

Technical Report: SPIDAL Summer REU 2021 Upgrading MDPOW and Adding Analysis Functionality

Alia Lescoulie

Department of Chemistry and Biochemistry, College of Science and Mathematics,
California Polytechnic State University, San Luis Obispo, CA, 93407

November 1, 2021

Abstract

MDPOW is an open source Python package for calculating water-octanol and water-cyclohexane partition coefficients using the molecular dynamics package GROMACS. During the SPIDAL REU, MDPOW was updated to support Python 3. A collection of classes were constructed, to simplify development of simulation analyses. These objects were designed for efficient and organized storage of a group of molecular dynamics systems, and provide a basic framework for developing analyses. All code is included in the latest release of MDPOW at <https://github.com/becksteinlab/mdpow>.

1 Introduction

Molecular dynamics (MD) is one method for computing water-octanol partition coefficients (P_{ow}) of drug-like molecules. This is accomplished through calculating the free energy of solvation from the simulation. With free energy values from water and octanol the partition coefficient can be obtained using $P_{ow} = (\Delta G_w - \Delta G_o)/(RT) \log(e)$ where $R = 8.31446261815 \times 10^{-3}$ kJ/mol · K is the universal gas constant and e is Euler's number [5].

A number of methods exist for obtaining free energies of systems in MD simulations. The two discussed here are Bennett Acceptance Ratio (BAR) and Thermodynamic Integration (TI) [9]. The free energy change ΔA from the canonical partition function Q is defined by the following [3]:

$$\Delta A = -\frac{\ln Q_1/Q_0}{\beta} \quad (1)$$

BAR overcomes the challenge of calculating partition functions by substituting the ratio of the partition functions Q_1/Q_0 for the ratio of probabilities (M_0/M_1) for a trial move which maintains the same configurational space but switches potential energies from the reference system U_0 to the unknown system U_1 and vice-versa [2]. The potential energy can be calculated with system positions.

Thermodynamic integration (TI) obtains changes in free energy by defining a parameter λ such that at $\lambda = 0$ $U = U_0$ and at $\lambda = 1$ $U = U_1$, where U_1 is the final state and U_0 is the initial state, and taking the partial derivative of A with respect to λ and integrating from 0 to 1 [3].

Both BAR and TI use λ as a parameter for the extent of solute-solvent interactions [3]. This requires running a MD simulation for each interaction at their respective λ values [7]. The process of initializing that number of simulations requires either the creation of a shell script running each or a tedious process of navigating directories, copying files, and running simulations. MDPOW makes this process easier with the user simply defining their settings and supplying a coordinate file, then running a few Python scripts.

MDPOW is not without its limitations. Before release 0.7.0 it was incompatible with anything newer than Python 2.7 and had few features for analysis beyond running BAR and TI calculations. Both of these features were added over the course of the SPIDAL REU.

Incompatibility with Python 3 created a number of challenges, in particular compatibility with newer versions of dependencies such as MDAnalysis, which dropped Python 2 compatibility [8, 6]. This made new development more difficult, and meant that users had to create an environment for MDPOW separate from the environment used for other portions of a project. For version 0.7.0,

retaining Python 2 cross-compatibility was necessary due to the long-term nature of scientific projects. This meant that many of the new features added in Python 3 such as f-strings, and in Python 3.6 and later, type annotations, were unavailable during the update process [5, 1].

An additional shortcoming of MDPOW addressed during the REU was its lack of a simple methodology for analyzing sets of alchemical free energy simulations. MDPOW, when running free energy perturbation (FEP) simulations, establishes a directory like the example in Fig 1, with simulation files sorted by solvent, interaction and λ -value. Running a simulation to get the partition coefficient can result in forty or more individual systems. For example using MDAnalysis Universe objects to load systems, one would have to write code to individually load each simulation each in its own directory. That's not to mention actually running analyses on that set of systems. Requiring users to each develop their own analysis tools for managing a larger number of systems would result in lost time and repeated work. Providing an analysis framework makes the process more user friendly.

