

Use cases of veriNA3d

Diego Gallego Perez^{1,2}

¹University of Barcelona - Dept. Biochemistry and Molecular Biomedicine

²Institute for Research in Biomedicine - (IRB Barcelona), Barcelona Institute of Science and Technology

2018-12-07

Contents

1	Introduction: Structural Bioinformatics in R	2
2	Parsing mmCIF files	2
2.1	Origin and standardization of mmCIF files	2
2.2	The CIF object	3
2.3	Bidirectional compatibility with bio3d	4
3	Main use case: Manage Nucleic Acid datasets	5
4	Use case: Querying the EMBL-EBI REST API	5
5	Use case: eRMSD	5
6	Use case: Generate substructures	5
	References	5

1 Introduction: Structural Bioinformatics in R

The R language provides an excellent interface for statistical analysis, which is also interesting from the point of view of structural data. This gap was filled in 2006 by the R package [bio3d](#) (Grant et al. 2006). It was presented as a suite of tools to handle PDB formatted structures, and trajectories. It integrates a variety of functions to analyse from sequence to 3D structure data (RMSD, NMA, PCA... see their [documentation](#) for details). As far as we know, bio3d represented the only structural package for R until now.

The R package presented in here, veriNA3d, does not replace bio3d at all. Rather, it was developed on top of it to cover additional necessities. The only common tool integrated in both packages is a parser for mmCIF files (see below). veriNA3d is mainly intended (but not limited) to the analysis of Nucleic Acids. It integrates a higher level of abstraction than bio3d since it also allows the analysis of datasets, in addition to analysis of single structures. The functions in the package could be divided in the following blocks:

- Dataset level: Functions to get and analyse lists of pdb IDs. This includes access to the [representative lists of RNA](#) by (Leontis and Zirbel 2012) and other analytical functions.
- Structure level: Functions to get data, parse mmCIF files and analyse these data.
- Plots: examples to show the results of the previous analysis.

The complete list of functions can be found in the README.md file within the package, also accessible on the gitlab [main page](#).

2 Parsing mmCIF files

2.1 Origin and standardization of mmCIF files

Atomic structural data of macromolecules has long been distributed in the PDB file format. However, one of its main limitations is the column size for the coordinates data, which didn't allowed to save molecules with more than 99999 atoms, more than 62 chains or more than 9999 residues (in a chain).

Given that the Protein Data Bank is continuously growing and accepting bigger structures (e.g. a whole *E.coli* ribosome has over 140000 atoms - pdbID 4V4S), an alternative file format became the standard: the mmCIF file format.

The mmCIFs are an evolution of the Crystallographic Information File (CIF), originally used for small molecule structures. It stands for **macromolecular CIF** file, and it has actually coexisted with the PDB format since the 1997. However, since the PDB is easier to parse and such big structures didn't populate the database at the time, most software has been developed for the PDB format.

The PDB format was definitely frozen in 2014. However, it will still coexist with the standard mmCIF format as long as all softwares evolve to accept mmCIFs. Following this trend, the bio3d R package integrated a read.cif function in their version 2.3. At that time, we had already started the development of our own cifParser function. Given that the mmCIF format is constantly evolving and that both functions take slightly different approaches, we decided to offer our own version of it, which might provide an useful and fast alternative for users working with mmCIF files.

2.2 The CIF object

Parsing a particular file format often involves creating a new class of object. In R, the principal [objects](#) are called S3, S4 and RC. Our container for mmCIF data is an S4 object called **CIF**, in contrast with the S3 object (called **pdb**) in `bio3d` - it is worth noting that this difference does **not** affect the compatibility between the two packages (see below for details).

Since different mmCIF files usually have different sections of data (in addition to the coordinates), we carried out an analysis that checked which ones are always present in all mmCIF files (this included all mmCIF files in the Protein Data Bank in March 2018), and reached a list of 14 items:

- Atom_site
- Atom_sites
- Atom_type
- Audit_author
- Audit_conform
- Chem_comp
- Database_2
- Entity
- Entry
- Exptl
- Pdbx_database_status
- Struct
- Struct_asym
- Struct_keywords

The detailed description of each these data sections can be found in the mmCIF [main site](#).

