

## Parallel Performance of Molecular Dynamics Trajectory Analysis

Journal:	<i>Concurrency and Computation: Practice and Experience</i>
Manuscript ID	Draft
Editor Selection:	Prof. David Walker
Wiley - Manuscript type:	Research Article
Date Submitted by the Author:	n/a
Complete List of Authors:	Khoshlessan, Mahzad; Arizona State University, Physics Paraskevagos, Ioannis; Rutgers University, Department of Electrical & Computer Engineering Fox, Geoffrey; Indiana University Bloomington, Digital Science Center Jha, Shantenu; Rutgers University, Department of Electrical & Computer Engineering Beckstein, Oliver; Arizona State University, Physics; Arizona State University, Center for Biological Physics
Keywords:	Python, Molecular Dynamics, Straggler, Trajectory Analysis, MPI, HPC

SCHOLARONE™  
Manuscripts

## RESEARCH ARTICLE

# Parallel Performance of Molecular Dynamics Trajectory Analysis

Mahzad Khoshlessan<sup>1</sup> | Ioannis Paraskevatos<sup>2</sup> | Geoffrey C. Fox<sup>3</sup> | Shantenu Jha<sup>2</sup> | Oliver Beckstein<sup>\*1,4</sup>

<sup>1</sup>Department of Physics, Arizona State University, Tempe, AZ 85281, USA

<sup>2</sup>Department of Electrical & Computer Engineering, Rutgers University, Piscataway, NJ 08854, USA

<sup>3</sup>Digital Science Center, Indiana University, Bloomington, IN 47405, USA

<sup>4</sup>Center for Biological Physics, Arizona State University, Tempe, AZ 85281, USA

**Correspondence**

\*Oliver Beckstein, Department of Physics, Arizona State University, Tempe, AZ 85281, USA. Email: [oliver.beckstein@asu.edu](mailto:oliver.beckstein@asu.edu)

**Summary**

The performance of biomolecular molecular dynamics (MD) simulations has steadily increased on modern high performance computing (HPC) resources but acceleration of the analysis of the output trajectories has lagged behind so that analyzing simulations is becoming a bottleneck. To close this gap, we studied the performance of parallel trajectory analysis with MPI and the Python MDAnalysis library on three different XSEDE supercomputers where trajectories were read from a Lustre parallel file system. Strong scaling performance was impeded by stragglers, MPI processes that were slower than the typical process. Stragglers were less prevalent for compute-bound workloads, thus pointing to file reading as a crucial bottleneck for scaling. However, a more complicated picture emerged in which both the computation and the ingestion of data exhibited close to ideal strong scaling behavior whereas stragglers were primarily caused by either large MPI communication costs or long times to open the single shared trajectory file. We improved overall strong scaling performance by either subfiling (splitting the trajectory into separate files) or MPI-IO with Parallel HDF5 trajectory files. We obtained near ideal strong scaling on up to 384 cores (16 nodes), thus reducing trajectory analysis times by two orders of magnitude compared to the serial approach.

**KEYWORDS:**

Python, MPI, HPC, MDAnalysis, MPI I/O, Global Arrays, HDF5, Straggler, Molecular Dynamics, Big Data, Trajectory Analysis

## 1 | INTRODUCTION

Molecular dynamics (MD) simulations are a powerful method to generate new insights into the function of biomolecules<sup>1–5</sup>. These simulations produce trajectories—time series of atomic coordinates—that now routinely include millions of time steps and can measure Terabytes in size. These trajectories need to be analyzed using statistical mechanics approaches<sup>6,7</sup> but because of the increasing size of data, trajectory analysis is becoming a bottleneck in typical biomolecular simulation scientific workflows<sup>8</sup>. Many data analysis tools and libraries have been developed to extract the desired information from the output trajectories from MD simulations<sup>9–22</sup> but few can efficiently use modern High Performance Computing (HPC) resources to accelerate the analysis stage. MD trajectory analysis primarily requires *reading* of data from the file system; the processed output data are typically negligible in size compared to the input data and therefore we exclusively investigate the reading aspects of trajectory I/O (i.e., the “I”). We focus on the *MDAnalysis* package<sup>17,18</sup>, which is an open-source object-oriented Python library for structural and temporal analysis of MD simulation trajectories and individual protein structures. Although *MDAnalysis* accelerates selected algorithms with OpenMP, it is not clear how to best use it for scaling up analysis on multi-node supercomputers. Here we discuss the challenges and lessons-learned for making parallel analysis on HPC resources feasible with *MDAnalysis*, which should also be broadly applicable to other general purpose trajectory analysis libraries.

Previously, we had used a parallel split-apply-combine approach to study the performance of the commonly performed “RMSD fitting” analysis problem<sup>23–25</sup>, which calculates the minimal root mean squared distance (RMSD) of the positions of a subset of atoms to a reference conformation under optimization of rigid body translations and rotations<sup>7,26,27</sup>. We had investigated two parallel implementations, one using *Dask*<sup>28</sup> and one using the message passing interface (MPI) with *mpi4py*<sup>29,30</sup>. For both *Dask* and MPI, we had previously only been able to obtain good strong scaling performance within a single node. Beyond a single node performance had dropped due to *straggler* tasks, a subset of tasks that had performed abnormally slower than the typical task execution times; the total execution time had become dominated by stragglers and overall performance had decreased. Stragglers are a well-known challenge to improving performance on HPC resources<sup>31</sup> but there has been little discussion of their impact in the biomolecular simulation community.

In the present study, we analyzed the MPI case in more detail to better understand the origin of stragglers with the goal to find parallelization approaches to speed up parallel post-processing of MD trajectories in the *MDAnalysis* library. We especially wanted to make efficient use of the resources provided by current supercomputers such as multiple nodes with hundreds of CPU cores and a Lustre parallel file system.

As in our previous study<sup>23</sup> we selected the commonly used RMSD algorithm implemented in *MDAnalysis* as a typical use case. We employed the single program multiple data (SPMD) paradigm to parallelize this algorithm on three different HPC resources (XSEDE’s *SDSC Comet*, *LSU SuperMic*, and *PSC Bridges*<sup>32</sup>). With SPMD, each process executes essentially the same operations on different parts of the data. The three clusters differed in their architecture but all used Lustre as their parallel file system. We used Python (<https://www.python.org/>), a machine-independent, byte-code interpreted, object-oriented programming language, which is well-established in the biomolecular simulation community with good support for parallel programming for HPC<sup>29,33</sup>. We found that communication and reading I/O were the two main scalability bottlenecks, with some indication that read I/O might have been interfering with the communications. We therefore focused on two different approaches to mitigate I/O bottlenecks: MPI parallel I/O (MPI-IO) with the HDF5 file format and subfilng (trajectory file splitting). For subfilng, we obtained good results with the *Global Arrays* package<sup>33,34</sup>, which provides a convenient layer to access and manage arrays over multiple MPI ranks. Both MPI-IO and subfilng eliminated stragglers and improved the performance with near ideal scaling,  $S(N) = N$ , i.e., the speed-up  $S$  scaled linearly with the number  $N$  of CPU cores while exhibiting a slope of one.

The paper is organized as follows: We first review stragglers and existing approaches to parallelizing MD trajectory analysis in section 2. We describe the software packages and algorithms in section 3 and the benchmarking environment in section 4. Section 5 explains how we measured performance. The main results are presented in section 6, with section 7 demonstrating reproducibility on different supercomputers. We provide general guidelines and lessons-learned in section 8 and finish with conclusions in section 9.

## 2 | BACKGROUND AND RELATED WORK

In our previous work, we found that straightforward implementation of simple parallelization with a split-apply-combine algorithm in Python failed to scale beyond a single compute node<sup>23</sup> because a few tasks (MPI-ranks or *Dask*<sup>28</sup> processes) took much longer than the typical task and so limited the overall performance. However, the cause for these *straggler* tasks remained obscure. Here, we studied the straggler problem in the context of an MPI-parallelized trajectory analysis algorithm in Python and investigated solutions to overcome it. We briefly review stragglers in section 2.1 and summarize existing approaches to parallel trajectory analysis in section 2.2.

### 2.1 | Stragglers

*Stragglers* or *outliers* were traditionally considered in the context of MapReduce jobs that consist of multiple tasks that all have to finish for the job to succeed: A straggler was a task that took an “unusually long time to complete”<sup>35</sup> and therefore substantially impeded job completion. In general, any component of a parallel workflow whose runtime exceeds a typical run time (for example, 1.5 times the median runtime) can be considered a straggler<sup>36</sup>. Stragglers are a challenge for improving performance on HPC resources<sup>31</sup>; they are a known problem in frameworks such as MapReduce<sup>35,36</sup>, Spark<sup>37–40</sup>, Hadoop<sup>35</sup>, cloud data centers<sup>31,41</sup>, and have a high impact on performance and energy consumption of big data systems<sup>42</sup>. Both internal and external factors are known to contribute to stragglers. Internal factors include heterogeneous capacity of worker nodes and resource competition due to other tasks running on the same worker node. External factors include resource competition due to co-hosted applications, input data skew, remote input or output source being too slow, faulty hardware<sup>35,43</sup>, and node mis-configuration<sup>35</sup>. Competition over scarce resources<sup>36</sup>, in particular the network bandwidth, was found to lead to stragglers in writing on Lustre file systems<sup>44</sup>. Garbage collection<sup>37,38</sup>, Java virtual machine (JVM) positioning to cores<sup>37</sup>, delays introduced while the tasks move from the scheduler to execution<sup>39</sup>, disk I/O during shuffling, Java’s just-in-time compilation<sup>38</sup>, output skew<sup>38</sup>, high CPU utilization, disk utilization, unhandled I/O access requests, and network package loss<sup>31</sup> were also among other external factors that might introduce stragglers. A wide variety of approaches have been investigated for detecting and mitigating stragglers, including tuning resource allocation and parallelism such as breaking the workload into many small tasks

that are dynamically scheduled at runtime<sup>45</sup>, slow Node-Threshold<sup>35</sup>, speculative execution<sup>35</sup> and cause/resource-aware task management<sup>36</sup>, sampling or data distribution estimation techniques, SkewTune to avoid data imbalance<sup>46</sup>, dynamic work rebalancing<sup>41</sup>, blocked time analysis<sup>47</sup>, and intelligent scheduling<sup>48</sup>.

In the present study, we analyzed large MD trajectories in parallel with MPI and Python and observed large variations in the completion time of individual MPI ranks. These variations bore some similarity to the straggler tasks observed in MapReduce frameworks so we approached analyzing and eliminating them in a similar fashion by systematically looking at different components of the problem, including read I/O from the shared Lustre file system and MPI communication. Even though we specifically worked in with the *MDAnalysis* package, all these principles and techniques are potentially applicable to MPI-parallelized data analysis in other Python-based libraries.

## 2.2 | Other Packages with Parallel Analysis Capabilities

Different approaches to parallelizing the analysis of MD trajectories have been proposed. HiMach<sup>14</sup> introduces scalable and flexible parallel Python framework to deal with massive MD trajectories, by combining and extending Google's MapReduce and the VMD analysis tool<sup>11</sup>. HiMach's runtime is responsible to parallelize and distribute Map and Reduce classes to assigned cores. HiMach uses parallel I/O for file access during map tasks and MPI\_Allgather in the reduction process. HiMach, however, does not discuss parallel analysis of analysis types that cannot be implemented via MapReduce. Furthermore, HiMach is not available under an open source license, which makes it difficult to integrate community contributions and add new state-of-the-art methods.

Wu et. al.<sup>49</sup> present a scalable parallel framework for distributed-memory post-simulation data analysis. This work consists of an interface that allows a user to write analysis programs sequentially, and the machinery that ensures these programs execute in parallel automatically. The main components of the proposed framework are (1) domain decomposition that splits computational domain into blocks with specified boundary conditions, (2) HDF5 based parallel I/O (3) data exchange that communicates ghost atoms between neighbor blocks, and (4) parallel analysis implementation of a real-world analysis application. This work does not discuss analysis methods which cannot be implemented using MapReduce and is limited to HDF5 file format.

Zazen<sup>50</sup> is a novel task-assignment protocol to overcome the I/O bottleneck for many I/O bound tasks. This protocol caches a copy of simulation output files on the local disks of the compute nodes of a cluster, and uses co-located data access with computation. Zazen is implemented in a parallel disk cache system and avoids the overhead associated with querying metadata servers by reading data in parallel from local disks. This approach has also been used to improve the performance of HiMach<sup>14</sup>. It, however, advocates a specific architecture where a parallel super-computer, which runs the simulations, immediately pushes the trajectory data to a local analysis cluster where trajectory fragments are cached on node-local disks. In the absence of such a specific workflow, one would need to stage the trajectory across nodes, and the time for data distribution is likely to reduce any gains from the parallel analysis.

VMD<sup>11,51</sup> provides molecular visualization and analysis tool through algorithmic and memory efficiency improvements, vectorization of key CPU algorithms, GPU analysis and visualization algorithms, and good parallel I/O performance on supercomputers. It is one of the most advanced programs for the visualization and analysis of MD simulations. It is, however, a large monolithic program, that can only be driven through its built-in Tcl interface and thus is less well suited as a library that allows the rapid development of new algorithms or integration into workflows.

CPPTraj<sup>19</sup> offers multiple levels of parallelization (MPI and OpenMP) in a monolithic C++ implementation. CPPTraj allows parallel reads between frames of the same trajectory but is especially geared towards processing an ensemble of many trajectories in parallel.

pyPczip<sup>52</sup> is a suite of software tools written in Python for compression and analysis of MD simulation data, in particular ensembles of trajectories. pyPczip is MPI parallelized and is specific to PCA-based investigations of MD trajectories and supports a wide variety of trajectory file formats (based on the capabilities of the underlying MDTraj package<sup>20</sup>). pyPczip can take one or many input MD trajectory files and convert them into a highly compressed, HDF5-based pcz format with insignificant loss of information. However, the package does not support general purpose analysis.

*In situ* analysis is an approach to execute analysis simultaneously with the running MD simulation so that I/O bottlenecks are mitigated<sup>53,54</sup>. Malakar et al.<sup>53</sup> studied the scalability challenges of time and space shared modes of analyzing large-scale MD simulations through a topology-aware mapping for simulation and analysis using the LAMMPS code. Similarly, Taufer and colleagues<sup>54</sup> presented their own framework for *in situ* analysis, which is based on the fast on-the-fly calculation of metadata that characterizes protein substructures via maximum eigenvalues of distance matrices. These metadata are used to index trajectory frames and enable targeted analysis of trajectory subsets. Both studies provide important ideas and approaches towards moving towards online-analysis in conjunction with a running simulation but are limited in generality.

All of the above frameworks provide tools for parallel analysis of MD trajectories. These frameworks, however, tend to fall short in providing parallelism in the context of a general and flexible library for the analysis of MD trajectories. Although straggler tasks are a common challenge arising in parallel analysis and are well-known in the data analysis community (see Section 2.1), there is, to our knowledge, little discussion about this problem in the biomolecular simulation community. Our own experience with a MapReduce approach in *MDAnalysis*<sup>23</sup> suggested that stragglers

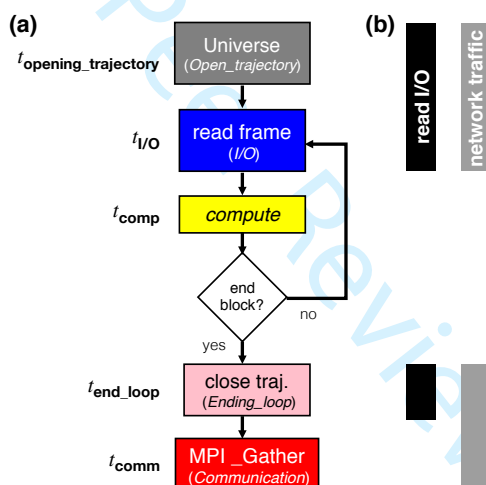
might be a somewhat under-appreciated problem. Therefore, in the present work we want to better understand requirements for efficient parallel analysis of MD trajectories in *MDAnalysis*, but to also provide more general guidance that could benefit developments in other libraries inside and outside of the scope of analysis of MD simulations.

### 3 | ALGORITHMS AND SOFTWARE PACKAGES

For our investigation of parallel trajectory analysis we focus on using MPI as the standard approach to parallelization in HPC. We employ the Python language, which is widely used in the scientific community because it facilitates rapid development of small scripts and code prototypes as well as development of large applications and highly portable and reusable modules and libraries. We use the *MDAnalysis* library to calculate a “RMSD time series” (explained in section 3.1) as a representative use case. Further details on the software packages are provided in sections 3.2–3.4.

#### 3.1 | RMSD Calculation with MDAnalysis

Simulation data exist in trajectories in the form of time series of atom positions and sometimes velocities. Trajectories come in a plethora of different and idiosyncratic file formats. *MDAnalysis*<sup>17,18</sup> is a widely used open source library to analyze trajectory files with an object oriented interface. The library is written in Python, with time critical code in C/C++/Cython. *MDAnalysis* supports most file formats of simulation packages including CHARMM<sup>55</sup>, Gromacs<sup>56</sup>, Amber<sup>57</sup>, and NAMD<sup>58</sup> and the Protein Data Bank<sup>59</sup> format. At its core, it reads trajectory data in different formats and makes them available through a uniform API; specifically, coordinates are represented as standard NumPy arrays<sup>60</sup>.



**FIGURE 1** Flow chart of the MPI-parallelized RMSD algorithm, Algorithm 1. (a) Each MPI process performs the same steps but reads trajectory frames from different blocks of the trajectory. The color scheme and labels in italics correspond to the colors and labels for measured timing quantities in the following graphs (e.g., Figs. 2c and 2d). The names of the corresponding timing quantities from Table 3 are listed next to each step. (b) Steps that access the shared Lustre file system with read I/O are included in the black bars; steps that communicate via the shared InfiniBand network are included in the gray bars. The Lustre file system is accessed through the network and hence all I/O steps also use the network.

As a test case that is representative of a common task in the analysis of biomolecular simulation trajectories we calculated the time series of the minimal structural root mean square distance (**RMSD**) after rigid body superposition<sup>7,27</sup>. The RMSD is used to show the rigidity of protein domains and more generally characterizes structural changes. It is calculated as a function of time  $t$  as

$$\text{RMSD}(t) = \min_{R, t} \sqrt{\frac{1}{N} \sum_{i=1}^N [(R \cdot \mathbf{x}_i(t) + \mathbf{t}) - \mathbf{x}_i^{\text{ref}}]^2} \quad (1)$$

where  $\mathbf{x}_i(t)$  is the position of atom  $i$  at time  $t$ ,  $\mathbf{x}_i^{\text{ref}}$  is its position in a reference structure and the distance between these two is minimized by finding the optimum  $3 \times 3$  rotation matrix  $R$  and translation vector  $\mathbf{t}$ . The optimum rigid body superposition was calculated with the QCPROT algorithm<sup>26,61</sup> (implemented in Cython and available through the *MDAnalysis.analysis.rms* module<sup>18</sup>).

The RMSD trajectory analysis was parallelized as outlined in the flow chart in Figure 1, with further details available in Algorithm 1. Each MPI process loads the core MDAnalysis structure (called the `Universe`), which includes loading a shared “topology” file with the simulation system information and opening the shared trajectory file. Each process operates on a different block of frames and iterates through them by reading the coordinates of a single frame into memory and performing the RMSD computation with them. Once all frames in the block are processed, the trajectory file is closed and results are communicated to MPI rank 0 using `MPI_Gather()`.

