

# MPI-parallel Molecular Dynamics Trajectory Analysis with the H5MD Format in the MDAnalysis Python Package

Edis Jakupovic<sup>‡</sup>, Oliver Beckstein<sup>‡\*</sup>

**Abstract**—Fill here

**Index Terms**—Molecular Dynamics Simulations, High Performance Computing, Python, MDAnalysis, HDF5, H5MD, MPI I/O

## Introduction

As HPC resources continue to increase, the size of molecular dynamics (MD) simulation files are now commonly terabytes in size, making serial analysis of these trajectory files impractical. Parallel analysis is a necessity for the efficient use of both HPC resources and a scientist's time. MDAnalysis is a widely used Python library that can read and write over 20 popular MD file formats while providing the same user-friendly interface [MADWB11], [GLB<sup>+</sup>16]. Previous work that focused on developing a task-based approach to parallel analysis found that an IO bound task only scaled to 12 cores due to a file IO bottleneck [SFMLIP<sup>+</sup>19]. Our previous feasibility study suggested that parallel reading via MPI-IO and HDF5 can lead to good scaling although it only used a reduced size custom HDF5 trajectory and did not provide a usable implementation of a true MD trajectory reader [KPF<sup>+</sup>20].

H5MD, or "HDF5 for molecular data", is an HDF5-based file format that is used to store MD simulation data, such as particle coordinates, box dimensions, and thermodynamic observables [dBCH14]. HDF5 is a structured, binary file format that organizes data into 2 objects: groups and datasets, which follows a hierarchical, tree-like structure, where groups represent nodes of the tree, and datasets represent the leaves [Col14]. The HDF5 library can be built on top of a message passing interface (MPI) implementation so that a file can be accessed in parallel on a parallel filesystem such as Lustre or BeeGFS. We implemented a parallel MPI-IO capable HDF5-based file format trajectory reader into MDAnalysis, H5MDReader, that adheres to H5MD specifications. H5MDReader interfaces with h5py, a high level Python package that provides a Pythonic interface to the HDF5 format such that accessing a file in parallel is as easy as passing

a keyword argument into h5py.File, and all of parallel disk access occurs under the hood.

We benchmarked H5MDReader's parallel reading capabilities with MDAnalysis on three HPC clusters: ASU Agave, SDSC Comet, and PSC Bridges. The benchmark consisted of a simple split-apply-combine scheme of an IO-bound task that split a 90k frame (113GB) trajectory into  $n$  chunks for  $n$  processes, where each process a task on their chunk of data, and then gathered the results back to the root process. For the computational task, we computed the time series root mean squared distance (RMSD) of the positions of the alpha carbons in the protein to their initial coordinates at the first frame of the trajectory. The RMSD calculation is not only a very common task performed to analyze the dynamics of the structure of a protein, but more importantly is a very fast computation that is heavily bounded by how quickly data can be read from the file. Therefore it provided an excellent analysis candidate to test the I/O capabilities of H5MDReader.

Across the three HPC clusters tested, the benchmarks were done on both a BeeGFS and Lustre parallel filesystem which is highly suited for multi-node MPI parallelization. We tested various algorithmic optimizations for our benchmark, including altering the stripe count, loading only necessary coordinate information with numpy.Masked\_arrays, and front loading all IO by loading the entire trajectory into memory prior to the RMSD calculation.

## BRIEFLY DISCUSS RESULTS AND CHUNKING

## Methods

We implemented a simple split-apply-combine parallelization algorithm that divides the number of frames in the trajectory evenly among all available processes. Each process receives a unique start and stop for which to iterate through their section of the trajectory and compute the RMSD at each frame. The data files used in our benchmark included a topology file YiiP\_system.pdb and a trajectory file YiiP\_system\_9ns\_center100x.h5md with 90100 frames. The trajectory file was converted on the fly with MDAnalysis to several different file formats. Table 1 gives all of these formats with how they are identified in this paper as well as their corresponding file size.

