

Class 9: Structural Bioinformatics pt. 1

Ellice Wang (PID: A16882742)

2025-02-06

Table of contents

The PDB database	1
2. Using Mol*	5
3. Introduction to Bio3D in R	8
4. Predicting functional dynamics	10
Setup	12
Search and retrieve ADK structures	13

The PDB database

The main repository of biomolecular structure data is called the PDB found at: <https://rcsb.org>

Let's see what this database contains. I went to PDB > Analyze > PDB Statistics > By Exp method and molecular type.

```
# load in database
pdb_stats <- read.csv("pdb_data_dist.csv")
pdb_stats
```

	Molecular.Type	X.ray	EM	NMR	Multiple.methods	Neutron	Other
1	Protein (only)	169,563	16,774	12,578	208	81	32
2	Protein/Oligosaccharide	9,939	2,839	34	8	2	0
3	Protein/NA	8,801	5,062	286	7	0	0
4	Nucleic acid (only)	2,890	151	1,521	14	3	1
5	Other	170	10	33	0	0	0
6	Oligosaccharide (only)	11	0	6	1	0	4
	Total						

```

1 199,236
2 12,822
3 14,156
4 4,580
5 213
6 22

```

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

82.83549% of structures in the PDB are solved by X-Ray while 10.75017% are solved by Electron microscopy.

```
pdb_stats$X.ray
```

```
[1] "169,563" "9,939" "8,801" "2,890" "170" "11"
```

The comma in these numbers is causing them to be read as character rather than numeric.

```

x <- pdb_stats$X.ray
sum(as.numeric(sub(",", "", x)))

```

```
[1] 191374
```

Or I can use the **readr** package and the `read_csv()` function.

```

library(readr)
pdb_stats <- read_csv("pdb_data_dist.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.
```

```
pdb_stats
```

```
# 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)    169563 16774 12578      208      81    32 199236
2 Protein/Oligosacc~ 9939  2839   34       8       2     0  12822
3 Protein/NA        8801  5062  286       7       0     0  14156
4 Nucleic acid (onl~ 2890   151 1521      14       3     1   4580
5 Other             170    10   33       0       0     0    213
6 Oligosaccharide (~ 11     0    6       1       0     4    22
```

```
colnames(pdb_stats)
```

```
[1] "Molecular Type"  "X-ray"           "EM"              "NMR"
[5] "Multiple methods" "Neutron"         "Other"           "Total"
```

```
library(janitor)
```

Warning: package 'janitor' was built under R version 4.4.1

Attaching package: 'janitor'

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

chisq.test, fisher.test

```
df <- clean_names(pdb_stats)
df
```

```
# 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)    169563 16774 12578      208      81    32 199236
2 Protein/Oligosacchar~ 9939  2839   34       8       2     0  12822
3 Protein/NA        8801  5062  286       7       0     0  14156
4 Nucleic acid (only)  2890   151 1521      14       3     1   4580
5 Other             170    10   33       0       0     0    213
6 Oligosaccharide (onl~ 11     0    6       1       0     4    22
```

Total number of X-ray structures

```
x_ray_tot <- sum(df$x_ray)
em_tot <- sum(df$em)
```

Total number of structures

```
tot_structure <- sum(df$total)
```

Proportion of X-ray structures and electron microscopy

```
x_ray_tot/tot_structure * 100
```

```
[1] 82.83549
```

```
em_tot/tot_structure * 100
```

```
[1] 10.75017
```

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

86.23852% of structures in the PDB are protein.

```
df[df$molecular_type == "Protein (only)",]$total/(tot_structure) * 100
```

```
[1] 86.23852
```

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?

There are 4,683 structures in the current PDB.

2. Using Mol*

The main Mol* homepage at: <https://molstar.org/viewer/> We can input our own PDB files or just give it a PDB database accession code (4 letter PDB code)

