

Class 9: Structural bioinformatics 1

Derek Chang (PID: 16942232)

The main database for structural data is called PDB (Protein Data Bank). Let's see what it contains:

Data: <https://tinyurl.com/pdbstats24>

```
pdbdb <- read.csv("pdb_stats.csv")
pdbdb
```

	Molecular.Type	X.ray	EM	NMR	Multiple.methods	Neutron	Other
1	Protein (only)	167,192	15,572	12,529	208	77	32
2	Protein/Oligosaccharide	9,639	2,635	34	8	2	0
3	Protein/NA	8,730	4,697	286	7	0	0
4	Nucleic acid (only)	2,869	137	1,507	14	3	1
5	Other	170	10	33	0	0	0
6	Oligosaccharide (only)	11	0	6	1	0	4
	Total						
1		195,610					
2		12,318					
3		13,720					
4		4,531					
5		213					
6		22					

Questions:

Q1: What percentage of structures in the PDB are solved by X-Ray and Electron Microscopy.

```
pdbdb$Total
```

```
[1] "195,610" "12,318" "13,720" "4,531" "213" "22"
```

I need to remove the commas and convert to numeric to do math:

```
as.numeric(sub(",", "", pdbdb$Total))
```

```
[1] 195610 12318 13720 4531 213 22
```

I could turn this snippet into a function to fix any future table i read like this:

```
x <- pdbdb$Total
as.numeric(sub(",", "", x))
```

```
[1] 195610 12318 13720 4531 213 22
```

```
comma2numeric <- function(x) {
  as.numeric(sub(",", "", x))
}
```

Test it

```
pdbdbnew<- comma2numeric(pdbdb$X.ray)
```

##Or try a different read/import function:

```
library(readr)
pdbdb <- read_csv("pdb_stats.csv")
```

```
Rows: 6 Columns: 8
```

```
-- Column specification -----
```

```
Delimiter: ","
```

```
chr (1): Molecular Type
```

```
dbl (3): Multiple methods, Neutron, Other
```

```
num (4): X-ray, EM, NMR, Total
```

```
i Use `spec()` to retrieve the full column specification for this data.
```

```
i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

```
sum(pdbdb$Total)
```

```
[1] 226414
```

```
sum(pdbdb$`X-ray`)/sum(pdbdb$Total) * 100
```

```
[1] 83.30359
```

```
sum(pdbdb$EM)/sum(pdbdb$Total) * 100
```

```
[1] 10.18091
```

83.3% of the structures are solved by X-ray and 10.18% of the structures are solved by EM.

Q2: What proportion of structures in the PDB are protein?

```
pdbdb
```

```
# A tibble: 6 x 8
  `Molecular Type`  `X-ray`    EM    NMR `Multiple methods` Neutron Other  Total
  <chr>             <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1 Protein (only)    167192 15572 12529      208      77      32 195610
2 Protein/Oligosacc~ 9639 2635 34      8        2        0 12318
3 Protein/NA        8730 4697 286      7        0        0 13720
4 Nucleic acid (onl~ 2869 137 1507     14        3        1 4531
5 Other             170 10 33      0        0        0 213
6 Oligosaccharide (~ 11 0 6      1        0        4 22
```

```
library(dplyr)
```

Attaching package: 'dplyr'

The following objects are masked from 'package:stats':

```
filter, lag
```

The following objects are masked from 'package:base':