The RMSD was determined for a subset of protein atoms, the  $N = 214 C_{\alpha}$  atoms of our test system (see section 4.3 for details). The time complexity for the RMSD Algorithm 1 is  $\mathcal{O}(T \times N^2)^{26}$  where  $T$  is the number of frames in the trajectory and  $N$  the number of particles included in the RMSD calculation.

---

**Algorithm 1** MPI-parallel Multi-frame RMSD Algorithm
 

---

**Input:** size: Total number of frames  
 ref: mobile group in the initial frame which will be considered as reference  
 start & stop: Starting and stopping frame index  
 topology & trajectory: files to read the data structure from  
**Output:** Calculated RMSD arrays

```

1: procedure Block_RMSD(topology, trajectory, ref, start, stop)
2:   u ← Universe(topology, trajectory)                                ▷ u hold all the information of the physical system
3:   g ← u.frames[start:stop]
4:   for  $\forall$ frame in g do
5:     results[frame] ← RMSD(g, ref)
6:   end for
7:   return results
8: end procedure
9:
10: MPI Init
11: rank ← rank ID
12: index ← indices of mobile atom group
13: xref0 ← Reference atom group,s position
14: out ← Block_RMSD(topology, trajectory, xref0, start=start, stop=stop)
15:
16: Gather(out, RMSD_data, rank_ID=0)
17: MPI Finalize
  
```

---

### 3.2 | MPI for Python (*mpi4py*)

MPI for Python (*mpi4py*) is a Python wrapper for the Message Passing Interface (MPI) standard and allows any Python program to employ multiple processors<sup>29,30</sup>. Performance degradation due to using *mpi4py* is not prohibitive<sup>29,30</sup> and the overhead is far smaller than the overhead associated with the use of interpreted versus compiled languages<sup>33</sup>. Overheads in *mpi4py* are small compared to C code if efficient raw memory buffers are used for communication<sup>29</sup>, as used in the present study.

### 3.3 | Global Arrays Toolkit

The *Global Arrays* (GA) toolkit provides users with a language interface that allows them to distribute data while maintaining the type of global index space and programming syntax similar to what is available when programming on a single processor<sup>34</sup>. *Global Arrays* is implemented with Fortran-77 and C bindings and provides C++ and Python interfaces. It allows manipulating physically distributed dense multi-dimensional arrays without explicitly defining communication and synchronization between processes. The underlying communication is determined by a runtime environment, which defaults to the *Communication runtime for Extreme Scale* (ComEx)<sup>62</sup>. ComEx uses shared memory for intra-node communication and inter-node communication employs ComEx with MPI. *Global Arrays in NumPy* (GAI<sub>N</sub>) extends GA to Python through Numpy<sup>33</sup>. The *Global Arrays* toolkit provides functions to create global arrays (`ga_create()`) and to copy data to (`ga_put()`) and from (`ga_get()`) such a global array, as well as additional functions for copying between arrays and freeing them<sup>33</sup>. When a global array is created (`ga_create()`) each process will create an array of the same shape and size, physically located in the local memory space of that process<sup>34</sup>. The GA library maintains a list of all these memory locations, which can be queried with the `ga_access()` function. Using a pointer returned by `ga_access()`, one can directly modify the data that is local to each process. When a process tries to access a block of data the request is first decomposed into individual blocks representing the contribution to the total request from the data held locally on each process (B. J. Palmer and J. Daily, personal communication). The requesting process then makes individual requests to each of the remote processes.

GA allows independent, asynchronous, and efficient access to logical blocks of physically distributed arrays, with no need for explicit cooperation by other processes; in particular, it allows data locality to be explicitly specified and used<sup>63</sup>. We investigated if communication cost could be reduced by using *Global Arrays*. Algorithm 2 describes the RMSD algorithm with *Global Arrays* instead of MPI.

---

**Algorithm 2** MPI-parallel Multi-frame RMSD using Global Arrays
 

---

**Input:** size: Total number of frames assigned to each rank  $N_b$   
**g\_a:** Initialized Global Arrays  
**xref0:** mobile group in the initial frame which will be considered as reference  
**start & stop:** that tell which block of trajectory (frames) is assigned to each rank  
**topology & trajectory:** files to read the data structure from  
**Include:** Block\_RMSD() from Algorithm 1

```

1: bsize  $\leftarrow$  ceil(trajectory.number_frames / size)
2: g_a  $\leftarrow$  ga.create(ga.C_DBL, [bsize*size,2], "RMSD")
3: buf  $\leftarrow$  np.zeros([bsize*size,2], dtype=float)
4: out  $\leftarrow$  Block_RMSD(topology, trajectory, xref0, start=start, stop=stop)
5: ga.put(g_a, out, (start,0), (stop,2))
6: if rank == 0 then
7:   buf  $\leftarrow$  ga.get(g_a, lo=None, hi=None)
8: end if

```

---

### 3.4 | MPI and Parallel HDF5

HDF5 is a structured self-describing hierarchical data format which is the standard mechanism for storing large quantities of numerical data in Python (<http://www.hdfgroup.org/HDF5>,<sup>64</sup>). Parallel HDF5 (*PHDF5*) typically sits on top of a MPI-IO layer and can use MPI-IO optimizations. In *PHDF5*, all file access is coordinated by the MPI library; otherwise, multiple processes would compete over accessing the same file on disk. MPI-based applications launch multiple parallel instances of the Python interpreter that communicate with each other via the MPI library. Implementation is straightforward as long as the user supplies a MPI communicator and takes into account some constraints required for data consistency<sup>64</sup>. *HDF5* itself handles nearly all the details involved with coordinating file access when the shared file is opened through the *mpio* driver.

MPI has two flavors of operation: collective (all processes have to participate in the same order) and independent (processes can perform the operation in any order or not at all)<sup>64</sup>. With *PHDF5*, modifications to file metadata must be performed collectively and although all processes perform the same task, they do not need to be synchronized<sup>64</sup>. Other tasks and any type of data operations can be performed independently by processes. In the present study, we use independent operations.

## 4 | BENCHMARK ENVIRONMENT

Our benchmark environment consisted of three different XSEDE<sup>32</sup> HPC resources (described in section 4.1), the software stack used (section 4.2), which had to be compiled for each resource, and the common test data set (section 4.3).

### 4.1 | HPC Resources

The computational experiments were executed on standard compute nodes of three XSEDE<sup>32</sup> supercomputers, *SDSC Comet*, *PSC Bridges*, and *LSU SuperMIC* (Table 1). *SDSC Comet* is a 2 PFlop/s cluster with 2,020 compute nodes in total. It is optimized for running a large number of medium-size calculations (up to 1,024 cores) to support the most prevalent type of calculation on XSEDE resources. *PSC Bridges* is a 1.35 PFlop/s cluster with different types of computational nodes, including 16 GPU nodes, 8 large memory and 2 extreme memory nodes, and 752 regular nodes. It was designed to flexibly support both traditional (medium scale calculations) and non-traditional (data analytics) HPC uses. *LSU SuperMIC* offers 360 standard compute nodes with a peak performance of 557 TFlop/s. The parallel file system on all three machines is Lustre (<http://lustre.org/>) and is shared between the nodes of each cluster.



Name	Nodes	Number of Nodes	CPUs	RAM	Network Topology	Scheduler and Resource Manager	parallel file system
<i>SDSC Comet</i>	Compute	6400	2 Intel Xeon (E5-2680v3) 12 cores/CPU, 2.5 GHz	128 GB DDR4 DRAM	56 Gbps IB	SLURM	Lustre
<i>PSC Bridges</i>	RSM	752	2 Intel Haswell (E5-2695 v3) 14 cores/CPU, 2.3 GHz	128 GB, DDR4-2133MHz	12.37 Gbps OPA	SLURM	Lustre
<i>LSU SuperMIC</i>	Standard	360	2 Intel Ivy Bridge (E5-2680) 10 cores/CPU, 2.8 GHz	64 GB, DDR3-1866MHz	56 Gbps IB	PBS	Lustre

**TABLE 1** Configuration of the HPC resources that were benchmarked. Only a subset of the total available nodes were used. IB: InfiniBand; OPA: Omni-Path Architecture.

4.2 | Software

Table 2 lists the tools and libraries that were required for our computational experiments. Many domain specific packages are not available in the standard software installation on supercomputers. We therefore had to compile them, which in some cases required substantial effort due to non-standard building and installation procedures or lack of good documentation. Because this is a common problem that hinders reproducibility we provide detailed version information, notes on the installation process, as well as comments on the ease of installation and the quality of the documentation in Table 2. For the MPI implementation we used Open MPI release 1.10.7 (<https://www.open-mpi.org/>) consistently everywhere. Detailed instructions to create the computing environments together with the benchmarking code can be found in the GitHub repository. Carefully setting up the same software stack on the three different supercomputers allowed us to clearly demonstrate the reproducibility of our results and showed that our findings were not dependent on machine specifics.

4.3 | Data Set

The test system contained the protein adenylate kinase with 214 amino acid residues and 3341 atoms in total<sup>65</sup> and the topology information (atoms types and bonds) was stored in a file in CHARMM PSF format. The test trajectory was created by concatenating 600 copies of a MD trajectory with 4,187 time frames (saved every 240 ps for a total simulated time of 1.004  $\mu$ s) in CHARMM DCD format<sup>66</sup> and converting to Gromacs XTC format trajectory, as described for the “600x” trajectory in Khoshlessan et al.<sup>23</sup>. The trajectory had a file size of about 30 GB and contained 2,512,200 frames (corresponding to 602.4  $\mu$ s simulated time). The file size was relatively small because water molecules that were also part of the original MD simulations were stripped to reduce the original file size by a factor of about 10; such preprocessing is a common approach if one is only interested in the protein behavior. Thus, the trajectory represents a small to medium system size in the number of atoms and coordinates that have to be loaded into memory for each time frame. The XTC format is a format with lossy compression<sup>67,68</sup>, which also contributed to the compact file size. XTC trades lower I/O demands for higher CPU demands during decompression and therefore performed well in our previous study<sup>23</sup>. Although 2,512,200 frames represents a long simulation for current standards, such trajectories will become increasingly common due to the use of special hardware<sup>69,70</sup> and GPU-acceleration<sup>56,71,72</sup>.

5 | METHODS

Documentation and benchmark codes are made available in the code repository <https://github.com/hpcanalytics/supplement-hpc-py-parallel-mdanalysis> under the GNU General Public License v3.0 (code) and the Creative Commons Attribution-ShareAlike (documentation) and are archived under DOI 10.5281/zenodo.3351616. These materials should enable users to recreate the computational environment on the tested XSEDE HPC resources (*SDSC Comet*, *PSC Bridges*, *LSU SuperMIC*), prepare data files, and run the computational experiments.

In the following we define the quantities and approach used for our performance measurements, with a full summary of all definitions in Table 3. We evaluated MPI performance of the parallel RMSD time series algorithm 1 by timing the total time to solution as well as the execution time for different parts of the code for individual MPI ranks with the help of the Python `time.time()` function.



Package	Version	Description	Ease of Installation	Documentation	Installation	Dependencies
GCC	4.9.4	GNU Compiler Collection	0	++	via configuration files, environment or command line options, minimal configuration is required	-
Open MPI	1.10.7	MPI Implementation	0	++	via configuration files, environment or command line options, minimal configuration is required	-
Global Arrays	5.6.1	Global Arrays	-	+	via configuration files, environment or command line options, several optional configuration settings available	MAMA, ARMC I MPI 1.x/2.x/3.x implementation like Open MPI built with shared/dynamic libraries, GCC
Python	2.7.13	Python language	+	++	Conda Installation	-
MPI4py	3.0.0	MPI for Python	+	++	Conda Installation	Python 2.7 or above, MPI 1.x/2.x/3.x implementation like Open MPI built with shared/dynamic libraries, Cython
GA4py	1.0	Global Arrays for Python	0	0	Python setuptools	Global Arrays, Python 2 only, MPI 1.x/2.x/3.x implementation like Open MPI built with shared/dynamic libraries, Cython, MPI4py, Numpy
PHDF5	1.10.1	Parallel HDF5	-	++	via configuration files, environment or command line options, several optional configuration settings available	MPI 1.x/2.x/3.x implementation like Open MPI GNU, MPFI90, MPICXX
H5py	2.7.1	Pythonic wrapper around the HDF5	+	++	Conda Installation	Python 2.7, or above, PHDF5, Cython
MDAnalysis	0.17.0	Python library to analyze trajectories from MD simulations	+	++	Conda Installation	Python >=2.7, Cython, GNU, Numpy

**TABLE 2** Detailed comparison on the dependencies and installation of different software packages used in the present study. Software was built from source or obtained via a package manager and installed on the multi-user HPC systems in Table 1. Evaluation of ease of installation and documentation uses a subjective scale with “++” (excellent), “+” (good), “0” (average), and “-” (difficult/lacking) and reflects the experience of a typical domain scientist at the graduate/post-graduate level in a discipline such as computational biophysics or chemistry.

## 5.1 | Timing Observables

We abbreviate the timings in the following as variables  $t_{L_n}$  where  $L_n$  refers to the line number in algorithm 1. We measured in the function `block_rmsd()` the *read I/O time* for ingesting the data of one trajectory frame from the file system into memory,  $t_{I/O}^{\text{frame}} = t_{L4}$ , and the *compute time* per trajectory frame to perform the computation,  $t_{\text{comp}}^{\text{frame}} = t_{L5}$ . The *total read I/O time for a MPI rank*,  $t_{I/O} = \sum_{\text{frame}=1}^{N_b} t_{I/O}^{\text{frame}}$ , is the sum over all I/O times for all the  $N_{\text{frames}}$  frames assigned to the rank; similarly, the *total compute time for a MPI rank* is  $t_{\text{comp}} = \sum_{\text{frame}=1}^{N_b} t_{\text{comp}}^{\text{frame}}$ . The time delay between the end of the last iteration and exiting the `for` loop is  $t_{\text{end\_loop}} = t_{L6}$ . The time  $t_{\text{opening\_trajectory}} = t_{L2} + t_{L3}$  measures the problem setup, which includes data structure initialization and opening of topology and trajectory files. The *communication time*,  $t_{\text{comm}} = t_{L16}$ , is the time to gather

Quantity	Definition
$N_b$	$N_{frames}^{total}/N$
$t_{end\_loop}$	$t_{L6}$
$t_{opening\_trajectory}$	$t_{L2} + t_{L3}$
$t_{comp}$	$\sum_{frame=1}^{N_b} t_{comp}^{frame}$
$t_{I/O}$	$\sum_{frame=1}^{N_b} t_{I/O}^{frame}$
$t_{all\_frame}$	$t_{L4} + t_{L5} + t_{L6}$
$t_{RMSD}$	$t_{L1} + \dots + t_{L8}$
$t_{comm/MPI}$	$t_{L16}$
$t_{comm/GA}$	$t_{L5} + t_{L6} + t_{L7} + t_{L8}$
$t_{comm}$	$t_{comm/MPI}$ (Alg. 1) or $t_{comm/GA}$ (Alg. 2)
$t_{Overhead1}$	$t_{all\_frame} - t_{I/O} - t_{comp} - t_{end\_loop}$
$t_{Overhead2}$	$t_{RMSD} - t_{all\_frame} - t_{opening\_trajectory}$
$t_N$	$t_{RMSD} + t_{comm}$
$\overline{t_{comp}}$	$\frac{1}{N} \sum_{rank=1}^N t_{comp}$
$\overline{t_{I/O}}$	$\frac{1}{N} \sum_{rank=1}^N t_{I/O}$
$\overline{t_{comm}}$	$\frac{1}{N} \sum_{rank=1}^N t_{comm}$
$t_{total}$	$\max t_N$

**TABLE 3** Summary of measured timing quantities. Timings are collected for the specified line numbers in the code, labeled as  $t_{L_n}$  where  $L_n$  refers to the line number in the corresponding algorithm.  $t_{comm/MPI}$  (in Algorithm 1) and  $t_{comm/GA}$  (in Algorithm 2) are both referred to as  $t_{comm}$  in the text. Variables in the top half of the table refer to measurements of an individual MPI rank. Variables in the bottom half are aggregates such as averages over all ranks or the total time to solution.

all data from all processor ranks to rank zero. The total time (for all frames) spent in `block_rmsd()` is  $t_{RMSD} = \sum_{i=1}^8 t_{L_i}$ . There are parts of the code in `block_rmsd()` that are not covered by the detailed timing information of  $t_{comp}$  and  $t_{I/O}$ . Unaccounted time is considered as *overhead*. We define  $t_{Overhead1}$  and  $t_{Overhead2}$  as the overheads of the calculations (see Table 3 for the definitions); both are expected to be negligible, which was the case in all our measurements. Finally, the *total time to completion of a single MPI rank*, when utilizing  $N$  cores for the execution of the overall experiment, is  $t_N$ , and as a result  $t_{RMSD} + t_{comm} \equiv t_N$ .

## 5.2 | Performance Parameters

We measured the *total time to solution*  $t_{total}(N)$  with  $N$  MPI processes on  $N$  cores, which is effectively  $t_{total}(N) \approx \max(t_N)$ . Strong scaling was quantified by the speed-up

$$S(N) = \frac{t_{total}(1)}{t_{total}(N)}, \quad (2)$$

relative to performance on a single core ( $t_{total}(1)$ ), and the efficiency

$$E(N) = \frac{S(N)}{N}. \quad (3)$$

Averages over ranks were calculated as

$$\overline{t_{comp}} = \frac{1}{N} \sum_{rank=1}^N t_{comp} = \frac{1}{N} \sum_{rank=1}^N \sum_{frame=1}^{N_b} t_{comp}^{frame}, \quad (4)$$

$$\overline{t_{I/O}} = \frac{1}{N} \sum_{rank=1}^N t_{I/O} = \frac{1}{N} \sum_{rank=1}^N \sum_{frame=1}^{N_b} t_{I/O}^{frame}, \quad (5)$$

and

$$\overline{t_{comm}} = \frac{1}{N} \sum_{rank=1}^N t_{comm}. \quad (6)$$

Additionally, we introduced two performance parameters that we found to be indicative of the occurrence of stragglers. We defined the ratio of compute time to read I/O time for the serial code as

$$R_{comp/I/O} = \frac{t_{comp}}{t_{I/O}} = \frac{t_{comp}/N_{frames}^{total}}{t_{I/O}/N_{frames}^{total}} = \frac{\overline{t_{comp}^{frame}}}{\overline{t_{I/O}^{frame}}} \quad (7)$$

where the last equality shows that the ratio can also be computed from the average times per frame,  $\overline{t_{\text{comp}}^{\text{frame}}}$  and  $\overline{t_{\text{I/O}}^{\text{frame}}}$ .  $R_{\text{comp}/\text{IO}}$  was calculated with the serial versions of our algorithms (on a single CPU core) in order to characterize the computational problem in the absence of parallelization. The ratio of compute to communication time was defined by the ratio of average total compute time to the average total communication time

$$R_{\text{comp}/\text{comm}} = \frac{\overline{t_{\text{comp}}}}{\overline{t_{\text{comm}}}}. \quad (8)$$

Because  $t_{\text{comm}}$  cannot be measured for a serial code, we estimated  $R_{\text{comp}/\text{comm}}$  from the rank-averages (Eqs. 4 and 6) for a given number of MPI ranks.

## 6 | COMPUTATIONAL EXPERIMENTS

We had previously measured the performance of the MPI-parallelized RMSD analysis task on two different HPC resources (*SDSC Comet* and *TACC Stampede*) and had found that it only scaled well up to a single node due to high variance in the runtime of the MPI ranks, similar to the straggler phenomenon observed in big-data analytics<sup>23</sup>. However, the ultimate cause for this high variance could not be ascertained. We therefore performed more measurements with more detailed timing information (see section 5) on *SDSC Comet* (described in this section) and two other supercomputers (summarized in section 7) in order to better understand the origin of the stragglers and find solutions to overcome them.