In order to obtain detailed timing information we instrumented code as follows:

```
1 class timeit(object):
2     def __enter__(self):
3         self._start_time = time.time()
```

<sup>‡</sup> Arizona State University

\* Corresponding author: obeckste@asu.edu

name	format	file size (GB)
H5MD_default	H5MD	113
H5MD_chunked	H5MD	113
H5MD_contiguous	H5MD	113
H5MD_gzipx1	H5MD	77
H5MD_gzipx9	H5MD	75
XTC	XTC	35
DCD	DCD	113
TRR	TRR	113

TABLE 1

```

4         return self
5
6     def __exit__(self, exc_type, exc_val, exc_tb):
7         end_time = time.time()
8         self.elapsed = end_time - self._start_time
9         # always propagate exceptions forward
10        return False

```

The `timeit` class was used as a context manager to record how long our benchmark spent on particular lines of code. Below, we give example code of how each benchmark was performed:

```

1 import MDAnalysis as mda
2 from MDAnalysis.analysis.rms import rmsd
3 from mpi4py import MPI
4 import numpy as np
5
6 def benchmark(topology, trajectory):
7     with timeit() as init_top:
8         u = mda.Universe(topology)
9     with timeit() as init_traj:
10        u.load_new(trajectory,
11                  driver="mpio",
12                  comm=MPI.COMM_WORLD)
13    t_init_top = init_top.elapsed
14    t_init_traj = init_traj.elapsed
15    CA = u.select_atoms("protein and name CA")
16    x_ref = CA.positions.copy()
17
18    total_io = 0
19    total_rmsd = 0
20    rmsd_array = np.empty(bsize, dtype=float)
21    for i, frame in enumerate(range(start, stop)):
22        with timeit() as io:
23            ts = u.trajectory[frame]
24            total_io += io.elapsed
25        with timeit() as rms:
26            rmsd_array[i] = rmsd(CA.positions,
27                                x_ref,
28                                superposition=True)
29        total_rmsd += rms.elapsed
30
31    with timeit() as wait_time:
32        comm.Barrier()
33    t_wait = wait_time.elapsed
34
35    with timeit() as comm_gather:
36        rmsd_buffer = None
37        if rank == 0:
38            rmsd_buffer = np.empty(n_frames,
39                                  dtype=float)
40        comm.Gatherv(sendbuf=rmsd_array,
41                    recvbuf=(rmsd_buffer,
42                             sendcounts),
43                    root=0)
44    t_comm_gather = comm_gather.elapsed

```

The time  $t_{\text{init\_top}}$  records the time it takes to load a universe from the topology file.  $t_{\text{init\_traj}}$  records the time it takes to open the trajectory file. The HDF5 file is opened with the `mpio` driver and the `MPI.COMM_WORLD` communicator to ensure the file is accessed in parallel via MPI I/O. It's important to separate

the topology and trajectory initialization times, as the topology file is not opened in parallel and represents a fixed cost each process must pay to open the file.  $t_{\text{I/O}}$  represents the time it takes to read the data for each frame into the corresponding `MDAnalysis.Universe.trajectory.ts` attribute. MDAnalysis reads data from MD trajectory files one frame, or "snapshot" at a time. Each time the `u.trajectory[frame]` is iterated through, MDAnalysis reads the file and fills in numpy arrays corresponding to that timestep. Each MPI process runs an identical copy of the script, but receives a unique `start` and `stop` variable such that the entire file is read in parallel.  $t_{\text{compute}}$  gives the total RMSD computation time.  $t_{\text{wait}}$  records how long each process waits before the results are gathered with `comm.Gather()`. Gathering the results is done collectively by MPI, which means all processes must finish their iteration blocks before the results can be returned. Therefore, it's important to measure  $t_{\text{wait}}$  as it represents the existence of "straggling" processes. If one process takes substantially longer than the others to finish its iteration block, all processes are slowed down.  $t_{\text{comm\_gather}}$  measures the time MPI spends communicating the results from each process back to the root process.

