



SOFTWARE TOOL ARTICLE

# pdb-tools: a swiss army knife for molecular structures [version 1; referees: 2 approved]

João P. G. L. M. Rodrigues <sup>1</sup>, João M.C. Teixeira <sup>2</sup>, Mikaël Trellet<sup>3</sup>, Alexandre M. J. J. Bonvin <sup>3</sup>

<sup>1</sup>Structural Biology, Stanford University School of Medicine, Stanford, California, 94305, USA

<sup>2</sup>Independent Researcher, Salamanca, Spain

<sup>3</sup>Bijvoet Center for Biomolecular Research, Utrecht University, Utrecht, 3584CH, The Netherlands

**V1** First published: 20 Dec 2018, 7:1961 (<https://doi.org/10.12688/f1000research.17456.1>)  
 Latest published: 20 Dec 2018, 7:1961 (<https://doi.org/10.12688/f1000research.17456.1>)

**Abstract**

The pdb-tools are a collection of Python scripts for working with molecular structure data in the Protein Data Bank (PDB) format. They allow users to edit, convert, and validate PDB files, from the command-line, in a simple but efficient manner. The pdb-tools are implemented in Python, without any external dependencies, and are freely available under the open-source Apache License at <https://github.com/haddocking/pdb-tools/> and on [PyPI](#).

**Keywords**

bioinformatics, chemistry, macromolecules, PDB, protein structure, Python, structural biology



This article is included in the [Python Collection](#) collection.

**Open Peer Review**

Referee Status:

	Invited Referees	
	1	2
<b>version 1</b> published 20 Dec 2018	 report	 report

- 1 **Bruno L. Victor** , University of Lisbon, Portugal
- 2 **Peter W. Rose** , University of California, San Diego, USA

**Discuss this article**

Comments (0)

**Corresponding authors:** João P. G. L. M. Rodrigues ([j.p.g.l.m.rodrigues@gmail.com](mailto:j.p.g.l.m.rodrigues@gmail.com)), Alexandre M. J. J. Bonvin ([a.m.j.j.bonvin@uu.nl](mailto:a.m.j.j.bonvin@uu.nl))

**Author roles:** **Rodrigues JPGLM:** Conceptualization, Software, Supervision, Validation, Writing – Original Draft Preparation, Writing – Review & Editing; **Teixeira JMC:** Software, Validation, Writing – Original Draft Preparation, Writing – Review & Editing; **Trellet M:** Software, Validation, Writing – Original Draft Preparation, Writing – Review & Editing; **Bonvin AMJJ:** Conceptualization, Software, Writing – Original Draft Preparation, Writing – Review & Editing

**Competing interests:** No competing interests were disclosed.

**Grant information:** J.P.G.L.M.R acknowledges funding from a Niels Stensen Fellowship and NIH grant R35GM122543. M.T. and A.M.J.J.B were supported by the European H2020 e-Infrastructure grant BioExcel (grant no. 675728).

*The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.*

**Copyright:** © 2018 Rodrigues JPGLM *et al.* This is an open access article distributed under the terms of the [Creative Commons Attribution Licence](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