```
intersect, setdiff, setequal, union
```

```
(sum(pdbdb$Total[1])/sum(pdbdb$Total)) *100
```

```
[1] 86.39483
```

86.39483% of the structures in the PDB are protein.

Q3: Type HIV in the PDB website search box on the home page and determine how many HIV-1 protease structures are in the current PDB?

In the current PDB there are 226414 HIV-1 Protease Structures.

2. Mol*

Mol* is a new web-based molecularviewer that we will use to learn the basics of here.

Q4: Water molecules normally have 3 atoms. Why do we see just one atom per water molecule in this structure?

We see just one atom per water molecule as this allows for better clarity of the water molecules and allows us to see other things better.

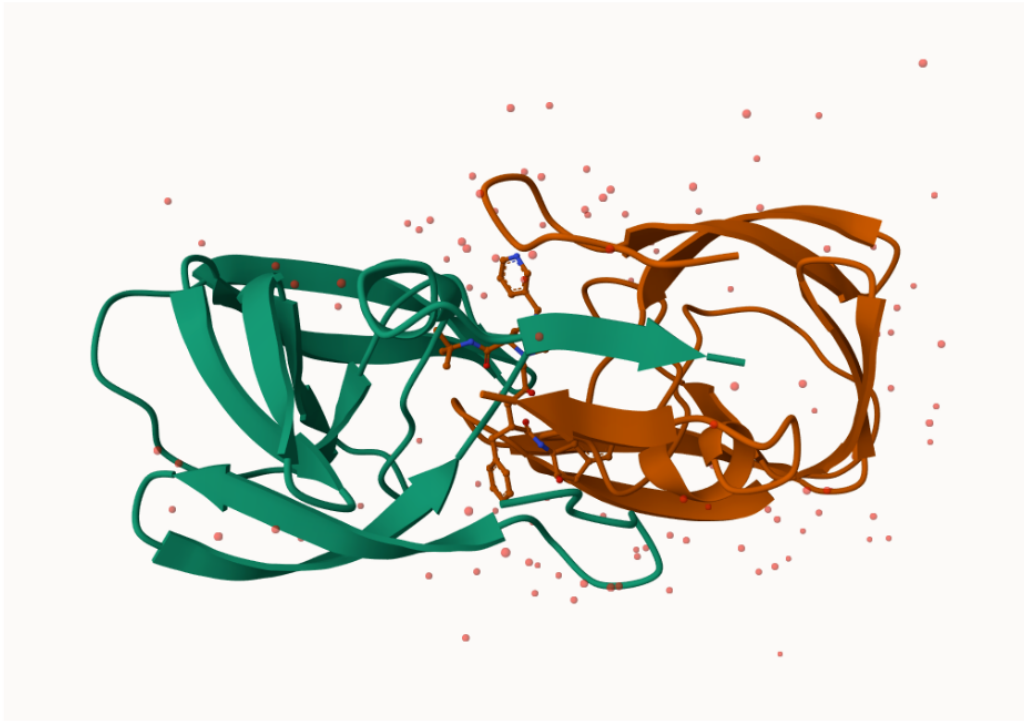


Figure 1: A first image from molstar

Q5: There is a critical “conserved” water molecule in the binding site. Can you identify this water molecule? What residue number does this water molecule have

