Diego Gallego Perez 1,2,3

¹University of Barcelona - Department of Biochemistry and Molecular Biomedicine

2018-12-18

Contents

1	Introduction: Structural Bioinformatics in R						
2	Parsing mmCIF files	3					
	2.1 Origin and standardization of mmCIF files	3					
	2.2 The CIF object	3					
	2.3 Bidirectional compatibility with bio3d	5					
3	Getting data from Application Programming Interfaces (API)						
	3.1 Querying the EMBL-EBI REST API	6					
4	The ϵ RMSD to compare structures						
5	Generate substructures						
6	Calculate phosphate pair-wise distances						
7	Manage Nucleic Acid datasets	15					
	7.1 Non-redundant lists of RNA	15					
	7.2 Queries	16					
	7.3 Analysing datasets	16					
	References	20					

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

³Joint BSC-IRB Program in Computational Biology

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 formated 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. This
 includes a wrapper of DSSR (Lu, Bussemaker, and Olson 2015) and a function to
 calculate the εRMSD (Bottaro, Di Palma, and Bussi 2014).
- 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 **m**acro**m**olecular **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 definetely 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 analisis 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 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")</pre>
## To download from PDB directly:
cif <- cifParser("1bau")</pre>
cif
#>
#> -- mmCIF with ID: 1BAU -----
#> Author description: NMR STRUCTURE OF THE DIMER INITIATION COMPLEX OF HIV-1
#> GENOMIC RNA, MINIMIZED AVERAGE STRUCTURE
#>
#> 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)</pre>
head(coords)
#> group_PDB id type_symbol label_atom_id label_alt_id label_comp_id
#> 1 ATOM 1 0 05' .
                  C C5'
C C4'
O O4'
C
#> 2
      ATOM 2
                                                    G
    ATOM 3
ATOM 4
#> 3
                                                     G
#> 4
                                                     G
      ATOM 5
#> 5
#> 6 ATOM 6 0 03'
#> label_asym_id label_entity_id label_seq_id pdbx_PDB_ins_code Cartn_x
1 1 1 1 1 1 1 1 1
           Α
                                                ? 24.965
            Α
#> 3
                                                ? 25.937
           Α
#> 4
#> 5
            Α
                1 1
            Α
#> Cartn_y Cartn_z occupancy B_iso_or_equiv pdbx_formal_charge auth_seq_id
#> 1 8.289 -15.135 1 0
                                       ? 1
                                0
                                                ?
#> 2 9.100 -14.503
                     1
                                                         1
#> 3 9.725 -15.512
                     1
                                0
                                               ?
                                                         1
#> 4 8.744 -16.162
                     1
                                0
                                                ?
                                                          1
#> 4 8.744 -16.162 1
#> 5 10.527 -16.627 1
#> 6 11.838 -16.252 1
                                0
                                                ?
                                                         1
                                0
                                                         1
#> auth_comp_id auth_asym_id auth_atom_id pdbx_PDB_model_num
```

#> 1	G	А	05'	1
#> 2	G	А	C5 '	1
#> 3	G	A	C4'	1
#> 4	G	A	04'	1
#> 5	G	Α	C3 '	1
#> 6	G	Α	03'	1

2.3 Bidirectional compatibility with bio3d

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

```
pdb <- cifAsPDB(cif)</pre>
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:
#>
        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 Getting data from Application Programming Interfaces (API)

Getting data is the first step of any pipeline, and parsing files is just one of the many ways data can be accessed. Application Programming Interfaces (API) are an intermediate point of access to a remote database. APIs offer the users a series of *endpoints* or *calls*, which are just **links**. Thus, instead of dealing directly with the database with SQL or other language, the user can just use the appropriate *call* to send a query to the API, and it will return the desired data.

In addition to parsing mmCIF files, veriNA3d also offers a series of functions to send queries to different APIs. Since sending queries to remote APIs requires Internet access, the full functionality of veriNA3d might depend on a good connection.