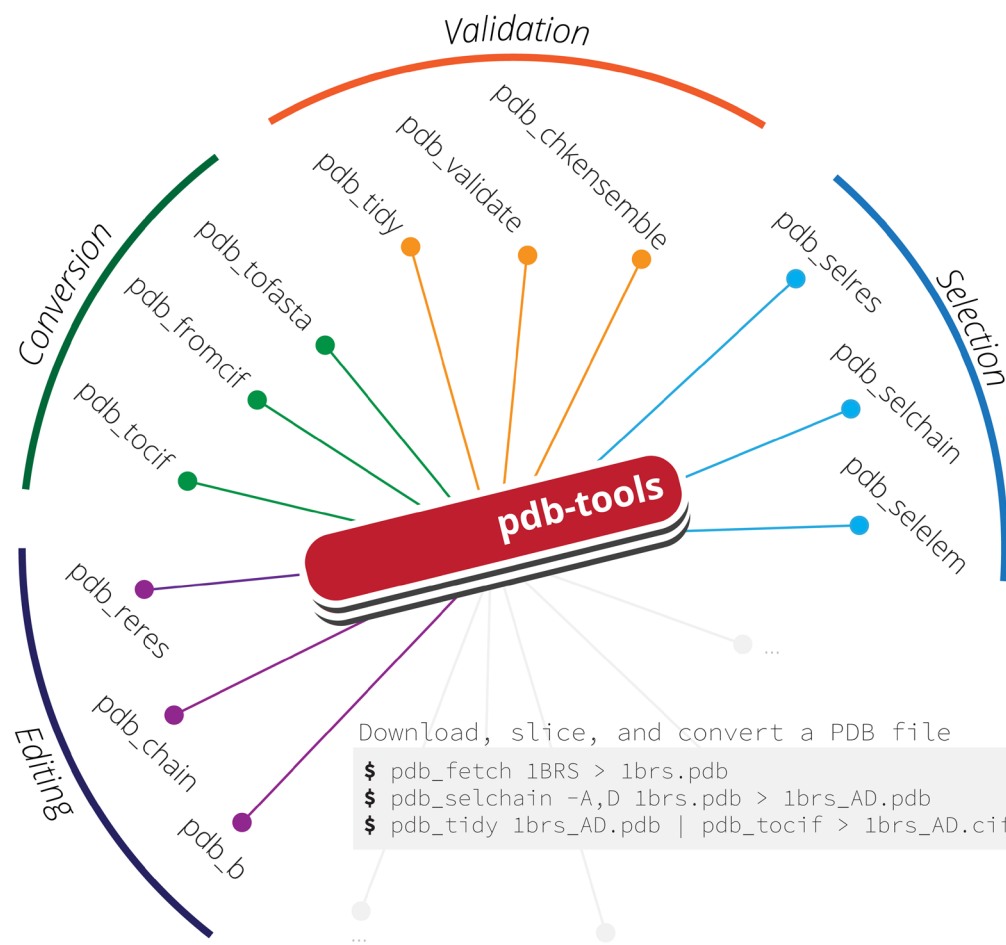
**How to cite this article:** Rodrigues JPGLM, Teixeira JMC, Trellet M and Bonvin AMJJ. **pdb-tools: a swiss army knife for molecular structures [version 1; referees: 2 approved]** *F1000Research* 2018, 7:1961 (<https://doi.org/10.12688/f1000research.17456.1>)

**First published:** 20 Dec 2018, 7:1961 (<https://doi.org/10.12688/f1000research.17456.1>)

## Introduction

Obtaining and analyzing three-dimensional structures of biological macromolecules, such as proteins or nucleic acids, is often a key step towards understanding their biological function. Of the many file formats used in structural biology to store three-dimensional coordinate data, the Protein Data Bank (PDB) format remains one of the most widely adopted, despite the introduction of a new standard - the mmCIF format - in recent years<sup>1,2</sup>. The PDB format encodes 44 possible different record types in a human-readable flat-text format, with each record having a number of fields with stringent spacing rules in a total of 80 characters. This strictness and complexity of the PDB format make manual editing difficult. As a result, researchers use molecular viewers, e.g. PyMOL<sup>3</sup> and VMD<sup>4</sup>, to perform basic editing tasks, such as selecting chains from a structure. More advanced edits, such as changing chain identifiers, deleting specific atoms, or renumbering residues require expertise with scripting languages, whose syntax varies from viewer to viewer. Other edits, such as selecting atomic positions based on occupancy values, are even more challenging to accomplish in molecular viewers, if possible at all. The alternative is for researchers to become proficient in a programming language, such as Python, and use one of its structural bioinformatics libraries, such as BioPython<sup>5</sup> or ProDy<sup>6</sup>.

Here, we present `pdb-tools`, a modular toolkit providing several everyday tasks when handling PDB files, namely downloading, editing, filtering, merging, sorting, and validation, as well as conversion to and from the more recent PDBx/mmCIF file format<sup>2</sup> (Figure 1). In order to shorten the learning curve for new users, all `pdb-tools` implement a unified command-line interface. Moreover, to support complex editing operations, we designed the toolkit to allow the serial concatenation of several tools in a pipeline without the need to read or write intermediary files. For developers, the coherent architecture simplifies maintenance and development of new tools.



**Figure 1. Overview of the tasks supported in `pdb-tools` and example tools.** A common usage example is shown on the gray box insert.

The pdb-tools are written in [Python](#) and are currently developed on GitHub under the open-source Apache License version 2.0. In addition, to simplify the installation procedure for the end-user, the toolkit is also available for download through [PyPI](#). Finally, pdb-tools are also part of the SBGrid initiative<sup>7</sup>.

## Methods

### Implementation

The pdb-tools toolkit is implemented in [CPython](#), using only modules of the standard library, which allows us to support a wide range of Python versions (2.7 and all of the 3.x series). The tools make use of Python generator expressions to improve performance and memory usage. This feature, along with support for streaming input and output data, allows users to create complex pipelines by chaining together different tools (e.g. format conversion followed by chain selection followed by editing followed by validation).

