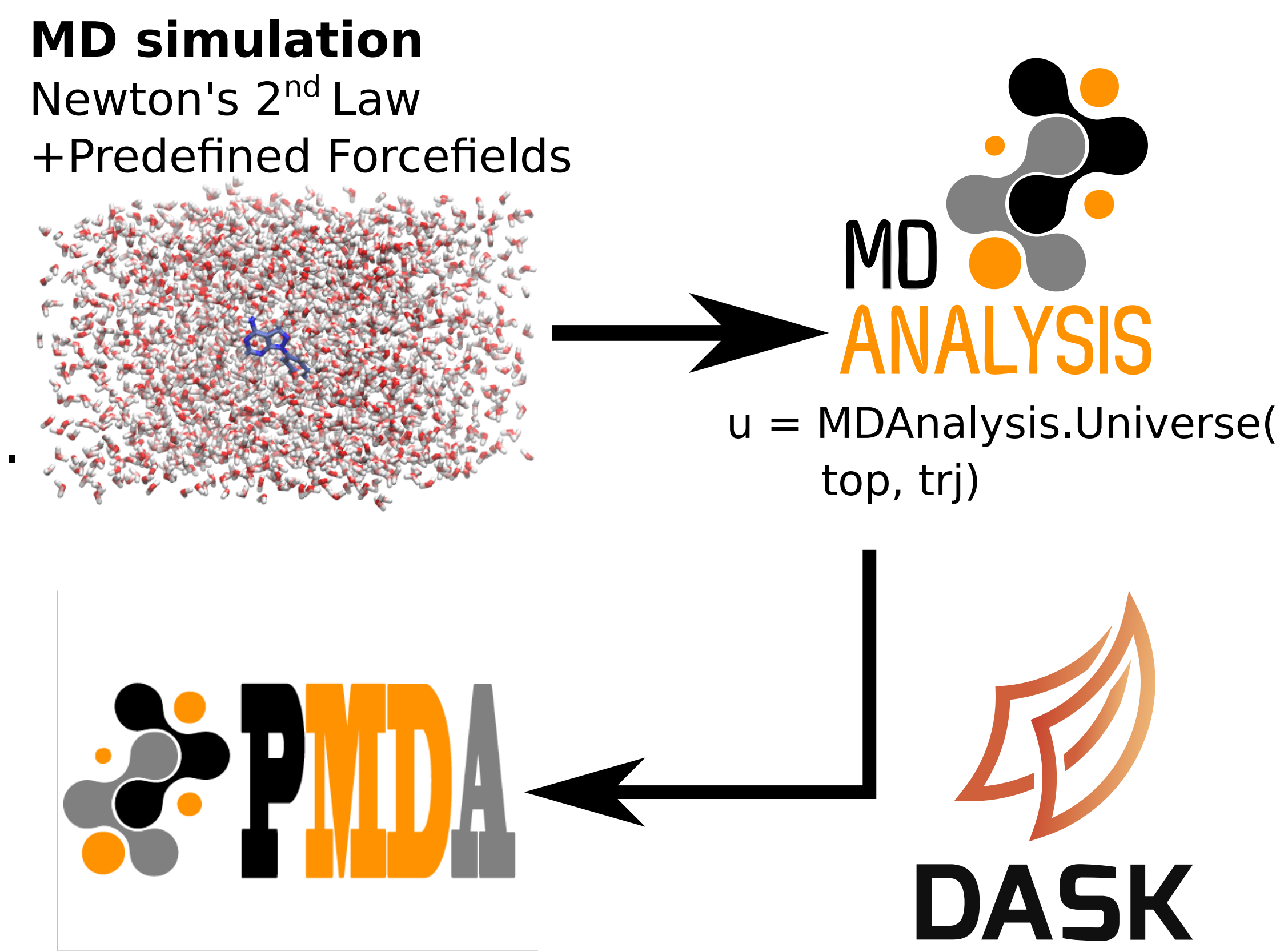


## Introduction

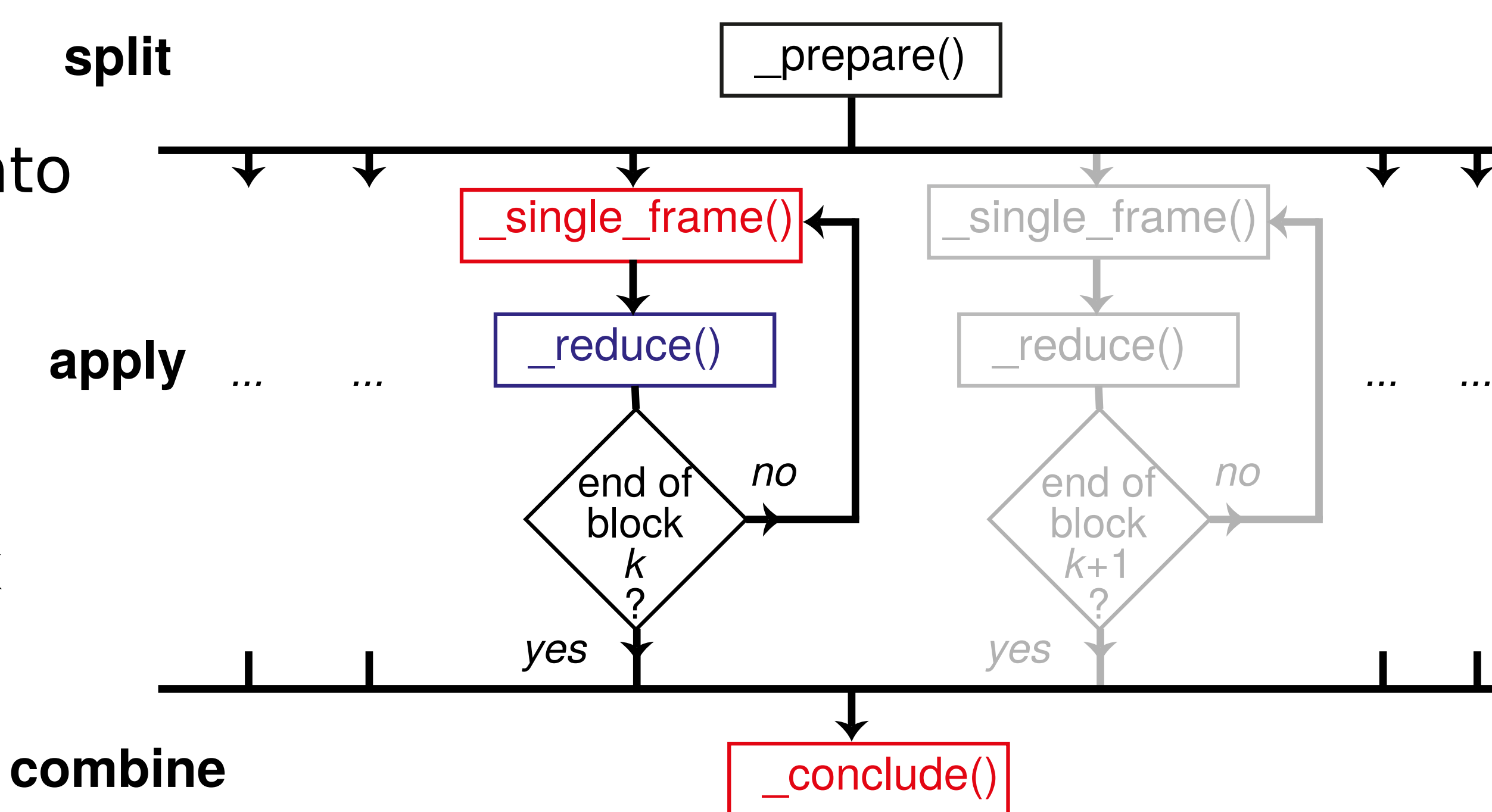
PMDA is a Python library that builds upon MDAAnalysis[1] and Dask[2] to provide parallel analysis algorithms for molecule dynamics (MD) simulations. At the core of PMDA is the idea that a common interface makes it easy to create code that can be easily parallelized.



## Methods

### split-apply-combine approach[3]:

The trajectory is split into blocks, analysis is performed separately and in parallel on each block ("apply"), then results from each block are gathered and combined.



## Acknowledgments

We would like to thank reviewer Cyrus Harrison for the idea to plot the fractional time spent on different stages of the program. This work was supported by the National Science Foundation under grant numbers ACI-1443054 and used the Extreme Science and Engineering Discovery Environment (XSEDE) supported by National Science Foundation grant number ACI-1548562. The SDSC Comet computer at the San Diego Supercomputer Center was used under allocation TG-MCB130177. Max Linke was supported by NumFOCUS under a small development grant.