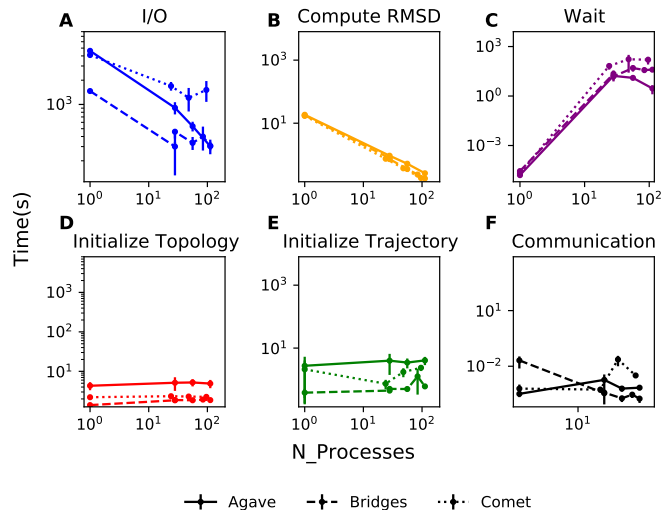
We applied this benchmark scheme to H5MD test files on Agave, Bridges, and Comet. We also tested 3 algorithmic optimizations: Lustre file striping, loading the entire trajectory into memory, and using Masked Arrays to only load the alpha carbon coordinates required for the RMSD calculation. For striping, we ran the benchmark on Bridges and Comet with a file stripe count of 48 and 96. For the into memory optimization, we used `MDAnalysis.Universe.transfer_to_memory()` to read the entire file in one go and pass all file I/O to the HDF5 library. For the masked array optimization, we allowed `u.load_new()` to take a list or array of atom indices as an argument, `sub`, so that the `MDAnalysis.Universe.trajectory.ts` arrays are instead initialized as `ma.masked_array`'s and only the indices corresponding to `sub` are read from the file.

Performance was quantified by measuring the I/O timing returned from the benchmarks, and strong scaling was assessed by calculating the speedup  $S(N) = t_1/t_N$  and the efficiency  $E(N) = S(N)/N$ .

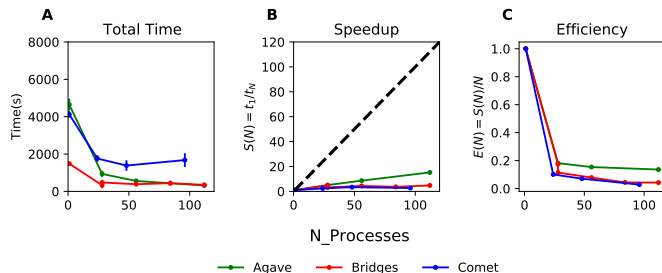
## Results and Discussion

### Default Benchmark Results

We first ran benchmarks with the simplest parallelization scheme of splitting the frames of the trajectory evenly among all participating processes. The H5MD file involved in the benchmarks was written with `pyh5md`, a python library that can easily read and write H5MD files. The datasets in the data file were chunked automatically by `h5py`'s auto-chunking algorithm. File I/O remains the largest contributor to the total benchmark time, as shown by Figure 1 (A). Figure 1 (B, D-F) also shows the initialization times, computation, MPI communication are negligible with regards to overall the overall analysis time.  $t_{\text{wait}}$ , however, becomes increasingly relevant as the number of processes increases (Figure 1 C), indicating a growing variance in the iteration block time across all processes. Although the total benchmark time continues to decrease as the number of processes increases, the total speedup observed is 15x (Figure 2 A,B).



**Fig. 1:** Benchmark timings breakdown for the ASU Agave, PSC Bridges, and SDSC Comet HPC clusters. The benchmark was run on up to 4 full nodes on each HPC, where  $N_{\text{processes}}$  was 1, 28, 56, and 112 for Agave and Bridges, and 1, 24, 48, and 96 on Comet. The `H5MD_default` file was used in the benchmark, where the trajectory was split in  $N$  chunks for each corresponding  $N$  process benchmark. Points represent the mean over three repeats with the standard deviation shown as error bars.



**Fig. 2:** Strong scaling performance of the RMSD analysis task of the `H5MD_default` data file on Agave, Bridges, and Comet.  $N_{\text{Processes}}$  ranged from 1 core, to 4 full nodes on each HPC, and the number of trajectory blocks was equal to the number of processes involved.