As a result, using 'pdb\_reres' to renumber the residues of a structure with 64,606 atoms takes ~0.17 seconds (on an Intel i7-7500U CPU) and ~10 MB of memory. Merging eight copies of the same structure using 'pdb\_merge' took ~0.74 seconds and ~18 MB of memory. These performance indicators make pdb-tools particularly useful when combined with shell scripting for batch processing entire collections of PDB files.

To aid the maintenance and development of pdb-tools, we wrote approximately 250 separate unit tests that provide 81% coverage for the tools and their possible usage options. In addition, each commit is checked for coding style against the PEP8 style guide, using flake8. Besides hosting the source code on GitHub, we use [setuptools](#) to package and distribute pdb-tools through [PyPI](#). Further, we use semantic versioning to indicate compatibility between releases with the same major version. The code is tested continuously using [Travis CI](#) and [AppVeyor](#) webhooks on GitHub, which monitor incoming pull requests and any commits to the master branch. We test on virtual instances of Windows 10 and Linux 16.04 LTS, under Python 2.7, 3.6, and 3.7. Nevertheless, the pdb-tools should run on most UNIX derivatives and Windows versions, provided a supported version of Python is installed.

Concerning documentation, each tool contains a self-contained help string, describing its purpose, command-line interface, and usage examples. This help text is available by running the tool without any argument. Additional documentation on the usage and installation of pdb-tools is available online on the project's web-page and in the README file accompanying the source code.

Finally, we develop pdb-tools as a community open-source project on GitHub and encourage contributions of any kind. To date, the project has more than 30 forks and has been starred by over a dozen users on GitHub. To help and promote these contributions, and based on our experience with building pdb-tools collaboratively, we have documents setting coding best practices, contribution guidelines, and a code of conduct for developers on our GitHub repository. In addition, we make use of and encourage users to contribute to our issue tracker, which documents our discussions and development decisions publicly.

## Use Cases

The pdb-tools follow a 'one tool, one job' design philosophy, where each tool performs a simple operation on an input file, while more complex operations can be achieved by chaining different tools together. Over the years, the pdb-tools have been used in both research and educational contexts, some of which we highlight below.

### Making selections in PDB files

Molecular viewers, such as [PyMOL](#) or [VMD](#), are the tools of choice for making selections of specific chains or residue ranges in PDB files. When handling more than one structure, programming libraries such as [Biopython](#) or [Prody](#) are better suited for the task. The selection tools included in pdb-tools offer a simple solution to make complex selections.

```
$ pdb_fetch lctf | pdb_keepcoord | pdb_reres > lctf.pdb
$ pdb_selres -1:30 lctf.pdb | pdb_selatom -CA > lctf_CA.pdb
$ head lctf_CA.pdb
ATOM      2  CA  GLU  A   1       17.706 17.982 -14.905 1.00 16.74      C
ATOM     11  CA  PHE  A   2       17.509 14.262 -14.184 1.00 13.24      C
... [truncated]
```

Importantly, users can use shell scripting to efficiently process large collections of structures. Those using UNIX-based systems can also make use of the [parallel](#) utility to distribute processing over several cores.

```
$ for pdbf in $( ls *.pdb )
$ do
$  pdb_selchain -A $pdbf | pdb_delelement -H | pdb_tidy > ${pdbf}_A_noH.pdb
$ done
```

### Preparing input files for molecular modelling using HADDOCK

HADDOCK is an integrative modeling software for protein interactions<sup>8</sup>, whose web server provides a user-friendly interface that oversaw over 200.000 job submissions to date. The simplest of submissions requires two structures, in PDB format, that cannot have more than one chain, cannot have overlapping residue numbers, and cannot have alternate locations for any atom. Often, users submitting jobs to the HADDOCK web server are faced with error messages because of such formatting requirements. Using pdb-tools, producing a HADDOCK-compliant PDB file is straightforward, as we demonstrate below.

```
$ pdb_fetch lbrs | pdb_keepcoord > lbrs.pdb # 6 chains
$ pdb_selchain -A,D lbrs.pdb | pdb_delhetatm | pdb_selaltloc > lbrs_AD.pdb
$ pdb_selchain -A lbrs_AD.pdb | pdb_tidy > lbrs_A.pdb
$ pdb_selchain -D lbrs_AD.pdb | pdb_tidy > lbrs_D.pdb
```

