

Class 10: Structural Bioinformatics Pt1

Jordan Prych (PID:A17080226)

Table of contents

1. The PDB Database	1
2. Using Mol*	5
3. Introduction to Bio3D in R	10
4. Predicting Functional Dynamics	13

1. The PDB Database

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

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

```
pdbstats <- read.csv("Data Export Summary.csv")
pdbstats
```

	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?

```
pdbstats$X.ray
```

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

These values are characters, not numeric, so you cannot do math with these characters. Commas make them characters.

I can fix this by replacing “,” for nothing “” with the `sub()` function:

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

```
[1] 191374
```

Use `install.packages()` to use the **readr** package and `read_csv()` function.

```
library(readr)
pdbstats <- read_csv("Data Export Summary.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.

```
pdbstats
```

A tibble: 6 x 8

	<code>`Molecular Type`</code> <chr>	<code>`X-ray`</code> <dbl>	EM <dbl>	NMR <dbl>	<code>`Multiple methods`</code> <dbl>	Neutron <dbl>	Other <dbl>	Total <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 (only)	2890	151	1521	14	3	1	4580
5 Other	170	10	33	0	0	0	213
6 Oligosaccharide (only)	11	0	6	1	0	4	22

I want to clean the column names as they are all lowercase and don't have spaces in them. Use the **janitor** package and `clean_names()` function.

```
library(janitor)
```

Attaching package: 'janitor'

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

`chisq.test`, `fisher.test`

```
df <- clean_names(pdbstats)
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 (only)  11     0     6           1       0     4    22
```

Total number of X-ray structures:

```
sum(df$x_ray)
```

```
[1] 191374
```

Total number of structures:

```
sum(df$total)
```

```
[1] 231029
```

percent:

```
sum(df$x_ray)/sum(df$total)*100
```

```
[1] 82.83549
```

percent of electron microscopy structures:

```
#total number of em structures  
sum(df$em)
```

```
[1] 24836
```

```
#percent  
sum(df$em)/sum(df$total)*100
```

```
[1] 10.75017
```

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

```
#total number of proteins  
sum(df[1:3, 8])
```

```
[1] 226214
```

```
#total number of structures  
sum(df$total)
```

```
[1] 231029
```

```
#proportion  
sum(df[1:3, 8])/sum(df$total)
```

```
[1] 0.9791585
```

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?

when searching the HIV-1 protease, 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).

The markdown code for inserting an image:



Figure 1: Molecular View of HSG

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

We only see one water molecule because in Mol*, hydrogen atoms and bonds with hydrogen are not represented. Therefore, since water has two hydrogen bonds, only the oxygen atom is shown 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.

This water molecule is at residue number 308. Shown in images below

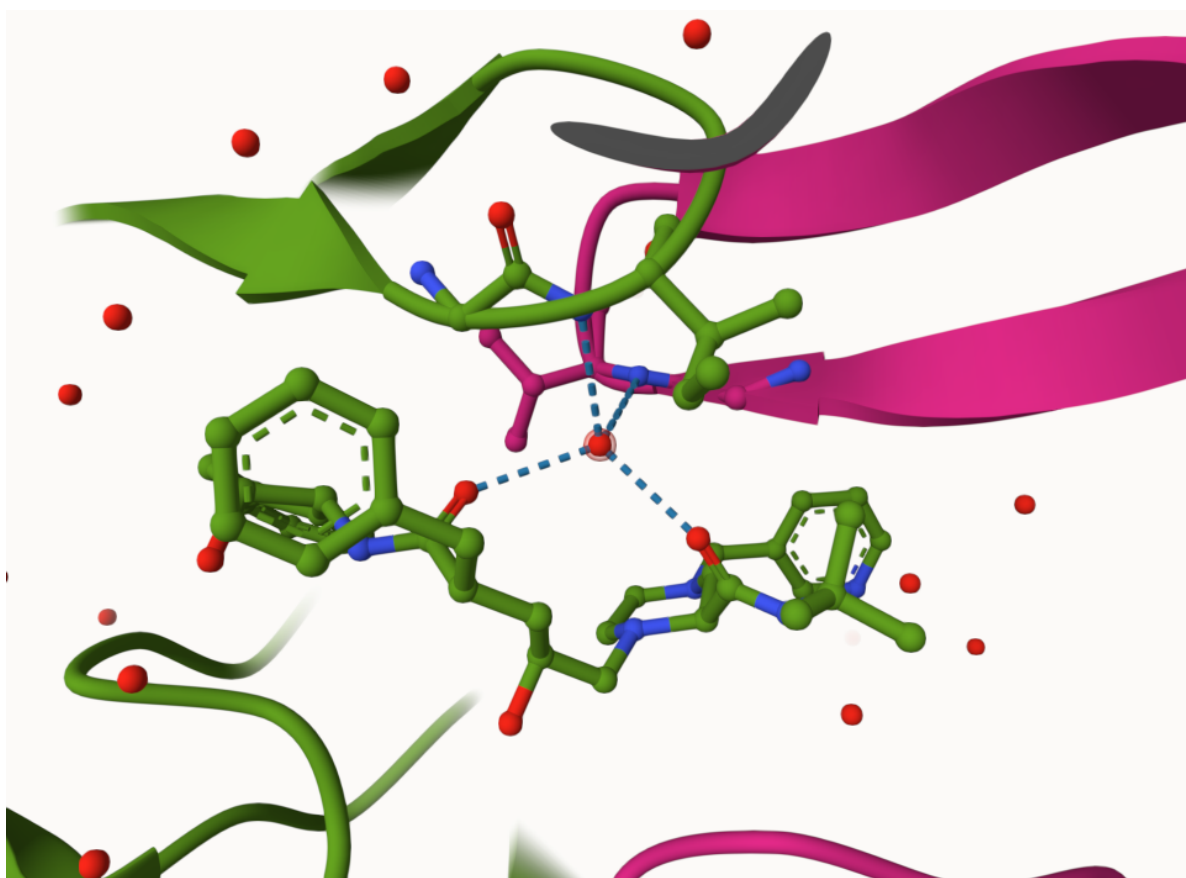


Figure 2: Water 308 in the Binding Site

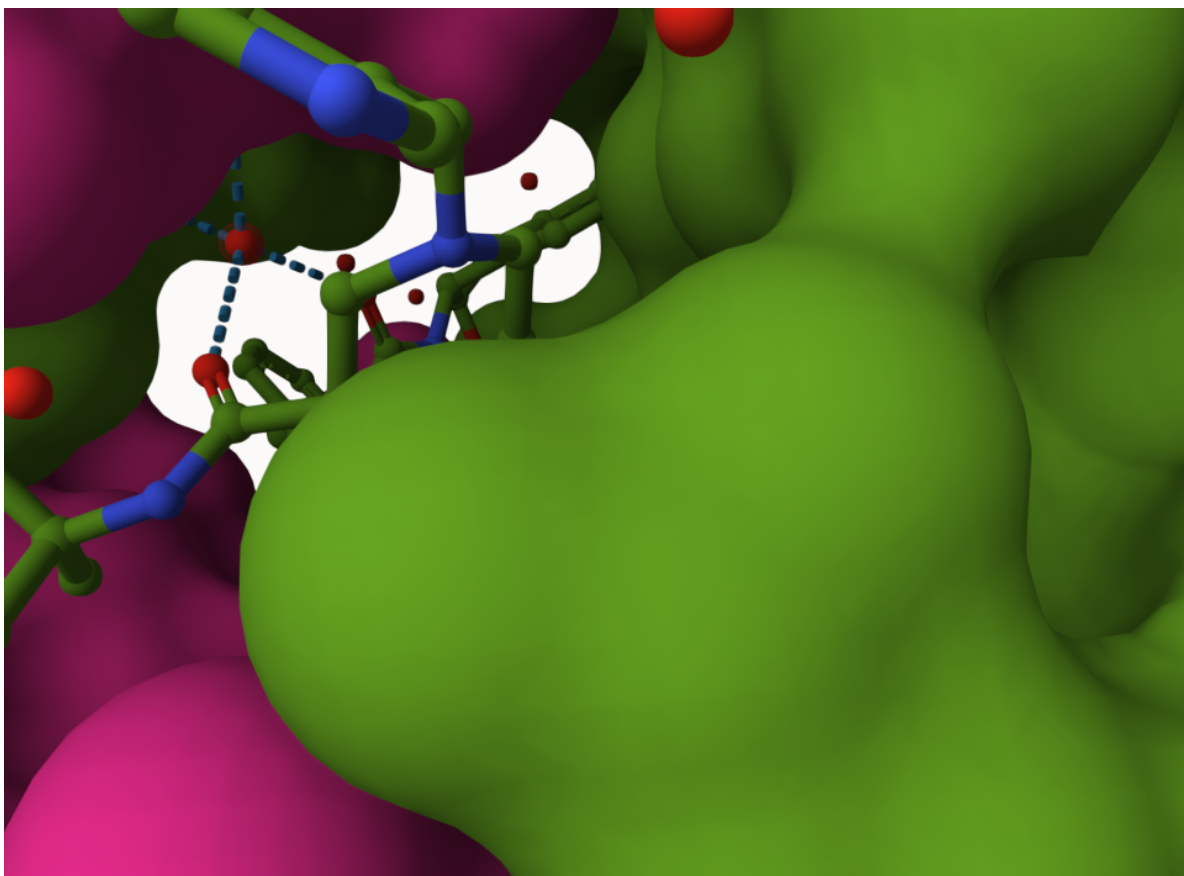


Figure 3: Surface Representation Showing Binding cavity of water 308

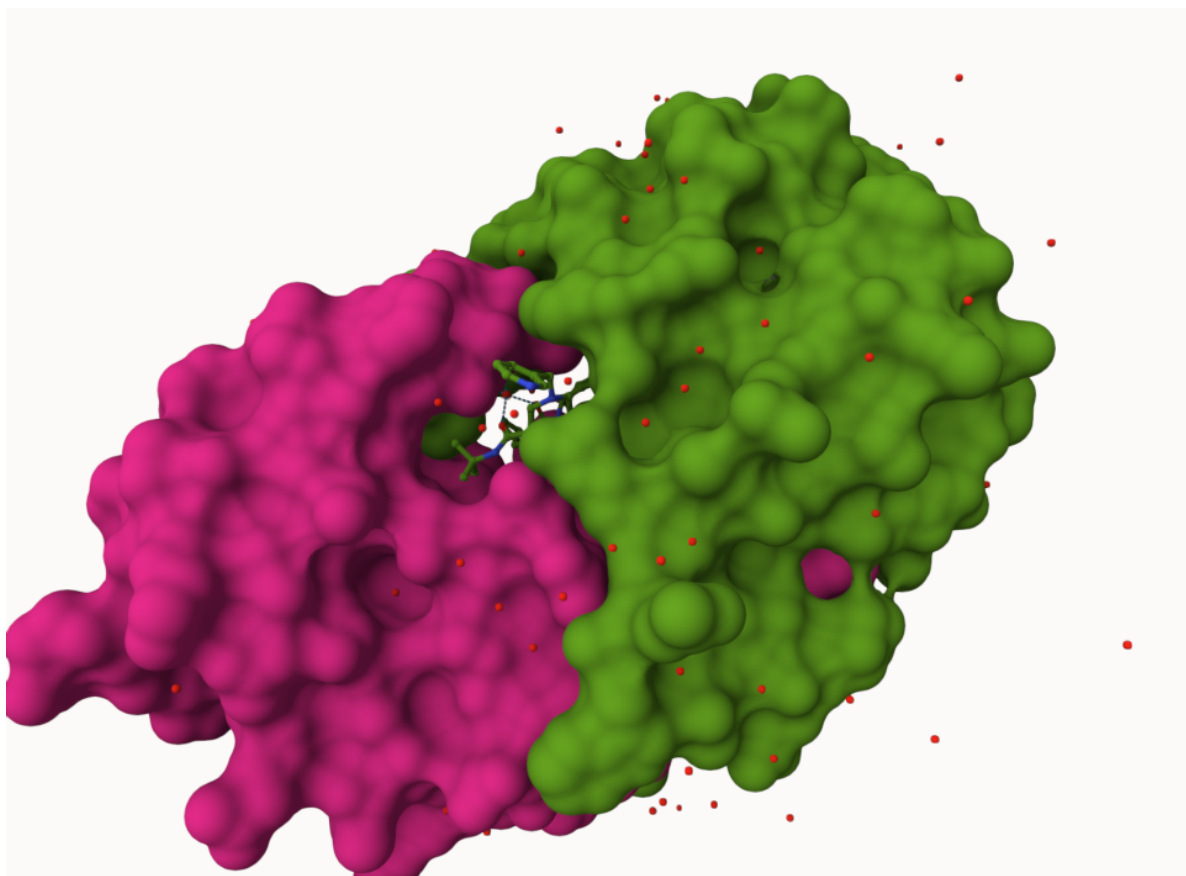


Figure 4: Overview of Surface representation of Water 308

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 5: Asp25 Amino Acid

3. Introduction to Bio3D in R

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