To see all the query functions, run:

?queryFunctions

IMPORTANT NOTE: The APIs accessed by veriNA3d are free of use with no limit of queries per user. However, this could change if the users of the APIs use them irresponsibly. Servers could eventually fall down if they receive more *calls* than they can actually process. To avoid that, veriNA3d actually saves in memory the result of any query, and any time that you use the same query again, it will take the cached result. To see this effect, run this test:

```
## Run a query for the first time, which will access the API
tech <- queryTechnique("4KQX", verbose=TRUE)

#> [1] "Querying: http://www.ebi.ac.uk/pdbe/api/pdb/entry/summary/4KQX"

#> [1] "Getting expType from API"

#> [1] "Saving expType in RAM"

## Run the same query for the second time, which will get it from memory tech <- queryTechnique("4KQX", verbose=TRUE)

#> [1] "Querying: http://www.ebi.ac.uk/pdbe/api/pdb/entry/summary/4KQX"

#> [1] "Getting expType from RAM"
```

However, this is only saved across the current session of R, and any script that uses the query functions will send them to the API every time it is run. VeriNA3d does not guarantee the correct service of the APIs and it does not monitor any of your processes. However, the API providers can known which IP address is querying their services at any time. To avoid overloading the servers, please make a responsible use (e.g. save locally the data that you use frequently).

3.1 Querying the EMBL-EBI REST API

An invaluable resource for structural & computational biologists is the PDBe REST API (Velankar et al. 2015). Around this technology, the package includes the following set of functions:

queryAuthors: List of authorsqueryReldate: Release datequeryDepdate: Deposition date

queryRevdate: Revision date

queryDescription: Author descriptionqueryEntities: Entitity information

queryFormats: File formats for the given ID

queryModres: Modified residuesqueryLigands: Ligands in structure

queryOrgLigands: Ligands in structure (substracting ions)

queryResol: Resolution (if applicable)queryTechnique: Experimental Technique

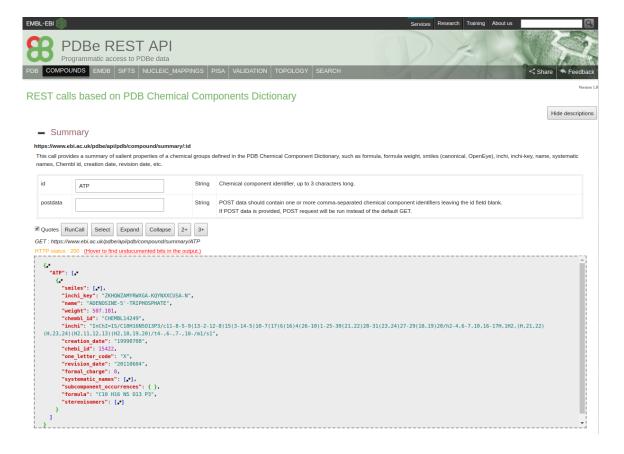
queryStatus: Released/Obsolete and related status information

The list of functions is **intendedly limited** in comparison with the dozens of *enpoints* of the REST API. Integrating them all would innecessarily increase the total amount of functions of the package. Moreover, the API might offer more and more *endpoints* with the time, and trying to keep them all would make the manteinance of this package more difficult. To allow the users to access their desired endpoints, an alternative method is provided.

The core of all the query functions is queryAPI, which integrates all the error-handling and cache functionalities. With the queryAPI function, any user can design their own queries, with a simple process. Herein a couple of examples.

3.1.1 Example 1

This snapshot shows the REST API website and a call that is not implemented in veriNA3d:



The link of this endpoint is: https://www.ebi.ac.uk/pdbe/api/pdb/compound/summary/ATP