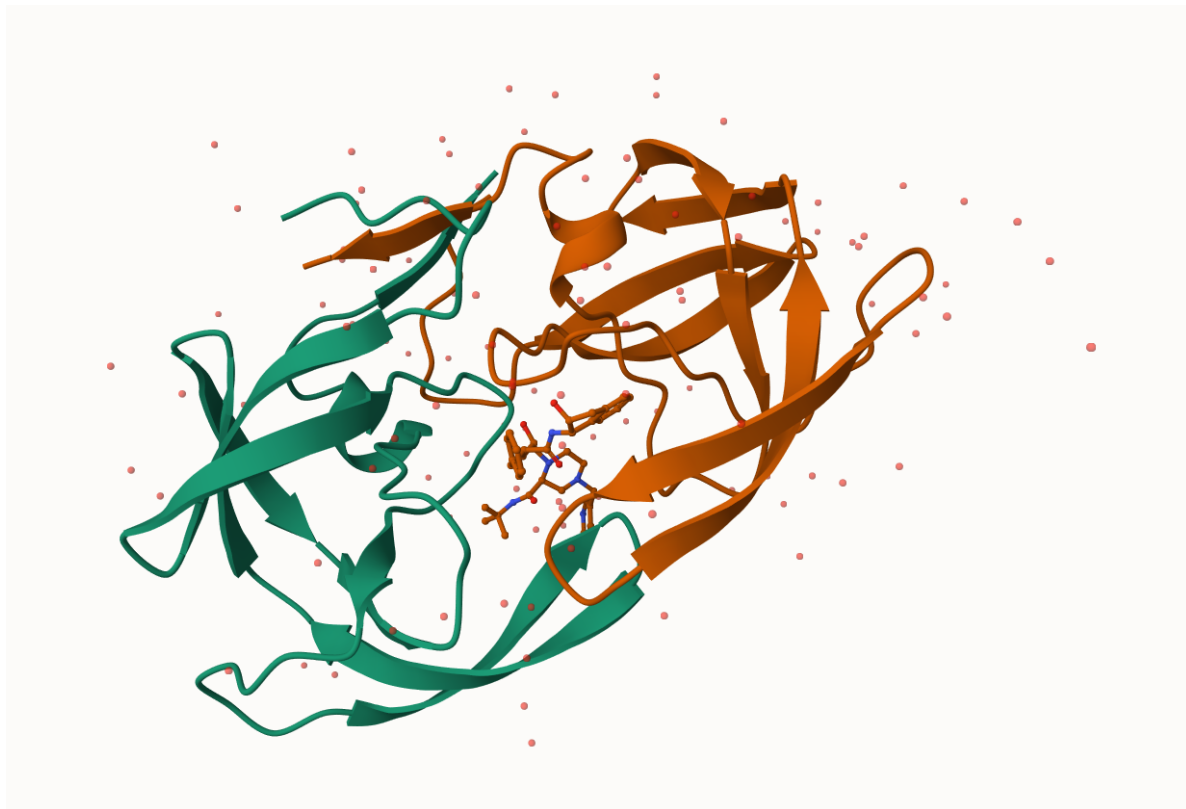


Figure 1: Molecular view of 1HSG

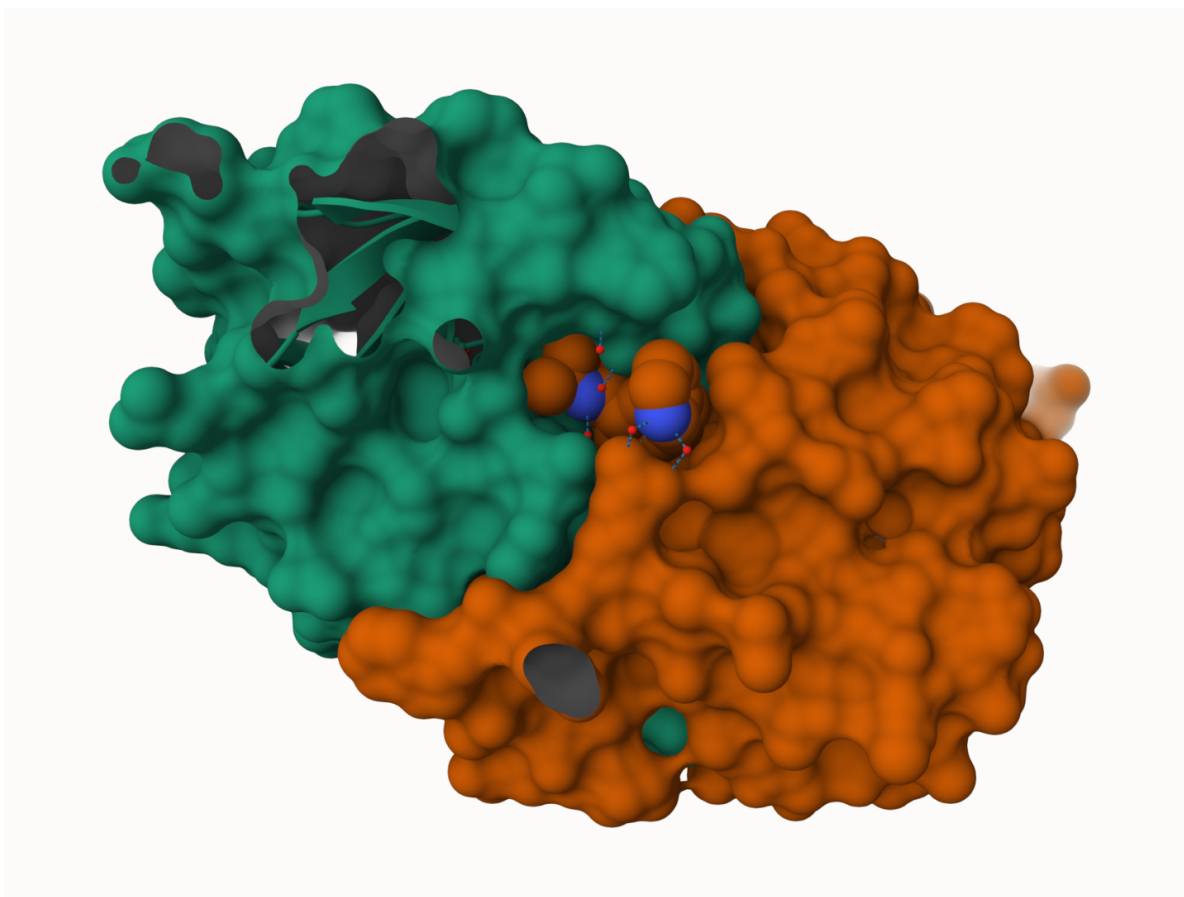


Figure 2: Surface representation

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

We only see 1 atom per water molecule in this structure because it is only showing the oxygen molecule of the water.

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

This water molecule has residue number 308.

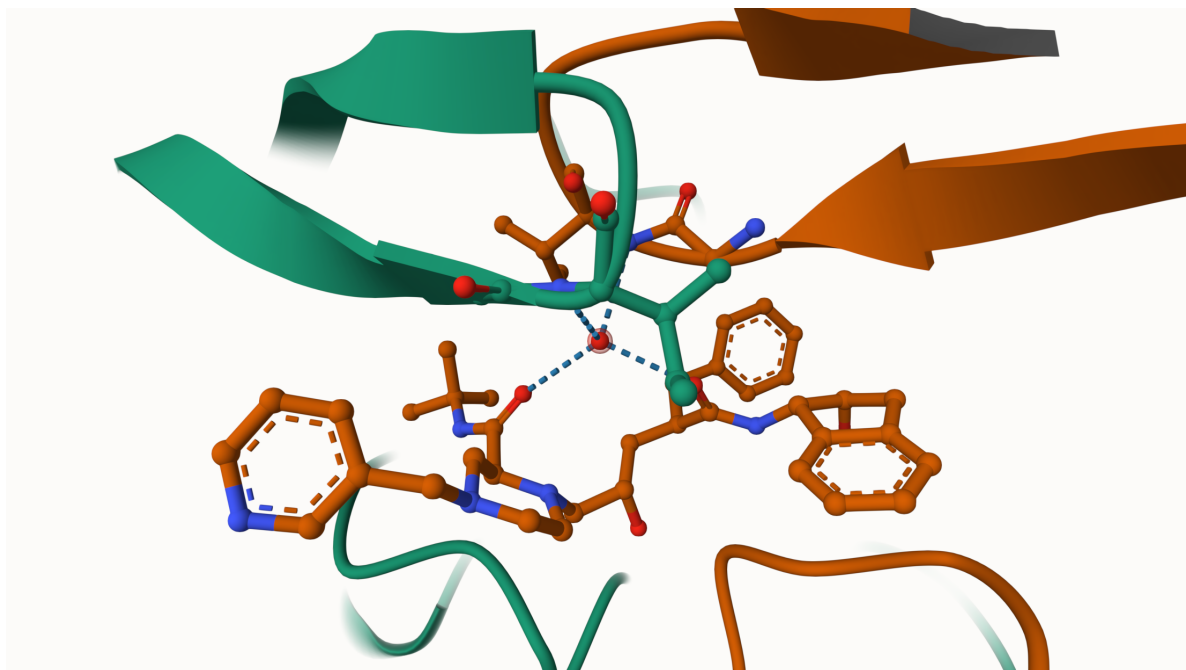


Figure 3: Water 308

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 4: The important D25 amino acids

3. Introduction to Bio3D in R

We can use the **bio3d** package for structural bioinformatics to read PDB data into R

```
library(bio3d)
```

Warning: package 'bio3d' was built under R version 4.4.1

```
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

```

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

```
length(pdbseq(pdb))
```

```
[1] 198
```

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

MK1

Q9: How many protein chains are in this structure?