#### Effects of Algorithmic Optimizations on File I/O

#### Conclusions

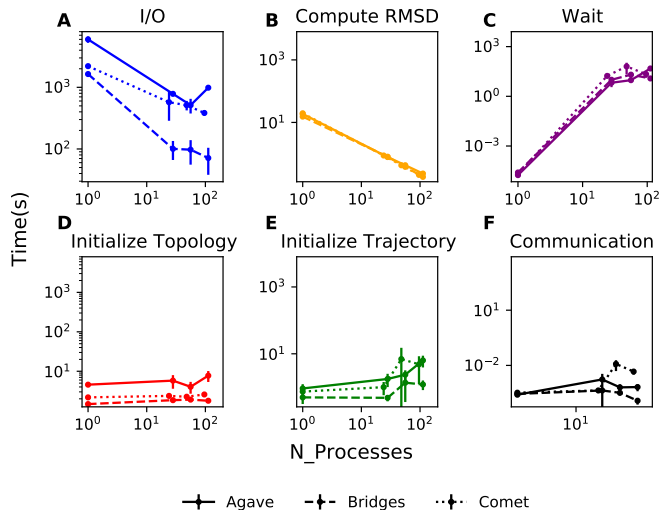
TODO

#### Acknowledgments

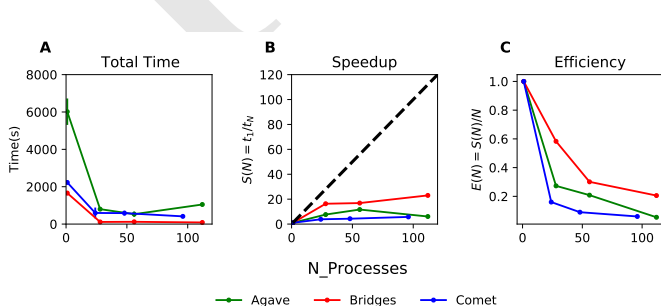
Funding was provided by the National Science Foundation for a REU supplement to award ACI1443054. The SDSC Comet computer at the San Diego Supercomputer Center was used under allocation TG-MCB130177. The authors acknowledge Research Computing at Arizona State University for providing HPC resources that have contributed to the research results reported within this paper. We would like to acknowledge Gil Speyer and Jason Yalim from the Research Computing Core Facilities at Arizona State University for advice and consultation.

#### REFERENCES

[Col14] Andrew Collette. Python and hdf5. In Meghan Blanchette and Rachel Roumeliotis, editors, *Python and HDF5*. O'Reilly

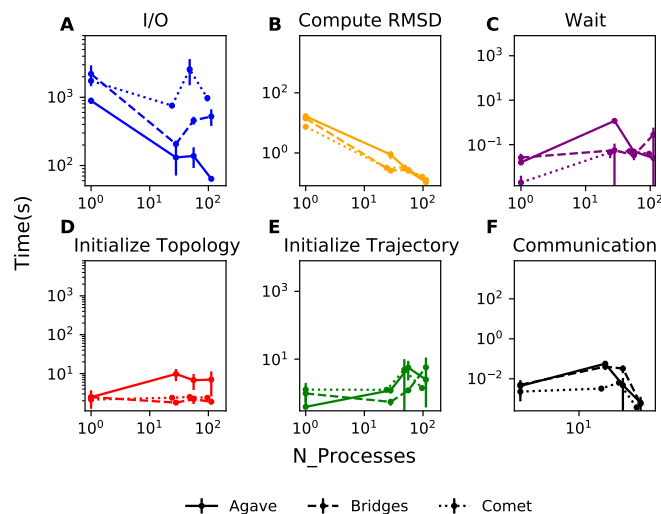


**Fig. 3:** Benchmark timings breakdown for the ASU Agave, PSC Bridges, and SDSC Comet HPC clusters for the `masked_array` optimization technique. The benchmark was run on up to 4 full nodes on each HPC, where  $N_{\text{processes}}$  was 1, 28, 56, and 112 for Agave and Bridges, and 1, 24, 48, and 96 on Comet. The `H5MD_default` file was used in the benchmark, where the trajectory was split in  $N$  chunks for each corresponding  $N$  process benchmark. Points represent the mean over three repeats with the standard deviation shown as error bars.