## References

- [1] Gowers, Richard J.; Linke, Max; Barnoud, Jonathan; Reddy, Tyler J. E.; Melo, Manuel N.; Seyler, Sean L.; Dotson, David L.; Domański, Jan; Buchoux, Sébastien; Kenney, Ian M.; and Beckstein, Oliver. MDAAnalysis: A Python package for the rapid analysis of molecular dynamics simulations. In S. Benthall and S. Rostrup, editors, Proceedings of the 15th Python in Science Conference, pages 102 – 109, Austin, TX, 2016. SciPy. URL: <https://www.mdanalysis.org/>.
- [2] Dask Development Team. Dask: Library for dynamic task scheduling, 2016, URL <https://dask.org>.
- [3] Hadley Wickham. The split-apply-combine strategy for data analysis. Journal of Statistical Software, 40(1), 2011. doi:10.18637/jss.v040.i01.

## Using PMDA

PMDA is released under the GNU General Public License, version 2

Source code is available in the public GitHub repository

<https://github.com/MDAnalysis/pmda/>.

### Installation

#### Install with conda:

```
conda config --add channels conda-forge
conda install pmda
```

#### Install with pip:

```
pip install --upgrade pmda
```

#### Install from source:

```
git clone git@github.com:MDAnalysis/pmda.git
cd pmda
python setup.py install
```

### User-defined Analysis

**pmda.custom.AnalysisFromFunction():** **pmda.parallel.ParallelAnalysisBase:**

```
import MDAAnalysis as mda
```

```
u = mda.Universe(top, traj)
protein = u.select_atoms('protein')
```

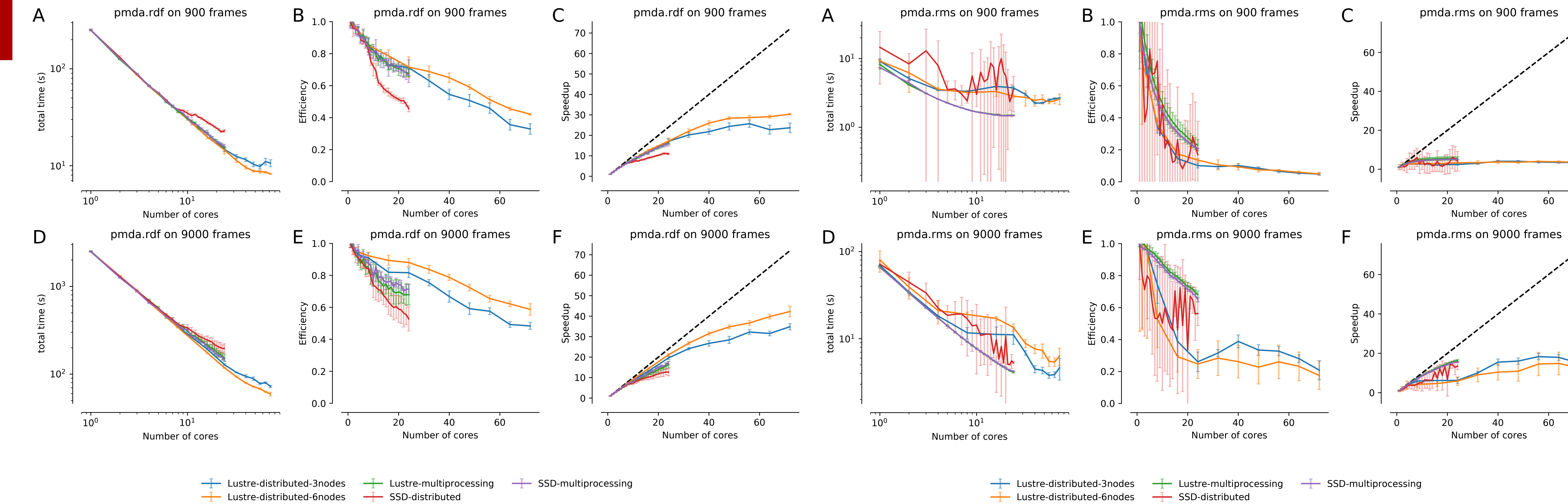
```
def rgyr(ag):
    return (ag.universe.trajectory.time,
            ag.radius_of_gyration())
```

```
import pmda.custom
parallel_rgyr =
    pmda.custom.AnalysisFromFunction(
        rgyr, u, protein)
parallel_rgyr.run(n_jobs=4, n_blocks=4)
print(parallel_rgyr.results)
```