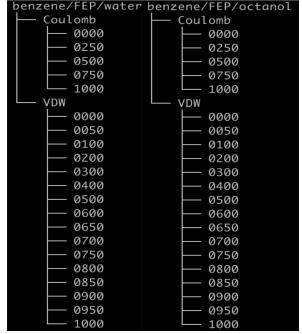


Figure 1: Example MDPOW FEP directory structure

2 Methods

2.1 Modernizing MDPOW

The differences between Python 2 and Python 3 presented compatibility issues in a number of key areas primarily: outdated tests, pickle compatibility issues, and changes to syntax. The first issue was ensuring that all tests ran in both Python versions so that cross compatibility could be validated. MDPOW testing is done in pytest, which eliminated yield tests in version 4.0 requiring their replacement. This was done using parameterization tests where the test parameter and a list of inputs are provided to a function. In the process of running tests, it became apparent reading pickle files was an incompatibility between versions. Pickle files, which save an instance of a Python object as a file, are used throughout MDPOW, particularly in testing. Pickles encoded in Python 2.7, when read normally using the six library in Python 3 raised byte errors. To solve this, pickle loading code was replaced with the block below throughout MDPOW.

```

1 if sys.version_info.major >= 3:
2     with open(gsolv, 'rb') as f:
3         G = pickle.load(f, encoding='latin1')
4 elif sys.version_info.major == 2:
5     G = pickle.load(gsolv.open())

```

By checking for the version, pickles can be loaded in byte mode when code is run in Python 3, ensuring compatibility. Patches like this, where separate blocks of code are executed based on version were avoided, but in this case was unavoidable. In MDPOW version 0.8.0 Python 2 support was dropped, and this style of patch has been replaced. The tests and pickle file represented the bulk of compatibility issues. Aside from that the remaining issues were simply the result of syntax differences present in older version of Python.

2.2 Implementation of Analysis Features

Addressing the aforementioned challenge of managing several simulations the `Ensemble` object was developed. The `Ensemble` object is a collection of MDAnalysis Universes stored in a Python dictionary. Python dictionaries, a type of hash table, have O(1) average search efficiency. This means that the time required to find an item is independent of the size of the table [4]. This search efficiency, and the organization created by the key-value pair structure of dictionaries, were the reasons for their selection when developing the `Ensemble` object. An `Ensemble` object is able to load the trajectory files present in a MDPOW molecules directory, automatically handling the process of directory navigation and system loading. The `Ensemble` object also replicates some of the functionality present in MDAnalysis Universe objects, namely it allows users to select atoms from the systems it stores. This returns a `EnsembleAtomGroup` object.

The `EnsembleAtomGroup` is to the MDAnalysis AtomGroup as the `Ensemble` is to the MDAnalysis Universe. It is initialized by the select atoms command of Ensemble. It, like the `Ensemble`

attempts to extend the functionality of MDAnalysis objects to collections of those objects. It has class methods to select atoms, return positions, and return the original `Ensemble` object. The code example below shows how an `Ensemble` can be initialized from a simulation directory, and how an `EnsembleAtomGroup` is created from an `Ensemble`.

```

1 benzene_dir = os.path.join('mdpow', 'tests', 'testing_resources',
2                             'states', 'benzene')
3 Benzene = mdpow.ensemble.Ensemble(dirname=benzene_dir)
4 DihedralGroup = Benzene.select_atoms('name C1 or name C2 or name C3 or name C4')
```

With the `Ensemble` and `EnsembleAtomGroup` objects an analysis framework analogous to `AnalysisBase` in MDAnalysis was developed for running calculations on collection of systems generated by MDPOW. Setting up an analysis for a MDPOW simulation requires simply sub classing `EnsembleAnalysis` as in the next listing. To implement an analysis, an `EnsembleAnalysisBase` subclass with a few basic methods must be defined. For an analysis to work it needs an `__init__` to accept parameters, a `_prepare` to create a place to store data, a `_single_universe` or `_single_frame` that runs an analysis generating data. A more detailed explanation of designing an analysis can be found in section 3.1.1.

```

1 class ExampleAnalysis(mdpow.ensemble.EnsembleAnalysis):
2     def __init__(self, ExampleEnsemble):
3         super(ExampleAnalysis, self).__init__(ExampleEnsemble.ensemble())
4         self._ensemble = ExampleEnsemble
5
6     def _prepare_ensemble(self):
7         self._cols = ['solvent', 'interaction', 'lambda', 'time', 'result']
8         self._result_dict = {x: [] for x in self._cols}
9         self.results = pd.DataFrame(columns=self._cols)
10
11    def _single_frame(self):
12        result = analysis_function(self._ensemble[self._key])
13        res_list = [self._key[0], self._key[1], self._key[2], self._ts.time, result]
14        for i in range(len(res_list)):
15            self._result_dict[self._col[i]].append(self.res_list[i])
16
17    def _conclude_ensemble(self):
18        for k in self._result_dict:
19            self.results[k] = self._result_dict[k]
```