The conserved water molecule in the binding site seems to be residue number 308.

Q6: Generate and save a figure clearly showing the two distinct chains of HIV-protease along with the ligand. You might also consider showing the catalytic residues ASP 25 in each chain and the critical water (we recommend “Ball & Stick” for these side-chains). Add this figure to your Quarto document.



Figure 2: A second image from molstar showing the important ASP25 amino acid molecules.

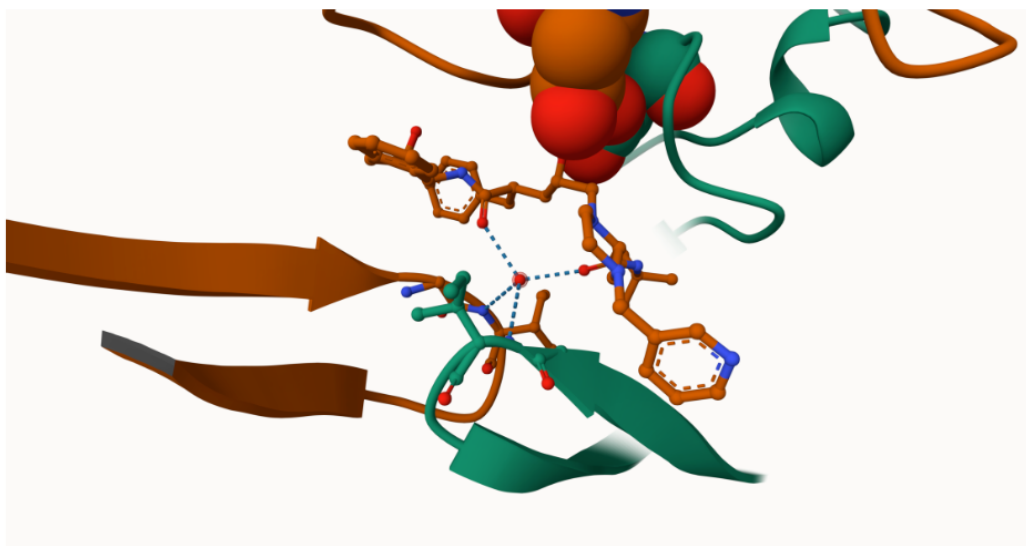


Figure 3: Image showing the two distinct chains of HIV-protease and the critical conserved water residue 308

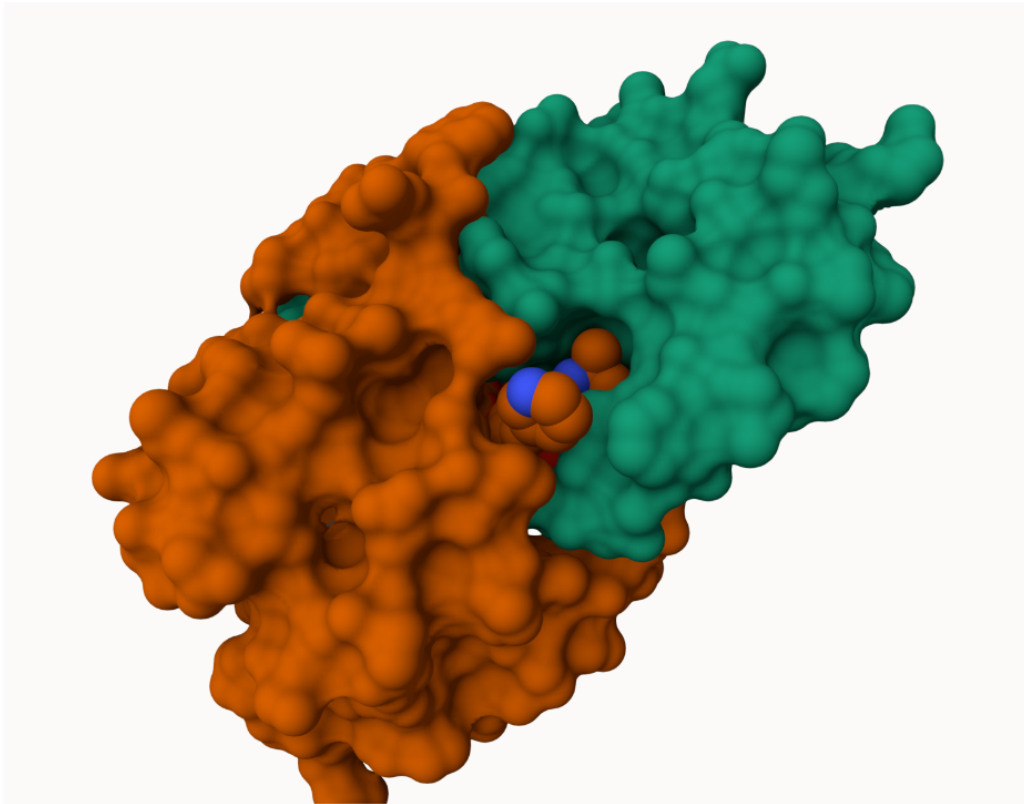


Figure 4: A third surface display image from molstar showing the Merk compound in the binding pocket

The Bio3D package

The bio3d package allows us to do all sorts of structural bioinformatics work in R.

Let's start with how it can read these PDB files:

```
library(bio3d)
pdb <- read.pdb("1HSG")
```

Note: Accessing on-line PDB file

`pdb`

```
Call: read.pdb(file = "1HSG")
```

```
Total Models#: 1
```

```
Total Atoms#: 1686, XYZs#: 5058 Chains#: 2 (values: A B)
```

```
Protein Atoms#: 1514 (residues/Calpha atoms#: 198)
```

```
Nucleic acid Atoms#: 0 (residues/phosphate atoms#: 0)
```

```
Non-protein/nucleic Atoms#: 172 (residues: 128)
```

```
Non-protein/nucleic resid values: [ HOH (127), MK1 (1) ]
```

```
Protein sequence:
```

```
PQITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWKPKMIGGIGGFIKVRQYD
QILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNFPQITLWQRPLVTIKIGGQLKE
ALLDTGADDTVLEEMSLPGRWKPKMIGGIGGFIKVRQYDQILIEICGHKAIGTVLVGPTP
VNIIGRNLLTQIGCTLNF
```

```
+ attr: atom, xyz, seqres, helix, sheet,
      calpha, remark, call
```

`attributes(pdb)`

```
$names
```

```
[1] "atom" "xyz" "seqres" "helix" "sheet" "calpha" "remark" "call"
```

```
$class
```

```
[1] "pdb" "sse"
```

`head(pdb$atom)`

	type	eleno	elety	alt	resid	chain	resno	insert	x	y	z	o	b
1	ATOM	1	N	<NA>	PRO	A	1	<NA>	29.361	39.686	5.862	1	38.10
2	ATOM	2	CA	<NA>	PRO	A	1	<NA>	30.307	38.663	5.319	1	40.62
3	ATOM	3	C	<NA>	PRO	A	1	<NA>	29.760	38.071	4.022	1	42.64
4	ATOM	4	O	<NA>	PRO	A	1	<NA>	28.600	38.302	3.676	1	43.40
5	ATOM	5	CB	<NA>	PRO	A	1	<NA>	30.508	37.541	6.342	1	37.87
6	ATOM	6	CG	<NA>	PRO	A	1	<NA>	29.296	37.591	7.162	1	38.40

```

      segid elesy charge
1  <NA>      N  <NA>
2  <NA>      C  <NA>
3  <NA>      C  <NA>
4  <NA>      O  <NA>
5  <NA>      C  <NA>
6  <NA>      C  <NA>

```

```

pdbseq(pdb)

```

```

  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20
"P" "Q" "I" "T" "L" "W" "Q" "R" "P" "L" "V" "T" "I" "K" "I" "G" "G" "Q" "L" "K"
21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40
"E" "A" "L" "L" "D" "T" "G" "A" "D" "D" "T" "V" "L" "E" "E" "M" "S" "L" "P" "G"
41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60
"R" "W" "K" "P" "K" "M" "I" "G" "G" "I" "G" "G" "F" "I" "K" "V" "R" "Q" "Y" "D"
61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80
"Q" "I" "L" "I" "E" "I" "C" "G" "H" "K" "A" "I" "G" "T" "V" "L" "V" "G" "P" "T"
81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 1
"P" "V" "N" "I" "I" "G" "R" "N" "L" "L" "T" "Q" "I" "G" "C" "T" "L" "N" "F" "P"
  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 21
"Q" "I" "T" "L" "W" "Q" "R" "P" "L" "V" "T" "I" "K" "I" "G" "G" "Q" "L" "K" "E"
22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41
"A" "L" "L" "D" "T" "G" "A" "D" "D" "T" "V" "L" "E" "E" "M" "S" "L" "P" "G" "R"
42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61
"W" "K" "P" "K" "M" "I" "G" "G" "I" "G" "G" "F" "I" "K" "V" "R" "Q" "Y" "D" "Q"
62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81
"I" "L" "I" "E" "I" "C" "G" "H" "K" "A" "I" "G" "T" "V" "L" "V" "G" "P" "T" "P"
82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99
"V" "N" "I" "I" "G" "R" "N" "L" "L" "T" "Q" "I" "G" "C" "T" "L" "N" "F"

```

Q7: How many amino acid residues are there in this pdb object?