### Pre-defined Analysis

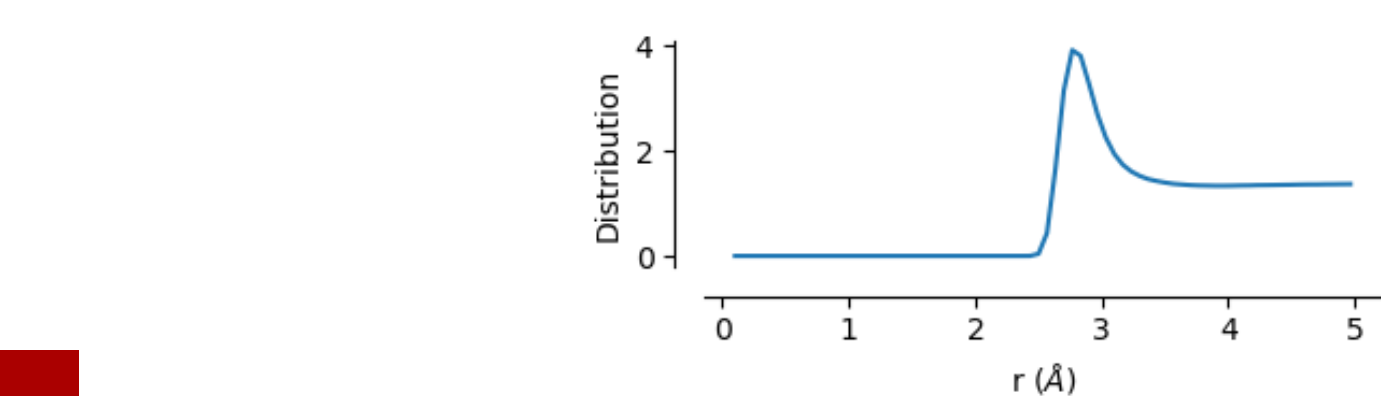
```
import MDAAnalysis as mda
from pmda import rms

u = mda.Universe(top, trj)
ca = u.select_atoms('name CA')
u.trajectory[0]
ref = u.select_atoms('name CA')
rmsd = rms.RMSD(ca, ref)
rmsd.run(n_jobs=4, n_blocks=4)
print(rmsd.rmsd)
```



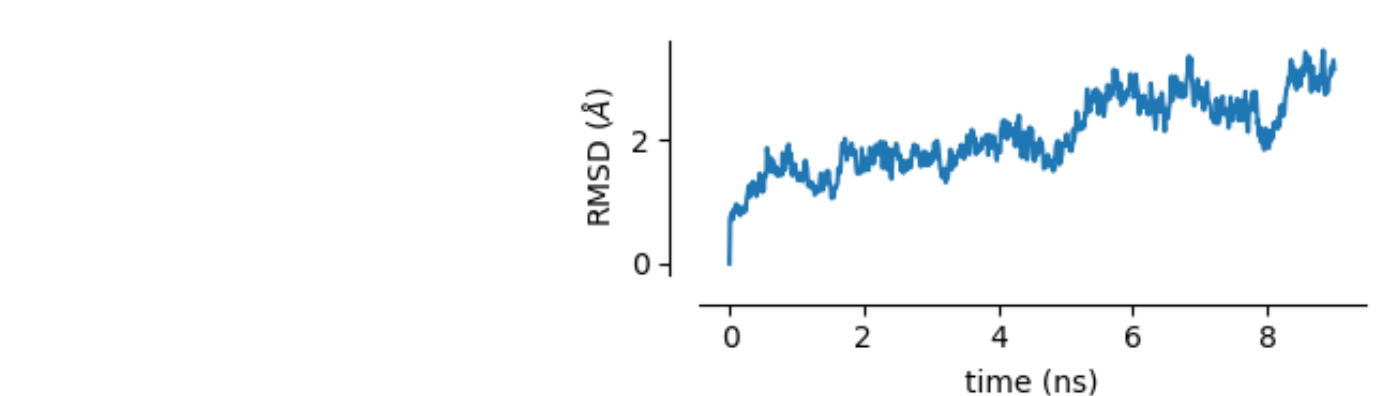
### RDF

Water oxygen-oxygen radial distribution function for all 24,239 oxygen atoms in the water molecules.