### 6.1 | RMSD Benchmark

We measured strong scaling for the RMSD analysis task (Algorithm 1) with the 2,512,200 frame test trajectory (section 4.3) on 1 to 72 cores (one to three nodes) of *SDSC Comet* (Figures 2a and 2b). We observed poor strong scaling performance beyond a single node (24 cores), comparable to our previous results<sup>23</sup>. A more detailed analysis showed that the RMSD computation, and to a lesser degree the read I/O, considered on their own, scaled well beyond 50 cores (yellow and blue lines in Figure 2c). But communication (sending results back to MPI rank 0 with `MPI_Gather()`; red line in Figure 2c) and the initial file opening (loading the system information into the `MDAnalysis.Universe` data structure from a shared “topology” file and opening the shared trajectory file; gray line in Figure 2c) started to dominate beyond 50 cores. Communication cost and initial time for opening the trajectory were distributed unevenly across MPI ranks, as shown in Figure 2d. The ranks that took much longer to complete than the typical execution time of the other ranks were the stragglers that hurt performance.

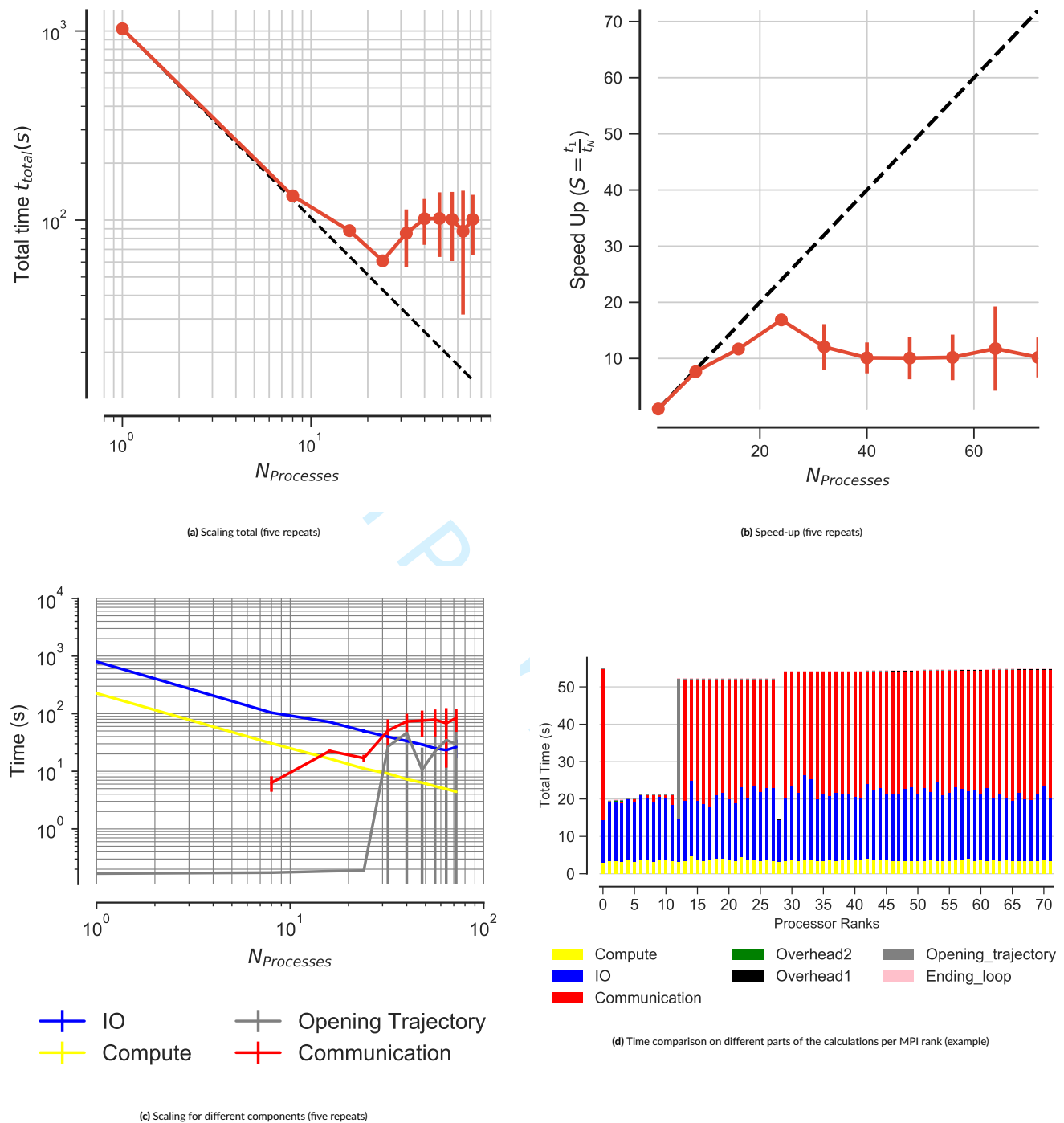
We qualitatively denoted by *straggler* any MPI rank that took at least about twice as long as the group of ranks that finished fastest, roughly following the original description of a straggler as a task that took an “unusually long time to complete”<sup>35</sup>. The fast-finishing ranks were generally clearly distinguishable in the per-rank timings such as in Figures 2d and A1d. Such a qualitative definition of stragglers was sufficient for our purpose to identify scalability bottlenecks, as shown in the following discussion.

### Identification of Scalability Bottlenecks

In the example shown in Figure 2d, 62 ranks out of 72 took about 60 s (the stragglers) whereas the remaining ranks only took about 20 s. In other instances, far fewer ranks were stragglers, as shown, for example, in Figure A1d. The detailed breakdown of the time spent on each rank (Figure 2d) showed that the computation,  $t_{\text{comp}}$ , was relatively constant across ranks. The time spent on reading data from the shared trajectory file on the Lustre file system into memory,  $t_{\text{I/O}}$ , showed variability across different ranks. The stragglers, however, appeared to be defined by occasionally much larger *communication* times,  $t_{\text{comm}}$  (line 16 in Algorithm 1), which were on the order of 30 s, and by larger times to initially open the trajectory (line 2 in Algorithm 1).  $t_{\text{comm}}$  varied across different ranks and was barely measurable for a few of them. Although the data in Figure 2d represented one run and in other instances different number of ranks were stragglers, the averages over all ranks in five independent repeats (Figure 2c) showed that increased  $t_{\text{comm}}$  were generally the reason for large variations in the run time for each rank. This initial analysis indicated that communication was a major issue that prevented good scaling beyond a single node but the problems related to file I/O also played an important role in limiting scaling performance.

### Influence of Hardware

We ran the same benchmarks on multiple HPC systems that were equipped with a Lustre parallel file system [XSEDE’s *PSC Bridges* (Fig. A1) and *LSU SuperMIC* (Fig. A2)], and observed the occurrence of stragglers, in a manner very similar to the results described for *SDSC Comet*. There was no clear pattern in which certain MPI ranks would always be a straggler, and neither could we trace stragglers to specific cores or nodes. Therefore, the phenomenon of stragglers in the RMSD case was reproducible on different clusters and thus appeared to be independent from the underlying hardware.



**FIGURE 2** Performance of the RMSD task parallelized with MPI on *SDCS Comet*. Results were communicated back to rank 0. Five independent repeats were performed to collect statistics. (a-c) The error bars show standard deviation with respect to the mean. In serial, there is no communication and no data points are shown for  $N = 1$  in (c). (d) Compute  $t_{\text{comp}}$ , read I/O  $t_{\text{I/O}}$ , communication  $t_{\text{comm}}$ , ending the for loop  $t_{\text{end\_loop}}$ , opening the trajectory  $t_{\text{opening\_trajectory}}$ , and overheads  $t_{\text{overhead1}}$ ,  $t_{\text{overhead2}}$  per MPI rank; see Table 3 for definitions. These are data from one run of the five repeats. MPI ranks 0, 12–27 and 29–72 are stragglers.

## 6.2 | Effect of Compute to I/O Ratio on Performance

The results in section 6.1 indicated opening the trajectory, communication, and read I/O to be important factors that appeared to correlate with stragglers. In order to better characterize the RMSD task, we computed the ratio between the time to complete the computation and the time spent on I/O per frame. The average values were  $\overline{t_{\text{comp}}^{\text{frame}}} = 0.09$  ms,  $\overline{t_{\text{I/O}}^{\text{frame}}} = 0.3$  ms, resulting in a compute-to-I/O ratio  $R_{\text{comp/I/O}} \approx 0.3$  (Eq. 7). Because  $R_{\text{comp/I/O}} \ll 1$ , the RMSD analysis task was characterized as I/O bound.

As we were not able to achieve good scaling beyond a single node, we hypothesized that decreasing the I/O load relative to the compute load would interleave read I/O with longer periods of computation, thus reducing the impact of I/O contention and the impact of stragglers. We therefore set out to measure compute bound tasks, i.e. ones with  $R_{\text{comp/I/O}} \gg 1$ . To measure the effect of the  $R_{\text{comp/I/O}}$  ratio on performance but leaving other parameters the same, we artificially increased the computational load by repeating the same RMSD calculation (line 10, algorithm 1) 40, 70 and 100 times in a loop, resulting in forty-fold ("40×"), seventy-fold ("70×"), and one hundred-fold ("100×") load increases.

### 6.2.1 | Effect of Increased Compute Workload

For an X-fold increase in workload, we expected the workload for the computation to scale with X as  $t_{\text{comp}}(X) = N_{\text{frames}}^{\text{total}} X \overline{t_{\text{comp}}^{\text{frame}}}$  while the read I/O workload  $t_{\text{I/O}}(X) = N_{\text{frames}}^{\text{total}} \overline{t_{\text{I/O}}^{\text{frame}}}$  (number of frames times the average time to read a frame) should remain independent of X. Therefore, the ratio for any X should be  $R_{\text{comp/I/O}}(X) = t_{\text{comp}}(X)/t_{\text{I/O}}(X) = X R_{\text{comp/I/O}}(X = 1)$ , i.e.,  $R_{\text{comp/I/O}}$  should just linearly scale with the workload factor X. The measured  $R_{\text{comp/I/O}}$  ratios of 11, 19, 27 for the increased computational workloads agreed with this theoretical analysis, as shown in Table 4.

Workload X	$t_{\text{comp}}$ (s)	$t_{\text{I/O}}$ (s)	$R_{\text{comp/I/O}}$	
			measured	theoretical
1×	226	791	0.29	
40×	8655	791	11	11
70×	15148	791	19	20
100×	21639	791	27	29

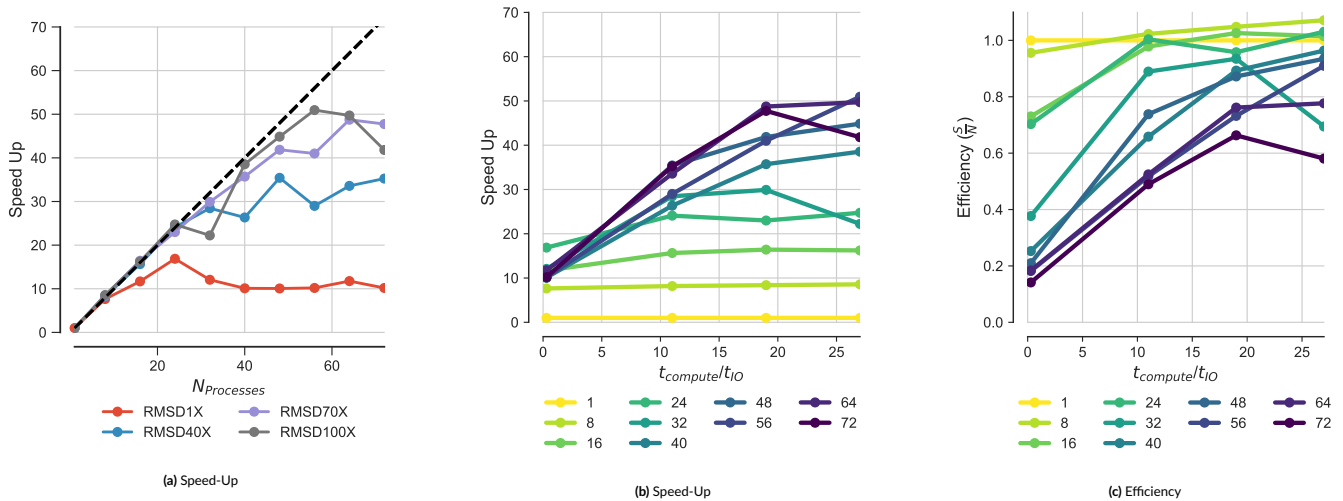
**TABLE 4** Change in  $R_{\text{comp/I/O}}$  ratio with change in the RMSD workload X. The RMSD workload was artificially increased in order to examine the effect of compute to I/O ratio on the performance. The reported compute and I/O time were measured based on the serial version using one core. The theoretical  $R_{\text{comp/I/O}}$  (see text) is provided for comparison.

We performed the experiments with increased workload to measure the effect of the  $R_{\text{comp/I/O}}$  ratio on performance (Figure 3). The strong scaling performance as measured by the speed-up  $S(N)$  improved with increasing  $R_{\text{comp/I/O}}$  ratio (Figure 3a). The calculations consistently showed better scaling performance to larger numbers of cores for higher  $R_{\text{comp/I/O}}$  ratios, e.g.,  $N = 56$  cores for the 70× RMSD task. The speed-up and efficiency approached their ideal value for each processor count with increasing  $R_{\text{comp/I/O}}$  ratio (Figures 3b and 3c). Even for moderately compute-bound workloads, such as the 40× and 70× RMSD tasks, increasing the computational workload over I/O reduced the impact of stragglers even though they still contributed to large variations in timing across different ranks and thus irregular scaling.

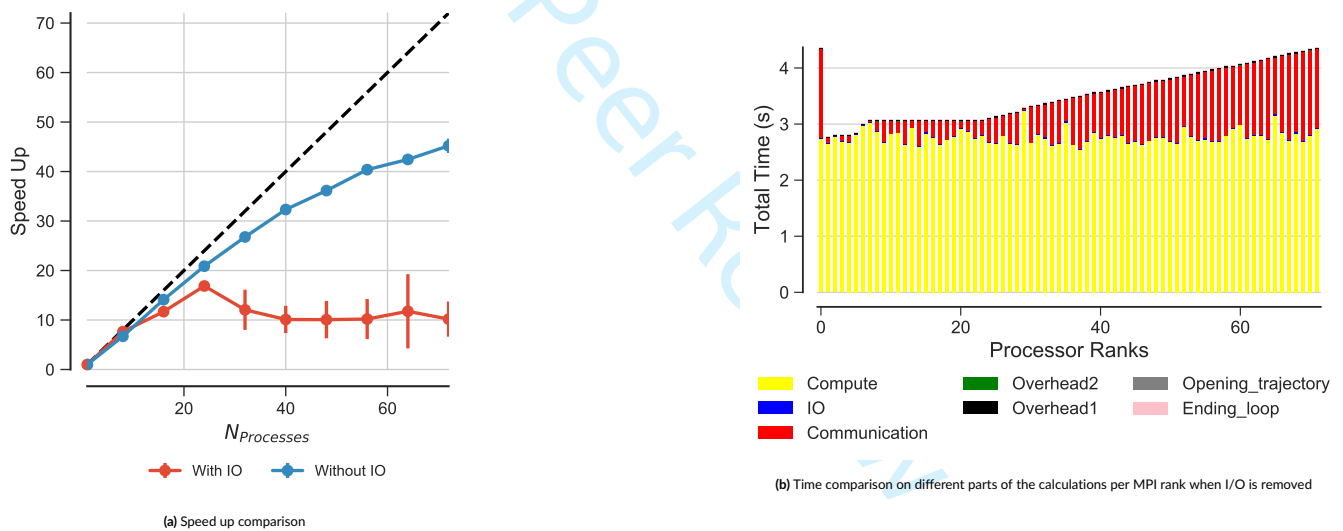
We also investigated the influence of the ratio of compute to communication costs ( $R_{\text{comp/comm}}$ , Eq. 8) on performance in B. We found evidence to support the hypothesis that a larger ratio was beneficial, provided I/O costs could also be reduced. However, read I/O ultimately seemed to be the key determinant for performance, as discussed in the next sections.

### 6.2.2 | Effect of Absence of Read I/O on Communication

In order to study an extreme case of a compute-bound task, we eliminated all I/O from the RMSD task by randomly generating artificial trajectory data in memory; the data had the same size as if they had been obtained from the trajectory file. The time for the data generation was excluded and no file access was necessary. Without any I/O, performance improved markedly (Figure 4), with reasonable scaling up to 72 cores (3 nodes). No stragglers were observed because overall communication time decreased and showed less variability; nevertheless, an increase in communication time prevented ideal scaling performance. Although in practice I/O cannot be avoided, this experiment demonstrated that accessing the trajectory file on the Lustre file system is at least one cause for the observed stragglers.



**FIGURE 3** Effect of  $R_{comp/IO}$  ratio on performance of the RMSD task on *SDSC Comet*. We tested performance for  $R_{comp/IO}$  ratios of 0.3, 11, 19, 27, which correspond to  $1\times$  RMSD,  $40\times$  RMSD,  $70\times$  RMSD, and  $100\times$  RMSD respectively. (a) Effect of  $R_{comp/IO}$  on the speed-up. (b) Change in speed-up with respect to  $R_{comp/IO}$  for different processor counts. (c) Change in the efficiency with respect to  $R_{comp/IO}$  for different processor counts.



**FIGURE 4** Comparison of the performance of the RMSD task with I/O ( $R_{comp/IO} \approx 0.3$ ) and without I/O ( $R_{comp/IO} = +\infty$ ) on *SDSC Comet*. Five repeats were performed to collect statistics. (a) Speed-up. The error bars show standard deviation with respect to the mean. (b) Compute  $t_{comp}$ , read I/O  $t_{I/O} = 0$ , communication  $t_{comm}$ , ending the for loop  $t_{end\_loop}$ , opening the trajectory  $t_{opening\_trajectory}$ , and overheads  $t_{overhead1}$ ,  $t_{overhead2}$  per MPI rank. (See Table 3 for definitions.)