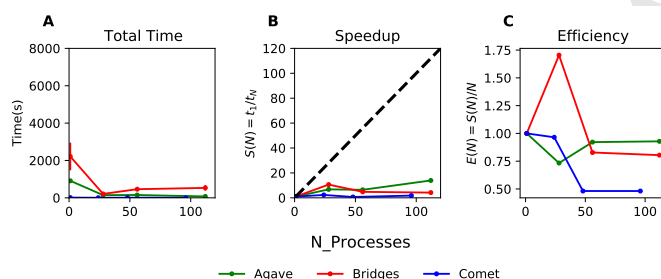


**Fig. 4:** Strong scaling performance of the RMSD analysis task with the `masked_array` optimization technique. The benchmark used the `H5MD_default` data file on Agave, Bridges, and Comet.  $N_{\text{Processes}}$  ranged from 1 core, to 4 full nodes on each HPC, and the number of trajectory blocks was equal to the number of processes involved.

- Media, Inc., 1005 Gravenstein Highway North, Sebastopol, CA 95472., 2014.
- [dBCH14] Pierre de Buyl, Peter H. Colberg, and Felix Höfling. H5MD: A structured, efficient, and portable file format for molecular data. *Computer Physics Communications*, 185(6):1546 – 1553, 2014. doi:10.1016/j.cpc.2014.01.018.
- [GLB<sup>+</sup>16] Richard J. Gowers, Max Linke, Jonathan Barnoud, Tyler J. E. Reddy, Manuel N. Melo, Sean L. Seyler, David L. Dotson, Jan Domański, Sébastien Buchoux, Ian M. Kenney, and Oliver Beckstein. MDAnalysis: A Python package for the rapid analysis of molecular dynamics simulations. In Sebastian Benthall and Scott Rostrup, editors, *Proceedings of the 15th Python in Science Conference*, pages 98–105, Austin, TX, 2016. SciPy. URL: <https://www.mdanalysis.org>, doi:10.25080/Majora-629e541a-00e.
- [KPF<sup>+</sup>20] Mahzad Khoshlessan, Ioannis Paraskevatos, Geoffrey C. Fox, Shantenu Jha, and Oliver Beckstein. Parallel performance of molecular dynamics trajectory analysis. *Concurrency and Computation: Practice and Experience*, 32:e5789, 2020. doi:10.1002/cpe.5789.
- [MADWB11] Naveen Michaud-Agrawal, Elizabeth Jane Denning,



**Fig. 5:** Benchmark timings breakdown for the ASU Agave, PSC Bridges, and SDSC Comet HPC clusters for the loading-into-memory optimization technique. The benchmark was run on up to 4 full nodes on each HPC, where  $N_{\text{processes}}$  was 1, 28, 56, and 112 for Agave and Bridges, and 1, 24, 48, and 96 on Comet. The `H5MD_default` file was used in the benchmark, where the trajectory was split in  $N$  chunks for each corresponding  $N$  process benchmark. Points represent the mean over three repeats with the standard deviation shown as error bars.



**Fig. 6:** Strong scaling performance of the RMSD analysis task with the `masked_array` optimization technique. The benchmark used the `H5MD_default` data file on Agave, Bridges, and Comet.  $N_{\text{Processes}}$  ranged from 1 core, to 4 full nodes on each HPC, and the number of trajectory blocks was equal to the number of processes involved.

Thomas B. Woolf, and Oliver Beckstein. MDAnalysis: A toolkit for the analysis of molecular dynamics simulations. *J Comp Chem*, 32:2319–2327, 2011. [doi:10.1002/jcc.21787](https://doi.org/10.1002/jcc.21787).

[SFMLIP<sup>+</sup>19] Shujie Fan, Max Linke, Ioannis Paraskevatos, Richard J. Gowers, Michael Gecht, and Oliver Beckstein. PMDA - Parallel Molecular Dynamics Analysis. In Chris Calloway, David Lippa, Dillon Niederhut, and David Shupe, editors, *Proceedings of the 18th Python in Science Conference*, pages 134 – 142, Austin, TX, 2019. SciPy. [doi:10.25080/Majora-7ddc1dd1-013](https://doi.org/10.25080/Majora-7ddc1dd1-013).