`EnsembleAnalysis` was developed to be similar to `AnalysisBase` in MDAnalysis. When developing an analysis several aspects must be defined; an `__init__`, which accepts parameters in this case lines 2-4. Among the parameters must be at least one `Ensemble`, which must be passed back to the parent class using the `super` function from Python as seen in line 3. Next `_prepare_ensemble` and `_prepare_universe` can be used to establish data structures used in the overall analysis and the individual systems. In lines 6-9 on the example a results dictionary and `DataFrame` are set up for organizing results by solvent, interaction, lambda, and time. Next are the methods responsible for generating results from the MD simulations, `_single_universe` which runs on each system in the `Ensemble`, and `_single_frame` which runs on each frame of each `Universe`. `_conclude_universe` and `_conclude_ensemble` are run after each `Universe` and the `Ensemble` respectively. When an `EnsembleAnalysis` based object is run, each system stored in the provided `Ensemble` is iterated over with `_single_universe` ran at that point, and `_single_frame` run on each frame of the `Universe`. In this case `_single_frame` consists of an example `example_function` which returns a result on line 12. In the proceeding lines that result is stored in `_result_dict` which was created in `_prepare_ensemble`. All of the information from each frame is recorded to ensure that the data is "tidy" [10]. Afterwards there is the option to define `_conclude_universe` which runs each time a new system is completed allowing for processing to occur after each system. Finally after iteration is complete `_conclude_ensemble` is run. In this case it saves the contents of the results dictionary into the results `DataFrame`, this is done at the end because dictionaries are far simpler to add data to within the program, but `DataFrame` is more convenient for saving and plotting data. In the final `DataFrame` each index has all the information, 'solvent', 'interaction', 'lambda', 'time', and 'result'. This allows for easy organization of data by multiple parameters for example returning all the results for water and VDW.

3 Results and Discussion

With the establishment of a framework for analyzing collections of simulations, such as those generated by FEP calculations, the process of developing analyses for MDPOW simulations is simplified. These are developed using objects created as subclasses of `EnsembleAnalysis`.

3.1 Implemented EnsembleAnalysis Methods

Using the `EnsembleAnalysis` framework two new methods were developed as part of the new MDPOW analysis submodule. The first to be discussed is `SolvationAnalysis` which quantifies solvent molecules within the given cutoff distances. It is discussed in section 3.2.1 specifically focusing on how to design an `EnsembleAnalysis` based class. The second is `DihedralAnalysis` in section 3.2.2 which discusses considerations for efficient design of `EnsembleAnalysis` classes.

3.1.1 Solvation Shell: Example for the EnsembleAnalysis Framework

One of the analyses developed with the `EnsembleAnalysis` framework, `SolvationAnalysis`, which returns the number of solvents within the given distances, can also serve as an example of how to develop an analysis more generally.

`SolvationAnalysis` quantifies the number of solvent atom within a given distance of the solute. It accomplishes this thought calculating the distances between and each molecule using the `capped_distance` function from MDAnalysis to get the indices for pairs which fall within the given cutoff [8, 6].