2 protein chains: Chain A & B

Looking at the 'pdb' object in more detail

```
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>

Let's try a new function not yet in the bio3d package. It requires the **r3dmol** package that we need to install with 'install.packages("r3dmol")'

```
library(r3dmol)
source("https://tinyurl.com/viewpdb")
# view.pdb(pdb, backgroundColor = "white")
```

4. Predicting functional dynamics

We can use the 'nma()' function in bio3d to predict the large-scale functional motions of biomolecules.

```
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:

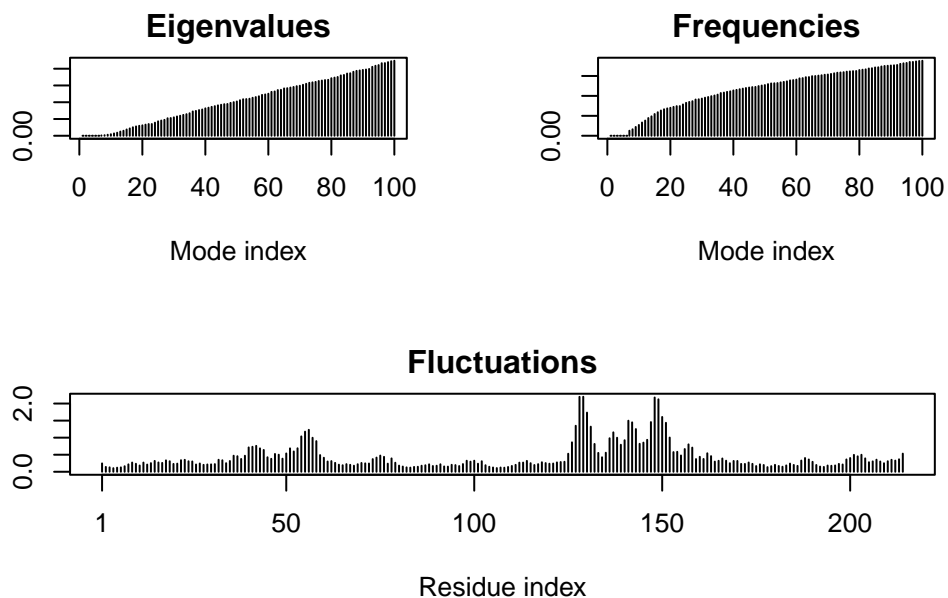
```
MRILLGAPGAGKGTQAQFIMEKYGIPQISTGDMRLRAAVKSGSELGKQAKDIMDAGKLV
DELVIALVKERIAQEDCRNGFLDGFPRTPQADAMKEAGINVDYVLEFDVPDELIVDKI
VGRRVHAPSGRVYHVKFNPPKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQMTAPLIG
YYSKEAEAGNTKYAKVDGTPVAEVRADLEKILG
```

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

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

```
Building Hessian...      Done in 0.057 seconds.
Diagonalizing Hessian... Done in 0.6 seconds.
```

```
plot(m)
```



Write out a trajectory of predicted molecule motion:

```
mktrj(m, file="adk_m7.pdb")
```

Load file into Mol*

Setup

Q10. Which of the packages above is found only on BioConductor and not CRAN?

msa is not found on CRAN

Q11. Which of the above packages is not found on BioConductor or CRAN?

bio3d-view is not found on either BioConductor or CRAN

Q12. True or False? Functions from the devtools package can be used to install packages from GitHub and BitBucket?

TRUE

Search and retrieve ADK structures

```
aa <- get.seq("lake_A")
```

Warning in get.seq("lake_A"): Removing existing file: seqs.fasta

Fetching... Please wait. Done.

```
aa
```

```
      1      .      .      .      .      .      60
pdb|1AKE|A  MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMLRAAVKSGSELGKQAKDIMDAGKLV
      1      .      .      .      .      .      60

      61      .      .      .      .      .      120
pdb|1AKE|A  DELVIALVKERIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFDVPDELIVDRI
      61      .      .      .      .      .      120

     121      .      .      .      .      .      180
pdb|1AKE|A  VGRRVHAPSGRVYHVKNPPKVEGKDDVTGEELTRKDDQEETVRKRLVEYHQMTPALIG
     121      .      .      .      .      .      180

     181      .      .      .      214
pdb|1AKE|A  YYSKEAEAGNTKYAKVDGTPVAEVRADLEKILG
     181      .      .      .      214
```

Call:

```
read.fasta(file = outfile)
```

Class:

```
fasta
```

Alignment dimensions:

```
1 sequence rows; 214 position columns (214 non-gap, 0 gap)
```

```
+ attr: id, ali, call
```

Q13. How many amino acids are in this sequence, i.e. how long is this sequence?

There are 214 amino acids in this sequence.