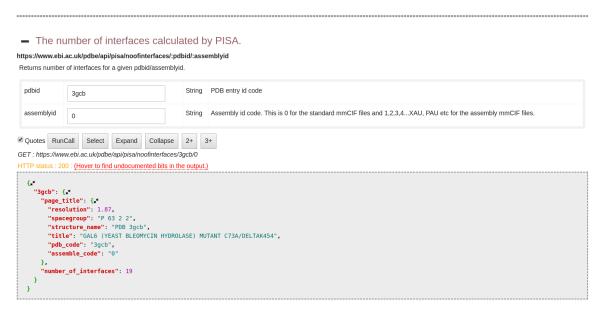
The queryAPI function can understand and send this query using the arguments 'ID', 'API', 'string1', and 'string2' properly:

```
atpsummary <- queryAPI(ID="ATP", API="ebi",</pre>
                       string1="pdb/compound/summary/", string2="")
str(atpsummary$ATP)
#> 'data.frame': 1 obs. of 15 variables:
#> $ smiles
                           :List of 1
#> ..$:'data.frame': 2 obs. of 3 variables:
    ....$ program: chr "CACTVS" "OpenEye OEToolkits"
#> ....$ version: chr "3.341" "1.5.0"
#> ....$ name : chr "Nc1ncnc2n(cnc12)[C@@H]30[C@H](C0[P@](0)(=0)0[P@@](0)(=0)0[P](0)(0)=0)[C@@H](0)[C@
#> $ inchi_key
                           : chr "ZKHQWZAMYRWXGA-KQYNXXCUSA-N"
#> $ name
#> $ weight
                           : chr "ADENOSINE-5'-TRIPHOSPHATE"
                           : num 507
#> $ chembl_id
                           : chr "CHEMBL14249"
#> $ inchi
                           : chr "InChI=1S/C10H16N5013P3/c11-8-5-9(13-2-12-8)15(3-14-5)10-7(17)6(16)4(26-
#> $ creation_date
                          : chr "19990708"
#> $ chebi_id
                           : int 15422
#> $ one_letter_code
                        : chr "X"
#> $ revision_date
                           : chr "20110604"
#> $ formal_charge
                          : int 0
#> $ rotmat_charge : Int o
#> $ systematic_names :List of 1
#> ..$:'data.frame': 2 obs. of 3 variables:
    ....$ program: chr "ACDLabs" "OpenEye OEToolkits"
#> ....$ version: chr "10.04" "1.5.0"
#> .... name : chr "adenosine 5'-(tetrahydrogen triphosphate)" "[[(2R,3S,4R,5R)-5-(6-aminopurin-9-yl)
#> $ subcomponent_occurrences:'data.frame': 1 obs. of 0 variables
#> $ formula
                           : chr "C10 H16 N5 013 P3"
#> $ stereoisomers
                           :List of 1
#> ..$:'data.frame': 1 obs. of 2 variables:
    ....$ name : chr "9-{5-0-[(S)-hydroxy{[(R)-hydroxy(phosphonooxy)phosphoryl]oxy}phosphoryl]-beta
   ....$ chem_comp_id: chr "HEJ"
```

- The common root in all the REST API *endpoints* is "https://www.ebi.ac.uk/pdbe/api/", which is internally managed by the function by using *API="ebi"*.
- The string1="pdb/compound/summary/" indicates everything that comes after the root and before the ID.
- The *ID="ATP"* obviously represents the desired structure, either a 4 character string for a pdbID or <= 3 character string for compounds.
- The *string2*="" is also a necessary argument that reflects that nothing comes after the ID.

3.1.2 Example 2

This snapshot shows a second call not implemented in veriNA3d:



The link of this endpoint is: https://www.ebi.ac.uk/pdbe/api/pisa/noofinterfaces/3gcb/0

The proper call with queryAPI would be:

This second example shows a case in which the "string2" argument is necessary. If you are unsure about the real link that is actually being constructed, you can always use *verbose=TRUE* to see it printed.

4 The ϵ RMSD to compare structures

A new interesting metric to compare Nucleic Acid structures is the ϵ RMSD (Bottaro, Di Palma, and Bussi 2014), currently available in the BaRNAba python package (Bottaro et al. 2018). The metric was implented in the eRMSD function and completely reproduces BaRNAba results for the structures tested.