```

sum(pdb$calpha)

```

```

[1] 198

```

There are 198 amino acid residues in this pdb object

Q8: Name one of the two non-protein residues?

The two non-protein residues are HOH and Mk1.

Q9: How many protein chains are in this structure?

There are two protein chains, chain A and B

```
unique(pdb$atom$chain)
```

```
[1] "A" "B"
```

##Predicting functional motions of a single structure

Let's do bioinformatics prediction of functional motion - i.e. the movements that one of these molecules needs to make to do its stuff.

```
adk <- read.pdb("6s36")
```

Note: Accessing on-line PDB file

PDB has ALT records, taking A only, rm.alt=TRUE

```
adk
```

```
Call: read.pdb(file = "6s36")
```

```
Total Models#: 1
```

```
Total Atoms#: 1898, XYZs#: 5694 Chains#: 1 (values: A)
```

```
Protein Atoms#: 1654 (residues/Calpha atoms#: 214)
```

```
Nucleic acid Atoms#: 0 (residues/phosphate atoms#: 0)
```

```
Non-protein/nucleic Atoms#: 244 (residues: 244)
```

```
Non-protein/nucleic resid values: [ CL (3), HOH (238), MG (2), NA (1) ]
```

```
Protein sequence:
```

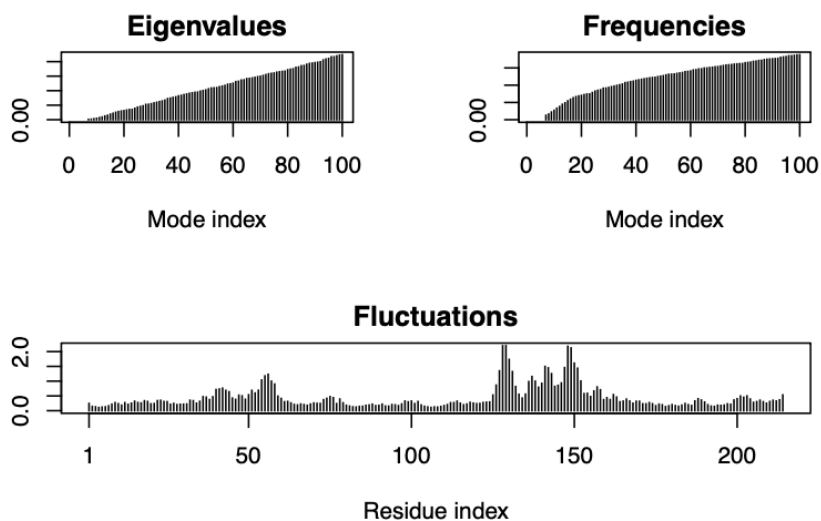
```
MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGMLRAAVKSGSELGKQAKDIMDAGKLV  
DELVIALVKERIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFDVPDELIVDKI  
VGRRVHAPSGRVYHVKNPPKVEGKDDVTGEELTRKDDQEETVRKRLVEYHQMTAPLIG  
YYSKEAEAGNTKYAKVDGTPVAEVRADLEKILG
```

```
+ attr: atom, xyz, seqres, helix, sheet,  
      calpha, remark, call
```

```
# Perform flexibility prediction  
m <- nma(adk)
```

```
Building Hessian...      Done in 0.01 seconds.  
Diagonalizing Hessian... Done in 0.229 seconds.
```

```
plot(m)
```



We will write out multi-model trajectory PDB files that we can make animations of predicted motions from.

```
mktrj(m, file = "adk.pdb")
```

We can open this file in mol* to play the trajectory.