build_topology()
307 {
308     int nmissing_bond, nmissing_angle, nmissing_dihedral, nmissing_improper, all;
309     #pragma omp parallel
310     #pragma omp single
311     {
312         #pragma omp task depend(out:nmissing_bond)
313         if (force->bond) {
314             nmissing_bond = neigh_bond->build();
315             ...
316         }
317         #pragma omp task depend(out:nmissing_angle)
318         if (force->angle) {
319             nmissing_angle = neigh_angle->build();
320             ...
321         }
322         #pragma omp task depend(out: nmissing_dihedral)
323         if (force->dihedral) {
324             nmissing_dihedral = neigh_dihedral->build();
325             ...
326         }
327         #pragma omp task depend(out: nmissing_improper)
328         if (force->improper) {
329             nmissing_improper = neigh_improper->build();
330             ...
331         }
332         #pragma omp task depend(in: nmissing_bond, nmissing_angle, nmissing_dihedral, \
333         nmissing_improper)
334         {
335             MPI_Allreduce() x 4
336             ...
337         } //end task
338     } //end single and parallel
339     #pragma omp taskwait
340 } //end single and parallel

```


341 }

Listing 2: Modified void Neighbor::build_topology() code including OpenMP taskification

343 The present modifications significantly reduce the execution time of the biggest sequential part and work
344 efficiently for a small number of threads. Note, however, that we are generating only 5 OpenMP tasks and
345 only 4 of them can run in parallel as the communications one needs to wait for the execution of the others.
346 So, when moving to the extreme case of using 48 threads (the number of cores on the target machine), more
347 parallelism is needed. To accomplish that, the loops that make the calculations inside each of the 4 build
348 functions have also been taskified: this allows the generation of sufficient tasks to feed all threads.

349 The results of these modifications can be appreciated in Figure 3 where the red lines mark explicitly
350 the region of code affected by these changes and the reduction in execution time (upper part of the figure
351 corresponds to the original LAMMPS OpenMP implementation and the bottom timeline corresponds to our
352 improved OpenMP version).

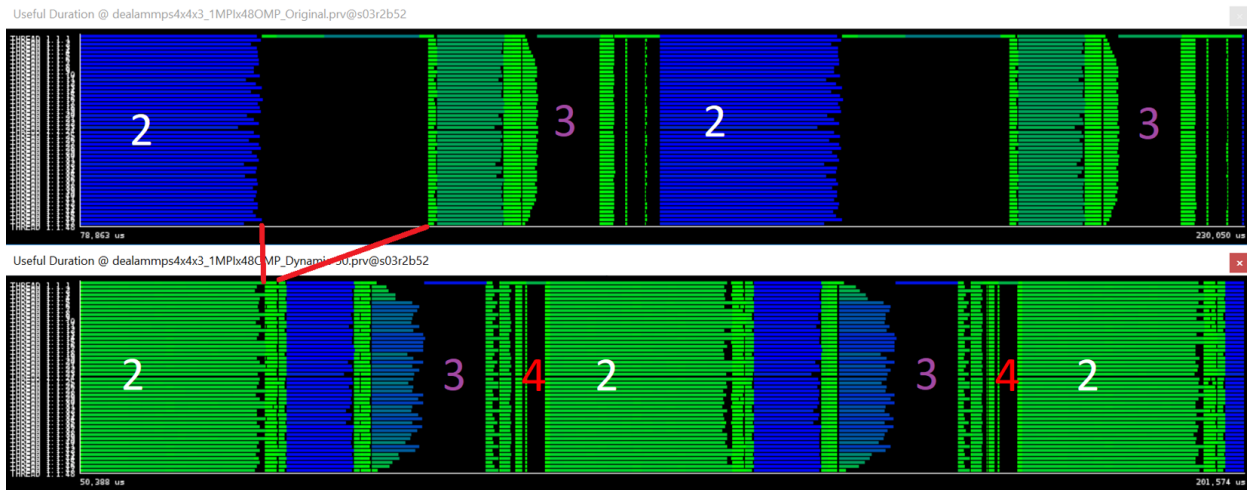


Figure 3: Improved version compared to the LAMMPS original OpenMP implementation.

353 4.2. Use of OpenMP dynamic scheduler

354 The second code modification has been done in the previous mentioned function NPairHalfBinNewton-
355 Tri::build(NeighList *list). It has been shown how by simply using the LAMMPS OpenMP implementation
356 the detected imbalance was alleviated. A close look to the source code shows, however, that the LAMMPS
357 OpenMP implementation does a static partition of the workload (like MPI does) so there is still some room
358 for improvement in this part of the code. This static partition of the workload is done through the use of
359 3 macros defined in npair_omp.h (see Listing 3). These 3 macros are widely used along all the LAMMPS
360 OpenMP code so the same code refactoring done in this section could be done in many other parts of the
361 code.

```
362 // get access to number of threads and per-thread data structures via FixOMP
363
364 #define NPAIR_OMP_INIT \
365     const int nthreads = comm->nthreads; \
366     const int ifix = modify->find_fix("package_omp")
367
368 // get thread id and then assign each thread a fixed chunk of atoms
369 #define NPAIR_OMP_SETUP(num) \
370     { \
371         const int tid = omp_get_thread_num(); \
372         const int idelta = 1 + num/nthreads; \
373         const int ifrom = tid*idelta; \

```

```

374     const int ito = ((ifrom + idelta) > num) \
375         ? num : (ifrom+idelta);           \
376     FixOMP *fix = static_cast<FixOMP *>(modify->fix [ ifix ] ); \
377     ThrData *thr = fix->get_thr(tid);     \
378     thr->timer(Timer::START);
379
380 #define NPAIR_OMP_CLOSE                \
381     thr->timer(Timer::NEIGH);          \
382 }
383

```

Listing 3: Macros defined in npair_omp.h

Listing 3 shows that the macro actually performing the workload partition is NPAIR_OMP_SETUP(num). It does so by dividing num among the number of available threads (i.e. in a loop of num iterations, defines the starting and end iteration that must be executed by each thread by setting ifrom and ito variables).

Once understood how these macros work, it is quite straightforward to implement the proposed approach. As it can be seen in Listing 4 the proposed change simply consists in substituting the original for that uses ifrom and ito variables by another one that, instead, starts at 0 and ends at nlocal (i.e. the value used in NPAIR_OMP_SETUP in this case). Of course, the new for is surrounded by a #pragma omp for schedule(dynamic) to do the worksharing (note that a #pragma omp parallel is not needed as it is already present at the beginning of the function in the original code, see Listing 4).

```

393 void NPairHalfBinNewtonTriOmp::build(NeighList *list)
394 {
395     ...
396     NPAIR_OMP_INIT;
397 #if defined(_OPENMP)
398 #pragma omp parallel default(none) shared(list)
399 #endif
400     NPAIR_OMP_SETUP(nlocal);
401     ...
402     #pragma omp for schedule(dynamic,50)
403     for (i = 0; i < nlocal; i++) {
404 // for (i = ifrom; i < ito; i++) {
405         ... (Computation)
406     }
407 }
408 NPAIR_OMP_CLOSE;
409 }
410

```

Listing 4: Sketch of NPairHalfBinNewtonTriOmp::build(NeighList *list) code including the dynamic OpenMP schedule

4.3. Enabling more OpenMP parallelism

Looking into the generated traces showed that there was a relatively large portion of code not OpenMP parallelized just before the execution of void NPairHalfBinNewtonTriOmp::build(NeighList *list). This is shown in Figures 3 and 4 marked with a red 4. The bottom timeline of Figure 4 represents (at the same timescale) the same part of the execution after parallelizing some functions in this region of code. The red lines going from one timeline to the other show the reduction in execution time achieved when using the new implemented parallel regions. Three different functions have been parallelized in this section of code but let us focus on the most important in terms of execution time: void NBinStandard::bin_atoms() (see Listing 5).