Building `EnsembleAnalysis` using inheritance, a concept from object oriented programming, simplifies the process of programming an analysis. Inheriting from a superclass gives a subclass the same methods and attributes, meaning that some methods common to all `EnsembleAnalysis` objects like `run` can be retained; other methods the user can override to run their own analysis code. To develop `SolvationAnalysis`, only three methods were overridden in the subclass.

```
1  class SolvationAnalysis(EnsembleAnalysis):
2      def __init__(self, solute: EnsembleAtomGroup, solvent: EnsembleAtomGroup,
3                  distances: List[float]):
4          self.check_groups_from_common_ensemble([solute, solvent])
5          super(SolvationAnalysis, self).__init__(solute.ensemble)
6          self._solute = solute
7          self._solvent = solvent
8          self._dists = distances
9
9      def _prepare_ensemble(self):
10         self._col = ['distance', 'solvent', 'interaction',
11                     'lambda', 'time', 'N_solvent']
12         self.results = pd.DataFrame(columns=self._col)
13         self._res_dict = {key: [] for key in self._col}
14
15     def _single_frame(self):
16         solute = self._solute[self._key]
17         solvent = self._solvent[self._key]
18         pairs, distances = capped_distance(solute.positions, solvent.positions,
19                                              max(self._dists), box=self._ts.dimensions)
20         solute_i, solvent_j = np.transpose(pairs)
21         for d in self._dists:
22             close_solv_atoms = solvent[solvent_j[distances < d]]
23             result = [d, self._key[0], self._key[1], self._key[2],
24                       self._ts.time, close_solv_atoms.n_atoms]
25             for i in range(len(self._col)):
26                 self._res_dict[self._col[i]].append(result[i])
27
28     def _conclude_ensemble(self):
29         for k in self._col:
30             self.results[k] = self._res_dict[k]
```

The first method `__init__` accepts arguments needed to generate an instance of the object. In this case the user selected solute and solvent, contained in `EnsembleAtomGroups`, and a list of the distances to be measured. On line 2 the `__init__` ensures that the two `EnsembleAtomGroup` objects originate from the same `Universe` using a built in method of `EnsembleAnalysis` called `check_groups_from_common_ensemble`. This is important as the groups in `_single_frame` must both

exist in the same Universe. Next the `Ensemble` from the aforementioned groups is passed into the parent class with the `super` function in line 4. Finally the items are saved as class attributes (saved in the class, and accessible to other methods). The underscore in front of method and attribute names by Python convention indicates that it is protected.

The next method `_prepare_ensemble` establishes data structures needed in the analysis. The final results are available to the user in a `DataFrame`, but for the sake of efficiency the results are stored in arrays contained in a dictionary with the same keys as the columns in the `DataFrame` as the analysis is run. On line 13 for brevity dictionary interpretation is give each key a blank array. With structures established for storing the results the next methods can focus on obtaining information from the simulation.

With data structures lined up, the actual number of solvent molecules can be counted with `_single_frame`. This is done using the the capped distances function in MDAnalysis on line 18. This computes a distance array for the given group positions and returns distances in that range [8, 6].

Finally the `_res_dict` data is inputted into the results `DataFrame`. This gives data that is neatly organized and easy to generate figures from. The data obtained from `SolvationShell` gives an idea of how solvent molecules are distributed around the solute. This data can also be compared between different λ values in a alchemical free energy simulation, offering incite into the role of different intermolecular interactions in the solvation process.

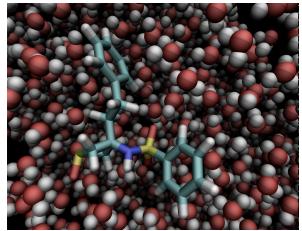


Figure 2: Solute in surrounded by water molecules.