The following example shows how to get the $\epsilon RMSD$ between two models of the same structure:

```
## Parse cif file
cif <- cifParser("2d18")

## Select a couple of models
model1 <- selectModel(cif=cif, model=1)
model3 <- selectModel(cif=cif, model=3)

## Calculate the eRMSD
eRMSD <- eRMSD(cif1=model1, cif2=model3)
eRMSD
#> [1] 0.3000208
```

5 Generate substructures

In many cases you might be interested on a particular region of a structure (e.g. a peptide from a complex, or a ligand and it's binding site). For a given structure, trimSphere can generate a smaller pdb object and save it to a PDB file if desired. The region of interest can be selected by using the chain identifier and the residue index, or with the atom.select function from bio3d (which in turn allows you to select regions of the structure in a variety of ways). In addition to the region of interest, the function can also include the surrounding region by seting a cutoff distance.

5.0.1 Example 1

```
## Parse human ribosome - takes around 12 seconds in R-3.5
cif <- cifParser("6ek0")</pre>
## Query entities and check them
ent <- queryEntities("6ek0")</pre>
head(ent[, c("entity_id", "molecule_name", "in_chains")])
  entity_id
                         molecule_name in_chains
#> 1
                    28S ribosomal RNA
         1
#> 2
           2
                                              L7
                     5S ribosomal RNA
#> 3
           3
                    5.8S ribosomal RNA
                                              L8
           4 60S ribosomal protein L8
#> 4
                                              LA
#> 5
           5 60S ribosomal protein L3
                                              LB
           6 60S ribosomal protein L4
                                              LC
## Generate a smaller pdb with the 60S ribosomal protein L8
chain <- "LA"
protL8 <- trimSphere(cif, chain=chain, cutoff=0)</pre>
protL8
#>
#> Call: trim.pdb(pdb = cif, eleno = outeleno)
     Total Models#: 1
#>
       Total Atoms#: 1898, XYZs#: 5694 Chains#: 1 (values: A)
#>
       Protein Atoms#: 1898 (residues/Calpha atoms#: 248)
#>
       Nucleic acid Atoms#: 0 (residues/phosphate atoms#: 0)
#>
       Non-protein/nucleic Atoms#: 0 (residues: 0)
#>
       Non-protein/nucleic resid values: [ none ]
#>
#>
#>
     Protein sequence:
#>
        GRVIRGQRKGAGSVFRAHVKHRKGAARLRAVDFAERHGYIKGIVKDIIHDPGRGAPLAKV
        VFRDPYRFKKRTELFIAAEGIHTGOFVYCGKKAOLNIGNVLPVGTMPEGTIVCCLEEKPG
#>
#>
        DRGKLARASGNYATVISHNPETKKTRVKLPSGSKKVISSANRAVVGVVAGGGRIDKPILK
        AGRAYHKYKAKRNCWPRVRGVAMNPVEHPFGGGNHQHIGKPSTIR...<cut>...LRGT
#>
#> + attr: atom, helix, sheet, segres, xyz,
```

```
#> calpha, call

## The same command with the argument file would save it directly:
protL8 <- trimSphere(cif, chain=chain, cutoff=0, file="output.pdb")</pre>
```

5.0.2 Example 2

To get the desired region of interest and its sorroundings, let's see a second example using the same structure:

```
## Load bio3d library
library(bio3d)
## Get pdb object from CIF
pdb <- cifAsPDB(cif)</pre>
## Get list of ligands in the human ribosome 6EK0
queryLigands("6ek0", onlyligands=T)
#> [1] "MG" "HMT" "ZN" "HYG"
## Get the atomic index for a desired ligand
HMTligand_inds <- which(pdb$atom$resid == "HMT")</pre>
## Use bio3d function to select the ligand using its atom indices
sel <- atom.select(pdb, eleno=HMTligand_inds)</pre>
## Get substructure and sorroundings at 10 Angstroms
HTMligand <- trimSphere(pdb, sel=sel, cutoff=5)</pre>
#> [1] "Computing distances ..."
#> [1] " ... done"
#> [1] "Finding the atom details ..."
#> [1] " ... done, the output is coming"
```