```

421 void NBinStandard::bin_atoms()
422 {
423     ...
424 #if defined(_OPENMP)
425 #pragma omp parallel for
426 #endif
427     for (i = 0; i < mbins; i++) binhead[i] = -1;
428 }
429

```

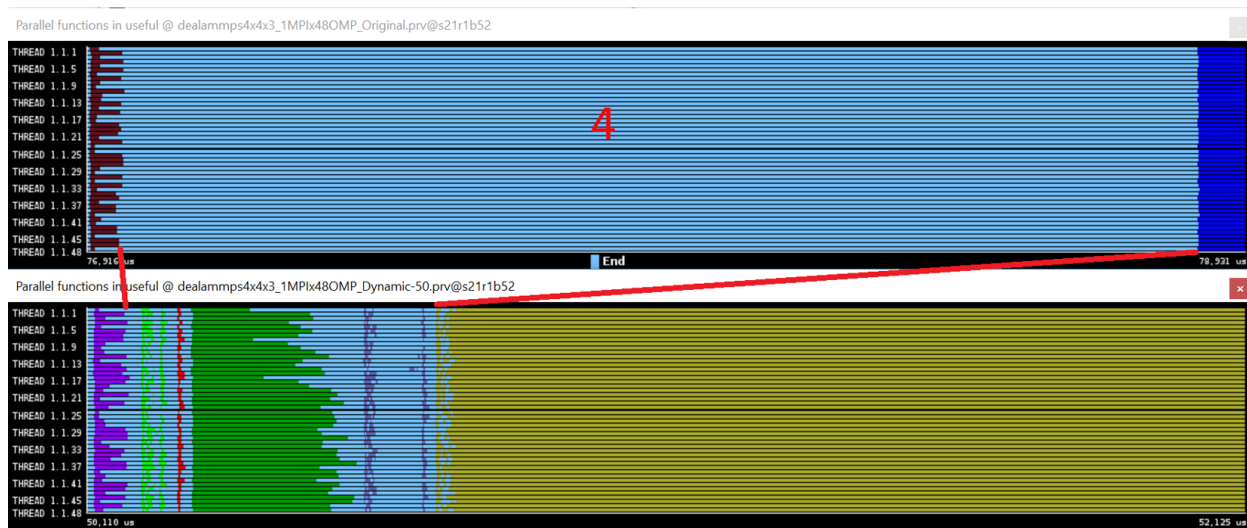


Figure 4: New implemented parallel regions.

```

430 ...
431
432 if (includegroup) {
433     int bitmask = group->bitmask[includegroup];
434     for (i = nall-1; i >= nlocal; i--) {
435         if (mask[i] & bitmask) {
436             ...
437         }
438     }
439     for (i = atom->nfirst-1; i >= 0; i--) {
440         ...
441     }
442 } else {
443 #if defined(_OPENMP)
444 #pragma omp parallel for private(ibin)
445 #endif
446     for (i = nall-1; i >= 0; i--) {
447         ibin = coord2bin(x[i]);
448         atom2bin[i] = ibin;
449         bins[i] = binhead[ibin];
450         binhead[ibin] = i;
451     }
452 }
453 }
454 }
455

```

Listing 5: Sketch of void NBinStandard::bin_atoms() code including the new added parallelization

Two parallel for worksharing constructs have been defined: the first one is not very relevant in terms of execution time and it corresponds to the red area in the bottom timeline in Figure 4. The important one is the located at the bottom of the Listing 5 which corresponds to the green region. As it can be seen it is a very simple parallel for that just needs to privatize ibin to work correctly.

Another important region of code candidate for enabling more OpenMP parallelism is the one marked with a purple 3 in Figure 3. This area corresponds to the execution of PPPM::poisson_ik_triclinic() method. A sketch of the original code can be seen in Listing 6.

```

463
464 void PPPM::poisson_ik_triclinic()
465 {
466

```

```

467 int i,j,k,n;
468
469 // x direction gradient
470
471 n = 0;
472 for (i = 0; i < nfft; i++) {
473     work2[n] = fxx[i]*work1[n+1];
474     work2[n+1] = -fxx[i]*work1[n];
475     n += 2;
476 }
477
478 fft2->compute(work2,work2,-1);
479
480 n = 0;
481 for (k = nzlo_in; k <= nzhi_in; k++)
482     for (j = nylo_in; j <= nyhi_in; j++)
483         for (i = nxlo_in; i <= nxhi_in; i++) {
484             vdx_brick[k][j][i] = work2[n];
485             n += 2;
486         }
487
488 // y direction gradient
489
490 ... //(same code for y direction)
491
492 // z direction gradient
493
494 ... //(same code for z direction)
495 }

```

Listing 6: Sketch of void PPM::poisson_ik_triclinic()) original code

497 The method consists on three differentiated parts (one for each x, y and z direction) with identical
498 structure: an initial loop, a call to fft2->compute() and, last, three nested loops. Unfortunately, variable n
499 prevents a direct parallelization of both the initial and the three nested loops.

500 The solution for the first loop is to incorporate the management of variable n (initialization and incre-
501 ment) to the control structure of the loop (see Listing 7). Once this is done, a simple pragma omp parallel
502 for will do the work.

```

503
504
505 void PPM::poisson_ik_triclinic()
506 {
507     int i,j,k,n;
508
509     // x direction gradient
510     #pragma omp parallel for
511     for (i = 0, n = 0; i < nfft; i++, n = n+2) {
512         work2[n] = fxx[i]*work1[n+1];
513         work2[n+1] = -fxx[i]*work1[n];
514     }
515
516     fft2->compute(work2,work2,-1);
517
518     int BS = (nyhi_in - nylo_in + 1) * (nxhi_in - nxlo_in + 1) * 2;
519     #pragma omp parallel
520     {
521         #pragma omp for private(n,j,i) nowait
522         for (k = nzlo_in; k <= nzhi_in; k++) {
523             n = (k - nzlo_in) * BS;
524             for (j = nylo_in; j <= nyhi_in; j++) {
525                 for (i = nxlo_in; i <= nxhi_in; i++) {
526                     vdx_brick[k][j][i] = work2[n];
527                     n += 2;
528                 } //i

```

```

529     } //j
530 } //k
531
532 // y direction gradient
533 #pragma omp for
534 for (i = 0, n = 0; i < nfft; i++, n=n+2) {
535     work2[n] = fky[i]*work1[n+1];
536     work2[n+1] = -fky[i]*work1[n];
537 }
538 } //parallel
539
540 ... //(rest of code omitted)
541
542 }

```

Listing 7: Sketch of parallelized version void PPPM::poisson_ik_triclinic() method

544 The solution for the second case (the three nested loops) is trickier: we need to privatize variable n and
545 to do that we need to manually calculate the initial value for n at each iteration of the outer-most loop.
546 This is easily done through the use of the added variable BS that stores the increments of the variable n in
547 the two inner-most loops. Once all of these is done (see Listing 7), the outer-most loop can be parallelized
548 by simply privatizing n , j and i .

549 As a last comment, the parallel region opened for the nested loops of x direction is used for the first loop
550 of y direction as depicted in Listing 7. The same is done between y direction and z direction (not shown in
551 the Listing).

552 4.4. Overlapping computation and communication

553 Last, but not least, we show here how to effectively overlap computation with MPI communication. More
554 precisely, we have worked in `remap_3d` function, which is called several times in the same region of code
555 mentioned on the last part of previous subsection. As it can be seen in Listing 8, the function consists in 4
556 differentiated parts:

- 557 1. A sequence of `MPI_Irecv` calls to receive data from other processes.
- 558 2. A sequence of (`pack`, `MPI_Send`) calls that packs and sends data to other processes.
- 559 3. A call to `pack` and `unpack` to manage the data of the calling process.
- 560 4. A sequence of (`MPI_Waitany`, `unpack`) to wait for the corresponding `MPI_Irecv` to get the data and
561 put it in the required memory location.