### RMSD

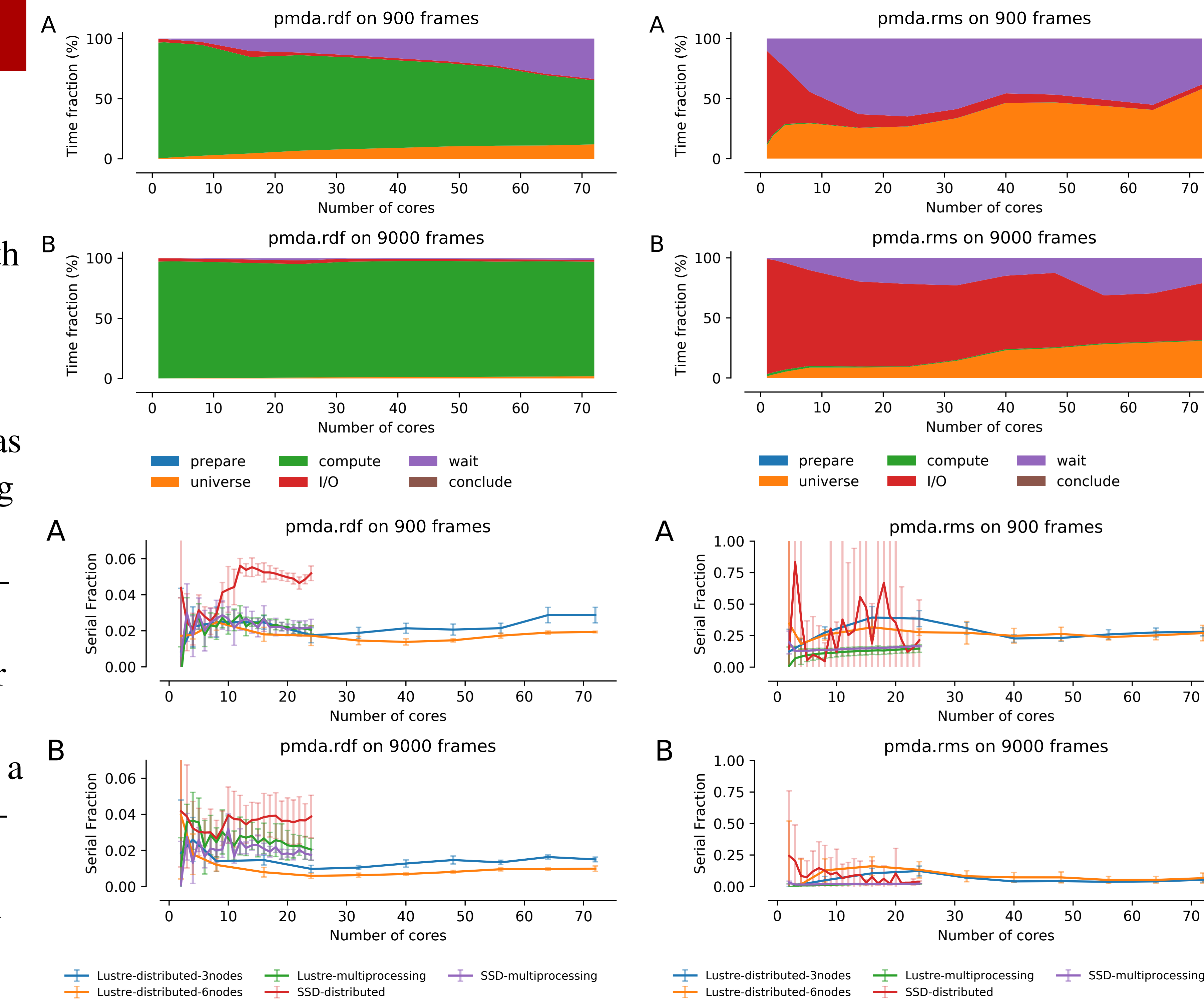
Time series of root mean square distance after optimum superposition (RMSD) of all 564 Cα atoms of a protein.



## Conclusion

The PMDA Python package provides a framework to parallelize analysis of MD trajectories with a simple split-apply-combine approach by combining Dask with MDAAnalysis.

We showed that performance depends on the type of analysis that is being performed. Compute-intensive tasks such as the RDF calculation can show good strong scaling up to about a hundred cores on a typical supercomputer and a sizable speed-up that approached 40. Such performance should make this an attractive solution for many users. For tasks such as the RMSD calculation, whose speed-up is limited by a considerable serial fraction, a single multi-core workstation seems sufficient to achieve speed-ups on the order of 10 and HPC resources would not be useful.



## Performance Evaluation

We tested two tasks(RDF and RMSD) with different combinations of Dask schedulers with different means to read the trajectory data as shown in the Table.

configuration label	file storage	scheduler	max nodes	max processes
Lustrre-distributed-3nodes	Lustrre	<i>distributed</i>	3	72
Lustrre-distributed-6nodes	Lustrre	<i>distributed</i>	6	72
Lustrre-multiprocessing	Lustrre	<i>multiprocessing</i>	1	24
SSD-distributed	SSD	<i>distributed</i>	1	24
SSD-multiprocessing	SSD	<i>multiprocessing</i>	1	24