5.0.3 Example 3

A third useful example would be the generation of a pdb with the interacting region between two molecules in the structure. To achieve the goal, veriNA3d also counts with the findBind ingSite function, as shown below:

```
## Parse another pdb for this example
pdb <- cifAsPDB("lnyb")

## Find region of interaction between RNA and protein
data <- findBindingSite(pdb, select="RNA", byres=TRUE)

## Get atom indices from interacting region molecules
eleno <- append(data$eleno_A, data$eleno_B)

## Select using bio3d</pre>
```

```
sel <- atom.select(pdb, eleno=eleno)

## Get substructure
trimSphere(pdb, sel=sel, file="interacting_site.pdb", verbose=FALSE)</pre>
```

6 Calculate phosphate pair-wise distances

Atomic distances are calculated internally in the trimSphere function, a resource that is also made available to the user through the functions measureEntityDist (returns all the distances between the atoms of two entities) and measureElenoDist (returns all the distances between a set of atoms).

Herein it is shown how to calculate the pair-wise distances between all the phosphate atoms of an RNA structure, which could be easily adapted for alpha carbons or other elements.

7 Manage Nucleic Acid datasets

VeriNA3d also provides a set of tools to work at a higher level of abstraction, which allow the user to manage datasets of structures at once. Such approaches are often found in the literature in a variety of studies, either with descriptive (Murray et al. 2003, Darré et al. (2016)), validation (Jain, Richardson, and Richardson 2015) or modeling (Z. Wang and Xu 2011, J. Wang et al. (2015)) purposes. However, as far as we know, no software with dataset analysis purposes has ever been freely distributed before.

7.1 Non-redundant lists of RNA

VeriNA3d integrates the function <code>getLeontisList</code>, which returns the desired non-redundant list of RNA structures, as found in their website (Leontis and Zirbel 2012). Here it is shown an example of how to make a non-redundant dataset of protein-RNA complexes with resolutions better than 2 Angstroms.

```
## Get non-redundant list from Leontis website
rnalist <- getLeontisList(release="3.47", threshold="2.0A")</pre>
head(rnalist)
     Equivalence_class Representative
                                                 Class_members
#> 1
         NR_2.0_29848.1
                                  3DIL | 1 | A
                                                       3DIL | 1 | A
#> 2
         NR_2.0_47162.1
                                  3PDR | 1 | X 3PDR | 1 | X; 3PDR | 1 | A
#> 3
         NR_2.0_39627.3
                                  2R8S | 1 | R
                                                       2R8S | 1 | R
         NR_2.0_35133.2
#> 4
                                  2Z75 | 1 | B
                                                       2Z75 | 1 | B
#> 5
         NR_2.0_66115.1
                                  6CB3 | 1 | B 6CB3 | 1 | B; 6CB3 | 1 | A
         NR_2.0_31222.1
                                  5FJC | 1 | A
                                                       5FJC | 1 | A
```

The previous object is a data-frame with three columns: * Equivalence class * Representative structure * Members of the Equivalence class

Note that the output is formated according with Leontis&Zirbel nomenclature (AAAA|M|C), where "AAAA" is the PDB accession code, "M" is the model and "C" is the Chain.

The representative structures might not necessarily have the desired properties (e.g. being a protein-RNA complex). To get only representatives which are complexed with a protein, all the members of the class should be checked, which can be done with getAltRepres.

```
## Set progressbar=TRUE to see the progress
protrna <- getAltRepres(rnalist=rnalist, type="protRNA", progressbar=FALSE)</pre>
head(protrna)
     Representative Equivalence_class
                < NA >
                        NR_2.0_29848.1
#> 2
                <NA>
                        NR_2.0_47162.1
#> 3
           2R8S | 1 | R
                        NR_2.0_39627.3
                <NA>
                        NR_2.0_35133.2
                <NA>
                        NR_2.0_66115.1
#> 5
                        NR_2.0_31222.1
                < NA >
```