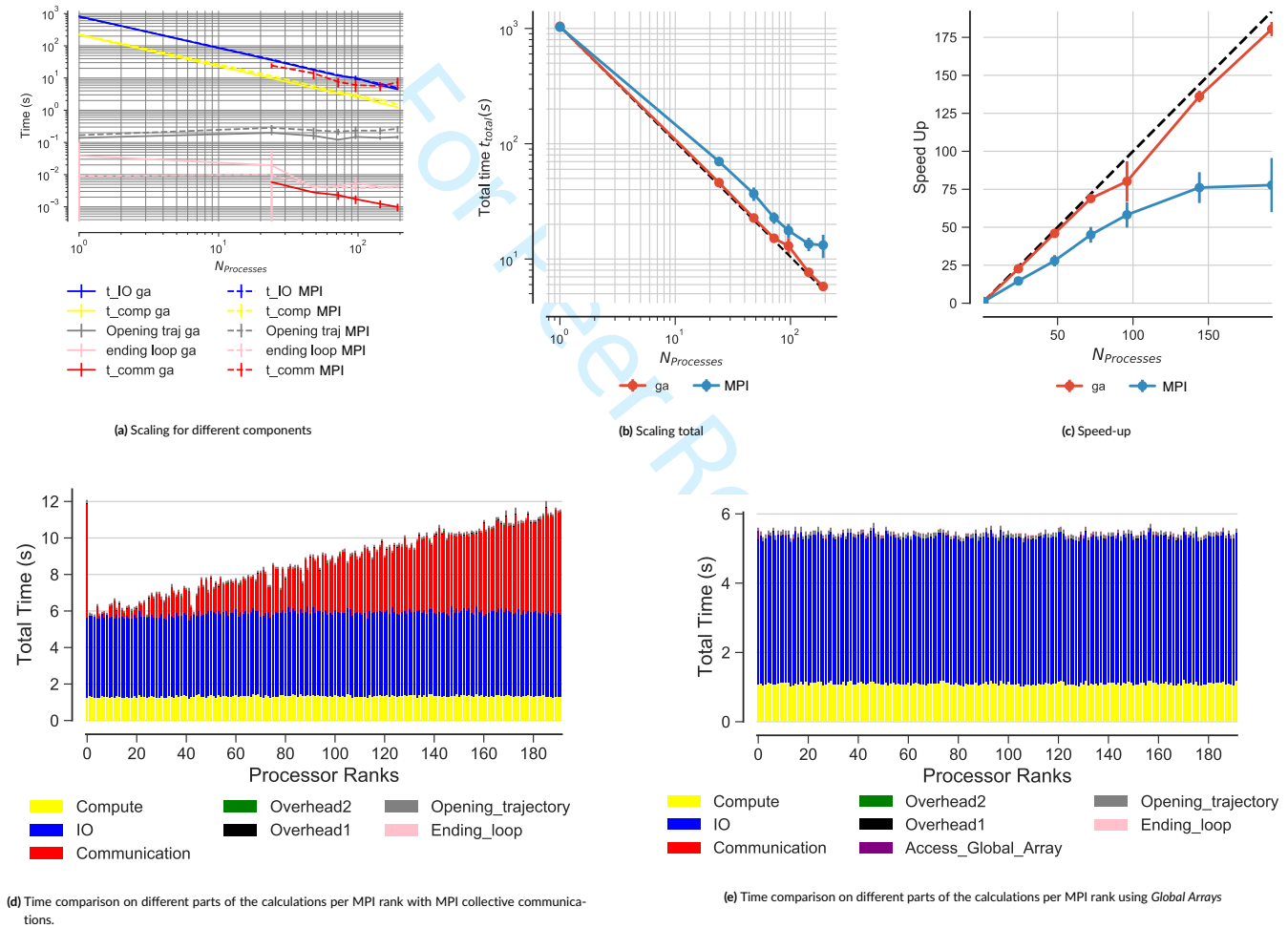
### 6.3 | Reducing I/O Cost

In order to improve performance we needed to employ strategies to avoid the competition over file access across different ranks when the  $R_{comp/IO}$  ratio was small. To this end, we experimented with two different ways for reducing the I/O cost: 1) splitting the trajectory file into as many segments as the number of processes, thus using file-per-process access, and 2) using the HDF5 file format together with MPI-IO parallel reads instead of the XTC trajectory format. We discuss these two approaches and their performance improvements in detail in the following sections.

### 6.3.1 | Splitting the Trajectories (“subfiling”)

Subfiling is a mechanism previously used for splitting a large multi-dimensional global array to a number of smaller subarrays in which each smaller array is saved in a separate file. Subfiling reduces the file system control overhead by decreasing the number of processes concurrently accessing a shared file<sup>73,74</sup>. Because subfiling is known to improve programming flexibility and performance of parallel shared-file I/O, we investigated splitting our trajectory file into as many trajectory segments as the number of processes. The trajectory file was split into  $N$  segments, one for each process, with each segment having  $N_b$  frames. This way, each process would access its own trajectory segment file without competing for file accesses.

We ran a benchmark up to 8 nodes (192 cores) and observed rather better scaling behavior with efficiencies above 0.6 (Figure 5b and 5c) with the delay time for stragglers reduced from 65 s to about 10 s for 72 processes. However, scaling was still far from ideal due to the MPI communication costs. Although the delay due to communication was much smaller than compared to parallel RMSD with shared-file I/O [compare Figure 5d ( $t_{\text{comm}}^{\text{Straggler}} \gg t_{\text{comp}} + t_{\text{I/O}}$ ) to Figure 2d ( $t_{\text{comm}}^{\text{Straggler}} \approx t_{\text{comp}} + t_{\text{I/O}}$ )], it was still delaying several processes and resulted in longer job completion times (Figure 5d). These delayed tasks impacted performance so that speed-up remained far from ideal (Figure 5c).



**FIGURE 5** Comparison of the performance of the RMSD task on SDSC Comet when the trajectories are split (*subfiling*). The communication step used either MPI collective communications (“MPI”, with `MPI_Gather()`) or *Global Arrays* (“ga”, as described in Section 6.4). In the case of *Global Arrays*, all ranks updated the global RMSD array (`ga_put()`) and rank 0 accessed the whole RMSD array through the global memory address (`ga_get()`). Five repeats were performed to collect statistics. (a) Compute and I/O scaling versus number of processes. In serial, there is no communication and no data points are shown for  $N = 1$ . (b) Total time scaling versus number of processes. (c) Speed-up. (a-c) The error bars show standard deviation with respect to the mean. (d-e) Compute  $t_{\text{comp}}$ , read I/O  $t_{\text{I/O}}$ , communication  $t_{\text{comm}}$ , access to the whole global RMSD array by rank 0,  $t_{\text{Access\_Global\_Array}}$ , ending the for loop  $t_{\text{end\_loop}}$ , opening the trajectory  $t_{\text{opening\_trajectory}}$ , and overheads  $t_{\text{overhead1}}$ ,  $t_{\text{overhead2}}$  per MPI rank; see Table 3 for the definitions.



The subfiling approach appeared promising but it required preprocessing of trajectory files and additional storage space for the segments. We benchmarked the necessary time for splitting the trajectory in parallel using different number of MPI processes (Table 5); in general the operation scaled well, with efficiencies  $> 0.8$  although performance fluctuated, as seen for the case on six nodes where the efficiency dropped to 0.34 for the run. These preprocessing times were not included in the estimates because we focused on better understanding the principal causes of stragglers and we wanted to make the results directly comparable to the results of the previous sections. Nevertheless, from an end user perspective, preprocessing of trajectories can be integrated in workflows (especially as the data in Table 5 indicated good scaling) and the preprocessing time can be quickly amortized if the trajectories are analyzed repeatedly.

$N_{\text{seg}}$	$N_p$	nodes	time (s)	S	E
24	24	1	89.9	1.0	1.0
48	48	2	46.8	1.9	0.96
72	72	3	33.7	2.7	0.89
96	96	4	25.1	3.6	0.89
144	144	6	43.7	2.1	0.34
192	192	8	13.5	6.7	0.83

**TABLE 5** The wall-clock time spent for writing  $N_{\text{seg}}$  trajectory segments using  $N_p$  processes using MPI on SDSC Comet. One set of runs was performed for the timings. Scaling S and efficiency E are relative to the 1 node case (24 MPI processes).

Often trajectories from MD simulations on HPC machines are produced and kept in smaller files or segments that can be concatenated to form a full continuous trajectory file. These trajectory segments could be used for the subfiling approach. However, it might not be feasible to have exactly one segment per MPI rank, with all segments of equal size, which constitutes the ideal case for subfiling. MDAnalysis can create virtual trajectories from separate trajectory segment files that appear to the user as a single trajectory. In C we investigated if this so-called *ChainReader* functionality could benefit from the subfiling approach. We found some improvements in performance but discovered limitations in the design of the ChainReader (namely that all segment files are initially opened) that will have to be addressed before equivalent performance can be reached.

6.3.2 | MPI-based Parallel HDF5

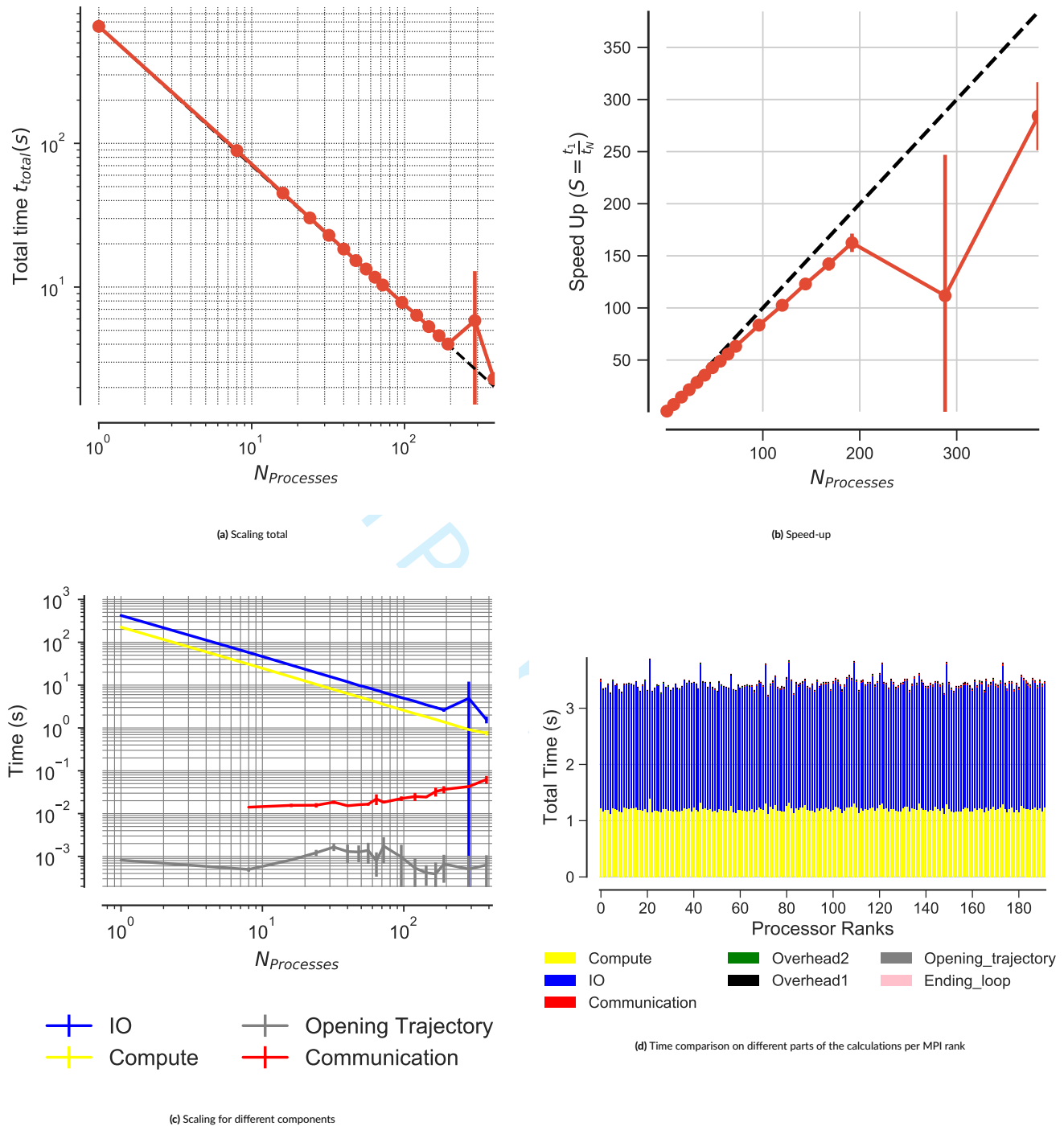
In the HPC community, parallel I/O with MPI-IO is widely used in order to address I/O limitations. We investigated MPI-based Parallel HDF5 to improve I/O scaling. We converted our XTC trajectory file into a simple custom HDF5 format so that we could test the performance of parallel I/O with the HDF5 file format. The code for this file format conversion can be found in the GitHub repository. The time it took to convert our XTC file with 2, 512, 200 frames into HDF5 format was about 5,400 s on a local workstation with network file system (NFS).

We ran our benchmark on up to 16 nodes (384 cores) on SDSC Comet and we observed near ideal scaling behavior Figures 6a and 6b) with parallel efficiencies above 0.8 on up to 8 nodes (Figure A3a) with no straggler tasks (Figure 6d). The trajectory reading I/O ( $t_{I/O}$ ) was the dominant contribution, followed by compute ( $t_{\text{comp}}$ ), but because both contributions scaled well, overall scaling performance remained good, even for 384 cores. We observed a low-performing outlier for 12 nodes (288 cores) with slower I/O than typical but did not further investigate. Importantly, the trajectory opening cost remained negligible (in the millisecond range) and the cost for MPI communication also remained small (below 0.1 s, even for 16 nodes). Overall, parallel MPI with HDF5 appeared to be a robust approach to obtain good scaling, even for I/O-bound tasks.

6.4 | Testing the Global Arrays Toolkit

The *Global Arrays* (GA) toolkit<sup>34</sup> is a convenient layer to represent and access arrays across multiple MPI ranks and nodes. Because of its convenience and possibly reduced communications overhead due to its use of shared memory on a physical node and MPI for inter-node communication (see Section 3.3) we wanted to compare parallel trajectory analysis with GA to the MPI-based implementation that was discussed in the previous sections.

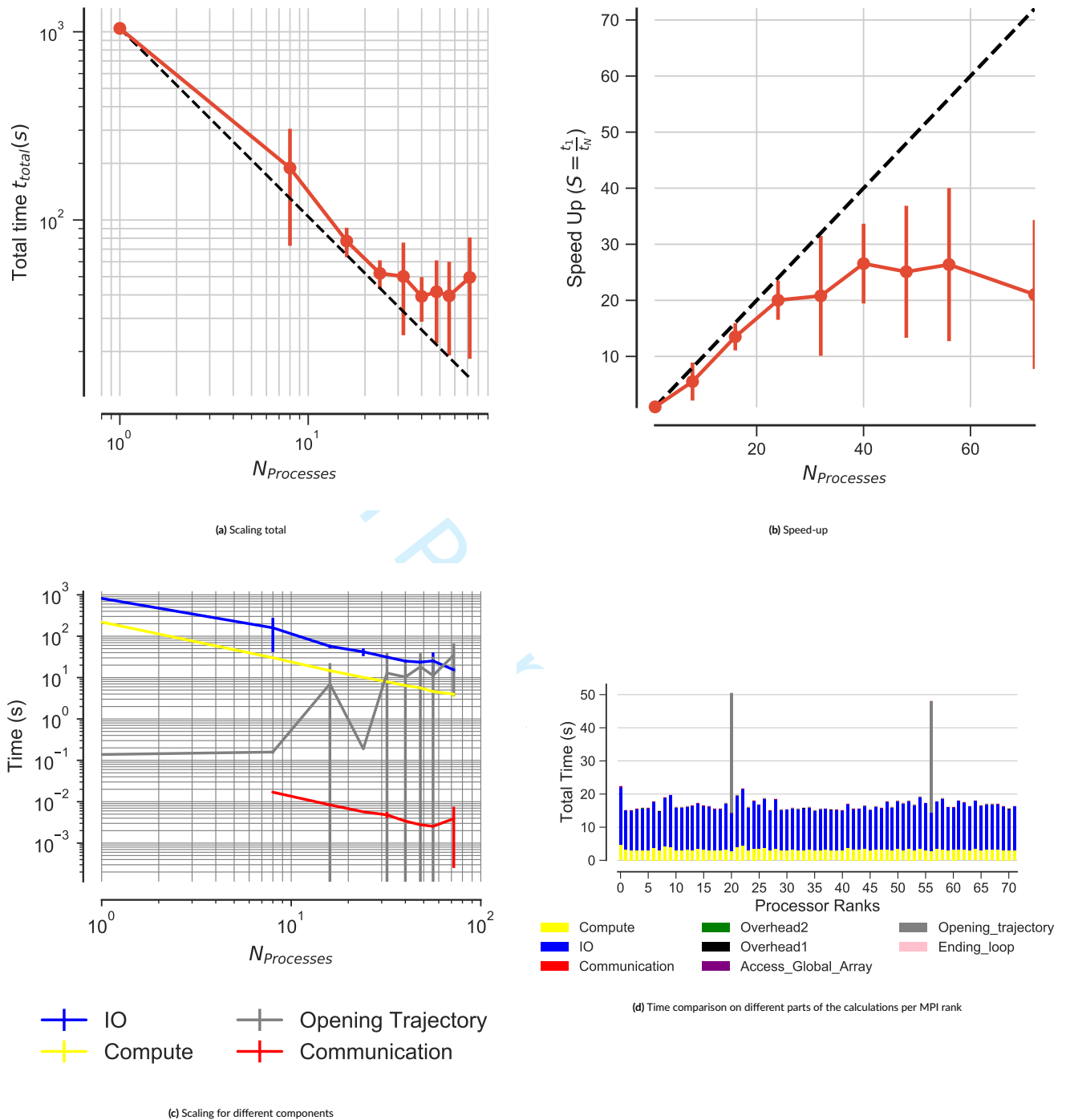
With GA, one large RMSD array called the *global array* was defined and each MPI rank updated its associated block in the global RMSD array using `ga_put()` (Algorithm 2). At the end, when all the processes exited the `block_rmsd()` function and updated their local block in the global array, rank 0 accessed the whole global array using `ga_access()`. In GA, the time for communication is  $t_{\text{ga\_put}} + t_{\text{ga\_access}}$ . We tested that the approach with GA (Algorithm 2) gave the same results as the previously discussed approach with `MPI_Gather()` (Algorithm 1).



**FIGURE 6** Performance of the RMSD task with MPI-based parallel HDF5 (MPI-IO) on SDSC Comet. Data are read from the file system from a shared HDF5 file format instead of XTC format (independent I/O) and results are communicated back to rank 0. Five repeats were performed to collect statistics. (a-c) The error bars show standard deviation with respect to the mean. In serial, there is no communication and no data points are shown for  $N = 1$  in (c). (d) Compute  $t_{\text{comp}}$ , read I/O  $t_{\text{I/O}}$ , communication  $t_{\text{comm}}$ , ending the for loop  $t_{\text{end\_loop}}$ , opening the trajectory  $t_{\text{opening\_trajectory}}$ , and overheads  $t_{\text{overhead1}}$ ,  $t_{\text{overhead2}}$  per MPI rank; see Table 3 for definitions. These are typical data from one run of the five repeats.

### Shared file

Using GA improved the strong scaling performance (Figures 7a and 7b) by reducing the communication time. Nevertheless, the remaining variation in the trajectory I/O part of the calculation and in particular the initial opening of the trajectory prevented ideal scaling (Figure 7c). Stragglers were



**FIGURE 7** Performance of the RMSD task using *Global Arrays* on *SDSC Comet*. All ranks updated the global RMSD array (`ga_put()`) and rank 0 accessed the whole RMSD array through the global memory address (`ga_get()`). Five repeats were performed to collect statistics. (a-c) The error bars show standard deviation with respect to the mean. In serial, there is no communication and not data points are shown for  $N = 1$  in (c). In (d), compute  $t_{comp}$ , read I/O  $t_{I/O}$ , communication  $t_{comm}$ , access to the whole global array by rank 0,  $t_{Access\_Global\_Array}$ , ending the for loop  $t_{end\_loop}$ , opening the trajectory  $t_{opening\_trajectory}$ , and overheads  $t_{overhead1}$ ,  $t_{overhead2}$  per MPI rank are shown; see Table 3 for definitions. These are typical data from one run of the five repeats. MPI ranks 20 and 56 were stragglers.

primarily due to the fact that all ranks had to open the same trajectory file at the beginning of the execution (Figure 7d). In this case, these slow processes took about 50 s, which was slower than the mean execution time of all other ranks of 17 s. Trajectory opening was already problematic

in the initial test (Figure 2c), which was still dominated by the communication cost. By substantially reducing communication cost, the simultaneous trajectory opening by multiple ranks emerged as the next dominant cause for stragglers.

