

Class 9: Structural Bioinformatics Pt. 1

Youn Soo Na (PID: A17014731)

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

Data from: <https://www.rcsb.org/stats> Or from alternate link: <https://tinyurl.com/pdbstats24>

```
pdb24 <- read.csv("pdb_stats.csv", row.names=1)
head(pdb24)
```

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

```
# Some of the "numeric" values are actually characters
# We need to change them to numeric values.
# as.numeric( sub(",", "", pdb24$Total))
# I could run this into a function to fix the whole table or any future table
# I read like this:
# x <- pdb24$Total
# as.numeric( sub(",", "", x) )
```

```

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

# Test it.
# comma2numeric(pdb24$X.ray)
# head(pdb24)

pdb24test <- apply(pdb24, 2, comma2numeric)
head(pdb24test)

```

	X.ray	EM	NMR	Multiple.methods	Neutron	Other	Total
[1,]	167192	15572	12529	208	77	32	195610
[2,]	9639	2635	34	8	2	0	12318
[3,]	8730	4697	286	7	0	0	13720
[4,]	2869	137	1507	14	3	1	4531
[5,]	170	10	33	0	0	0	213
[6,]	11	0	6	1	0	4	22

```

## 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
```

And answer the following questions:

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

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

```
[1] 93.4845
```

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

```
# library(dplyr)
# pdbdb %>%
#   filter(rowSums(sapply(., function(x) grepl("protein", x, ignore.case = TRUE))) > 0)
pdbdb$Total[1]/sum(pdbdb$Total) * 100
```

```
[1] 86.39483
```

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?