Many equivalence classes will lack a proper representative, so a NA value is returned. To get the definite list in a more friendly format and removing missing data, use represAsDataFrame.

```
nrlist <- represAsDataFrame(protrna)</pre>
head(nrlist)
    Equivalence_class pdb model chain
#> 1
       NR_2.0_55702.1 1C9S
                            1
#> 2 NR_2.0_14136.11 1DFU
                                   Μ
#> 3 NR_2.0_14136.12 1DFU
                              1
                                    Ν
#> 4
     NR_2.0_18996.1 1DI2
                             1
                                   Ε
#> 5 NR_2.0_50004.1 1FXL
                            1
                                   B
       NR_2.0_61174.1 1H2C
#> 6
```

7.2 Queries

When analysing a dataset, a particular piece of data might be useful, such as the file formats available or if a particular ligand is present or not in all the pdb files. A handy function to approach this need is applyToPDB. Here it's shown how to get the structures in the previous list that contain Mg ions:

```
## Set progressbar=TRUE to see the progress
nrlist_mg <- applyToPDB(FUN=hasHetAtm, listpdb=nrlist$pdb,</pre>
                      hetAtms="MG", progressbar=FALSE)
nrlist <- cbind(nrlist, Mg=nrlist_mg[, 2])</pre>
head(nrlist)
   Equivalence_class pdb model chain
      NR_2.0_55702.1 1C9S
                           1
                                   W FALSE
#> 2 NR_2.0_14136.11 1DFU
                             1
                                   M TRUE
#> 3 NR_2.0_14136.12 1DFU
                             1
                                   N TRUE
#> 4 NR_2.0_18996.1 1DI2 1
                                  E FALSE
     NR_2.0_50004.1 1FXL
#> 5
                            1
                                   B FALSE
#> 6 NR_2.0_61174.1 1H2C
                                   R FALSE
                             1
## To see only the structures containing Mg use
head(nrlist[nrlist$Mg == TRUE, ])
     Equivalence_class pdb model chain
       NR_2.0_14136.11 1DFU
                           1 M TRUE
#> 2
#> 3
       NR_2.0_14136.12 1DFU
                           1
                                   N TRUE
#> 7
      NR_2.0_70420.1 1J1U
                             1
                                   B TRUE
#> 9
       NR_2.0_26264.1 1JID
                              1
                                    B TRUE
#> 19
       NR_2.0_18961.1 1WMQ
                              1
                                    D TRUE
        NR_2.0_19039.1 1WPU
                              1
```

In the example, the argument hetAtms is passed to the hasHetAtm function.