```
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

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

198 residues

```
#returns amino acids
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"		

```
#how many?
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 chains A and 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 new function not yet in the bio3d package. It requires the **r3dmol** package that we need to install with `install.packages("rd3mol")` and `install.packages("shiny")`

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

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:

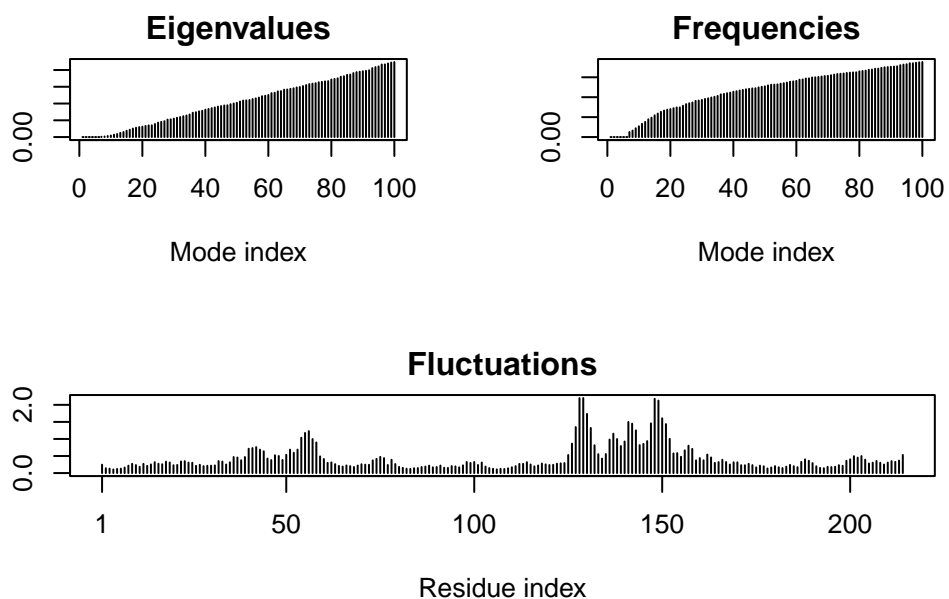
```
MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMRLRAAVKSGSELGKQAKDIMDAGKLV
DELVIALVKERIAQEDCRNGFLLDGFPRTPQADAMKEAGINVDYVLEFDVPDELIVDKI
VGRRVHAPSGRVYHVKNPPKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQMTAPLIG
YYSKEAEAGNTKYAKVDGTPVAEVRADLEKILG
```

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

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

```
Building Hessian...      Done in 0.05 seconds.  
Diagonalizing Hessian... Done in 0.28 seconds.
```

```
plot(m)
```



Peaks are functional spots predicted to move in the molecule.