```

562
563
564 void remap_3d(FFT_SCALAR *in , FFT_SCALAR *out , FFT_SCALAR *buf ,
565             struct remap_plan_3d *plan)
566 {
567     ... // (omitted code)
568
569     // post all recvs into scratch space
570     for (irecv = 0; irecv < plan->nrecv; irecv++) {
571         MPI_Irecv(&scratch[plan->recv_bufloc[irecv]], plan->recv_size[irecv],
572                 MPI_FFT_SCALAR, plan->recv_proc[irecv], 0,
573                 plan->comm, &plan->request[irecv]);
574     }
575
576     // send all messages to other procs
577     for (isend = 0; isend < plan->nsend; isend++) {
578         plan->pack(&in[plan->send_offset[isend]],
579                 plan->sendbuf, &plan->packplan[isend]);
580         MPI_Send(plan->sendbuf, plan->send_size[isend], MPI_FFT_SCALAR,
581                 plan->send_proc[isend], 0, plan->comm);
582     }
583

```

```

584 // copy in -> scratch -> out for self data
585 if (plan->self) {
586     isend = plan->nsend;
587     irecv = plan->nrecv;
588     plan->pack(&in[plan->send_offset[isend]],
589             &scratch[plan->recv_bufloc[irecv]],
590             &plan->packplan[isend]);
591     plan->unpack(&scratch[plan->recv_bufloc[irecv]],
592               &out[plan->recv_offset[irecv]],&plan->unpackplan[irecv]);
593 }
594
595 // unpack all messages from scratch -> out
596 for (i = 0; i < plan->nrecv; i++) {
597     MPI_Waitany(plan->nrecv, plan->request, &irecv, MPI_STATUS_IGNORE);
598     plan->unpack(&scratch[plan->recv_bufloc[irecv]],
599               &out[plan->recv_offset[irecv]], &plan->unpackplan[irecv]);
600 }
601
602 ... // (omitted code)
603 }

```

Listing 8: Sketch of the original code in `remap_3d()` method

The changes done in the code to allow computation and communication overlapping can be seen in Listing 9 and are summarized in the following items:

1. All code is wrapped by a parallel and a single constructs to create the parallel OpenMP region and allow only one thread to enter the code to create tasks.
2. The loop corresponding to point number 2 of the original code has been split into two loops: one loop doing all the packs and the other doing all the `MPI_Send`. The loop doing the packs has been moved to the very beginning of the function and each pack has been defined as a task.
3. To allow the previous taskification, `plan->sendbuf` has been redefined (not shown) and now is a buffer of buffers indexed by `isend`: this allows for all pack tasks to be independent.
4. As it is independent of the rest of the communications, the self-data management of point 3 of the original code has been moved next and taskified.
5. A `taskwait` is needed just after the loop with `MPI_Irecv` because the following `MPI_Send` needs the tasks with packs defined in point 2 to be finished.
6. Finally, the ununpacks of the last loop have been defined as tasks: in this way, the following `MPI_Waitany` does not need to wait for the previous unpack to finish.

```

620
621
622 void remap_3d(FFT_SCALAR *in, FFT_SCALAR *out, FFT_SCALAR *buf,
623             struct remap_plan_3d *plan)
624 {
625     ... // (omitted code)
626
627     #pragma omp parallel
628     #pragma omp single
629     {
630         for (isend = 0; isend < plan->nsend; isend++) {
631             #pragma omp task firstprivate(isend)
632             plan->pack(&in[plan->send_offset[isend]],
633                     &plan->sendbuf[isend*plan->sendbuf_size], &plan->packplan[isend]);
634         }
635
636         // copy in -> scratch -> out for self data
637         if (plan->self) {
638             isend = plan->nsend;
639             irecv = plan->nrecv;
640             #pragma omp task firstprivate(isend, irecv)
641             {

```