### Subfilig

We tested subfilig (see Section 6.3.1) with GA to reduce the initial delay due to trajectory opening. Under otherwise identical conditions as in the previous section we now observed near ideal scaling behavior with efficiencies above 0.9 (Figure 5b and 5c) without any straggler tasks (Figure 5e). Although the reason why in our case GA appeared to be more efficient than direct MPI-based communication remained unclear, these results showed that contention for file access clearly impacted performance. By removing the contention, near ideal scaling could be achieved.

## 6.5 | Likely Causes of Stragglers

The data indicated that an increase in the duration of both MPI communication and trajectory file access lead to large variability in the run time of individual MPI processes and ultimately poor scaling performance beyond a single node. A discussion of likely causes for stragglers begins with the observation that opening and reading a single trajectory file from multiple MPI processes appeared to be at the center of the problem.

In MDAnalysis, individual trajectory frames are loaded into memory, which ensures that even systems with tens of millions of atoms can be efficiently analyzed on resources with moderate RAM sizes. The test trajectory (file size 30 GB) had 2,512,200 frames in total so each frame was about 0.011 MB in size. With  $t_{I/O} \approx 0.3$  ms per frame, the data were ingested at a rate of about 40 MB/s for a single process. For 24 MPI ranks (one SDSC Comet node), the aggregated reading rate would have been about 1 GB/s, well within the available bandwidth of 56 Gb/s of the InfiniBand network interface that served the Lustre file system, but nevertheless sufficient to produce substantial constant network traffic.

Furthermore, in our study the default Lustre stripe size value was 1 MB, i.e., the amount of contiguous data stored on a single Lustre object storage target (OST). Each I/O request read a single Lustre stripe but because the I/O size (0.011 MB) was smaller than the stripe size, many of these I/O requests were likely just accessing the same stripe on the same OST but nevertheless had to acquire a new reading lock for each request. The reason for this behavior is related to ensuring POSIX consistency that relates to POSIX I/O API and POSIX I/O semantics, which can have adverse effects on scalability and performance. Parallel file systems like Lustre implement sophisticated distributed locking mechanisms to ensure consistency. For example, locking mechanisms ensures that a node can not read from a file or part of a file while it might be being modified by another node. In fact, when the application I/O is not designed in a way to avoid scenarios where multiple nodes are fighting over locks for overlapping extents, Lustre can suffer from scalability limitations<sup>75</sup>. Continuously keeping metadata updated in order to have fully consistent reads and writes (POSIX metadata management), requires writing a new value for the file's last-accessed time (POSIX atime) every time a file is read, imposing a significant burden on parallel file system<sup>76</sup>. It was observed that contention for the interconnect between OSTs and compute nodes due to MPI communication may lead to variable performance in I/O measurements<sup>77</sup>. Conversely, our data suggest that single-shared-file I/O on Lustre can negatively affect MPI communication as well, even at moderate numbers (tens to hundreds) of concurrent requests, similar to recent network simulations that predicted interference between MPI and I/O traffic<sup>78</sup>. This work indicated that MPI traffic (inter-process communication) can be affected by increasing I/O, and in particular, a few MPI processes were always delayed by one to two orders of magnitude more than the median time. In summary, these observations in the context of the work by Brown et al.<sup>78</sup> suggest that our observed stragglers with large variance in the communication step might be due to interference of MPI communications with the I/O requests on the same network.

## 7 | REPRODUCIBILITY AND PERFORMANCE COMPARISON ON DIFFERENT CLUSTERS

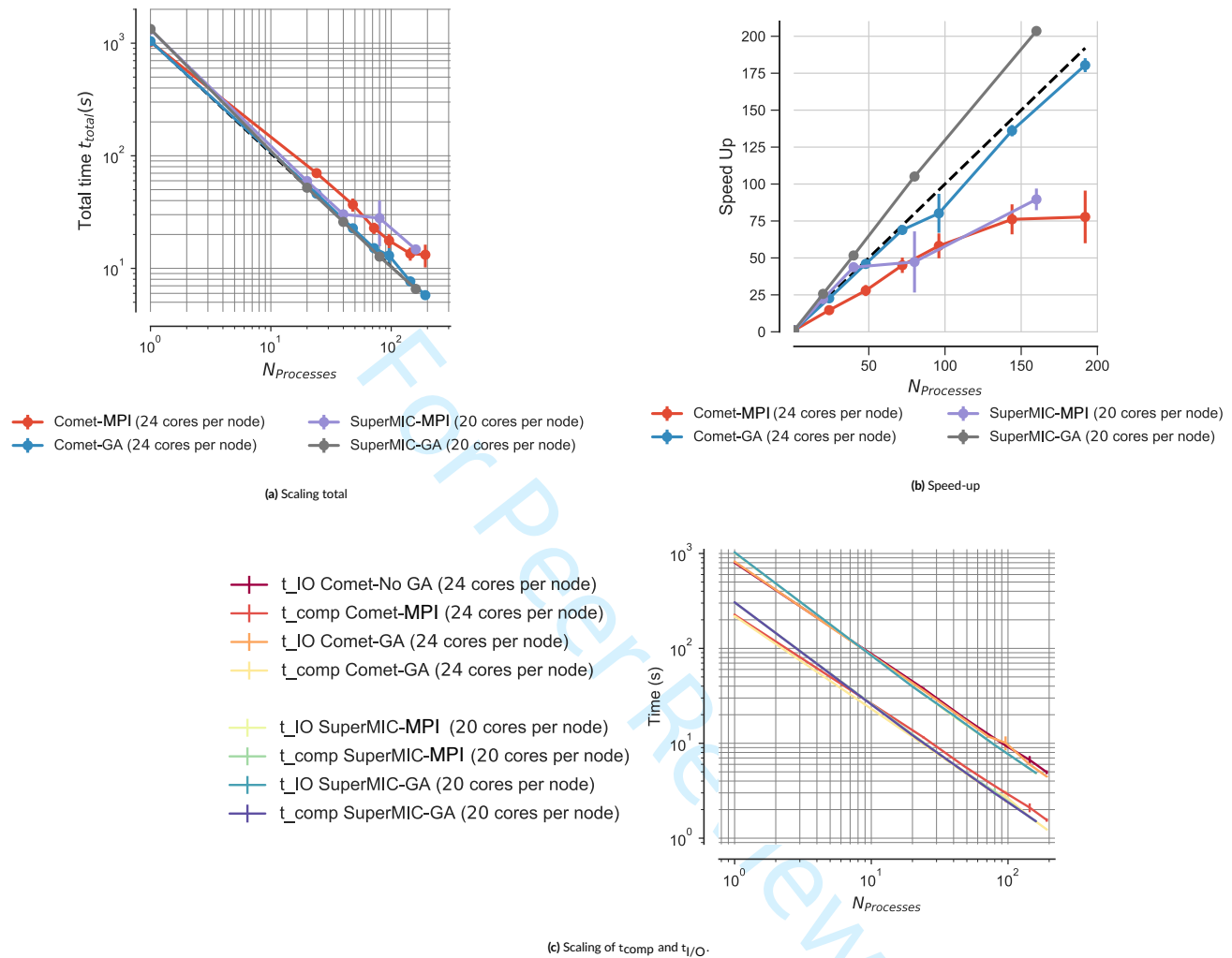
In this section we compare the performance of the RMSD task on different HPC resources (Table 1) to examine the robustness of the methods we used for our performance study and to ensure that the results are general and independent from the specific HPC system. Scripts and instructions to set up the computational environments and reproduce our computational experiments are provided in a git repository as described in section 5.

In A, we demonstrated that stragglers occur on PSC Bridges (Figure A1) and LSU SuperMIC (Figure A2) in a manner similar to the one observed on SDSC Comet (section 6.1). We performed additional comparisons for several cases discussed previously, namely (1) splitting the trajectories with blocking collective communications in MPI, (2) splitting the trajectories with Global Arrays for communications, and (3) MPI-based parallel HDF5.

### 7.1 | Splitting the Trajectories

Figure 8 shows the strong scaling of the RMSD task on different HPC resources. Splitting the trajectories with Global Arrays for communication resulted in very good scaling performance on LSU SuperMIC, similar to the results obtained on SDSC Comet. The results with MPI blocking collective communication (instead of Global Arrays) were also comparable between the two clusters, with scaling far from ideal due to the communication

cost (see section 6.3.1 and Figures 5d and A4). Overall, the scaling of the RMSD task is better on *LSU SuperMIC* than on *SDSC Comet* and the performance gap increased with increasing core number. The results on *LSU SuperMIC* confirmed the conclusion obtained on *SDSC Comet* that at least in this case Global Arrays performed better than MPI blocking collective communication.

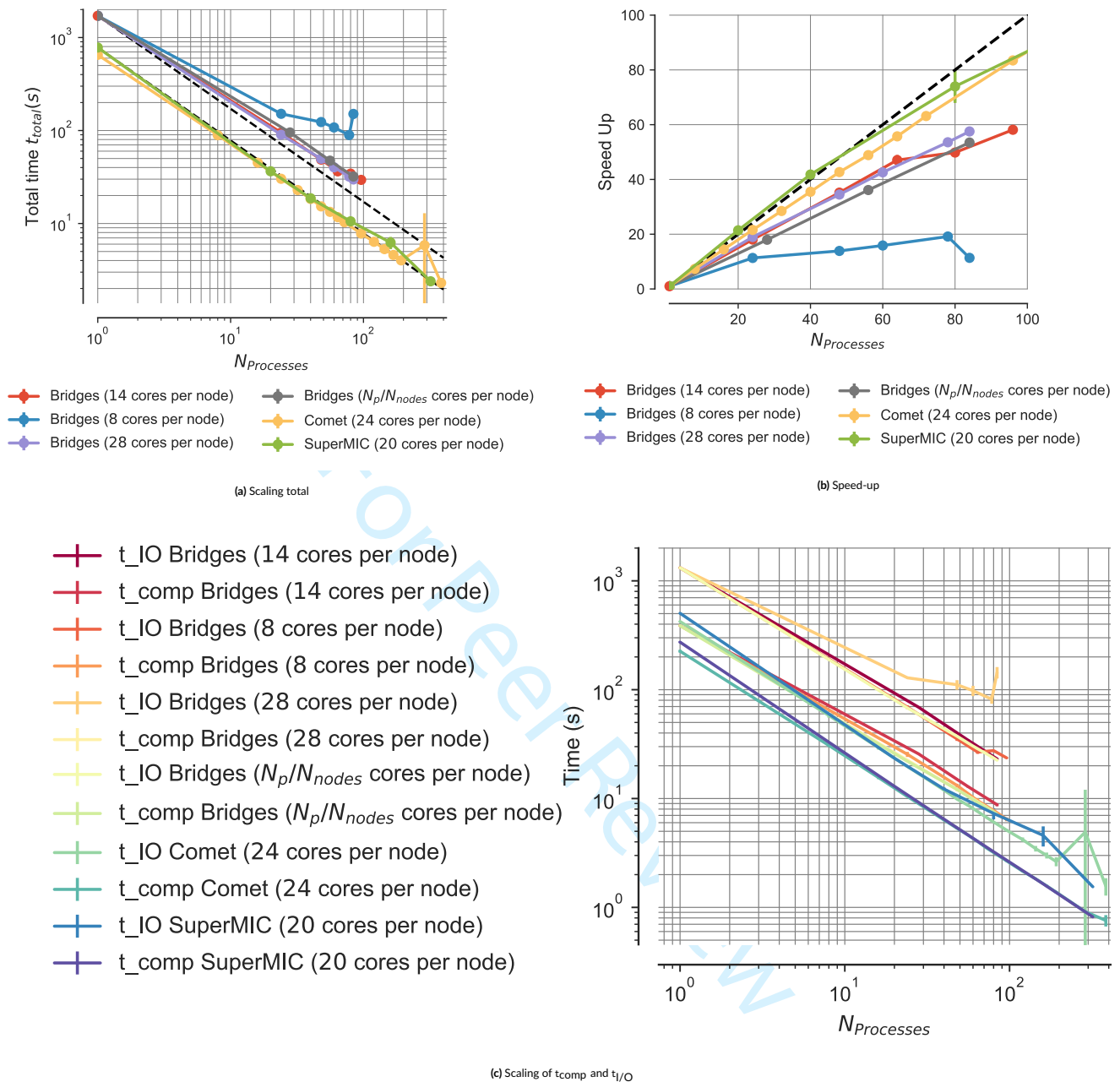


**FIGURE 8** Comparison of the performance of the RMSD task across different clusters (*SDSC Comet*, *LSU SuperMIC*) when the trajectories are split (*subfiling*). Results were communicated back to rank 0 either with MPI collective communications (label “MPI”) or using *Global Arrays* (label “GA”). Five repeats were performed to collect statistics. The error bars show the standard deviation with respect to the mean.

## 7.2 | MPI-based Parallel HDF5

Figure 9 shows the scaling on *SDSC Comet*, *LSU SuperMIC*, and *PSC Bridges* using MPI-based parallel HDF5. Performance on *SDSC Comet* and *LSU SuperMIC* was very good with near ideal linear strong scaling. The performance on *PSC Bridges* was sensitive to how many cores per node were used. Using all 28 cores in a node resulted in poor performance but decreasing the number of cores per node and equally distributing processes over nodes improved the scaling (Figure 9), mainly by reducing variation in the I/O times.

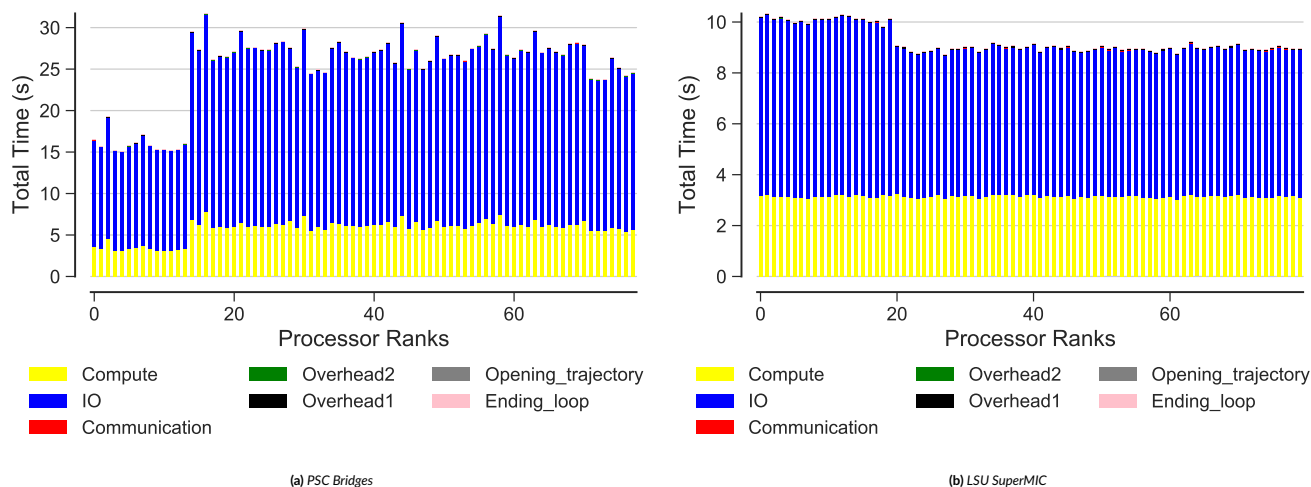
The main difference between the runs on *PSC Bridges* and *SDSC Comet/LSU SuperMIC* appeared to be the variance in  $t_{I/O}$  (Figure 9c). The I/O time distribution was fairly small and uniform across all ranks on *SDSC Comet* and *LSU SuperMIC* (Figures 10b and 6d). However, on *PSC Bridges* the I/O time was on average about two and a half times larger and the I/O time distribution was also more variable across different ranks (Figure 10a).



**FIGURE 9** Comparison of the performance of the RMSD task across different clusters (*SDSC Comet*, *PSC Bridges*, *LSU SuperMIC*) with MPI-IO. Data were read from a shared HDF5 file instead of an XTC file, using MPI independent I/O in the PHDF5 library. Results were communicated back to rank 0.  $N_p/N_{nodes}$  indicates that number of processes used for the task were equally distributed over all compute nodes. Five repeats were performed to collect statistics. The error bars show standard deviation with respect to mean. In (b) only results up to 100 cores are shown to simplify the comparison; see Fig. 6b for *SDSC Comet* and Fig. A3c for *LSU SuperMic* data.

### 7.3 | Comparison of Compute and I/O Scaling Across Different Clusters

A full comparison of compute and I/O scaling across different clusters for different test cases and algorithms is shown in Table 6. For MPI-based parallel HDF5, both the compute and I/O time on *Bridges* were consistently larger than their corresponding values on *SDSC Comet* and *LSU SuperMIC*. For example, with one core the corresponding compute and I/O time were  $t_{comp} = 387$  s,  $t_{I/O} = 1318$  s versus 225 s, 423 s on *SDSC Comet* and 273 s, 503 s on *LSU SuperMIC*. This performance difference became larger with increasing core number. When the trajectories were split and Global Arrays was used for communication both *SDSC Comet* and *LSU SuperMIC* showed similar performance.



**FIGURE 10** Examples of timing per MPI rank for the RMSD task with MPI-based parallel HDF5 on (a) *PSC Bridges* and (b) *LSU SuperMIC*. Five repeats were performed to collect statistics and these were typical data from one run of the five repeats. Compute  $t_{\text{comp}}$ , read I/O  $t_{\text{I/O}}$ , communication  $t_{\text{comm}}$ , ending the for loop  $t_{\text{end\_loop}}$ , opening the trajectory  $t_{\text{opening\_trajectory}}$ , and overheads  $t_{\text{overhead1}}$ ,  $t_{\text{overhead2}}$  per MPI rank; see Table 3 for definitions.

Overall, the results from *SDSC Comet* and *LSU SuperMIC* are consistent with each other. Performance on *PSC Bridges* seemed sensitive to the exact allocation of cores on each node but nevertheless the approaches that decreased the occurrence of stragglers on *SDSC Comet* and *LSU SuperMIC* also improved performance on *PSC Bridges*. Thus, the findings described in the previous sections are valid for a range of different HPC clusters with Lustre file systems.

## 8 | GUIDELINES FOR IMPROVING PARALLEL TRAJECTORY ANALYSIS PERFORMANCE

Although the performance measurements were performed with *MDAnalysis* and therefore capture some details of this library such as the specific timings for file reading, we believe that the broad picture is fairly general and applies to any Python-based approach that uses MPI for parallelizing trajectory access with a split-apply-combine approach. Based on the lessons that we learned, we suggest the following guidelines to improve strong scaling performance:

**Heuristic 1** Calculate compute to I/O ratio ( $R_{\text{comp}/\text{IO}}$ , Eq. 7) and compute to communication ratio ( $R_{\text{comp}/\text{comm}}$ , Eq. 8).  $R_{\text{comp}/\text{IO}}$  determines whether the task is compute bound ( $R_{\text{comp}/\text{IO}} \gg 1$ ) or IO bound ( $R_{\text{comp}/\text{IO}} \ll 1$ ).  $R_{\text{comp}/\text{comm}}$  determines whether the task is communication bound ( $\frac{t_{\text{comp}}}{t_{\text{comm}}} \ll 1$ ) or compute bound ( $R_{\text{comp}/\text{IO}} \gg 1$ ).