7.3 Analysing datasets

When you have your desired list of PDB IDs, the function pipeNucData and pipeProtNucData can extract multiple structural parameters for all the dataset.

```
## After the download is finished, this dataset is analised in less than 2 min
## (single core, intel i5 2.3Ghz). Set progressbar=TRUE to see the progress.
```

```
ntinfo <- pipeNucData(pdbID=nrlist$pdb,</pre>
                        model=nrlist$model,
                        chain=nrlist$chain,
                        progressbar=FALSE, cores=2)
  #> [1] "Download completed, saved in: /tmp/Rtmplr4rLQ"
  str(ntinfo)
  #> 'data.frame': 3130 obs. of 83 variables:
  #> $ ntID : int 1 2 3 4 5 6 7 8 9 10 ...
                        : chr "1C9S" "1C9S" "1C9S" "1C9S" ...
  #> $ pdbID
                        : chr "1" "1" "1" "1" ...
  #> $ model
  #> $ chain
                        : chr "W" "W" "W" "W" ...
                    : chr "101" "102" "103" "104" ...

: chr "?" "?" "?" ...

: Factor w/ 2 levels "pu","py": 1 1 2 1 1 1 1 2 1 1 ...

: chr "G" "A" "U" "G" ...
#> $ resno
  #> $ insert
```

```
#> $ chi_valid
                          : logi TRUE TRUE TRUE TRUE TRUE TRUE ...
#> $ kappa_valid
                          : logi FALSE FALSE FALSE FALSE FALSE ...
#> $ base_exists
                          : logi TRUE TRUE TRUE TRUE TRUE TRUE ...
#> $ eta_valid
                         : logi FALSE FALSE FALSE FALSE TRUE ...
#> $ theta_valid
                         : logi TRUE FALSE FALSE FALSE TRUE TRUE ...
#> $ dist.03p.post_P
                         : num 1.61 1.61 1.61 1.61 1.62 ...
#> $ dist.Clp.N9 : num 1.47 1.49 NA 1.49 1.49 ...
#> $ dist.C1p.N1
                        : num 5.08 5.18 1.49 5.08 5.19 ...
#> $ angle.P.05p.C5p : num 114 121 111 122 123 ...
#> $ angle.05p.C5p.C4p : num 108 112 112 106 105 ...
#> $ angle.C5p.C4p.C3p : num 115 119 114 115 113 ...
#> $ angle.C5p.C4p.O4p : num 107 106 108 106 108 ...
#> $ angle.C4p.C3p.O3p : num 110 111 109 112 111 ...
#> $ angle.C4p.C3p.C2p : num 102 104 103 103 104 ...
#> $ angle.C4p.O4p.C1p : num 109 110 109 109 111 ...
#> $ angle.C3p.O3p.post_P: num 120 122 119 122 118 ...
#> $ angle.C3p.C2p.C1p : num 103.2 99.7 102.6 101.2 102.6 ...
#> $ angle.C2p.C1p.O4p : num 109 108 105 107 107 ...
#> $ angle.03p.C3p.C2p : num 112 110 116 109 110 ...
#> $ angle.C3p.C2p.O2p : num 114 113 111 111 113 ...
#> $ angle.C1p.C2p.O2p : num 105 113 114 112 107 ...
                        : num NA 295 311 329 277 ...
#> $ alpha
#> $ beta
                        : num 161 190 175 172 158 ...
#> $ gamma
                        : num 52 38.2 51.5 22.4 24.4 ...
#> $ delta
                          : num 92.5 93.3 133.7 146.3 94.6 ...
                   279.2 300.9 141 85.2 271.8 ...

: num    0.863 12.774 333.108 338.676 5

: num    339.5 329.2 36.9 33.7 336 ...

: num    30.4 36.7 328.2 327.9 32.3 ...

: num    329 329.3 17 21.2 330.4 ...

: num    19.17 11.39 6.21 359.84 15
                      : num 211 226 209 245 215 ...
#> $ epsilon
#> $ zeta
                        : num 0.863 12.774 333.108 338.676 5.655 ...
#> $ nu0
#> $ nu1
#> $ nu2
#> $ nu3
#> $ nu4
                         : num 19.17 11.39 6.21 359.84 15.11 ...
#> $ kappa
                         : num NA NA NA NA NA NA NA NA NA ...
#> $ eta
                        : num NA 174 179 216 151 ...
#> $ theta
                        : num 217.1 240 113.3 91.8 218.9 ...
#> $ chi
                          : num 259 203 254 246 266 ...
#> $ pu_amp
                         : num 32.7 37.6 37 35 33.2 ...
                      : num 16.8 359.13 152.17 161.72 8.53 ...
#> $ pu_phase
                          : num 4.35 4.75 2.26 1.58 4.49 ...
#> $ Dp
```

In this case, pipeNucData is downloading each structure from internet, but you can also download them manually before execution, save them in a directory and provide it to the function. When reading files locally, the function is obviously faster. The file format can either be pdb or cif, being the cif files the recommended option (they can also be read when compressed in .gz). To make it work with local files, the extension should also be provided (e.g. ".cif.gz", ".cif" or ".pdb").

```
path="/your/path/to/the/dataset/", extension=".cif.gz")
```