```

642     plan->pack(&in[plan->send_offset[isend]],
643              &scratch[plan->recv_bufloc[irecv]],
644              &plan->packplan[isend]);
645     plan->unpack(&scratch[plan->recv_bufloc[irecv]],
646               &out[plan->recv_offset[irecv]],&plan->unpackplan[irecv]);
647 }
648 }
649
650 // post all recvs into scratch space
651 for (irecv = 0; irecv < plan->nrecv; irecv++) {
652     MPI_Irecv(&scratch[plan->recv_bufloc[irecv]], plan->recv_size[irecv],
653             MPI FFT_SCALAR, plan->recv_proc[irecv], 0,
654             plan->comm,&plan->request[irecv]);
655 }
656 #pragma omp taskwait
657 // send all messages to other procs
658 for (isend = 0; isend < plan->nsend; isend++) {
659     MPI_Send(&plan->sendbuf[isend*plan->sendbuf_size], plan->send_size[isend], MPI FFT_SCALAR,
660            plan->send_proc[isend], 0, plan->comm);
661 }
662
663 // unpack all messages from scratch -> out
664 for (i = 0; i < plan->nrecv; i++) {
665     MPI_Waitany(plan->nrecv, plan->request,&irecv, MPI_STATUS_IGNORE);
666     #pragma omp task firstprivate(irecv)
667     plan->unpack(&scratch[plan->recv_bufloc[irecv]],
668               &out[plan->recv_offset[irecv]],&plan->unpackplan[irecv]);
669 }
670 } //parallel and single
671 ... // (omitted code)
672 }
673 }

```

Listing 9: Sketch of the modified code in `remap_3d()` method, including the code reordering and the taskification of packs and unpacks

Figure 5 shows how computation and communication have been effectively overlapped. The timelines correspond to a trace of a run with 8 MPI processes with 6 OpenMP threads each. The upper timeline represents the MPI calls (being pink \rightarrow `MPI_Irecv`, blue \rightarrow `MPI_Send`, and green \rightarrow `MPI_Waitany`) and the bottom timeline represents task execution. Figure 5 is a zoom of just one invocation of `remap_3d` method for the first 2 processes (12 threads in total) used in the execution. In this case, `isend = irecv = 3` so each process executes 3 pack tasks + the pack/unpack task corresponding to the self data (the tasks on the left part), and 3 unpack tasks (on the right). It can be seen how the packs are now overlapped with the `MPI_Irecv` calls at the beginning and how the `MPI_Waitany`s at the end do not need to wait for the execution of the previous unpack.

5. Evaluation

5.1. Environment

The experiments have been performed on MareNostrum4 [26]. This supercomputer is based on Intel Xeon Platinum processors from the Skylake generation. It is a Lenovo system composed of SD530 Compute Racks, an Intel Omni-Path high performance network interconnect and running SuSE Linux Enterprise Server as operating system. Compute nodes are equipped with:

- 690 • 2 sockets Intel Xeon Platinum 8160 CPU with 24 cores each @ 2.10GHz for a total of 48 cores per
- 691 node.
- 692 • L1d 32K; L1i 32K; L2 cache 1024K; L3 cache 33729K.
- 693 • 96 GB of main memory 1.88 GB/core.

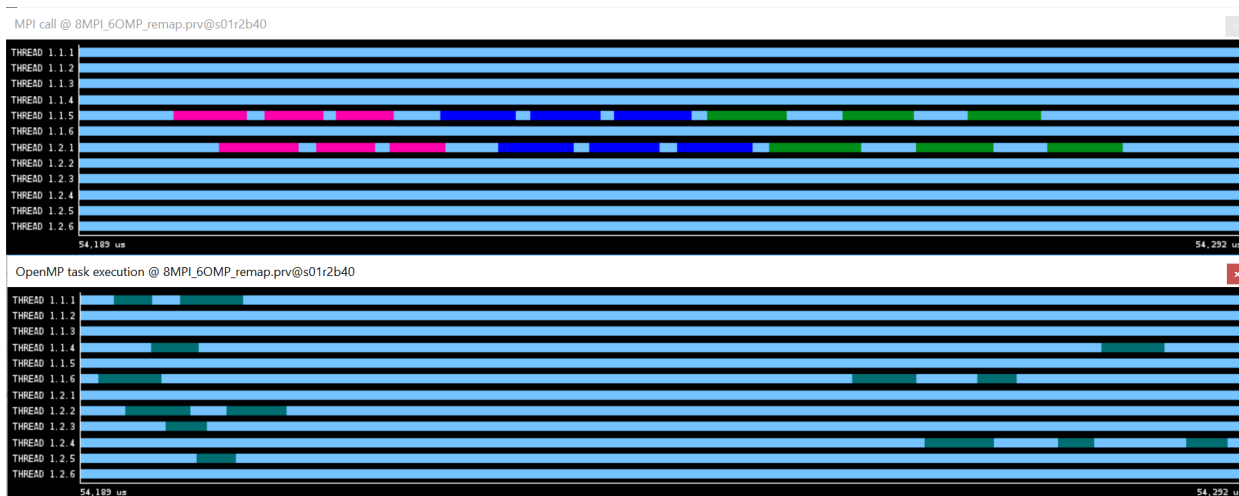


Figure 5: Overlapping computation and communication in remap_3d method.

- 694 • 100 Gbit/s Intel Omni-Path HFI Silicon 100 series PCI-E adapter.
- 695 • 10 Gbit Ethernet.
- 696 • 200 GB local SSD available as temporal storage during jobs.
- 697 • The processors support well-known vectorization instructions such as SSE, AVX up to AVX-512.

698 The software environment used is as follows:

- 699 • LAMMPS 20Nov19
- 700 • Intel 17.0.4 20170411 compiler

701 5.2. Benchmark Characterization

702 Figure 6 shows efficiency metrics for all the testcases and benchmarks studied in this work. These
 703 performance metrics correspond to MPI-only executions and using 48 MPI ranks in all cases. Most of them
 704 present a poor parallel efficiency, considering that only 48 MPI ranks are considered. Particularly bad are
 705 the cases of short-range interactions benchmarks and the CG-GO testcase. The reasons for poor parallel
 706 efficiency are diverse. While the main problem of short-range interactions benchmarks is Communication
 707 Efficiency, the most limiting factor of the CG-GO benchmark is Load Balance with an extremely low
 708 value (in contrast with the rest of the benchmarks). Mid-range interactions benchmarks have very similar
 709 characteristics: although with slightly different weights on the two components, they both have the same
 710 Parallel Efficiency value. Finally, Rhodopsin is the best performing benchmark. Note that, the Epoxy
 711 testcase presents quite different characteristics when compared with CG-GO as discussed later.

712 So with all these benchmarks we cover very different scenarios both in terms of type of simulation and
 713 in terms of performance metrics characteristics.

714 5.3. Execution Time

715 In this section we present the wall time execution of all testcases and benchmarks for different setups
 716 including the so-called "Vanilla" (i.e. the regular MPI-only execution), "Balance" (i.e. MPI-only execution