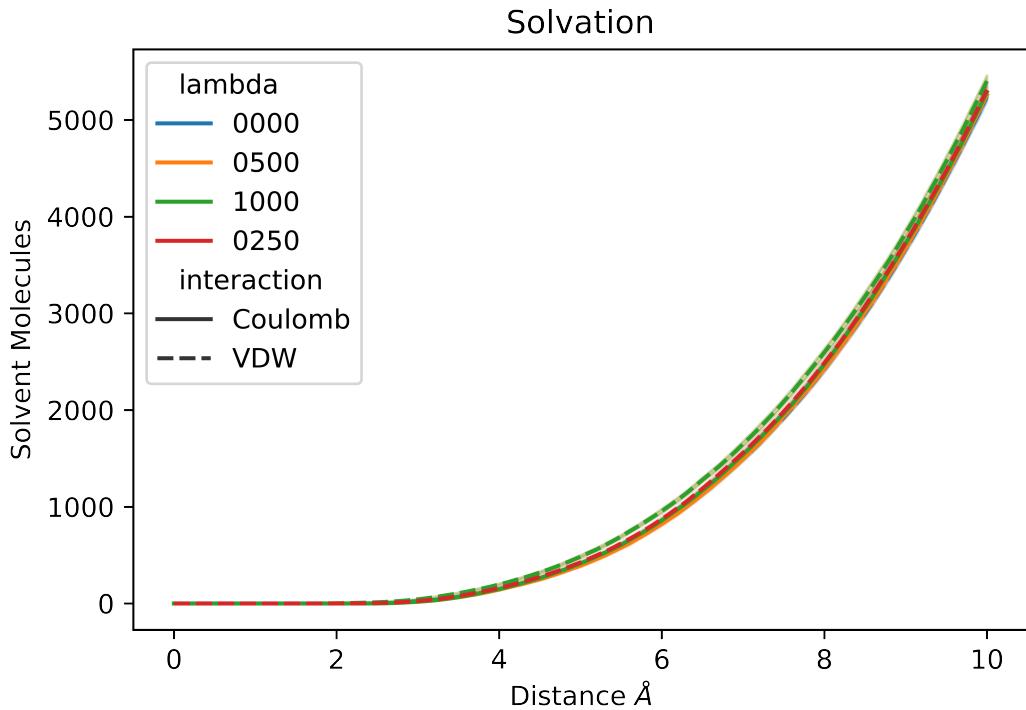


Figure 3: Example solvation shell plot

3.1.2 Dihedral Analysis: Building an EnsembleAnalysis Efficiently

Additionally `DihedralAnalysis` was developed using the `EnsembleAnalysis` framework. `DihedralAnalysis` accepts a list of `EnsembleAtomGroup` objects and returns the dihedral angles over the course of the trajectory. It accomplishes this using the `calc_dihedrals` function from MDAnalysis [8, 6]. This function is written in Cython, which translates python code to C, making it signifiganlty faster than code written in python [Mull2018].

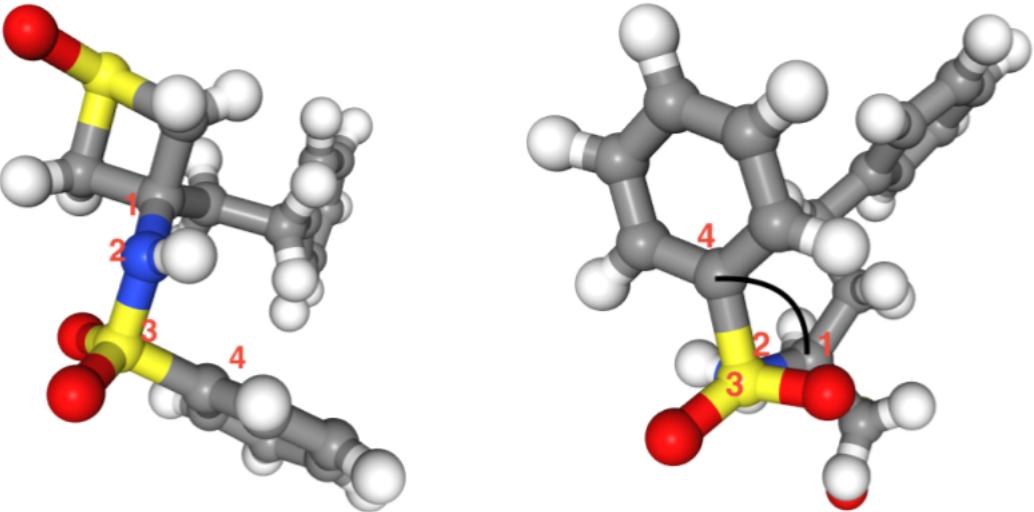


Figure 4: Example dihedral

Dihedral angles are calculated from the positions of four atoms, and describe the geometry of two portions of a molecule between a chemical bond. This gives a better understanding of the spatial relationship between groups in a molecule.

The implementation of `DihedralAnalysis` in MDPOW was based on the MDAnalysis version designed for one `Universe`, accepting a list of `AtomGroup` objects each containing the four atoms of a dihedral group. It then repackages that list of `AtomGroup` objects into four groups each containing one atom from the user given list [Mull2018]. The advantage of this approach is that for each frame of the trajectory `calc_dihedrals`, which can accept arrays, is only run a single time, calculating the list of dihedrals simultaneously. This is more efficiency due to the aforementioned efficiency of code run in C as compared to python.