### Extracting the aminoacid/nucleotide sequence of a PDB file

Although PDB files may include SEQRES records that store the sequence of the construct used for structure determination, the final structure does not necessarily include all the amino acids or nucleotides. This discrepancy is particularly important when building alignments for homology modeling. To this end, we include in pdb-tools a PDB to FASTA converter that extracts the sequence of the structure directly from the ATOM/HETATM records. The user can specify if the entire structure should be exported as a single FASTA sequence record, or divided by chain, using the **-multi** option.

```
$ pdb_fetch lbrs | pdb_tofasta
>PDB|ABCDEF
INTFDGVADYQLQTYHKLPDNYITKSEAQALGWVASKGNLADVAPGKSIGGDIFSNREGKL
PGKSGRTWREADINYTSGFNRNSDRILYSSDWLIYKTTDHYQTFTKIRAQVINTFDGVADY
... [truncated]
$
$ pdb_fetch lbrs | pdb_tofasta -multi
>PDB|A
VINTFDGVADYQLQTYHKLPDNYITKSEAQALGWVASKGNLADVAPGKSIGGDIFSNREGK
LPGKSGRTWREADINYTSGFNRNSDRILYSSDWLIYKTTDHYQTFTKIR
>PDB|B
AQVINTFDGVADYQLQTYHKLPDNYITKSEAQALGWVASKGNLADVAPGKSIGGDIFSNRE
GKLPKSGRTWREADINYTSGFNRNSDRILYSSDWLIYKTTDHYQTFTKIR
... [truncated]
```

### Converting to and from mmCIF format

Due to several limitations with the PDB file format, the Protein Data Bank recently introduced the new PDBx/mmCIF format. To support this decision, we include in pdb-tools converters to and from the mmCIF format. These converters allow users to implicitly use pdb-tools on mmCIF files. However, large structures that cannot abide by the PDB format specification (too many chains or residues or atoms) cannot be converted.

```
$ pdb_fromcif lbrs.cif | pdb_selchain -A | pdb_tocif > lbrs_A.cif
```

### Summary

The pdb-tools provide a command-line interface to edit three-dimensional molecular structures in the PDB format. Since it is written in Python, without any third-party dependencies, the toolkit is available on Linux, Mac OS, and Windows and supported on both Python 2 and 3. The source code is hosted on GitHub and licensed under the open-source Apache License version 2.0, to encourage contributions from the community, and available for installation through PyPI and SBGrid. As such, the pdb-tools are an efficient, portable, and extremely flexible toolkit for both starting and experienced researchers in structural biology.

### Software availability

pdb-tools is available from: <https://haddocking.github.io/pdb-tools/>, <https://pypi.org/project/pdb-tools/>, and <https://sbgrid.org/software/titles/pdb-tools>

Source code available from: <https://github.com/haddock/pdb-tools>, <https://github.com/haddock/pdb-tools/releases/tag/2.0.0-rc1>

Archived source code as at time of publication: <http://doi.org/10.5281/zenodo.2168070><sup>9</sup>

License: Apache License, version 2.0

### Grant information

J.P.G.L.M.R acknowledges funding from a Niels Stensen Fellowship and NIH grant R35GM122543. M.T. and A.M.J.J.B were supported by the European H2020 e-Infrastructure grant BioExcel (grant no. 675728).

*The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.*

### Acknowledgments

The authors acknowledge the members of the Bonvin group for fruitful discussions and suggestions of new tools over the years.

### References

- Westbrook JD, Fitzgerald PM: **The PDB format, mmCIF, and other data formats.** *Methods Biochem Anal.* 2003; **44**: 161–179.  
[PubMed Abstract](#) | [Publisher Full Text](#)
- wwPDB consortium: **Protein Data Bank: the single global archive for 3D macromolecular structure data.** *Nucleic Acids Res.* 2018; gky949.  
[PubMed Abstract](#) | [Publisher Full Text](#)
- Schrödinger LLC: **The PyMOL Molecular Graphics System, Version 1.8.** 2015.
- Humphrey W, Dalke A, Schulten K: **VMD: visual molecular dynamics.** *J Mol Graph.* 1996; **14**(1): 33–38, 27-8.  
[PubMed Abstract](#) | [Publisher Full Text](#)
- Cock PJ, Antao T, Chang JT, *et al.*: **Biopython: freely available Python tools for computational molecular biology and bioinformatics.** *Bioinformatics.* 2009; **25**(11): 1422–1423.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Bakan A, Meireles LM, Bahar I: **ProDy: protein dynamics inferred from theory and experiments.** *Bioinformatics.* 2011; **27**(11): 1575–1577.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Morin A, Eisenbraun B, Key J, *et al.*: **Collaboration gets the most out of software.** *eLife.* 2013; **2**: e01456.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- van Zundert GCP, Rodrigues JPGLM, Trellet M, *et al.*: **The HADDOCK2.2 Web Server: User-Friendly Integrative Modeling of Biomolecular Complexes.** *J Mol Biol.* 2016; **428**(4): 720–725.  
[PubMed Abstract](#) | [Publisher Full Text](#)
- Rodrigues J, Teixeira JMC, Trellet M, *et al.*: **haddock/pdb-tools: Release Candidate 1 (Version 2.0.0-rc1).** *Zenodo.* 2018.  
<http://www.doi.org/10.5281/zenodo.2168070>