 717 but including the balancing mechanisms provided by the LAMMPS implementation) and different hybrid
 718 configurations. 48 cores are used in all cases and, for the hybrid configurations, only configurations up to
 719 the point where using more OpenMP threads translates in worst performance than the MPI-only case are
 720 shown. In all benchmarks, the number of timesteps has been increased to have an execution wall time of
 721 at least 1 minute for the Vanilla case in order to better appreciate execution time differences (i.e. the wall
 722 times reported are the ones provided by LAMMPS and they only have a precision of seconds).

	Testcases:		Short-range:		Mid-range:		Long-range:
	Epoxy	CG-GO	Granular chute	Polymer chain	EAM	Lennard-Jones	Rhodopsin
Parallel efficiency	0,69	0,40	0,38	0,38	0,58	0,58	0,64
-- Load balance	0,85	0,42	0,85	0,80	0,81	0,95	0,94
-- Communication efficiency	0,81	0,95	0,45	0,48	0,72	0,62	0,69

Figure 6: Efficiency metrics of all testcases and benchmarks using 48 MPI ranks.

5.3.1. Testcases

Let us start by the CG-GO testcase as it is the one that shows the greatest load imbalance. Figure 7 (left part) presents the execution times for different setups of this benchmark. The bars represent the maximum time spent by an MPI rank on a given code section as reported by LAMMPS in its performance execution report while the blue line that traverses the figure represents the total wall time of the different executions. So the difference between both heights gives an idea of the load imbalance that affects a given configuration. As it can be seen, this is huge for the Vanilla case. The Balance version reduces this difference a lot (mainly by reducing the "Comm" maximum execution time). Note, however how a wide range of hybrid configurations (from 24 to 8 MPI processes) do this better achieving also lower (i.e. better) wall execution times.

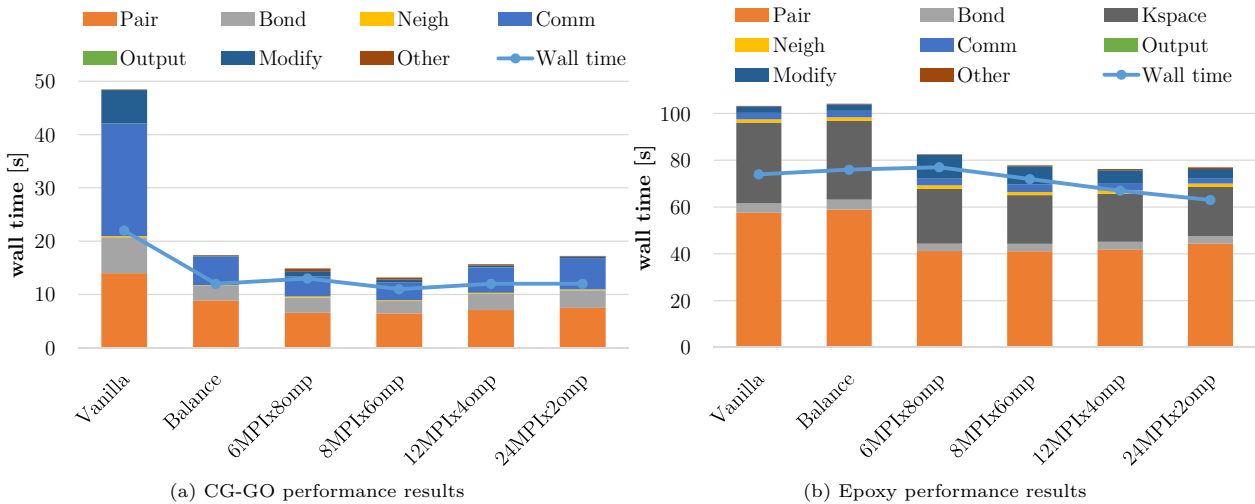


Figure 7: Performance results for test cases inputs

Figure 8 represent 10 timesteps of the executions using the Vanilla, Balance and Hybrid versions at the same timescale. As it can be seen, the Vanilla version presents a heavy load imbalance: this allows the Balance version to achieve a 43% improvement in execution time (timeline in the middle). But, more interestingly, the Hybrid version (bottom timeline) is 12% faster than the Balance version and it achieves a 50% improvement when compared with the Vanilla version.

Figure 9 shows the parallel functions executed by the Hybrid version. The time spent in OpenMP parallel regions in this case is 80,6% of the total execution time.

A very different scenario is shown in the right plot of the figure, that presents the execution time of the Epoxy resin testcase. Note the different characteristics of this testcase when compared with the CG-GO. First, the imbalance is not so relevant for the Vanilla scenario. In fact, it can be seen how the use of the balance mechanisms provided by LAMMPS actually make the execution slower (i.e. it adds overhead without any improvement). Hybrid configurations from 24 to 8 MPI ranks perform better in terms of both load balance and execution time, mainly due to the better performance of the pairing ("Pair", orange in the figure) of the hybrid cases, meaning that the OpenMP parallelization is more efficient than the MPI one.

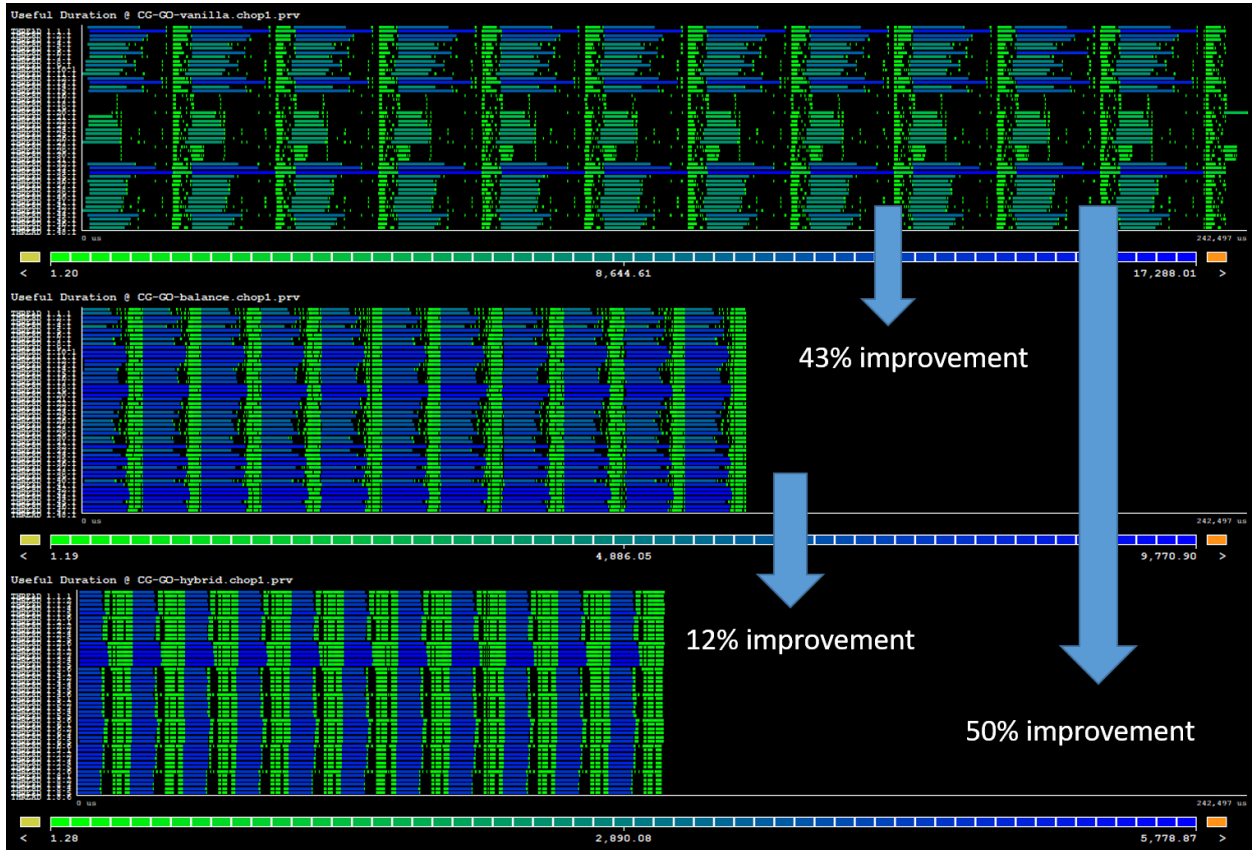


Figure 8: Useful duration timelines for the CG-GO testcase (upper: Vanilla, middle: Balance, lower:Hybrid).

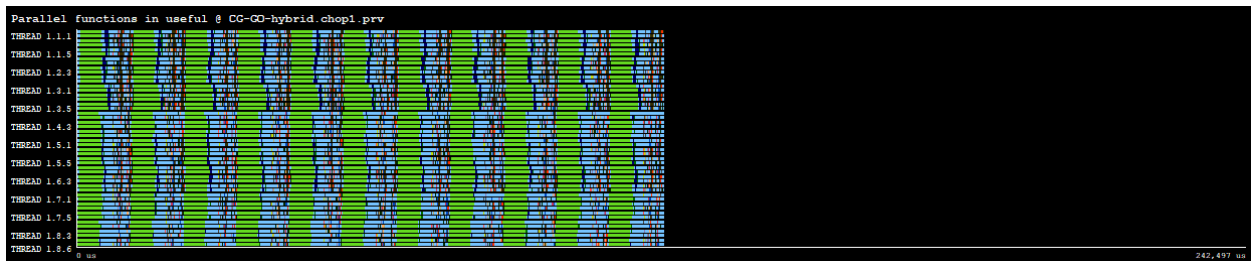


Figure 9: Parallel functions timeline for the CG-GO testcase (Hybrid version).

747 This testcase clearly shows how the Hybrid version is able to improve the Vanilla even when the Balance
 748 version is not, demonstrating that hybridization provides other benefits than just load balance.

749 5.3.2. Short-range interactions benchmarks

750 Figure 10 present the results for the short-range interactions benchmarks. For both benchmarks the
 751 analysis is actually the same: nor the Balance mechanism nor the Hybrid solution are able to improve the
 752 Vanilla setup.

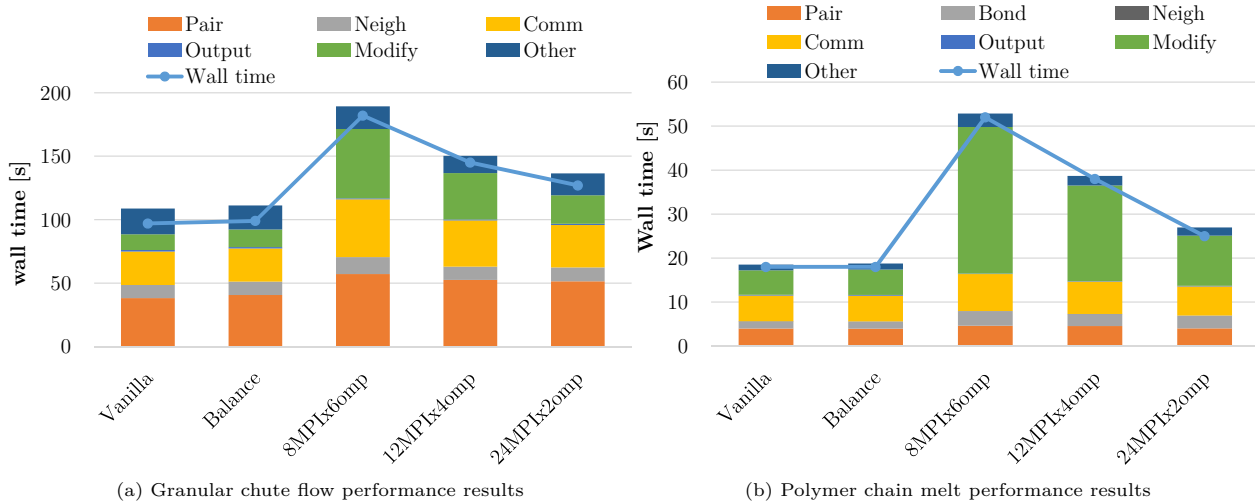


Figure 10: Short-range interactions benchmarks results

753 For the Balance mechanism, the explanation is clear: the problem of the Vanilla setup, if any, is not load
 754 imbalance. For the Hybrid configuration, we will take the Polymer chain benchmark as a representative but
 755 a similar analysis could be done for the Granular chute. Figure 11 represents 10 timesteps of the Polymer
 756 chain benchmark. The timelines, of course at the same timescale, show two main reasons that explain why
 757 this short-range interactions cases do not benefit from the use of a Hybrid implementation:

- 758 1. The most computational intensive part (dark blue) is not OpenMP parallelized so the execution time
 759 is increased.