As discussed in Section 6.2, for I/O bound problems the interference between MPI communication and I/O traffic can be problematic<sup>51,78</sup> and the performance of the task will be affected by the straggler tasks that delay job completion time.

**Heuristic 2** For  $R_{\text{comp}/\text{IO}} \geq 1$ , single-shared-file I/O can be used and will not decrease performance. One may consider the following cases:

**Heuristic 2.1** If  $R_{\text{comp}/\text{comm}} \gg 1$ , the task is compute bound and will scale well as shown in Figure 3.

**Heuristic 2.2** If  $R_{\text{comp}/\text{comm}} \ll 1$ , one might consider using *Global Arrays* to improve scaling by utilizing efficient distribution of data via the shared arrays (section 6.4).

**Heuristic 3** For  $R_{\text{comp}/\text{IO}} \leq 1$  the task is I/O bound and single-shared-file I/O should be avoided to remove unnecessary metadata operations. One might want to consider the following steps:

**Heuristic 3.1** If there is access to HDF5 format, use MPI-based Parallel HDF5 (Section 6.3.2).

**Heuristic 3.2** If the trajectory file is not in HDF5 format then one can consider subfiling and split the single trajectory file into as many trajectory segments as the number of processes. Splitting the trajectories can be easily performed in parallel and trajectory conversion may be integrated into the beginning of standard workflows for MD simulations. Alternatively, trajectories may be kept in smaller



Cluster	Gather	File Access	Time	Serial	N <sub>Processes</sub>									
					Comet: 24 Bridges: 24 SuperMIC: 20	Comet: 48 Bridges: 48 SuperMIC: 40	Comet: 72 Bridges: 60 SuperMIC: 80	Comet: 96 Bridges: 78 SuperMIC: 160	Comet: 144 Bridges: 84 SuperMIC: 320	Comet: 192 SuperMIC: 320	Comet: 384 SuperMIC: 320			
Comet	MPI	Single	t <sub>i/o</sub>	791 ± 5.22	49 ± 3.45	29 ± 1.3	26 ± 9.19	-	-	-	-			
			t <sub>comp</sub>	225 ± 5.4	11 ± 0.75	6 ± 0.35	4 ± 0.48	-	-	-	-			
Bridges	MPI	Single	t <sub>i/o</sub>	770 ± 10.8	38 ± 0.84	33 ± 19.4	15 ± 1.6	-	-	-	-			
			t <sub>comp</sub>	221 ± 3.9	11 ± 0.43	6 ± 0.32	4 ± 0.18	-	-	-	-			
SuperMIC	MPI	Single	t <sub>i/o</sub>	1014.51 ± 2.94	48.08 ± 0.35	24.5 ± 0.79	12 ± 0.31	-	6.24 ± 0.38	-	-			
			t <sub>comp</sub>	303.85 ± 2.3	14.56 ± 0.14	7.4 ± 0.25	3.7 ± 0.12	-	1.8 ± 0.04	-	-			
Comet	GA	Single	t <sub>i/o</sub>	820 ± 18.49	41 ± 8.99	23 ± 4.14	15 ± 2.06	-	-	-	-			
			t <sub>comp</sub>	219 ± 9.8	10 ± 0.3	5 ± 0.48	3 ± 0.54	-	-	-	-			
Comet	MPI	Splitting	t <sub>i/o</sub>	799 ± 5.22	37 ± 1.22	18 ± 0.18	12 ± 0.14	9 ± 0.3	6 ± 0.66	4 ± 0.23	-			
			t <sub>comp</sub>	225 ± 5.4	11 ± 0.31	5 ± 0.07	3 ± 0.04	3 ± 0.11	2 ± 0.23	1 ± 0.07	-			
SuperMIC	MPI	Splitting	t <sub>i/o</sub>	1013.75 ± 2.8	39.99 ± 0.36	19.18 ± 0.25	9.61 ± 0.28	-	4.83 ± 0.06	-	-			
			t <sub>comp</sub>	304.26 ± 2.55	12.41 ± 0.22	5.99 ± 0.09	3.08 ± 0.13	-	1.5 ± 0.01	-	-			
Comet	GA	Splitting	t <sub>i/o</sub>	820 ± 18.5	36 ± 0.78	17 ± 0.3	11 ± 0.23	10 ± 1.7	5 ± 0.14	4 ± 0.07	-			
			t <sub>comp</sub>	219 ± 9.5	9 ± 0.22	4 ± 0.07	3 ± 0.04	2 ± 0.4	1 ± 0.05	1 ± 0.02	-			
SuperMIC	GA	Splitting	t <sub>i/o</sub>	1027.62 ± 10.32	39.62 ± 0.2	19.66 ± 0.1	9.57 ± 0.1	-	4.86 ± 0.05	-	-			
			t <sub>comp</sub>	305.78 ± 3.47	12.16 ± 0.1	6.01 ± 0.007	2.97 ± 0.1	-	1.51 ± 0.03	-	-			
Comet	MPI	PHDF5	t <sub>i/o</sub>	423 ± 5.88	19 ± 0.3	9 ± 0.13	6 ± 0.06	5 ± 0.12	3 ± 0.2	3 ± 0.25	1.57 ± 0.29			
			t <sub>comp</sub>	225 ± 6.55	10 ± 0.12	5 ± 0.1	3 ± 0.04	2 ± 0.05	1 ± 0.04	1 ± 0.03	0.76 ± 0.09			
Bridges	MPI	PHDF5	t <sub>i/o</sub>	1318.87 ± 10.42	67.93 ± 0.52	37.37 ± 0.2	30.35 ± 0.15	24.16 ± 0.89	22.5 ± 0.17	-	-			
			t <sub>comp</sub>	387.8 ± 5.51	21.97 ± 0.38	12.12 ± 0.34	9.79 ± 0.24	7.72 ± 0.03	7.18 ± 0.08	-	-			
SuperMIC	MPI	PHDF5	t <sub>i/o</sub>	503.69 ± 2.57	12.96 ± 0.06	6.46 ± 0.02	3.2 ± 0.01	-	1.64 ± 0.01	-	0.82 ± 0.004			
			t <sub>comp</sub>	273.54 ± 4.7	23.44 ± 0.29	12.22 ± 0.43	7.3 ± 0.85	-	4.59 ± 0.96	-	1.55 ± 0.009			

**TABLE 6** Comparison of the compute and I/O scaling for different test cases and number of processes. Five repeats were performed to collect statistics. The mean value and the standard deviation with respect to mean are reported for each case.

chunks if they are already produced in batches; for instance, when running simulations with *Gromacs*<sup>56</sup>, the `gmx mdrun -noappend` option produces individual trajectory segments instead of extending an existing trajectory file.

**Heuristic 3.3** In case of  $R_{\text{comp/comm}} \ll 1$ , use of *Global Arrays* may be considered to potentially improve scaling (Section 6.3.1).

## 9 | CONCLUSIONS

We analyzed the strong scaling performance of a typical task when analyzing MD trajectories, the calculation of the time series of the RMSD of a protein, with the widely used Python-based *MDAnalysis* library. The task was parallelized with MPI following the *split-apply-combine* approach by having each MPI process analyze a contiguous segment of the trajectory. This approach did not scale beyond a single node because straggler MPI processes exhibited large upward variations in runtime. Stragglers were primarily caused by either increased MPI communication costs or increased time to open the single shared trajectory file whereas both the computation and the ingestion of data exhibited close to ideal strong scaling behavior. Stragglers were less prevalent for compute-bound workloads (i.e.,  $R_{\text{comp/IO}} \gg 1$ ), suggesting that file read I/O was responsible for poor MPI communication. In particular, artificially removing all I/O substantially improved performance of the communication step and thus brought overall performance close to ideal (i.e., linear increase in speed-up with processor count with slope one). By performing benchmarks on three different XSEDE supercomputers we showed that our results were independent from the specifics of the hardware and local environment. Our results hint at the possibility that stragglers might be due to the competition between MPI messages and the Lustre file system on the shared InfiniBand interconnect, which would be consistent with other similar observations<sup>51</sup> and theoretical predictions by Brown et al.<sup>78</sup>. One possible interpretation of our results is that for a sufficiently large per-frame compute workload, read I/O interferes much less with communication than for an I/O bound task that almost continuously accesses the file system. This interpretation suggested that we needed to improve read I/O to reduce interference.

We investigated subfiling (splitting of the trajectories into separate files, one for each MPI rank) and MPI-based parallel I/O. Subfiling improved scaling, especially when combined with the *Global Arrays* toolkit. *Global Arrays* reduced the communication cost compared to MPI collective communications even though it only acts as programming layer to access data across multiple nodes in a convenient array form and also uses MPI for its inter-node data exchange. Subfiling with *Global Arrays* achieved nearly ideal scaling up to 192 cores (8 nodes on *SDSC Comet*). When we used MPI-based parallel I/O through HDF5 together with MPI for communications we achieved nearly ideal performance up to 384 cores (16 nodes on *SDSC Comet*) and speed-ups of two orders of magnitude compared to the serial execution. The latter approach appears to be a promising way forward as it directly builds on very widely used technology (MPI-IO and HDF5) and echoes the experience of the wider HPC community that parallel file I/O is necessary for efficient data handling.

The biomolecular simulation community suffers from a large number of trajectory file formats with very few being based on HDF5, with the exception of the H5MD format<sup>79</sup> and the MDTraj HDF5 format<sup>20</sup>. Our work suggests that HDF5-based formats should be seriously considered as the default for MD simulations if users want to make efficient use of their HPC systems for analysis. Alternatively, enabling MPI-IO for trajectory readers in libraries such as *MDAnalysis* might also provide a path forward to better read performance.

We summarized our findings in a number of guidelines for improving the scaling of parallel analysis of MD trajectory data. We showed that it is feasible to run an I/O bound analysis task on HPC resources with a Lustre parallel file system and achieve good scaling behavior up to 384 CPU cores with an almost 300-fold speed-up compared to serial execution. Although we focused on the *MDAnalysis* library, similar strategies are likely to be more generally applicable and useful to the wider biomolecular simulation community.

## ACKNOWLEDGEMENTS

We are grateful to Sarp Oral for insightful comments on this manuscript. This work was supported by the National Science Foundation under grant numbers ACI-1443054 and ACI-1440677. This work used the Extreme Science and Engineering Discovery Environment (XSEDE), which is supported by National Science Foundation grant number ACI-1548562. *SDSC Comet* at the San Diego Supercomputer Center, *LSU SuperMic* at Louisiana State University, and *PSC Bridges* at the Pittsburgh Supercomputing Center were used under allocations TG-MCB090174 and TG-MCB130177.

## References

1. Borhani DW, Shaw DE. The future of molecular dynamics simulations in drug discovery. *J Comput Aided Mol Des* 2012; 26(1): 15–26. doi: 10.1007/s10822-011-9517-y.

2. Dror RO, Dirks RM, Grossman JP, Xu H, Shaw DE. Biomolecular simulation: a computational microscope for molecular biology. *Annu Rev Biophys* 2012; 41: 429–52. doi: 10.1146/annurev-biophys-042910-155245.
3. Orozco M. A theoretical view of protein dynamics. *Chem Soc Rev* 2014; 43: 5051–5066. doi: 10.1039/C3CS60474H.
4. Perilla JR, Goh BC, Cassidy CK, et al. Molecular dynamics simulations of large macromolecular complexes. *Current Opinion in Structural Biology* 2015; 31: 64 – 74. doi: 10.1016/j.sbi.2015.03.007.
5. Bottaro S, Lindorff-Larsen K. Biophysical experiments and biomolecular simulations: A perfect match? *Science* 2018; 361(6400): 355–360. doi: 10.1126/science.aat4010.
6. Tuckerman ME. *Statistical Mechanics: Theory and Molecular Simulation*. Oxford, UK: Oxford University Press, 2010.
7. Mura C, McAnany CE. An introduction to biomolecular simulations and docking. *Molecular Simulation* 2014; 40(10-11): 732–764. doi: 10.1080/08927022.2014.935372.
8. Cheatham T, Roe D. The impact of heterogeneous computing on workflows for biomolecular simulation and analysis. *Computing in Science Engineering* 2015; 17(2): 30–39. doi: 10.1109/MCSE.2015.7.
9. Kneller GR, Keiner V, Kneller M, Schiller M. nmoldyn: A program package for a neutron scattering oriented analysis of molecular dynamics simulations. *Computer Physics Communications* 1995; 91(1): 191 – 214. doi: 10.1016/0010-4655(95)00048-K.
10. Hinsen K, Pellegrini E, Stachura S, Kneller GR. nmoldyn 3: Using task farming for a parallel spectroscopy-oriented analysis of molecular dynamics simulations. *Journal of Computational Chemistry* 2012; 33(25): 2043–2048. doi: 10.1002/jcc.23035.
11. Humphrey W, Dalke A, Schulten K. VMD – Visual Molecular Dynamics. *J Mol Graph* 1996; 14: 33–38.
12. Hinsen K. The molecular modeling toolkit: a new approach to molecular simulations. *Journal of Computational Chemistry* 2000; 21(2): 79–85.
13. Grant BJ, Rodrigues APC, ElSawy KM, McCammon JA, Caves LSD. Bio3d: an R package for the comparative analysis of protein structures. *Bioinformatics* 2006; 22(21): 2695–6. doi: 10.1093/bioinformatics/btl461.
14. Tu T, Rendleman CA, Borhani DW, et al. A scalable parallel framework for analyzing terascale molecular dynamics simulation trajectories. In: 2008 SC - International Conference for High Performance Computing, Networking, Storage and Analysis. Austin, TX, USA: IEEE, 2008; 1–12. doi: 10.1109/SC.2008.5214715.
15. Romo TD, Grossfield A. LOOS: An extensible platform for the structural analysis of simulations. In: 31st Annual International Conference of the IEEE EMBS. Minneapolis, Minnesota, USA: IEEE, 2009; 2332–2335.
16. Romo TD, Leioatts N, Grossfield A. Lightweight object oriented structure analysis: Tools for building tools to analyze molecular dynamics simulations. *Journal of Computational Chemistry* 2014; 35(32): 2305–2318. doi: 10.1002/jcc.23753.
17. Michaud-Agrawal N, Denning EJ, Woolf TB, Beckstein O. MDAnalysis: A toolkit for the analysis of molecular dynamics simulations. *J Comp Chem* 2011; 32: 2319–2327. doi: 10.1002/jcc.21787.
18. Gowers RJ, Linke M, Barnoud J, et al. MDAnalysis: A Python package for the rapid analysis of molecular dynamics simulations. In: Benthall S, Rostrup S, eds., *Proceedings of the 15th Python in Science Conference*. Austin, TX: SciPy, 2016; 102 – 109. doi: 10.25080/Majora-629e541a-00e.
19. Roe DR, Thomas E Cheatham I. PTRAJ and CPPTRAJ: Software for processing and analysis of molecular dynamics trajectory data. *Journal of Chemical Theory and Computation* 2013; 9(7): 3084–3095. doi: 10.1021/ct400341p. PMID: 26583988.
20. McGibbon RT, Beauchamp KA, Harrigan MP, et al. MDTraj: A modern open library for the analysis of molecular dynamics trajectories. *Biophysical Journal* 2015; 109(8): 1528 – 1532. doi: 10.1016/j.bpj.2015.08.015.
21. Yesylevskyy SO. Pteros 2.0: Evolution of the fast parallel molecular analysis library for C++ and python. *Journal of Computational Chemistry* 2015; 36(19): 1480–1488. doi: 10.1002/jcc.23943.
22. Doerr S, Harvey MJ, Noé F, De Fabritiis G. HTMD: High-throughput molecular dynamics for molecular discovery. *Journal of Chemical Theory and Computation* 2016; 12(4): 1845–1852. doi: 10.1021/acs.jctc.6b00049.

23. Khoshlessan M, Paraskevatos I, Jha S, Beckstein O. Parallel analysis in MDAnalysis using the Dask parallel computing library. In: Katy Huff, David Lippa, Dillon Niederhut, Pacer M, eds., *Proceedings of the 16th Python in Science Conference*. Austin, TX: SciPy, 2017; 64–72. doi: 10.25080/shinma-7f4c6e7-00a.
24. Paraskevatos I, Luckow A, Khoshlessan M, et al. Task-parallel analysis of molecular dynamics trajectories. In: *ICPP 2018: 47th International Conference on Parallel Processing, August 13–16, 2018, Eugene, OR, USA*. Association for Computing Machinery, New York, NY, USA: ACM, 2018; Article No. 49.
25. Shujie Fan, Max Linke, Ioannis Paraskevatos, Richard J Gowers, Michael Gecht, Oliver Beckstein. PMDA - Parallel Molecular Dynamics Analysis. In: Chris Calloway, David Lippa, Dillon Niederhut, David Shupe, eds., *Proceedings of the 18th Python in Science Conference*. Austin, TX: SciPy, 2019; 134 – 142. doi: 10.25080/Majora-7ddc1dd1-013.
26. Liu P, Agrafiotis DK, Theobald DL. Fast determination of the optimal rotational matrix for macromolecular superpositions. *J Comput Chem* 2010; 31(7): 1561–3. doi: 10.1002/jcc.21439.
27. Leach AR. *Molecular Modelling. Principles and Applications*. Longman, 1996.
28. Rocklin M. Dask: Parallel computation with blocked algorithms and task scheduling. In: Huff K, Bergstra J, eds., *Proceedings of the 14th Python in Science Conference (SciPy 2015)*. Austin, TX: SciPy, 2015; 130–136. doi: 10.25080/Majora-7b98e3ed-013.
29. Dalcín LD, Paz RR, Kler PA, Cosimo A. Parallel distributed computing using python. *Advances in Water Resources* 2011; 34(9): 1124 – 1139. doi: 10.1016/j.advwatres.2011.04.013.
30. Dalcín LD, Paz R, Storti M. MPI for python. *Journal of Parallel and Distributed Computing* 2005; 65(9): 1108 – 1115. doi: 10.1016/j.jpdc.2005.03.010.
31. Garraghan P, Ouyang X, Yang R, McKee D, Xu J. Straggler root-cause and impact analysis for massive-scale virtualized cloud datacenters. *IEEE Transactions on Services Computing* 2016; 12: 91–104. doi: 10.1109/TSC.2016.2611578.
32. Towns J, Cockerill T, Dahan M, et al. XSEDE: Accelerating scientific discovery. *Computing in Science & Engineering* 2014; 16(5): 62–74. doi: 10.1109/MCSE.2014.80.
33. Daily JA. *GAiN: Distributed Array Computation with Python*. Master's thesis, School of Electrical Engineering and Computer Science, Washington State University, Pullman, WA, 2009.
34. Nieplocha J, Palmer B, Tipparaju V, Krishnan M, Trease H, Aprà E. Advances, applications and performance of the Global Arrays shared memory programming toolkit. *The International Journal of High Performance Computing Applications* 2006; 20(2): 203–231.
35. Dean J, Ghemawat S. MapReduce: simplified data processing on large clusters. *Communications of the ACM* 2008; 51(1): 107–113. doi: 10.1145/1327452.1327492.
36. Ananthanarayanan G, Kandula S, Greenberg A, et al. Reining in the outliers in map-reduce clusters using Mantri. In: *Proceedings of the 9th USENIX Conference on Operating Systems Design and Implementation, OSDI'10*. Berkeley, CA, USA: USENIX Association, 2010; 265–278.
37. Kyong J, Jeon J, Lim SS. Improving scalability of apache spark-based scale-up server through docker container-based partitioning. In: *Proceedings of the 6th International Conference on Software and Computer Applications - ICSCA '17*. New York, USA: ACM Press, 2017; 176–180. doi: 10.1145/3056662.3056686.
38. Ousterhout K. *Architecting for Performance Clarity in Data Analytics Frameworks*. Ph.D. thesis, EECS Department, University of California, Berkeley, Berkeley, CA, 2017. URL <https://www2.eecs.berkeley.edu/Pubs/TechRpts/2017/EECS-2017-158.html>.
39. Gittens A, Devarakonda A, Racah E, et al. Matrix factorizations at scale: A comparison of scientific data analytics in spark and C+MPI using three case studies. In: *IEEE International Conference on Big Data (Big Data)*. 2016; 204–213. doi: 10.1109/BigData.2016.7840606.
40. Yang H, Liu X, Chen S, Lei Z, Du H, Zhu C. Improving Spark performance with MPTE in heterogeneous environments. In: *2016 International Conference on Audio, Language and Image Processing (ICALIP)*. IEEE, 2016; 28–33. doi: 10.1109/ICALIP.2016.7846627.
41. Kirpichov E, Denielou M. No shard left behind: dynamic work rebalancing in Google Cloud Dataflow. Google Cloud Blog, 2016. URL <https://cloud.google.com/blog/products/gcp/no-shard-left-behind-dynamic-work-rebalancing-in-google-cloud-dataflow>. Accessed Aug 24, 2019.

