

Class 09: Structural Bioinformatics 1

Kaitlyn Powell

1: Introduction to the RCSB Protein Data Bank (PDB)

The RCSB Protein Data Bank (PDB)

Protein structures by x-ray crystallography dominate this database dominate this data base. We are skipping Q1 - Q3 as the website is too slow for us.

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

92.85% of structures in the PDB are solved by X-Ray and Electron Microscopy.

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

86.99% of 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?

There are 4,008 HIV-1 protease structures in the current PDB.

2. Visualizing the HIV-1 protease structure

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

We only see one atom per water molecule in this structure due to the fact that the hydrogen molecules are so small that the resolution of the image can not visualize them. However, the oxygen molecules are large enough to be seen, so the atom that is present for each water molecule is the oxygen.

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 is positioned right in between the ligand and the binding site, and plays a very important role in the binding of the ligand. The residue number that this water molecule has is 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 (we recommend “Ball & Stick” for these side-chains). Add this figure to your Quarto document.

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

A way for indinavir, or even larger ligands and substrates could enter the the binding site would be for the polymer to be more flexible in order to allow a larger space where the ligand binds. Bonds can also be broken within the polymer in order to make the binding site larger, therefore allowing larger substrates/ligands to bind.



Figure 1: HIV-Pr structure from 1hsg

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?

The secondary structure elements that are likely to only form in the dimer rather than the monomer are ligands that require two different binding sites. Since monomers only have one chain, it is not possible to have two binding sites. However, in a dimer, there are two chains and therefore a possibility for two binding sites.

#3. Introduction to Bio3D in R

Bio3D is an R package for structural bioinformatics. To use it we need to call it up with the `library()` function (just like any package).

```
library(bio3d)
```

To read a PDB file, we can use `read.pdb()`

```
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:

```
PQITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWPKMKMIGGIGGFIKVRQYD
QILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNFPQITLWQRPLVTIKIGGQLKE
ALLDTGADDTVLEEMSLPGRWPKMKMIGGIGGFIKVRQYDQILIEICGHKAIGTVLVGPTP
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"
```

The ATOM records of a PDB file are stored in `pdb$atom`

```
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>

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

There are 198 amino acid residues in this pdb object.

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

MK1 is one of the two non-protein residues.

Q9: How many protein chains are in this structure?

There are 2 chains in this structure.

4. Comparative structure analysis of Adenylate Kinase

Comparative analysis of Adenylate kinase (ADK)

Search and retrieve ADK structures

We will start our analysis with a single PDB id (code from the PDB database): 1AKE

First we get it's primary sequence:

```
aa <- get.seq("1ake_a")
```

Warning in get.seq("1ake_a"): Removing existing file: seqs.fasta

Fetching... Please wait. Done.

```
aa
```

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

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

121 . . . . 180
pdb|1AKE|A VGRRVHAPSGRVYHVKNPPKVEGKDDVTGEELTRKDDQEETVRKRLVEYHQM TAPLIG
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

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

msa is only found on BioConductor and not CRAN.

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

bio3d-11 is not found on BioConductor or CRAN.

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

TRUE

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

There are 214 amino acids in this sequence.

```
# Blast or hmmer search
#b <- blast.pdb(aa)

# Plot a summary of search results
#hits <- plot(b)
# List out some 'top hits'
#head(hits$ pdb.id)
```

Use these ADK structures for analysis:

```
hits <- NULL
hits$ pdb.id <- c('1AKE_A', '6S36_A', '6RZE_A', '3HPR_A', '1E4V_A', '5EJE_A', '1E4Y_A', '3X2S_A',
```

Download all these PDB files from the online database...

```
# Download related PDB files
files <- get.pdb(hits$ pdb.id, path="pdbs", split=TRUE, gzip=TRUE)
```

Warning in get.pdb(hits\$ pdb.id, path = "pdbs", split = TRUE, gzip = TRUE): pdbs/
1AKE.pdb.gz exists. Skipping download