- 760 2. Only very small parts of the less computational intensive part (light green) are parallelized with
 761 OpenMP, leading also to an increased execution time. This can be perfectly seen in Figure 12 where
 762 the OpenMP parallel regions are depicted (meaning light blue no parallel region at all, i.e. sequential
 763 execution). Actually, the percentage of time of the whole execution spent in OpenMP parallel regions
 764 is only of 23%. This suggests that there is a lot of room for improvement by parallelizing other parts
 765 of the code used by this benchmark in similar ways as explained in Section 3.3.

766 5.3.3. Mid-range interactions benchmarks

767 Figure 13 (left) presents the execution times of the EAM benchmark for different configurations. In this
 768 case, all the versions perform quite similar. It is noticeable, however that the best performing version is
 769 Hybrid for the 24MPIx2omp case: 70 seconds in contrast with the 76 seconds of the Vanilla or the 75 seconds
 770 of the Balance version. A similar analysis can be done for the Lennar-Jones benchmark (Figure 13 (right)):
 771 the only noticeable difference is that in this case Balance is a bit worst than Vanilla (just one second) while
 772 the 24MPIx2omp Hybrid configuration is still able to improve by 4 seconds the Vanilla case.

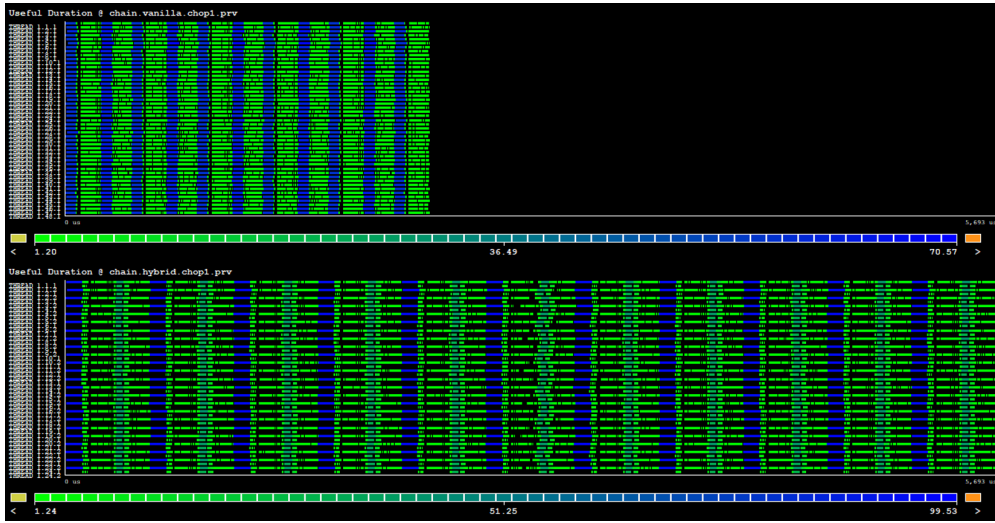


Figure 11: Useful duration timelines for the Polymer chain melt benchmark (Upper part: Vanilla execution, lower part: Hybrid version).

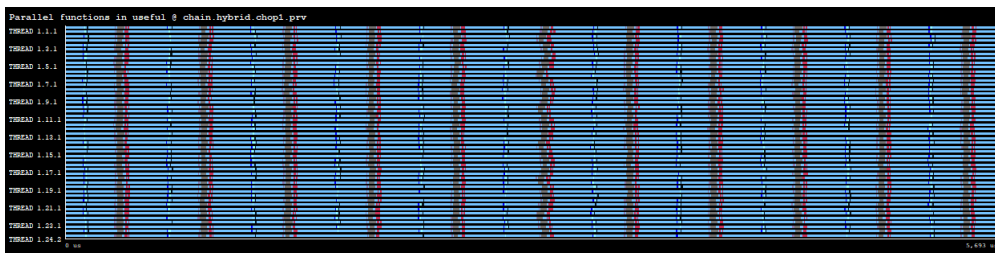


Figure 12: Parallel functions timeline for the Polymer chain melt benchmark (Hybrid version).

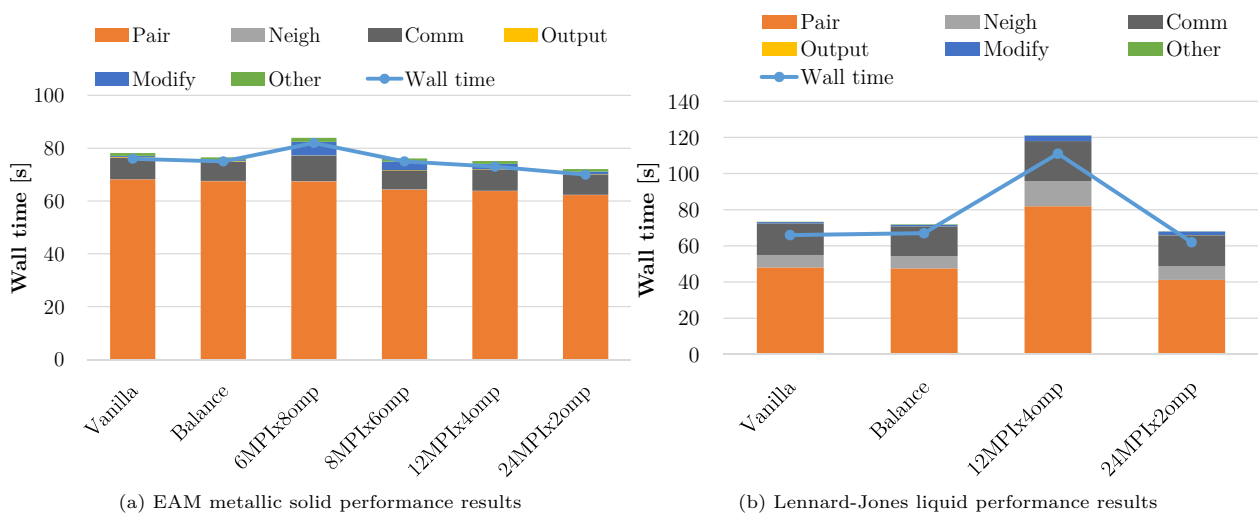


Figure 13: Mid-range interactions benchmarks results

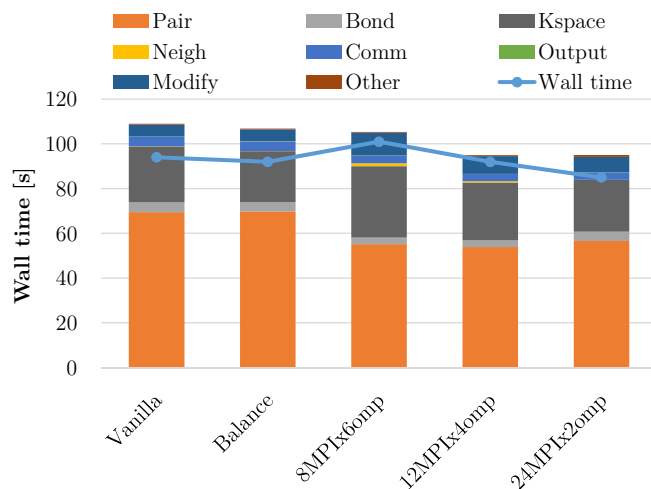


Figure 14: Rhodopsin protein benchmark results.

773 5.3.4. Long-range interactions benchmarks

774 Figure 14 presents the execution time of the Rhodopsin benchmark for different configurations. The
 775 three Hybrid configurations on the right are able to outperform both Vanilla and Balance versions. The
 776 pairing process (orange bar) is much faster in the Hybrid configurations.

777 Figure 15 represents two timelines of 10 timesteps of the Rhodopsin benchmark execution at the same
 778 timescale. As it can be visually noted, the execution of the Hybrid case is significantly faster. This gain
 779 in performance comes mainly from the pairing phase (blue sections in the timelines) done in the compute
 780 function in the PairLJCharmmCoulLongOMP module of LAMMPS which is faster for the Hybrid case.

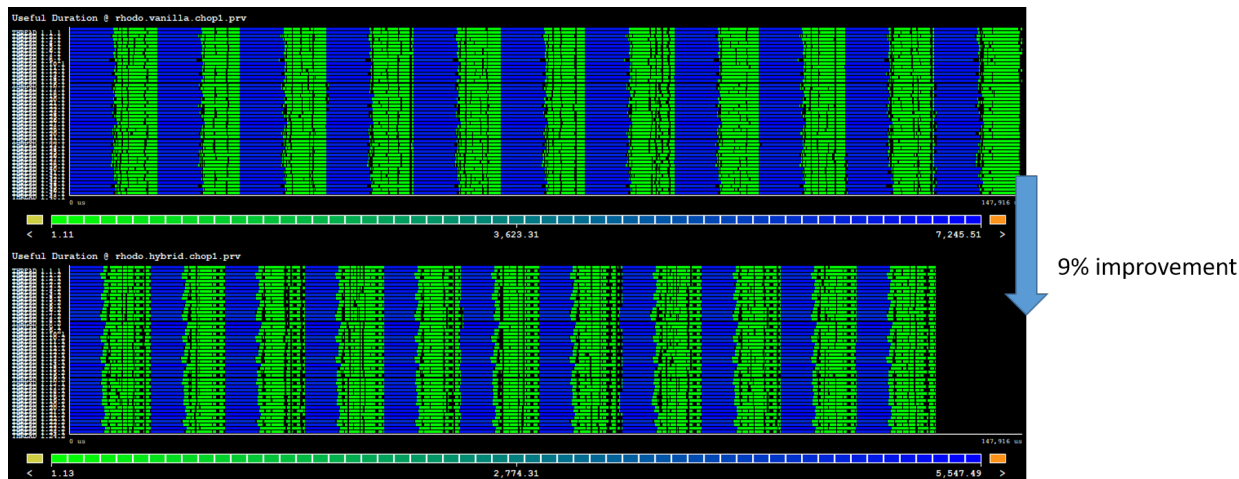


Figure 15: Useful duration timelines for the Rhodopsin protein benchmark (top: Vanilla, bottom: Hybrid).

781 As it can be seen on the red parts of Figure 16, this function is fully OpenMP parallelized. Note, however,
 782 that there are still other parts of the code not parallelized with OpenMP (light blue in Figure 16 and black
 783 in bottom timeline of Figure 15), making the execution of the green areas in Figure 15 being slightly faster
 784 for the Vanilla case. But even so, the improvement achieved by the OpenMP implementation of compute
 785 subroutine is able to compensate by far this loss in performance.

786 All in all, the ratio of time in parallel regions with respect to the whole execution time is pretty high:
 787 76,6%. We can now compare short-range with long-range interactions benchmarks and explain why Hybrid

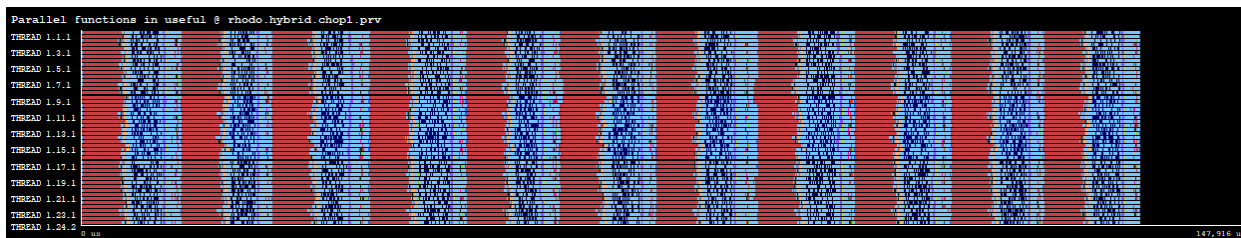


Figure 16: Parallel functions timeline for the Rhodopsin protein benchmark (Hybrid version).

788 implementation is able to improve the performance of the later but not of the formers. Comparing Figures 12
 789 and 16 one can see how the percentage of time outside any parallel OpenMP region (light blue areas) for
 790 both benchmarks is drastically different. Actually, the percentage of time inside OpenMP parallel regions
 791 of the Polymer chain benchmark is only 23,7%. This is exacerbated if we focus on the pairing phases: while
 792 in the case of Rhodopsin benchmark this corresponds to the red parts in Figure 16, so it represents a high
 793 percentage of the whole execution, this is not the case for the Polymer chain benchmark where the pairing
 794 (light brown in Figure 12) corresponds to a very residual part of the whole execution (in this benchmark,
 795 subroutine compute of PairLJCutOMP module).

796 So, at the end, the reason that ultimately explains the different behaviour of the two kind of benchmarks
 797 is that the weight of the pairing process (which is the most important part parallelized with OpenMP in all
 798 benchmarks) in the whole execution is very high for the long-range interactions benchmarks (45,6% in the
 799 Rhodopsin benchmark) while it is insignificant for the short-range interactions benchmarks (8,4% for the
 800 Polymer chain benchmark).

801 6. Conclusions