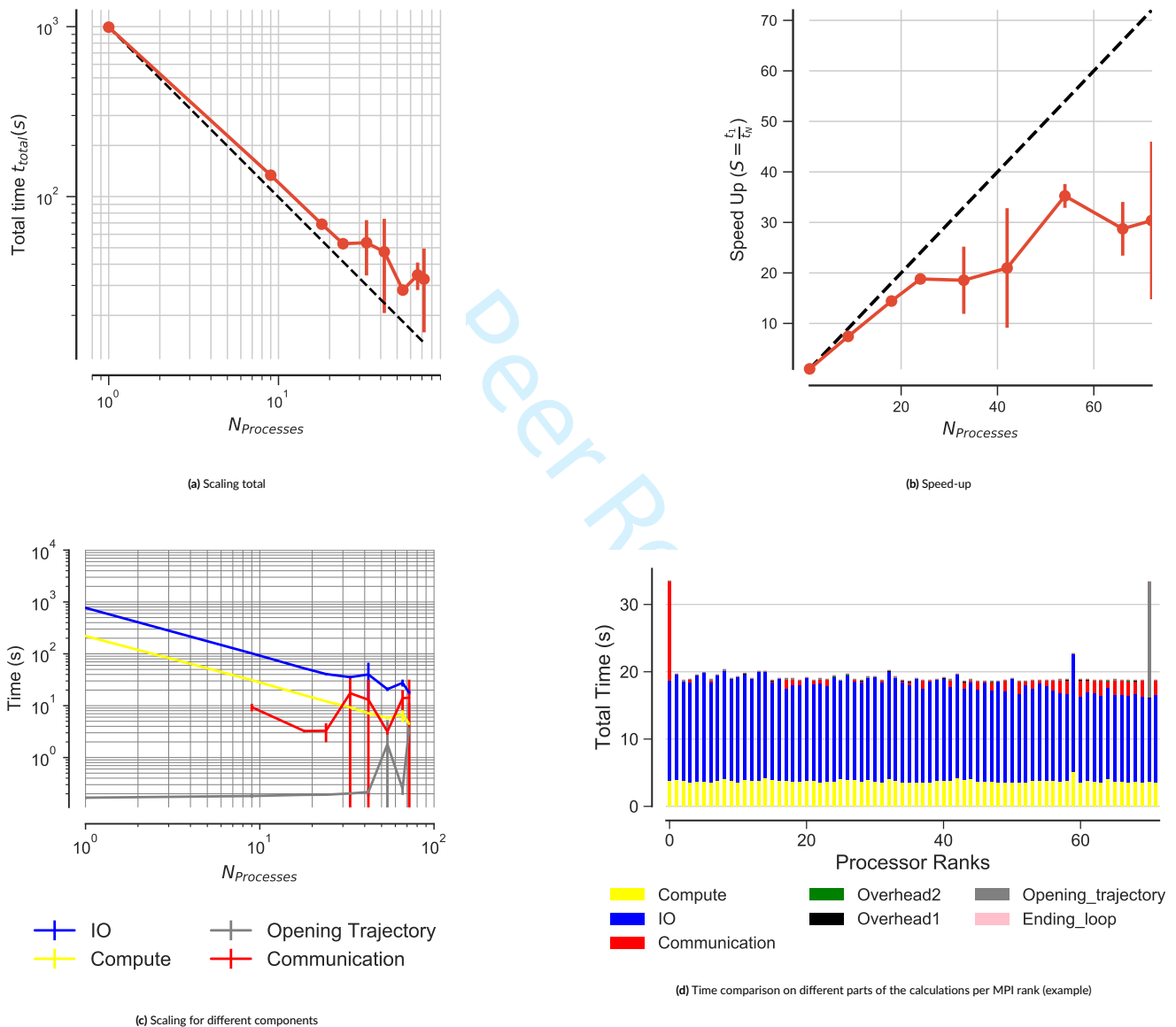
42. Phan TD. *Energy-efficient Straggler Mitigation for Big Data Applications on the Clouds*. Ph.D. thesis, École normale supérieure de Renne, 2017.
43. Chen Q, Liu C, Xiao Z. Improving mapreduce performance using smart speculative execution strategy. *IEEE Transactions on Computers* 2014; 63(4): 954–967. doi: 10.1109/TC.2013.15.
44. Xie B, Chase J, Dillow D, et al. Characterizing output bottlenecks in a supercomputer. In: *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis, SC '12*. Los Alamitos, CA, USA: IEEE Computer Society Press, 2012; 8:1–8:11.
45. Rosen J, Zhao B. Fine-grained micro-tasks for mapreduce skew-handling. Tech. rep., EECS, UC Berkeley, 2012. URL <https://pdfs.semanticscholar.org/3617/916adb83f33f8df7d0b3bfc23d0de80da9b7.pdf>.
46. Kwon Y, Balazinska M, Howe B, Rolia J. Skewtune: Mitigating skew in mapreduce applications, pages 25–36. In: *SIGMOD'12. SIGMOD '12 Proceedings of the 2012 ACM SIGMOD International Conference on Management of Data*, 2012; Pages 25–36. doi: 10.1145/2213836.2213840.
47. Ousterhout K, Rasti R, Ratnasamy S, Shenker S, Chun BG. Making sense of performance in data analytics frameworks. In: *NSDI'15 Proceedings of the 12th USENIX Conference on Networked Systems Design and Implementation*, ISBN: 978-1-931971-218. 2015; Pages 293–307.
48. Abdul-Wahid B, Feng H, Rajan D, et al. Awe-wq, fast-forwarding molecular dynamics using the accelerated weighted ensemble. *Journal of Chemical Information and Modeling* 2014; 54: 3033–3043.
49. Wu G, Song H, Lin D. A scalable parallel framework for microstructure analysis of large-scale molecular dynamics simulations data. *Computational Materials Science* 2018; 144: 322–330.
50. Tu T, Rendleman CA, Miller PJ, Sacerdoti F, Dror RO, Shaw DE. Accelerating parallel analysis of scientific simulation data via zazen. In: *8th USENIX Conference on File and Storage Technologies, San Jose, CA, USA*. 2010; 129–142. URL [http://www.usenix.org/events/fast10/tech/full\\_papers/tu.pdf](http://www.usenix.org/events/fast10/tech/full_papers/tu.pdf).
51. Stone JE, Israilewitz B, Schulten K. Early experiences scaling vmd molecular visualization and analysis jobs on blue waters. In: *Proceedings of the 2013 Extreme Scaling Workshop (Xsw 2013)*. Washington, DC, USA: IEEE Computer Society, 2013; 43–50.
52. Shkurtia A, Goni R, Andrio P, et al. pyPcazip: A PCA-based toolkit for compression and analysis of molecular simulation data. *SoftwareX* 2016; 5: 44–50.
53. Malakar P, Knight C, Munson T, Vishwanath V, Papka ME. Scalable in situ analysis of molecular dynamics simulations. In: *ISAV'17 Proceedings of the In Situ Infrastructures on Enabling Extreme-Scale Analysis and Visualization*. 2017; 1–6.
54. Johnston T, Zhang B, Liwo A, Crivelli S, Taufer M. In situ data analytics and indexing of protein trajectories. *J Comput Chem* 2017; 38(16): 1419–1430. doi: 10.1002/jcc.24729.
55. Brooks BR, Brooks III CL, Mackerell ADJ, et al. CHARMM: the biomolecular simulation program. *J Comp Chem* 2009; 30(10): 1545–1614. doi: 10.1002/jcc.21287.
56. Abraham MJ, Murtola T, Schulz R, et al. GROMACS: High performance molecular simulations through multi-level parallelism from laptops to supercomputers. *SoftwareX* 2015; 1–2: 19 – 25. doi: 10.1016/j.softx.2015.06.001.
57. Case DA, Cheatham TE 3rd, Darden T, et al. The amber biomolecular simulation programs. *J Comput Chem* 2005; 26(16): 1668–1688. doi: 10.1002/jcc.20290.
58. Phillips J, Braun R, Wang W, et al. Scalable molecular dynamics with NAMD. *J Comput Chem* 2005; 26: 1781–1802. doi: 10.1002/jcc.20289.
59. Burley SK, Berman HM, Bhikadiya C, et al. Protein Data Bank: the single global archive for 3D macromolecular structure data. *Nucleic Acids Research* 2018; 47(D1): D520–D528. doi: 10.1093/nar/gky949.
60. Van Der Walt S, Colbert SC, Varoquaux G. The numpy array: a structure for efficient numerical computation. *Computing in Science & Engineering* 2011; 13(2): 22–30. doi: 10.1109/MCSE.2011.37.
61. Theobald DL. Rapid calculation of RMSDs using a quaternion-based characteristic polynomial. *Acta Crystallogr A* 2005; 61(Pt 4): 478–80. doi: 10.1107/S0108767305015266.

62. Daily J, Vishnu A, Palmer B, van Dam H, Kerbyson D. On the suitability of MPI as a PGAS runtime. In: *2014 21st International Conference on High Performance Computing (HiPC)*. 2014; 1–10. doi: 10.1109/HiPC.2014.7116712.
63. Nieplocha J, Harrison RJ, Littlefield RJ. Global Arrays: A non-uniform-memory-access programming model for high-performance computers. *Journal of Supercomputing* 1996; 10(2): 169–189.
64. Collette A. Python and hdf5. In: Blanchette M, Roumeliotis R, eds., *Python and HDF5*. O'Reilly Media, Inc., 1005 Gravenstein Highway North, Sebastopol, CA 95472., 2014; .
65. Seyler SL, Beckstein O. Sampling of large conformational transitions: Adenylate kinase as a testing ground. *Molec Simul* 2014; 40(10–11): 855–877. doi: 10.1080/08927022.2014.919497.
66. Seyler S, Beckstein O. Molecular dynamics trajectory for benchmarking MDAnalysis. figshare, 2017. doi: 10.6084/m9.figshare.5108170.
67. Lindahl E, Hess B, van der Spoel D. Gromacs 3.0: A package for molecular simulation and trajectory analysis. *J Mol Mod* 2001; 7(8): 306–317. doi: 10.1007/s008940100045.
68. Spångberg D, Larsson DSD, van der Spoel D. Trajectory NG: portable, compressed, general molecular dynamics trajectories. *J Mol Model* 2011; 17(10): 2669–85. doi: 10.1007/s00894-010-0948-5.
69. Shaw DE, Dror RO, Salmon JK, et al. Millisecond-scale molecular dynamics simulations on anton. In: *SC '09: Proceedings of the Conference on High Performance Computing Networking, Storage and Analysis*. New York, NY, USA: ACM, 2009; 1–11. doi: 10.1145/1654059.1654099.
70. Shaw DE, Grossman JP, Bank JA, et al. Anton 2: Raising the bar for performance and programmability in a special-purpose molecular dynamics supercomputer. In: *SC '14: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*. 2014; 41–53. doi: 10.1109/SC.2014.9.
71. Salomon-Ferrer R, Götz AW, Poole D, Le Grand S, Walker RC. Routine microsecond molecular dynamics simulations with amber on gpus. 2. explicit solvent particle mesh ewald. *Journal of Chemical Theory and Computation* 2013; 9(9): 3878–3888. doi: 10.1021/ct400314y.
72. Glaser J, Nguyen TD, Anderson JA, et al. Strong scaling of general-purpose molecular dynamics simulations on gpus. *Computer Physics Communications* 2015; 192: 97–107. doi: 10.1016/j.cpc.2015.02.028.
73. Choudhary A, Liao Wk, Gao K, et al. Scalable I/O and analytics. *Journal of Physics: Conference Series* 2009; 180(012048).
74. Son SW, Sehrish S, keng Liao W, Oldfield R, Choudhary A. Reducing I/O variability using dynamic I/O path characterization in petascale storage systems. *Journal of Supercomputing* 2017; 73(5): pp 2069–2097.
75. Lin KW, Chou J, Byna S, Wu K. Optimizing fast query performance on Lustre file system. In: *SSDBM Proceedings of the 25th International Conference on Scientific and Statistical Database Management*, Article No. 29. 2013; .
76. Lockwood G. What is so bad about POSIX I/O? The Next Platform, 2017. URL <https://www.nextplatform.com/2017/09/11/whats-bad-posix-io/>. Accessed Aug 24, 2019.
77. Mache J, Lo V, Garg S. The impact of spatial layout of jobs on I/O hotspots in mesh networks. *Journal of Parallel and Distributed Computing* 2005; 65(10): 1190 – 1203. doi: 10.1016/j.jpdc.2005.04.020.
78. Brown KA, Jain N, Matsuoka S, Schulz M, Bhatele A. Interference between I/O and MPI traffic on fat-tree networks. In: *Proceedings of the 47th International Conference on Parallel Processing, ICPP 2018*. New York, NY, USA: ACM, 2018; 7:1–7:10. doi: 10.1145/3225058.3225144.
79. de Buyl P, Colberg PH, Höfling F. H5MD: A structured, efficient, and portable file format for molecular data. *Computer Physics Communications* 2014; 185(6): 1546 – 1553. doi: 10.1016/j.cpc.2014.01.018.

## APPENDIX

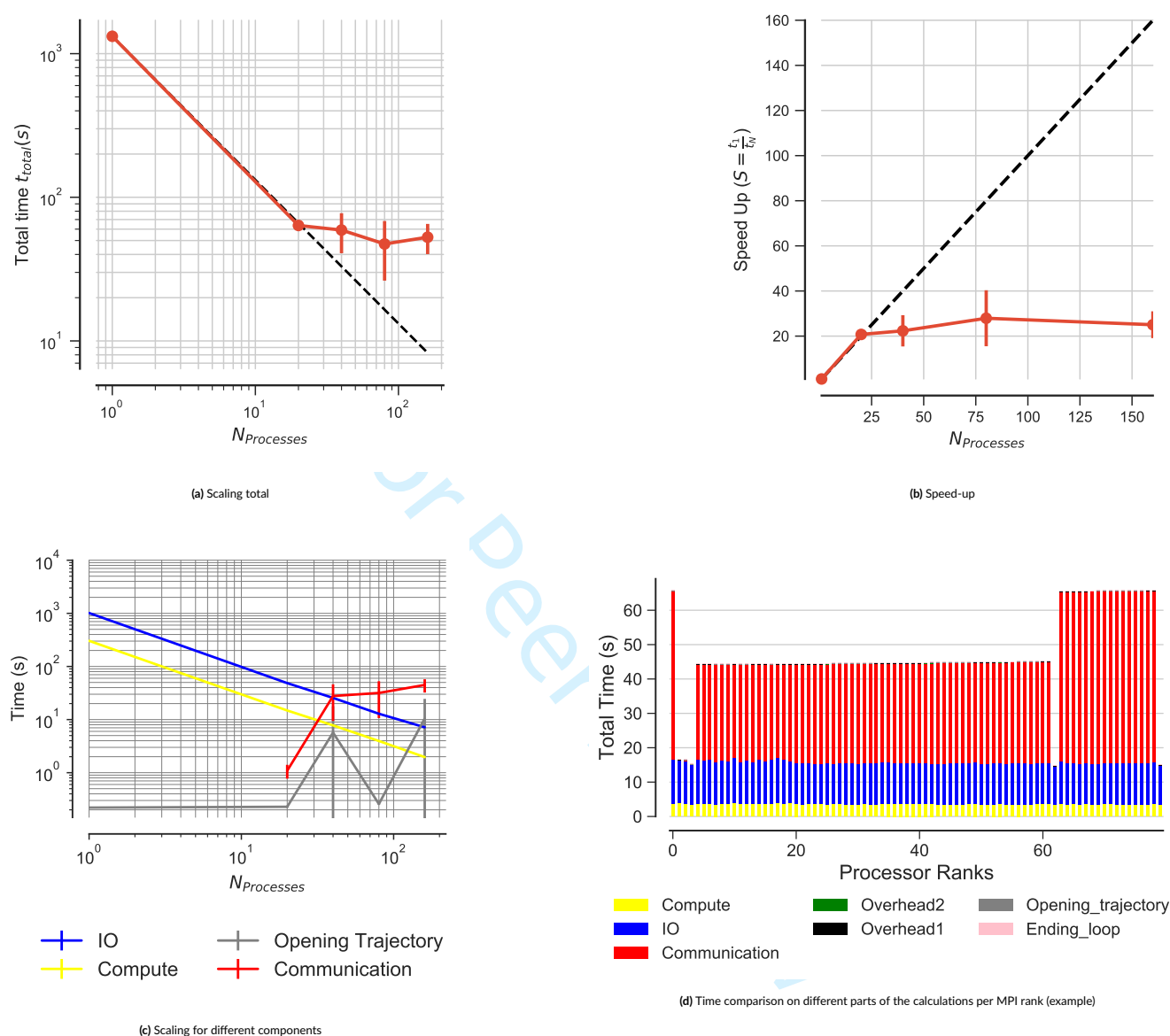
## A ADDITIONAL DATA

Figure A1 shows performance of the RMSD task on *PSC Bridges*.



**FIGURE A1** *PSC Bridges*: Performance of the RMSD task. Results are communicated back to rank 0. Five independent repeats were performed to collect statistics. (a-c) The error bars show standard deviation with respect to the mean. In serial, there is no communication and hence no data point is shown for  $N = 1$  in (c). (d) Compute  $t_{\text{comp}}$ , read I/O  $t_{\text{I/O}}$ , communication  $t_{\text{comm}}$ , ending the for loop  $t_{\text{end\_loop}}$ , opening the trajectory  $t_{\text{opening\_trajectory}}$ , and overheads  $t_{\text{overhead1}}$ ,  $t_{\text{overhead2}}$  per MPI rank; see Table 3 for definitions. These are data from one run of the five repeats. MPI ranks 0 and 70 are stragglers.