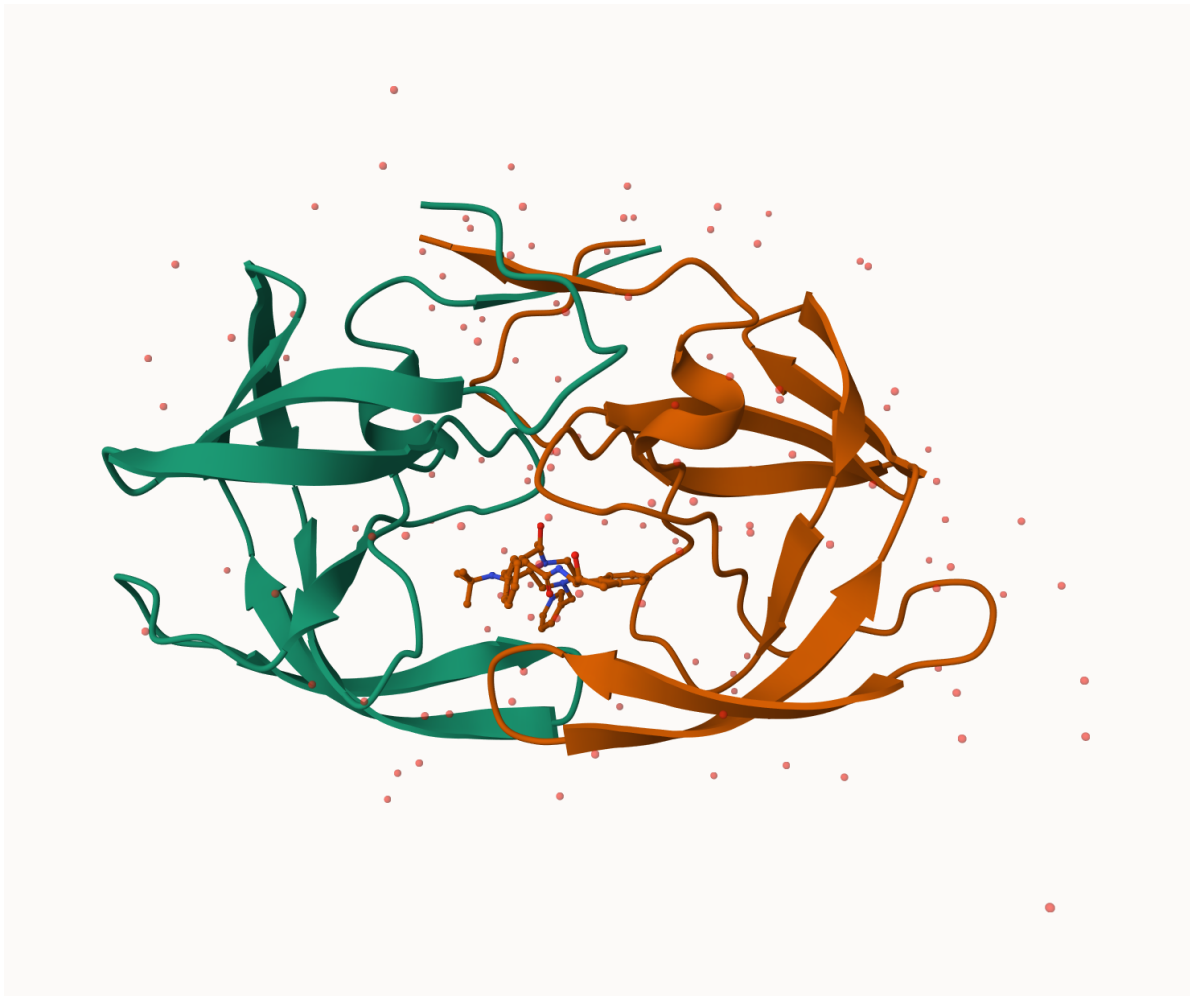
```
4553
```

Mol*

Mol* (pronounced “molstar”) is a new web-based molecular viewer that we will need to learn the basics of here.