`DihedralAnalysis` attempted to replicated this approach, but extended to lists of `EnsembleAtomGroup`. Prior to running calculations, the dihedral groups atoms are repacked into new `AtomGroups` each containing one atom from each provided dihedral group. Those `AtomGroups` are then repackaged into four `EnsembleAtomGroups`. That ensures that for, like in the MDAnalysis version, `calc_dihedrals` is only called once per frame rather than iteration over list of dihedral groups and calculating the dihedral angle each time.

When applied to `EnsembleAtomGroups` the added inefficiency of repackaging the dihedral groups outweighed the added efficiency of removing iteration from `_single_frame`. This was found by comparing the execution time of the two methods as additional dihedral groups were added. The iterative method outperformed the repackaging method by a relatively consistent margin each time.

4 Conclusion

The modernization of MDPOW and development of the analysis submodule over the course of the SIPDAL REU increased the usability and utility of the library. The update to Python 3 will ensure compatibility with future versions of dependency libraries, and ensure access to future features of those libraries. Additionally the `Ensemble` objects and `EnsembleAnalysis` framework simplify the development of analytical tools within MDPOW.

Acknowledgments

Funding was provided by the National Science Foundation for a REU supplement to award ACI1443054.

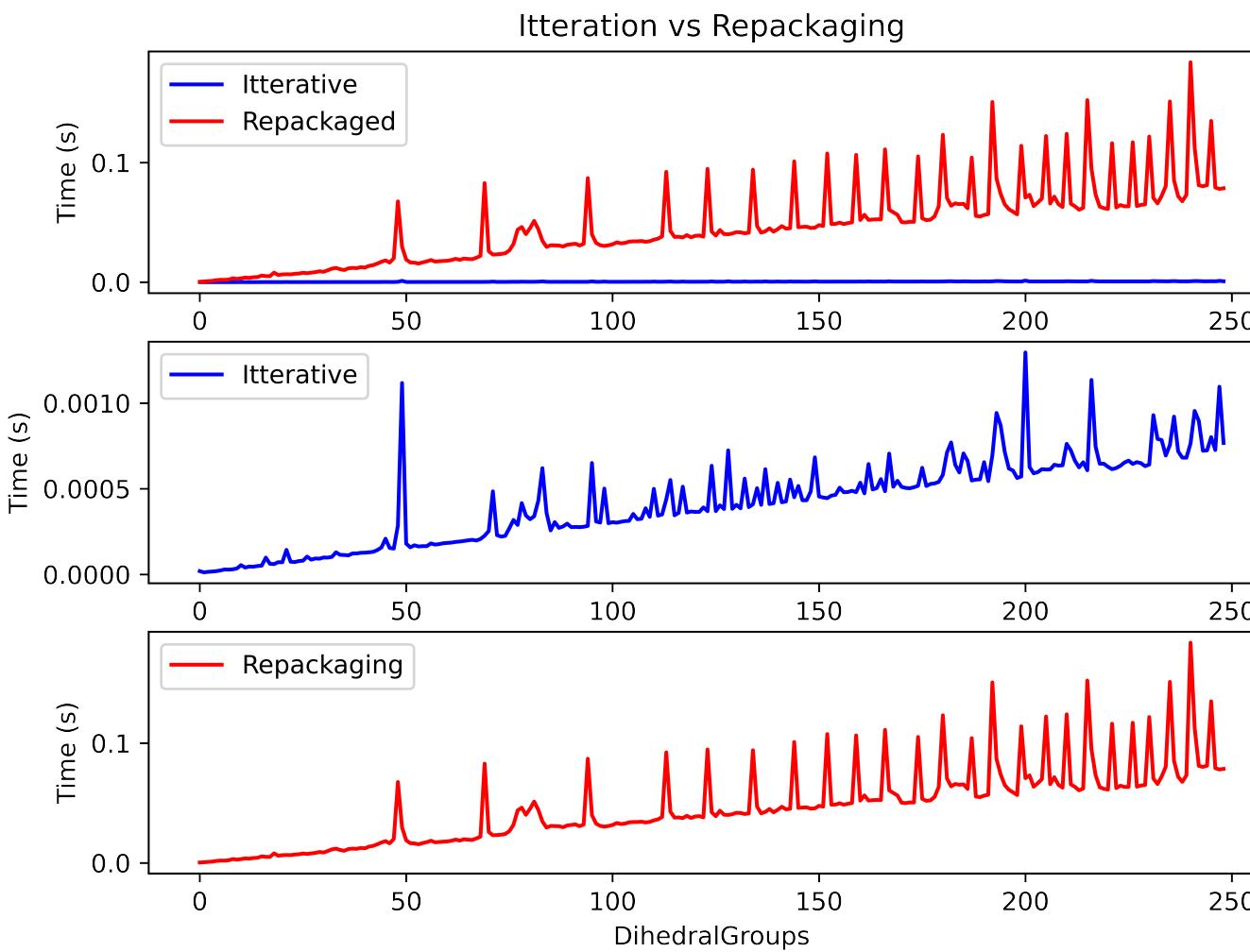


Figure 5: Execution time of Iterating compared to Repackaging

References

- [1] Oliver Beckstein and Bogdan I. Iorga. "Prediction of hydration free energies for aliphatic and aromatic chloro derivatives using molecular dynamics simulations with the OPLS-AA force field". In: *Journal of Computer-Aided Molecular Design* 26.5 (May 2012), pp. 635–645. ISSN: 1573-4951. DOI: 10.1007/s10822-011-9527-9. URL: <https://doi.org/10.1007/s10822-011-9527-9>.
- [2] Charles H Bennett. "Efficient estimation of free energy differences from Monte Carlo data". en. In: *Journal of Computational Physics* 22.2 (Oct. 1976), pp. 245–268. ISSN: 00219991. DOI: 10.1016/0021-9991(76)90078-4. URL: <https://linkinghub.elsevier.com/retrieve/pii/0021999176900784>.