Warning in get.pdb(hits\$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE): pdbs/
6S36.pdb.gz exists. Skipping download

Warning in get.pdb(hits\$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE): pdbs/
6RZE.pdb.gz exists. Skipping download

Warning in get.pdb(hits\$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE): pdbs/
3HPR.pdb.gz exists. Skipping download

Warning in get.pdb(hits\$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE): pdbs/
1E4V.pdb.gz exists. Skipping download

Warning in get.pdb(hits\$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE): pdbs/
5EJE.pdb.gz exists. Skipping download

Warning in get.pdb(hits\$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE): pdbs/
1E4Y.pdb.gz exists. Skipping download

Warning in get.pdb(hits\$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE): pdbs/
3X2S.pdb.gz exists. Skipping download

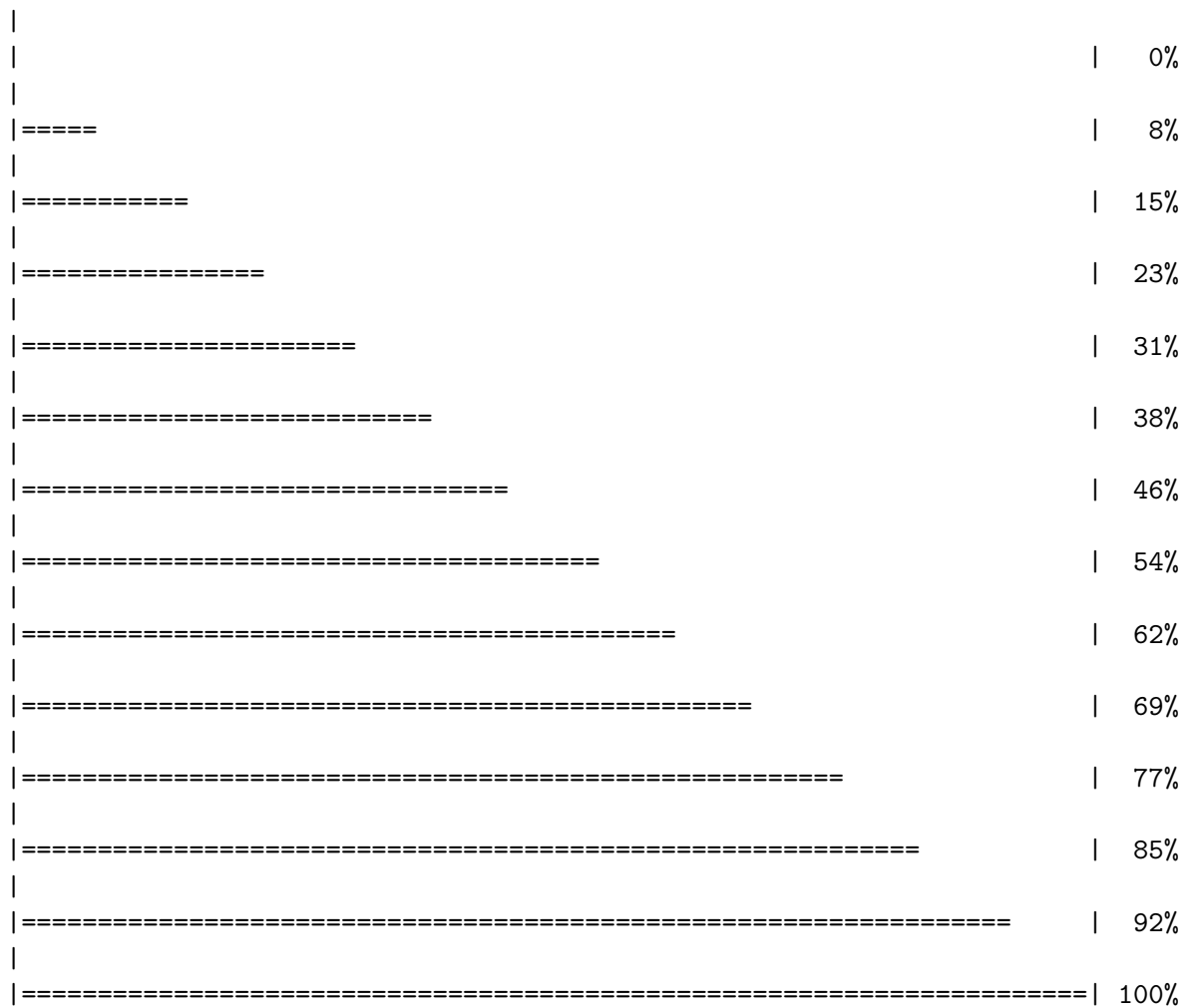
Warning in get.pdb(hits\$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE): pdbs/
6HAP.pdb.gz exists. Skipping download