<https://molstar.org/viewer/>

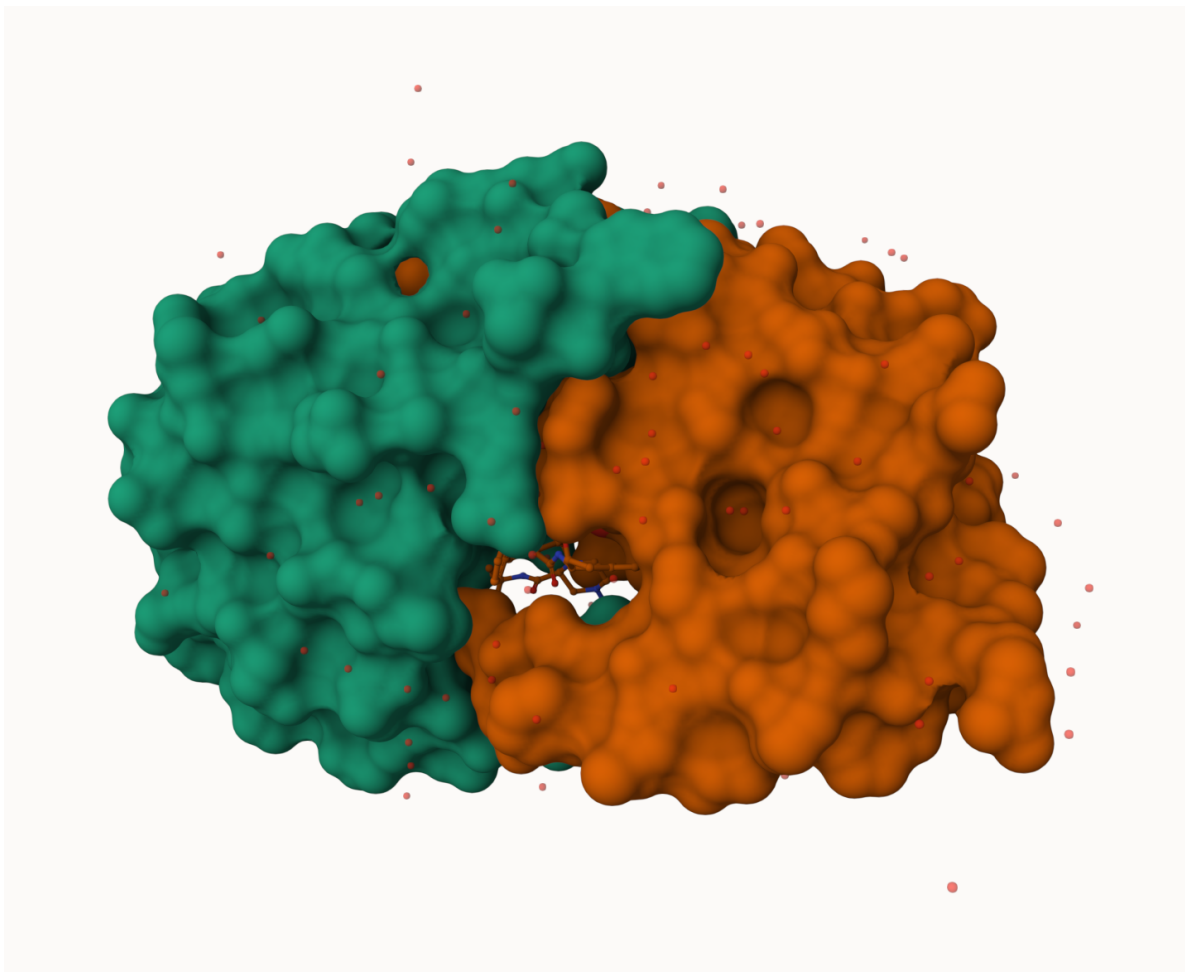
We will use PDB code: 1HSG

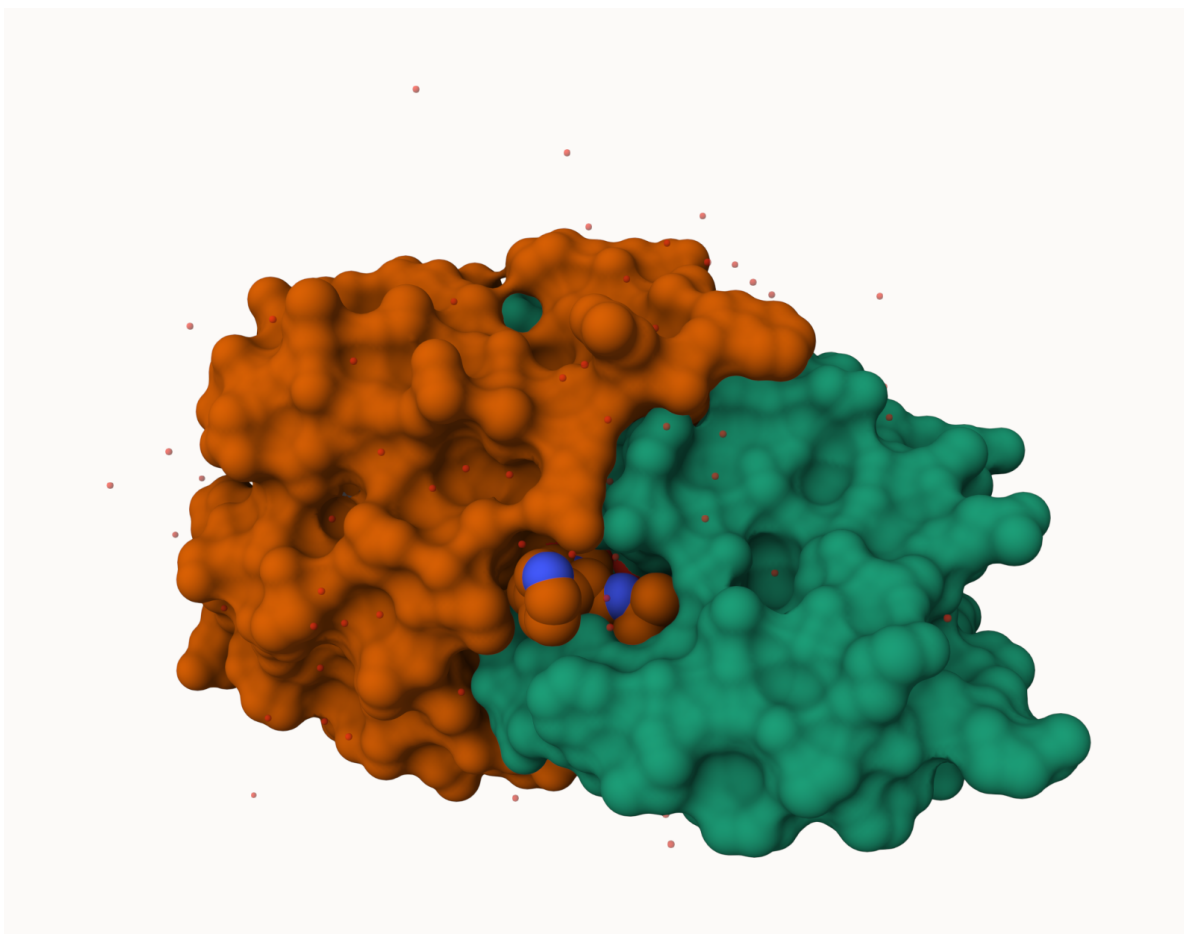


Note: This is an aspartic protease that uses 2 aspartic acid

Here are some more custome images:







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

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 water molecular is labeled HOH 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.