# Open Peer Review

Current Referee Status:  

Version 1

Referee Report 18 January 2019

<https://doi.org/10.5256/f1000research.19090.r42177>



**Peter W. Rose** 

Structural Bioinformatics Laboratory, San Diego Supercomputer Center, University of California, San Diego, La Jolla, CA, USA

While there are many, often complex, tools available to manipulate pdb files, pdb-tools stands out as a light-weight set of command line tools for common use cases. The requirement for each tool to perform one task makes these tools easy to use in shell scripts and to combine into a pipeline for batch processing a set of PDB structures.

On the software side, this is a well-designed open source tool kit available in GitHub that follows best software development practices, including documentation, checking for consistent coding style, high test coverage, using continuous integration, semantic versioning, and deployment on PyPi. A liberal open source license encourages community contributions. The modular design makes this project easy to maintain and expand compared to many of the legacy tools still in use.

Recommendation:

List all, not just a subset, of the available tools with an example on the documentation page. That way, a Google search will find those tools and examples are helpful for new users.

**Is the rationale for developing the new software tool clearly explained?**

Yes

**Is the description of the software tool technically sound?**

Yes

**Are sufficient details of the code, methods and analysis (if applicable) provided to allow replication of the software development and its use by others?**

Yes

**Is sufficient information provided to allow interpretation of the expected output datasets and any results generated using the tool?**

Yes

**Are the conclusions about the tool and its performance adequately supported by the findings presented in the article?**

Yes

**Competing Interests:** No competing interests were disclosed.

**Referee Expertise:** Structural Bioinformatics, Scientific Software Development

**I have read this submission. I believe that I have an appropriate level of expertise to confirm that it is of an acceptable scientific standard.**

Referee Report 04 January 2019

<https://doi.org/10.5256/f1000research.19090.r42178>



**Bruno L. Victor** 

BiolSI, Biosystems and Integrative Sciences Institute, Department of Chemistry and Biochemistry, Faculty of Sciences, University of Lisbon, Lisbon, Portugal

In this work, Rodrigues et al, present to the readers a collection of python based tools to edit, convert and validate PDB files, in a very simple and efficient way. The authors clearly describe in the methods section the procedure used to implement the pdb-tools toolkit, and they even provide some simple use cases where the toolkit can be applied. I have installed the toolkit in my workstation and I did not have any kind of problem. I have tested several of the available functionalities and I have found it to be very useful and easy-to-use on day-to-day work (as the authors clearly state in the article). This type of contribution software is very useful for any scientist (with beginner or advanced skills) dealing with structure data in PDB format.

**Is the rationale for developing the new software tool clearly explained?**

Yes

**Is the description of the software tool technically sound?**

Yes

**Are sufficient details of the code, methods and analysis (if applicable) provided to allow replication of the software development and its use by others?**

Yes

**Is sufficient information provided to allow interpretation of the expected output datasets and any results generated using the tool?**

Yes

**Are the conclusions about the tool and its performance adequately supported by the findings presented in the article?**

Yes

**Competing Interests:** No competing interests were disclosed.

**Referee Expertise:** Molecular modelling, computational chemistry, computational medicinal chemistry, software development applied to drug discovery projects



**I have read this submission. I believe that I have an appropriate level of expertise to confirm that it is of an acceptable scientific standard.**

---

The benefits of publishing with F1000Research:

- Your article is published within days, with no editorial bias
- You can publish traditional articles, null/negative results, case reports, data notes and more
- The peer review process is transparent and collaborative
- Your article is indexed in PubMed after passing peer review
- Dedicated customer support at every stage

For pre-submission enquiries, contact [research@f1000.com](mailto:research@f1000.com)

**F1000Research**