- [3] Christophe Chipot and Andrew Pohorille, eds. *Free energy calculations*. Springer Series in Chemical Physics 86. Berlin: Springer, 2007.
- [4] Thomas H. Cormen, ed. *Introduction to algorithms*. en. 3rd ed. OCLC: ocn311310321. Cambridge, Mass: MIT Press, 2009. ISBN: 978-0-262-03384-8 978-0-262-53305-8.

- [5] Shujie Fan, Hristo Nedev, Ranjit Vijayan, Bogdan I Iorga, and Oliver Beckstein. “Precise force-field-based calculations of octanol-water partition coefficients for the SAMPL7 molecules”. en. In: 35 (2021), pp. 853–870.
- [6] Richard Gowers, Max Linke, Jonathan Barnoud, Tyler Reddy, Manuel Melo, Sean Seyler, Jan Domański, David Dotson, Sébastien Buchoux, Ian Kenney, and Oliver Beckstein. “MDAnalysis: A Python Package for the Rapid Analysis of Molecular Dynamics Simulations”. en. In: Austin, Texas, 2016, pp. 98–105. doi: 10.25080/Majora-629e541a-00e. URL: https://conference.scipy.org/proceedings/scipy2016/oliver_beckstein.html.
- [7] Justin Lemkul. “From Proteins to Perturbed Hamiltonians: A Suite of Tutorials for the GROMACS-2018 Molecular Simulation Package [Article v1.0]”. en. In: *Living Journal of Computational Molecular Science* 1.1 (2019). ISSN: 25756524. doi: 10.33011/livecoms.1.1.5068. URL: <https://www.livecomsjournal.org/article/5068-from-proteins-to-perturbed-hamiltonians-a-suite-of-tutorials-for-the-gromacs-2018-molecular-simulation-package-article-v1-0>.
- [8] Naveen Michaud-Agrawal, Elizabeth J. Denning, Thomas B. Woolf, and Oliver Beckstein. “MDAnalysis: A toolkit for the analysis of molecular dynamics simulations”. en. In: *Journal of Computational Chemistry* 32.10 (July 2011), pp. 2319–2327. ISSN: 01928651. doi: 10.1002/jcc.21787.
- [9] Andrew Pohorille, Christopher Jarzynski, and Christophe Chipot. “Good Practices in Free-Energy Calculations”. In: *The Journal of Physical Chemistry B* 114.32 (2010), pp. 10235–10253. doi: 10.1021/jp102971x. URL: <http://pubs.acs.org/doi/abs/10.1021/jp102971x>.
- [10] Hadley Wickham. “Tidy data”. In: *The Journal of Statistical Software* 59 (10 2014). URL: <http://www.jstatsoft.org/v59/i10/>.