Write out a trajectory of the predicted molecular motion:

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

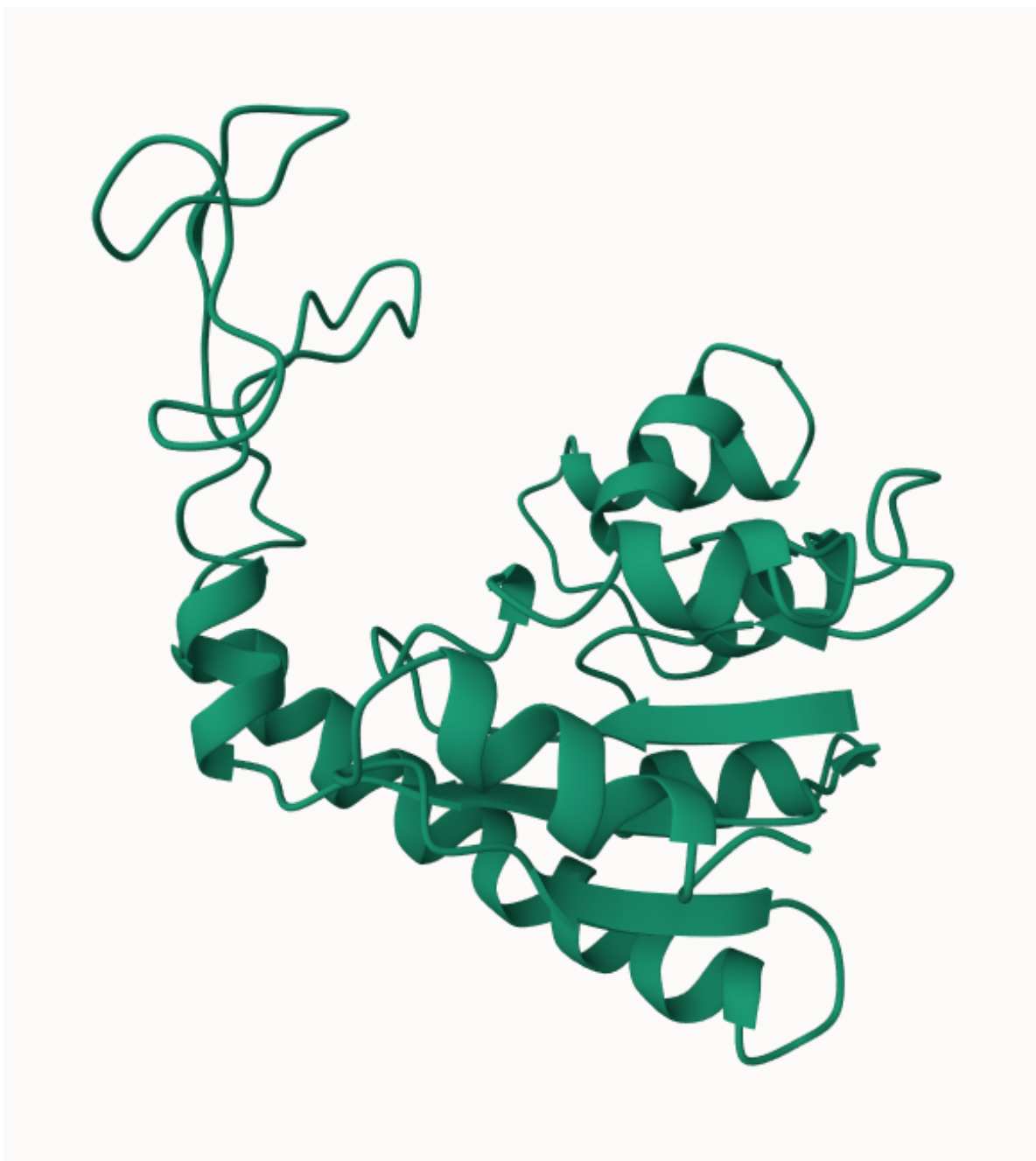


Figure 6: ADK Protein

We downloaded the animation trajectory, but this cannot be rendered into a PDF file.