The second dataset analysis available in the dataset is designed to study the interactions between proteins and nucleic acids, and would be executed as follows:

```
## After the download is finished, this dataset is analised in less than 1 min
## (single core, intel i5 2.3Ghz). Set progressbar=TRUE to see the progress.
aantinfo <- pipeProtNucData(pdbID=nrlist$pdb,</pre>
                          model=nrlist$model,
                          chain=nrlist$chain,
                          progressbar=FALSE, cores=2)
#> [1] "Download completed, saved in: /tmp/Rtmplr4rLQ"
str(aantinfo)
#> 'data.frame': 85187 obs. of 23 variables:
#> $ pdbID : Factor w/ 169 levels "1C9S", "1DFU",...: 1 1 1 1 1 1 1 1 1 1 1 ...
#> $ model : Factor w/ 1 level "1": 1 1 1 1 1 1 1 1 1 1 ...
#> $ eleno_A : num 1 2 3 4 5 6 7 8 9 10 ...
#> $ eleno_B : num 7303 7319 7303 7317 7303 ...
#> $ distance : num 8.57 8.41 9.13 7.25 6.36 ...
#> $ elety_A : chr "P" "OP1" "OP2" "05'" ...
#> $ resid_A : chr "G" "G" "G" "G" ...
#> $ resno_A : int 101 101 101 101 101 101 101 101 101 ...
#> $ chain_A : chr "W" "W" "W" "W" ...
#> $ insert_A : chr "?" "?" "?" "?" ...
#> $ alt_A : chr "." "." "." ...
#> $ entid_A : int 1 1 1 1 1 1 1 1 1 ...
#> $ b_A : num 45.1 44.4 46.2 46.6 46.5 ...
#> $ asym_id_A: chr "A" "A" "A" "A" ...
#> $ elety_B : chr "CB" "CE2" "CB" "CD2" ...
#> $ resid_B : chr "ARG" "PHE" "ARG" "PHE" ...
#> $ resno_B : int 31 32 31 31 31 32 32 32 ...
#> $ chain_B : chr "L" "L" "L" "L" ...
#> $ insert_B : chr "?" "?" "?" "?" ...
#> $ alt_B : chr "." "." "." ...
#> $ entid_B : int 2 2 2 2 2 2 2 2 2 2 ...
#> $ b_B : num 33.7 32.5 33.7 36.6 33.7 ...
#> $ asym_id_B: chr "M" "M" "M" "M" ...
```

References

Bottaro, S., G. Bussi, G. Pinamonti, S. Reiber, W. Boomsma, and K. Lindorff-Larsen. 2018. "Barnaba: Software for Analysis of Nucleic Acids Structures and Trajectories." *RNA*, no. Epub (November).

Bottaro, S., F. Di Palma, and G Bussi. 2014. "The Role of Nucleobase Interactions in RNA Structure and Dynamics." *Nucleic Acids Research* 42 (21): 13306–14.

Darré, L., I. Ivani, P.D. Dans, H. Gómez, A. Hospital, and M. Orozco. 2016. "Small Details Matter: The 2'-Hydroxyl as a Conformational Switch in RNA." *Journal of the American Chemical Society* 138 (50): 16355–63.

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.

Jain, S., D.C. Richardson, and J.S. Richardson. 2015. "Computational Methods for RNA Structure Validation and Improvement." In *Methods in Enzymology*, edited by S.A. Woodson and F.H.T. Allain, 558:181–212. Elsevier.

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.

Lu, X.J., H.J. Bussemaker, and W.K. Olson. 2015. "DSSR: An Integrated Software Tool for Dissection the Spatial Structure of RNA." *Nucleic Acids Research* 43 (21): e142.

Murray, L.J.W., W.B. Arendall III, D.C. Richardson, and J.S. Richardson. 2003. "RNA Backbone Is Rotameric." *Proceedings of the National Academy of Sciences of the United States of America* 100 (24): 13904–9.

Velankar, S., G. Van Ginkel, Y. Alhroub, G.M. Battle, J.M. Berrisford, M.J. Conroy, J.M. Dana, et al. 2015. "PDBe: Improved Accessibility of Macromolecular Structure Data from PDB and EMDB." *Nucleic Acids Research* 44 (D1): D385–395.

Wang, J., Y. Zhao, C. Zhu, and Y. Xiao. 2015. "3dRNAscore: A Distance and Torsion Angle Dependent Evaluation Function of 3D RNA Structures." *Nucleic Acids Research* 43 (10): e63.

Wang, Z., and J. Xu. 2011. "A Conditional Random Fields Method for RNA Sequence-Structure Relationship Modeling and Conformation Sampling." *Bioinformatics* 27 (13): i102–10.