Warning in get.pdb(hits\$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE): pdbs/
6HAM.pdb.gz exists. Skipping download

Warning in get.pdb(hits\$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE): pdbs/
4K46.pdb.gz exists. Skipping download

Warning in get.pdb(hits\$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE): pdbs/
3GMT.pdb.gz exists. Skipping download

Warning in get.pdb(hits\$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE): pdbs/
4PZL.pdb.gz exists. Skipping download



##Align and superpose structures

Align all these structures

```
# Align releated PDBs
pdbbs <- pdbaln(files, fit = TRUE, exefile="msa")
```

Reading PDB files:

```
pdbbs/split_chain/1AKE_A.pdb
pdbbs/split_chain/6S36_A.pdb
pdbbs/split_chain/6RZE_A.pdb
pdbbs/split_chain/3HPR_A.pdb
```



```

pdbs/split_chain/1E4V_A.pdb
pdbs/split_chain/5EJE_A.pdb
pdbs/split_chain/1E4Y_A.pdb
pdbs/split_chain/3X2S_A.pdb
pdbs/split_chain/6HAP_A.pdb
pdbs/split_chain/6HAM_A.pdb
pdbs/split_chain/4K46_A.pdb
pdbs/split_chain/3GMT_A.pdb
pdbs/split_chain/4PZL_A.pdb
    PDB has ALT records, taking A only, rm.alt=TRUE
.    PDB has ALT records, taking A only, rm.alt=TRUE
.    PDB has ALT records, taking A only, rm.alt=TRUE
.    PDB has ALT records, taking A only, rm.alt=TRUE
..   PDB has ALT records, taking A only, rm.alt=TRUE
....  PDB has ALT records, taking A only, rm.alt=TRUE
.    PDB has ALT records, taking A only, rm.alt=TRUE
...

```

Extracting sequences

```

pdb/seq: 1    name: pdbs/split_chain/1AKE_A.pdb
    PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 2    name: pdbs/split_chain/6S36_A.pdb
    PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 3    name: pdbs/split_chain/6RZE_A.pdb
    PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 4    name: pdbs/split_chain/3HPR_A.pdb
    PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 5    name: pdbs/split_chain/1E4V_A.pdb
pdb/seq: 6    name: pdbs/split_chain/5EJE_A.pdb
    PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 7    name: pdbs/split_chain/1E4Y_A.pdb
pdb/seq: 8    name: pdbs/split_chain/3X2S_A.pdb
pdb/seq: 9    name: pdbs/split_chain/6HAP_A.pdb
pdb/seq: 10   name: pdbs/split_chain/6HAM_A.pdb
    PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 11   name: pdbs/split_chain/4K46_A.pdb
    PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 12   name: pdbs/split_chain/3GMT_A.pdb
pdb/seq: 13   name: pdbs/split_chain/4PZL_A.pdb