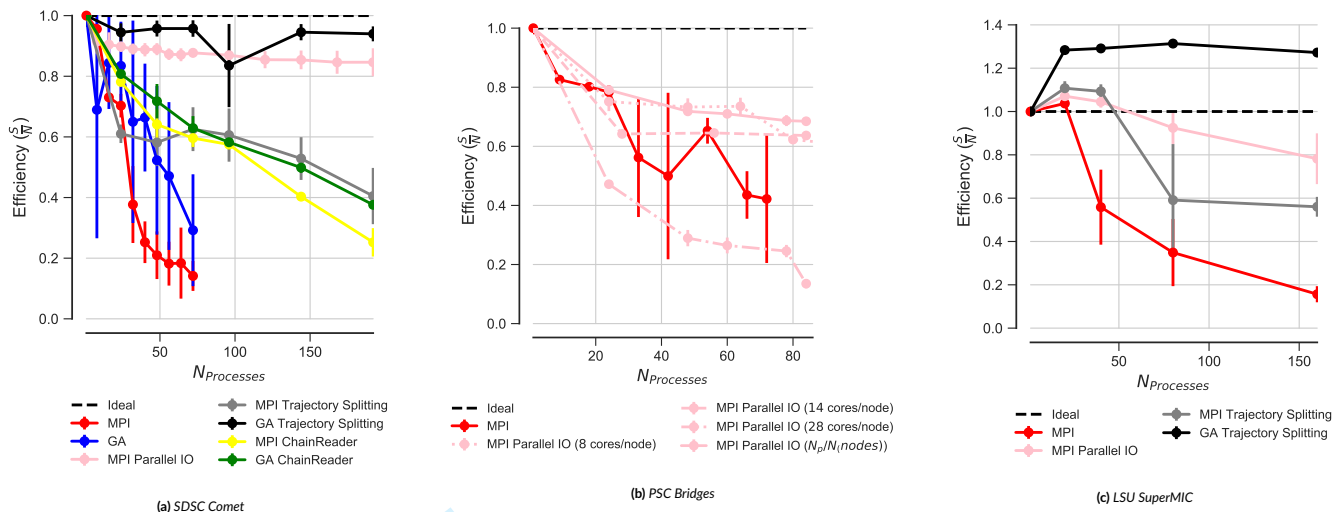


Figure A2 shows performance of the RMSD task on *LSU SuperMIC*.

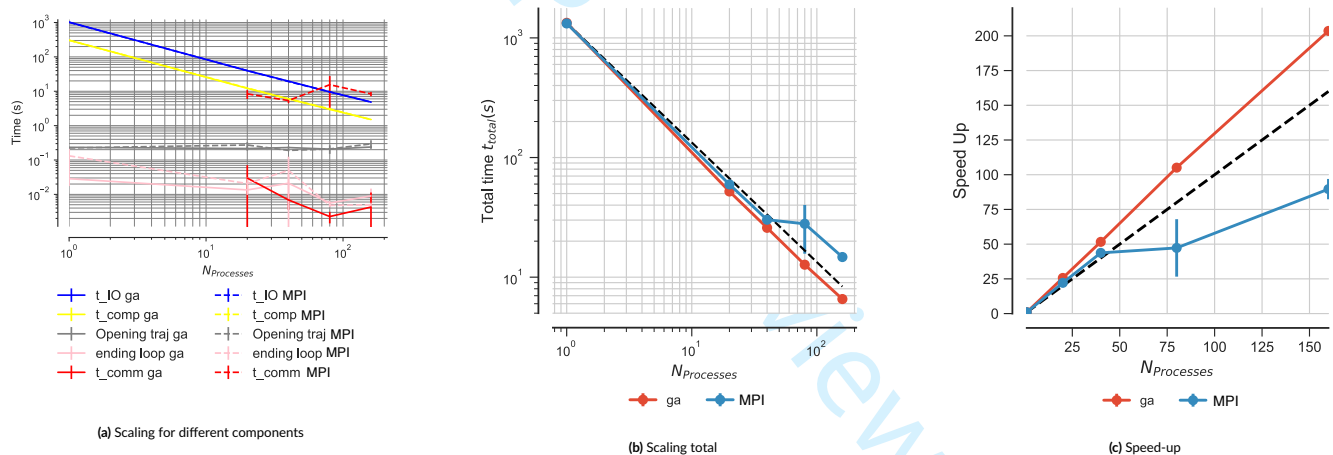
**FIGURE A2** *LSU SuperMIC*: Performance of the RMSD task with MPI. Results are communicated back to rank 0. Five independent repeats were performed to collect statistics. (a-c) The error bars show standard deviation with respect to mean. In serial, there is no communication and hence the data points for  $N = 1$  are not shown in (c). (d) Compute  $t_{\text{comp}}$ , read I/O  $t_{\text{I/O}}$ , communication  $t_{\text{comm}}$ , ending the for loop  $t_{\text{end\_loop}}$ , opening the trajectory  $t_{\text{opening\_trajectory}}$ , and overheads  $t_{\text{overhead1}}$ ,  $t_{\text{overhead2}}$  per MPI rank; see Table 3 for definitions. These are data from one run of the five repeats.

Figure A3 shows comparison of the parallel efficiency of the RMSD task between different test cases on *SDSC Comet*, *PSC Bridges*, and *LSU SuperMIC*.

Figure A4 shows how RMSD task scales with the increase in the number of cores when the trajectories are split using *Global Arrays* for communication compared to using MPI for communications on *LSU SuperMIC*.



**FIGURE A3** Comparison of the parallel efficiency between different test cases on (a) *SDSC Comet* (data for “MPI Parallel IO” are only shown up to 192 cores for better comparison across different scenarios, see Fig. 6b for equivalent scaling data up to 384 cores), (b) *PSC Bridges*, and (c) *LSU SuperMIC*. Five repeats were performed to collect statistics and error bars show standard deviation with respect to mean.



**FIGURE A4** *LSU SuperMIC*: Comparison of the performance of the RMSD task with subfiling and using either MPI (“MPI”) or *Global Arrays* (“ga”) for communication. For *Global Arrays*, all ranks update the global array (`ga_put()`) and rank 0 accesses the whole RMSD array through the global memory address (`ga_get()`). Five repeats were performed to collect statistics. (a) Compute and I/O scaling versus number of processes. (b) Total time scaling versus number of processes. (c) Speed-up. (a-c) The error bars show standard deviation with respect to mean.

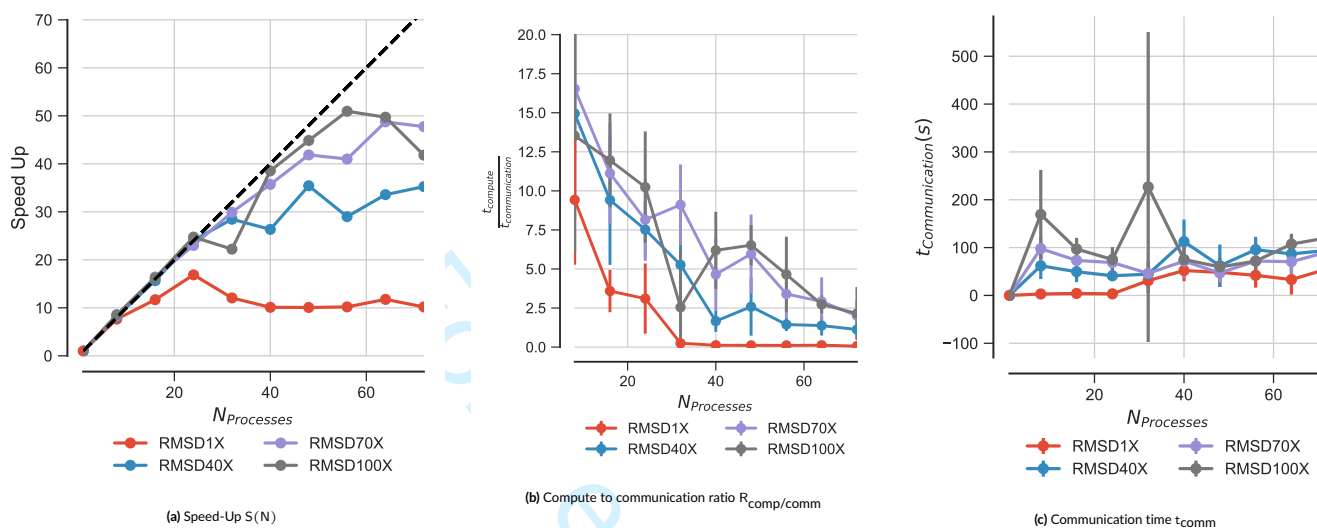
## B EFFECT OF $R_{comp/comm}$ ON PERFORMANCE

In Section 6.3, we improved scaling limitations due to read I/O by splitting the trajectory, but scaling remained far from ideal when MPI communication was used; the use of *Global Arrays* lead to better scaling (see Section 6.4) because the effective communication cost was reduced. Although we were not able to identify the reason for the better performance of *Global Arrays* (it still uses MPI as a communicator), the results motivated an analysis in terms of the communication costs. In addition to the compute to I/O ratio  $R_{comp/IO}$  discussed in Section 6.2 we defined another performance parameter called the compute to communication ratio  $R_{comp/comm}$  (Eq. 8).

We analyzed the data for variable workloads (see Section 6.2) in terms of the  $R_{comp/comm}$  ratio. The performance clearly depended on the  $R_{comp/comm}$  ratio (Figure B5). Performance improved with increasing  $R_{comp/comm}$  ratios (Figure B5b and B5a) even if the communication time was larger (Figure B5c). Although we still observed stragglers due to communication at larger  $R_{comp/comm}$  ratios (70× RMSD and 100× RMSD), their effect on performance remained modest because the overall performance was dominated by the compute load. Evidently, as long as overall performance

is dominated by a component such as compute that scales well, then performance problems with components such as communication will be masked and overall acceptable performance can still be achieved (Figures B5a and B5b).

Communication was usually not problematic within one node because of the shared memory environment. For less than 24 processes, i.e., a single compute node on *SDSC Comet*, the scaling was good and  $R_{\text{comp/comm}} \gg 1$  for all RMSD loads (Figures B5a and B5b). However, beyond a single compute node ( $> 24$  cores), scaling appeared to improve with increasing  $R_{\text{comp/comm}}$  ratio while the communication overhead decreased in importance (Figures B5a and B5b).



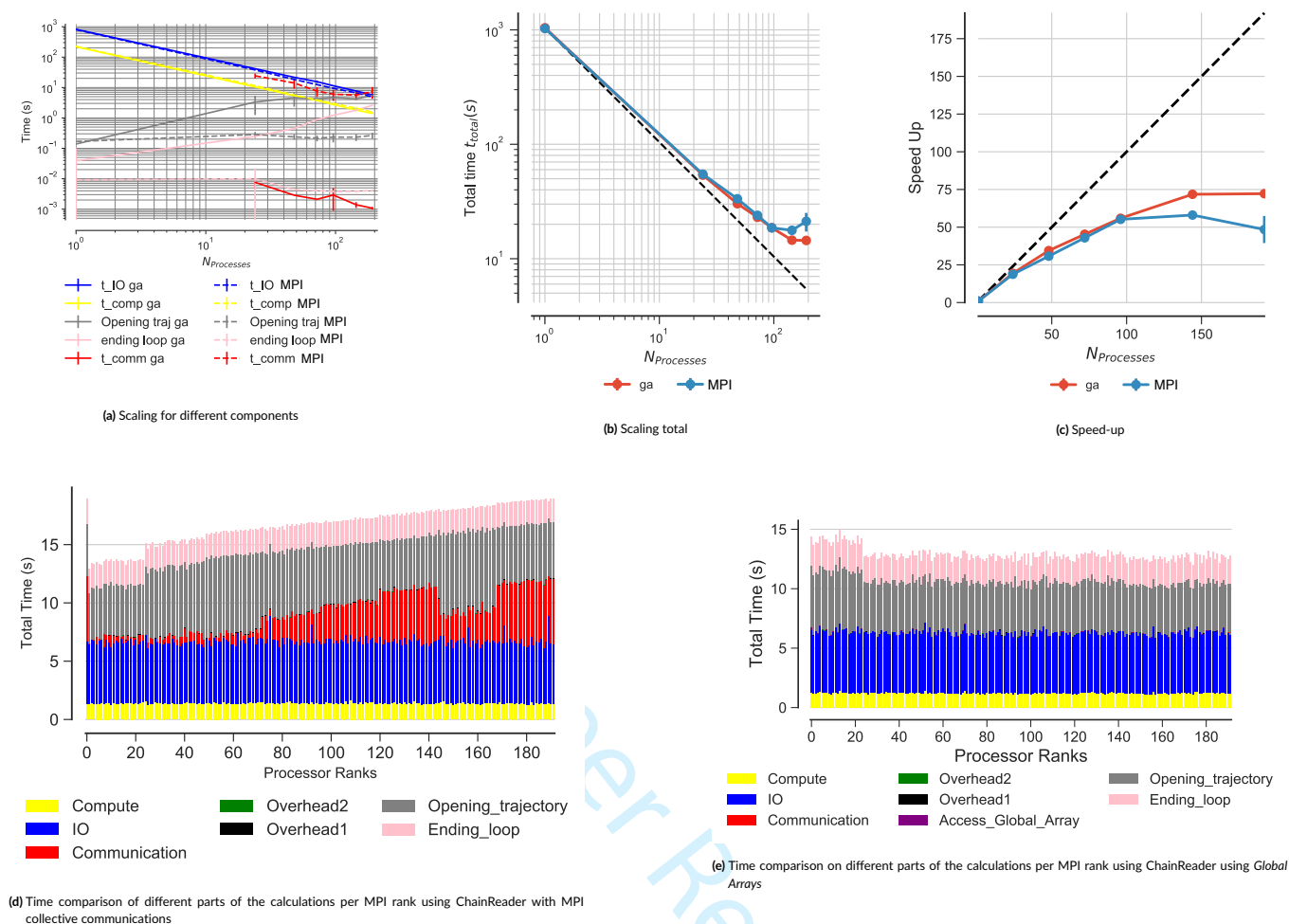
**FIGURE B5** Effect of the ratio of compute to communication time  $R_{\text{comp/comm}}$  on scaling performance on *SDSC Comet*. (a) Scaling for different computational workloads. (Same as Figure 3a.) (b) Change in  $R_{\text{comp/comm}}$  with the number of processes  $N$  for different workloads. (c) Comparison of communication time for different RMSD workloads. Five repeats were performed to collect statistics and error bars show standard deviation with respect to mean.

## C PERFORMANCE OF THE CHAINREADER FOR SPLIT TRAJECTORIES

In section 6.3.1 we showed how subfiling (splitting the trajectories) would help to overcome I/O limitations and improve scaling. However, the number of trajectories may not necessarily be equal to the number of processes. For example, trajectories from MD simulations on supercomputers are often kept in small segments in individual files that need to be concatenated later to form a trajectory that can be analyzed with common tools. Such segments might be useful for subfiling but making sure that the number of processes is equal to the number of trajectory files will not always be feasible. *MDAnalysis* can transparently represent multiple trajectories as one virtual trajectory using the *ChainReader*. This feature is convenient for serial analysis when trajectories are maintained as segments. In the current implementation of *ChainReader*, each process opens all the trajectory segment files but I/O will only happen from a specific block of the trajectory specific to that process only.

We wanted to test if the *ChainReader* would benefit from the gains measured for the subfiling approach. Specifically, we measured if the MPI-parallelized RMSD task (with  $N_p$  ranks) would benefit if the trajectory was split into  $N_{\text{seg}} = N_p$  trajectory segments, corresponding to an ideal scenario.

In order to perform our experiments we had to work around an issue with the XTC format reader in *MDAnalysis* that was related to the XTC random-access functionality that the *MDAnalysis.coordinates.XTC.XTCReader* class provides: The Gromacs XTC format<sup>67,68</sup> is a lossy-compression, XDR-based file format that was never designed for random access and the compressed format itself does not support fast random seeking. The *XTCReader* stores persistent offsets for trajectory frames to disk<sup>18</sup> in order to enable efficient access to random frames. These offsets will be generated automatically the first time a trajectory is opened and the offsets are stored in hidden `*.xtc_offsets.npz` files. The advantage of these persistent offset files is that after opening the trajectory for the first time, opening the same file will be very fast, and random access is immediately available. However, stored offsets can get out of sync with the trajectory they refer to. To prevent the use of stale offset data, trajectory file data (number of atoms, size of the file and last modification time) are also stored for validation. If any of these parameters change the



**FIGURE C6** Comparison on the performance of the MDAnalysis ChainReader for the RMSD task on *SDSC Comet* when the trajectories are split; for the communication step either collective MPI ("MPI") or *Global Arrays* ("ga") was used. In case of *Global Arrays*, all ranks update the global array (`ga_put()`) and rank 0 accesses the whole RMSD array through the global memory address (`ga_get()`). Five repeats were performed to collect statistics. (a) Compute and I/O scaling versus number of processes. (b) Total time scaling versus number of processes. (c) Speed-up. (a-c) The error bars show standard deviation with respect to the mean. (d-e) Compute  $t_{\text{comp}}$ , read I/O  $t_{\text{I/O}}$ , communication  $t_{\text{comm}}$ , access to the whole global array by rank 0  $t_{\text{Access\_Global\_Array}}$ , ending the for loop  $t_{\text{end\_loop}}$ , opening the trajectory  $t_{\text{opening\_trajectory}}$ , and overheads  $t_{\text{overhead1}}$ ,  $t_{\text{overhead2}}$  per MPI rank. (See Table 3 for the definitions.)

offsets are recalculated. If the XTC changes but the offset file is not updated then the offset file can be detected as invalid. With ChainReader in parallel, each process opens all the trajectories because each process builds its own `MDAnalysis.Universe` data structure. If an invalid offset file is detected (perhaps because of wrong file modification timestamps across nodes), several processes might want to recalculate these parameters and rebuild the offset file, which can lead to a race condition. In order to avoid the race condition, we removed this check from MDAnalysis for the purpose of the measurements presented here, but this comes at the price of not checking the validity of the offset files at all; future versions of MDAnalysis may lift this limitation. We obtained the results for the ChainReader with this modified version of MDAnalysis that eliminates the race condition by assuming that XTC index files are always valid.

Figure C6 shows the results for performance of the ChainReader for the RMSD task using either collective MPI or Global Arrays (GA) for the communication step. With GA good strong scaling performance was observable up to 144 cores (Figure C6c); without GA, strong scaling plateaued for more than 92 cores. In both cases, strong scaling performance was worse than what was achieved when each MPI process was assigned its own trajectory segment as described in Section 6.3.1. The strong scaling performance did not suffer because of the computation and the read I/O because both  $t_{\text{comp}}$  and  $t_{\text{I/O}}$  showed excellent strong scaling up to 196 cores (Figure C6a). Instead the time for ending the for loop  $t_{\text{end\_loop}}$ , which

includes the time for closing the trajectory file, and opening the trajectory  $t_{\text{opening\_trajectory}}$  appeared to be the scaling bottleneck. These results differed from the subfiling results (section 6.3.1) where neither  $t_{\text{end\_loop}}$  nor  $t_{\text{opening\_trajectory}}$  limited scaling (Figures 5d and 5e).

Although we did not further investigate the deeper cause for the reduced scaling performance of the ChainReader, we speculate that the primary problem is related to each MPI rank having to open all trajectory files in their ChainReader instance even though they will only read from a small subset. For  $N_p$  ranks and  $N_{\text{seg}}$  file segments, in total,  $N_p N_{\text{seg}}$  file opening/closing operations have to be performed. Each server that is part of a Lustre file system can only handle a limited number of I/O requests (read, write, stat, open, close, etc.) per second. A large number of such requests, from one or more users and one or more jobs, can lead to contention for storage resources. For  $N_p = N_{\text{seg}} = 100$ , the Lustre file system has to perform 10,000 of these operations almost simultaneously, which might degrade performance.

For Peer Review