802 This paper has shown the potential of the proposed hybrid MPI+OpenMP approach to effectively alle-
 803 viate performance problems such as load imbalance in LAMMPS simulations.

804 LAMMPS provides a ready-to-use balacing mechanism to partially solve load balancing problems in
 805 molecular dynamics simulations with non uniform atom densities. The LAMMPS balancing mechanism
 806 shows high efficiency compared to MPI-only simulation in case of high load imbalance (CG-GO testcase). We
 807 have introduced the use of an MPI+OpenMP hybrid implementation of LAMMPS as a third option. Fur-
 808 thermore, we have complemented the current partial OpenMP implementation of LAMMPS with additions
 809 and modifications, driven by the Epoxy testcase.

810 Our proposed modified version LAMMPS has been extensively compared against the baseline MPI case
 811 and against the use of the LAMMPS balance mechanism. Five benchmarks present in the LAMMPS
 812 distribution with varying range of interaction (short, mid and long-range) together with two testcases with
 813 very different characteristics were used for the comparison.

814 For the short-range interactions benchmarks, the regular MPI-only version was the best performing. As
 815 long as they do not present a load imbalance problem, the balancing mechanism does not provide any benefit
 816 in this case. The problem with the hybrid setup in this case is that only a small fraction of the simulations
 817 (~20%) runs in OpenMP parallel regions. This suggests that more additions similar to the ones proposed
 818 in Section 3.3 could be done to the OpenMP LAMMPS implementation.

819 For the mid-range interactions benchmarks, the hybrid option was the best in all cases. The balance
 820 mechanism only improves a bit the EAM simulation while it is the worst option for the Lennard-Jones
 821 benchmark. These results indicate that the hybrid implementation is able to improve performance metrics
 822 others than load balance, such as communication efficiency.

823 In the Rhodopsin benchmark (long-range interactions) the execution time is mainly dominated by the
 824 LAMMPS pairing process. The obtained results show that the OpenMP parallelization of the pairing is
 825 much faster than the MPI one, making the hybrid approach the best option also for this benchmark.

826 Regarding the highly-imbalanced testcase (CG-GO), the balancing mechanism shows its potential achiev-
827 ing a 43% improvement with respect of the regular MPI simulation. Interestingly, the hybrid version is able
828 to improve even further, up to 50%.

829 In the case of the Epoxy resin testcase (the one that motivated the implementations explained in Sec-
830 tion 3.3), the use of the balance mechanism only adds overhead (the execution is slower than for the regular
831 MPI version). The hybrid implementation, on the contrary, is the best option showing again that it is able
832 to improve simulations where the load balance is not the main problem.

833 So, the overall conclusion is that LAMMPS hybrid setups are able to handle scenarios with very high
834 load imbalance at least as well (if not better) as the LAMMPS balance mechanism while also providing
835 benefits in other scenarios where load balance is not the main performance bottleneck.

836 Our suggestion to LAMMPS developers, and MD in general, is to put effort into hybridizing the code
837 with an MPI+OpenMP strategy instead of implementing ad-hoc balancing methods. Because the hybrid
838 code not only can address more dynamic load imbalances but also improve the parallel efficiency reducing
839 the communication.

840 7. Acknowledgements

841 This work is partially supported by the Spanish Government through Programa Severo Ochoa (SEV-
842 2015-0493), by the Spanish Ministry of Science and Technology (TIN2015-65316-P), by the Generalitat
843 de Catalunya (2017-SGR-1414), and by the European POP CoE (GA n. 824080). This work is also
844 funded as part of the European Union Horizon 2020 research and innovation programme under grant agree-
845 ment nos. 800925 (VECMA project; www.vecma.eu) and 823712 (CompBioMed2 Centre of Excellence;
846 www.compbioimed2.eu), as well as the UK EPSRC for the UK High-End Computing Consortium (grant no.
847 EP/R029598/1).

848 References

- 849 [1] LAMMPS. [link].
850 URL <https://lammps.sandia.gov/>
- 851 [2] S. Plimpton, Fast parallel algorithms for short-range molecular dynamics, *Journal of computational physics* 117 (1) (1995)
852 1–19.
- 853 [3] K. D. Devine, E. G. Boman, R. T. Heaphy, B. A. Hendrickson, J. D. Teresco, J. Faik, J. E. Flaherty, L. G. Gervasio, New
854 challenges in dynamic load balancing, *Applied Numerical Mathematics* 52 (2-3) (2005) 133–152.
- 855 [4] Y. Deng, R. F. Peierls, C. Rivera, An adaptive load balancing method for parallel molecular dynamics simulations, *Journal*
856 *of Computational Physics* 161 (1) (2000) 250 – 263. doi:<https://doi.org/10.1006/jcph.2000.6501>.