The CIF object is created by the `cifParser` function and contains these 14 sections of data, which can be accessed with the CIF accessors. To see the accessor functions run:

```
library(veriNA3d)
?cif_accessors
```

To read a mmCIF file and access the coordinates data, use:

```
## To parse a local mmCIF file:
# cif <- cifParser("your-file.cif")
## To download from PDB directly:
cif <- cifParser("1bau")
cif
#>
#> -- mmCIF with ID: 1BAU -----
#>
#> Author description:  NMR STRUCTURE OF THE DIMER INITIATION COMPLEX OF HIV-1 GENOMIC RNA, MINIMIZED AVERAG
#> mmCIF version:      5.279
#>
#> To extract coordinates and other data use accessor functions
#> (type ?cif_accessors for details)
## To see the coordinates:
coords <- cifAtom_site(cif)
head(coords)
#>  group_PDB id type_symbol label_atom_id label_alt_id label_comp_id
```

Use cases of veriNA3d

```
#> 1      ATOM 1      0      05'      .      G
#> 2      ATOM 2      C      C5'      .      G
#> 3      ATOM 3      C      C4'      .      G
#> 4      ATOM 4      0      04'      .      G
#> 5      ATOM 5      C      C3'      .      G
#> 6      ATOM 6      0      03'      .      G
#> label_asym_id label_entity_id label_seq_id pdbx_PDB_ins_code Cartn_x
#> 1            A            1            1            ? 23.989
#> 2            A            1            1            ? 24.965
#> 3            A            1            1            ? 25.937
#> 4            A            1            1            ? 26.741
#> 5            A            1            1            ? 25.259
#> 6            A            1            1            ? 24.868
#> Cartn_y Cartn_z occupancy B_iso_or_equiv pdbx_formal_charge auth_seq_id
#> 1    8.289 -15.135      1            0            ?      1
#> 2    9.100 -14.503      1            0            ?      1
#> 3    9.725 -15.512      1            0            ?      1
#> 4    8.744 -16.162      1            0            ?      1
#> 5   10.527 -16.627      1            0            ?      1
#> 6   11.838 -16.252      1            0            ?      1
#> auth_comp_id auth_asym_id auth_atom_id pdbx_PDB_model_num
#> 1            G            A      05'            1
#> 2            G            A      C5'            1
#> 3            G            A      C4'            1
#> 4            G            A      04'            1
#> 5            G            A      C3'            1
#> 6            G            A      03'            1
```

2.3 Bidirectional compatibility with bio3d

Using the `cifParser` will often be the first step to analyse a structure. However, if the analysis requires any of the `bio3d` functions, then a conversion should be done with the `cifAsPDB` function.

```
pdb <- cifAsPDB(cif)
pdb
#>
#> Call: "1BAU"
#>
#> Total Models#: 1
#> Total Atoms#: 1486, XYZs#: 4458 Chains#: 2 (values: A B)
#>
#> Protein Atoms#: 0 (residues/Calpha atoms#: 0)
#> Nucleic acid Atoms#: 1486 (residues/phosphate atoms#: 46)
#>
#> Non-protein/nucleic Atoms#: 0 (residues: 0)
#> Non-protein/nucleic resid values: [ none ]
#>
#> Nucleic acid sequence:
```

Use cases of veriNA3d

```
#>      GGCAAUGAAGCGCGCACGUUGCCGGCAAUGAAGCGCGCACGUUGCC  
#>  
#> + attr: atom, xyz, calpha, model, flag, call
```

It takes a CIF object and generates an equivalent pdb object (as used by all bio3d functions). In addition, all veriNA3d functions are prepared to accept as input either the CIF or pdb objects. Therefore, the compatibility between the two packages is bidirectional.

3 Main use case: Manage Nucleic Acid datasets

Get Leontis list, change representative structures and analyse them with one of the pipelines

4 Use case: Querying the EMBL-EBI REST API

Query EBI API and how to construct new queries!

5 Use case: eRMSD

For two NMR models of RNA/DNA, compute the eRMSD

6 Use case: Generate substructures

For a given structure (CIF or PDB), generate a smaller PDB with the region of interest and surroundings

References

Grant, B.J., A.P.C. Rodrigues, K.M. ElSawy, J.A. McCammon, and L.S.D. Caves. 2006. "Bio3d: An R Package for the Comparative Analysis of Protein Structures." *Bioinformatics* 22 (21): 2695–6.

Leontis, N.B., and C.L. Zirbel. 2012. "Nonredundant 3D Structure Datasets for RNA Knowledge Extraction and Benchmarking." In *RNA 3D Structure Analysis and Prediction*, edited by N. Leontis and E. Westhof, 27:281–98. Springer Berlin Heidelberg.