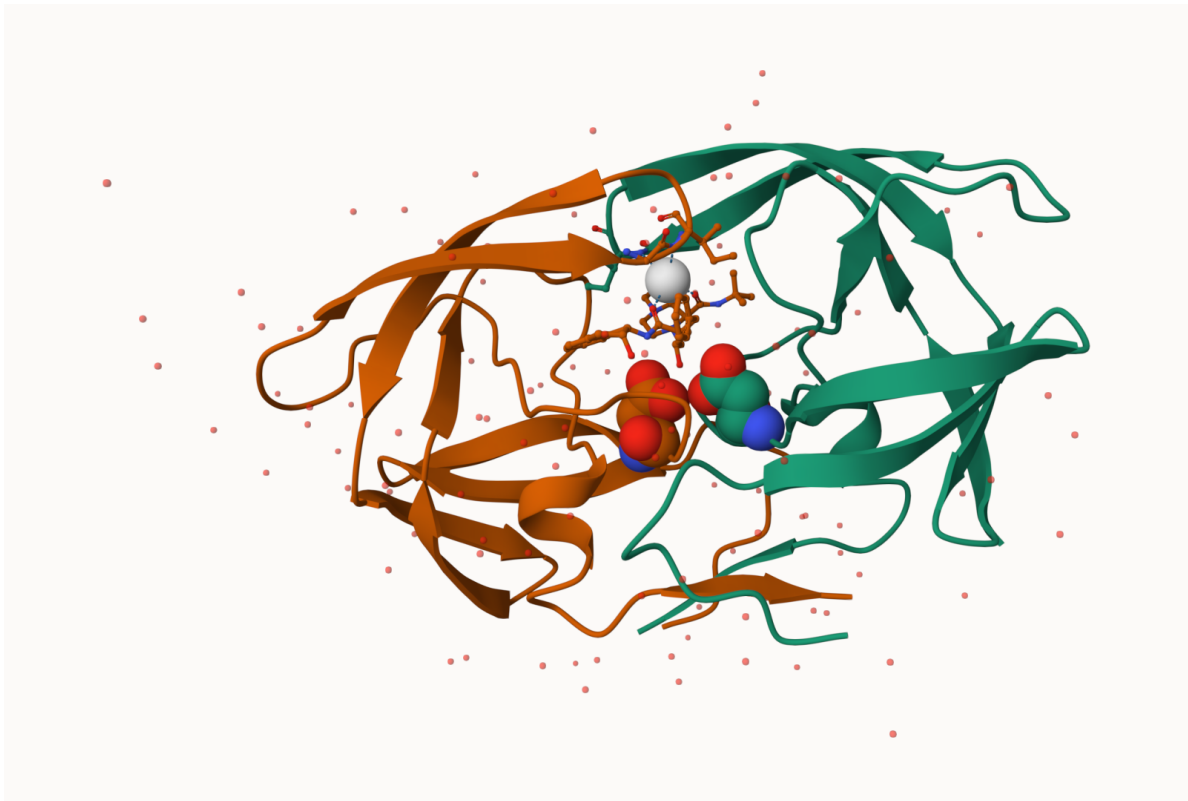


Figure 1: Water molecule is represented as a white ball

Discussion Topic: Can you think of a way in which indinavir, or even larger ligands and substrates, could enter the binding site?

Q7: [Optional] As you have hopefully observed HIV protease is a homodimer (i.e. it is composed of two identical chains). With the aid of the graphic display can you identify secondary structure elements that are likely to only form in the dimer rather than the monomer?

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

```
# get a quick summary
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
```

```
# MK1 is Merk1 ligand
```

```
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"

```

```

#What is position 25 of the protein sequence?

```

```

pdbseq(pdb)[25]

```

```

25
"D"

```

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

```
sum(pdb$calpha)
```

```
[1] 198
```

```
198
```

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

```
[1] 198
```

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

HOH or MK1

Q9: How many protein chains are in this structure?

2 chains

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

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

Predicting functional motions of a single structure

Let's do a bioinformatics prediction of functional motions - 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:
```

```
MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMLRAAVKSGSELGKQAKDIMDAGKLV  
TDELVIALVKERIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFDVPDELIVDKI  
VGRRVHAPSGRVYHVKFNPVKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQM  
TAPLIGYYSKEAEAGNTKYAKVDGTPVAEVRADLEKILG
```

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

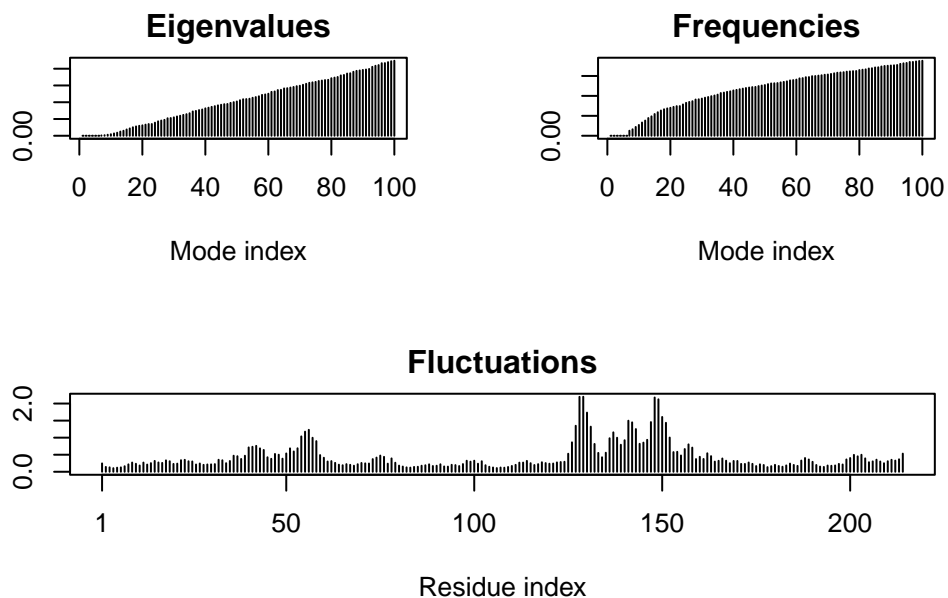
Perform a flexibility prediction

```
m <- nma(adk)
```

```
Building Hessian... Done in 0.014 seconds.
```

```
Diagonalizing Hessian... Done in 0.285 seconds.
```

```
plot(m)
```



Write out multi-model PDB file that we can use to make an animation of the predicted motions. To view a *movie* of these predicted motions we can generate a molecular trajectory using the `mktrj()` function

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

Now, I can open this in Mol* to play the movie.