857 URL <http://www.sciencedirect.com/science/article/pii/S002199910096501X>
- 858 [5] S. Plimpton, R. Pollock, M. Stevens, Particle-mesh ewald and rrespa for parallel molecular dynamics simulations, *Proc.*
859 *8th SIAM Conf. on Parallel Processing for Scientific Computing* (08 2000).
- 860 [6] C. Walshaw, M. Cross, Mesh partitioning: a multilevel balancing and refinement algorithm, *SIAM Journal on Scientific*
861 *Computing* 22 (1) (2000) 63–80.
- 862 [7] G. Karypis, V. Kumar, A fast and high quality multilevel scheme for partitioning irregular graphs, *SIAM Journal on*
863 *scientific Computing* 20 (1) (1998) 359–392.
- 864 [8] D. F. Harlacher, H. Klimach, S. Roller, C. Siebert, F. Wolf, Dynamic load balancing for unstructured meshes on space-
865 filling curves, in: *2012 IEEE 26th International Parallel and Distributed Processing Symposium Workshops & PhD Forum*,
866 *IEEE*, 2012, pp. 1661–1669.
- 867 [9] K. Schloegel, G. Karypis, V. Kumar, A unified algorithm for load-balancing adaptive scientific simulations, in: *SC’00:*
868 *Proceedings of the 2000 ACM/IEEE Conference on Supercomputing*, *IEEE*, 2000, pp. 59–59.
- 869 [10] LAMMPS-Balance. [link].
870 URL <https://lammps.sandia.gov/doc/balance.html>
- 871 [11] C. Huang, O. Lawlor, L. V. Kale, Adaptive mpi, in: *International workshop on languages and compilers for parallel*
872 *computing*, Springer, 2003, pp. 306–322.
- 873 [12] B. Acun, A. Gupta, N. Jain, A. Langer, H. Menon, E. Mikida, X. Ni, M. Robson, Y. Sun, E. Totonni, et al., Parallel
874 programming with migratable objects: Charm++ in practice, in: *SC’14: Proceedings of the International Conference for*
875 *High Performance Computing, Networking, Storage and Analysis*, *IEEE*, 2014, pp. 647–658.
- 876 [13] M. Etinski, J. Corbalan, J. Labarta, M. Valero, A. Veidenbaum, Power-aware load balancing of large scale mpi applications,
877 in: *2009 IEEE International Symposium on Parallel & Distributed Processing*, *IEEE*, 2009, pp. 1–8.
- 878 [14] M. Garcia, J. Corbalan, J. Labarta, LeWI: A Runtime Balancing Algorithm for Nested Parallelism, in: *Proceedings of the*
879 *International Conference on Parallel Processing (ICPP09)*, 2009.

- 880 [15] M. Garcia-Gasulla, F. Mantovani, M. Josep-Fabrego, B. Eguzkitza, G. Houzeaux, Runtime mechanisms to survive new
881 hpc architectures: a use case in human respiratory simulations, *The International Journal of High Performance Computing*
882 *Applications* 34 (1) (2020) 42–56.
- 883 [16] R. Rabenseifner, G. Hager, G. Jost, Hybrid mpi/openmp parallel programming on clusters of multi-core smp nodes, in:
884 2009 17th Euromicro international conference on parallel, distributed and network-based processing, IEEE, 2009, pp.
885 427–436.
- 886 [17] R. Rabenseifner, G. Wellein, Communication and optimization aspects of parallel programming models on hybrid archi-
887 tectures, *The International Journal of High Performance Computing Applications* 17 (1) (2003) 49–62.
- 888 [18] Barcelona Supercomputing Center, Paraver.
889 URL <https://tools.bsc.es/paraver>
- 890 [19] V. Pillet, J. Labarta, T. Cortes, S. Girona, Paraver: A tool to visualize and analyze parallel code, in: *Proceedings of*
891 *WoTUG-18: transputer and occam developments*, Vol. 44, 1995, pp. 17–31.
- 892 [20] Barcelona Supercomputing Center, Extrae.
893 URL <https://tools.bsc.es/extrae>
- 894 [21] H. Servat, et al., Framework for a productive performance optimization, *Parallel Computing* 39 (8) (2013) 336–353.
- 895 [22] D. Fincham, Parallel computers and molecular simulation, *Molecular Simulation* 1 (1-2) (1987) 1–45.
896 arXiv:<https://doi.org/10.1080/08927028708080929>, doi:10.1080/08927028708080929.
897 URL <https://doi.org/10.1080/08927028708080929>
- 898 [23] W. Smith, Molecular dynamics on hypercube parallel computers, *Computer Physics Communications* 62 (2) (1991) 229 –
899 248. doi:[https://doi.org/10.1016/0010-4655\(91\)90097-5](https://doi.org/10.1016/0010-4655(91)90097-5).
900 URL <http://www.sciencedirect.com/science/article/pii/0010465591900975>
- 901 [24] S. Plimpton, B. Hendrickson, A new parallel method for molecular dynamics simulation of macromolecular systems, *Journal*
902 *of Computational Chemistry* 17 (3) (1996) 326–337. doi:[https://doi.org/10.1002/\(SICI\)1096-987X\(199602\)17:3<326::AID-](https://doi.org/10.1002/(SICI)1096-987X(199602)17:3<326::AID-JCC7>3.0.CO;2-X)
903 [JCC7>3.0.CO;2-X](https://doi.org/10.1002/(SICI)1096-987X(199602)17:3<326::AID-JCC7>3.0.CO;2-X).
- 904 [25] R. Rabenseifner, G. Hager, G. Jost, Hybrid mpi/openmp parallel programming on clusters of multi-core smp nodes, in:
905 2009 17th Euromicro International Conference on Parallel, Distributed and Network-based Processing, 2009, pp. 427–436.
906 doi:10.1109/PDP.2009.43.
- 907 [26] Barcelona Supercomputing Center, Marenostum4.
908 URL <https://www.bsc.es/marenostum/marenostum>