```

```

[Truncated_Name:1] 1AKE_A.pdb
[Truncated_Name:2] 6S36_A.pdb
[Truncated_Name:3] 6RZE_A.pdb
[Truncated_Name:4] 3HPR_A.pdb
[Truncated_Name:5] 1E4V_A.pdb
[Truncated_Name:6] 5EJE_A.pdb
[Truncated_Name:7] 1E4Y_A.pdb
[Truncated_Name:8] 3X2S_A.pdb
[Truncated_Name:9] 6HAP_A.pdb
[Truncated_Name:10] 6HAM_A.pdb
[Truncated_Name:11] 4K46_A.pdb
[Truncated_Name:12] 3GMT_A.pdb
[Truncated_Name:13] 4PZL_A.pdb

1          .          .          .          40
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
-----MRIILLGAPVAGKGTQAQFIMEKYGIPQIS
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
-----MRIILLGALVAGKGTQAQFIMEKYGIPQIS
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
-----MRIILLGAPGAGKGTQAQFIMAKFGIPQIS
-----MRLILLGAPGAGKGTQANFIKEKFGIPQIS
TENLYFQSNAMRIILLGAPGAGKGTQAKIIEQKYNIAHIS
          **~*****  *****  *  *~ *  **
1          .          .          .          40

[Truncated_Name:1] 1AKE_A.pdb
[Truncated_Name:2] 6S36_A.pdb
[Truncated_Name:3] 6RZE_A.pdb
[Truncated_Name:4] 3HPR_A.pdb
[Truncated_Name:5] 1E4V_A.pdb
[Truncated_Name:6] 5EJE_A.pdb
[Truncated_Name:7] 1E4Y_A.pdb
[Truncated_Name:8] 3X2S_A.pdb
[Truncated_Name:9] 6HAP_A.pdb
[Truncated_Name:10] 6HAM_A.pdb
[Truncated_Name:11] 4K46_A.pdb
[Truncated_Name:12] 3GMT_A.pdb
[Truncated_Name:13] 4PZL_A.pdb

41          .          .          .          80
TGDMLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE
TGDMLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE
TGDMLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE
TGDMLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE
TGDMLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE
TGDMLRAAVKSGSELGKQAKDIMDACKLVDELVIALVKE
TGDMLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE
TGDMLRAAVKSGSELGKQAKDIMDCGKLVDELVIALVKE
TGDMLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVRE
TGDMLRAAIAKSGSELGKQAKDIMDAGKLVDEIIIALVKE
TGDMLRAAIAKAGTELGKQAKSVIDAGQLVSDDIILGLVKE
TGDMLRAAVKAGTPLGVEAKTYMDEGKLVPSLIIGLVKE
TGDMIRETIKSGSALGQELKKVLDAGELVSDEFIIVKVD
****~*  ~* *~ **  *  ~*  ** *  ^^ ~* ^^
41          .          .          .          80

[Truncated_Name:1] 1AKE_A.pdb
[Truncated_Name:2] 6S36_A.pdb
[Truncated_Name:3] 6RZE_A.pdb
[Truncated_Name:4] 3HPR_A.pdb
[Truncated_Name:5] 1E4V_A.pdb

81          .          .          .          120
RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD

```

[Truncated_Name:6] 5EJE_A.pdb	RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:7] 1E4Y_A.pdb	RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:8] 3X2S_A.pdb	RIAQEDSRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:9] 6HAP_A.pdb	RICQEDSRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:10] 6HAM_A.pdb	RICQEDSRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:11] 4K46_A.pdb	RIAQDDCAKGFLLDGFPR TIPQADGLKEVGVVVDYVIEFD
[Truncated_Name:12] 3GMT_A.pdb	RLKEADCANGYLF DGFPR TIAQADAMKEAGVAIDYVLEID
[Truncated_Name:13] 4PZL_A.pdb	RISKNDCNNGFLLDGVPR TIPQAQELDKLGVNIDYIVEVD
	*^ * *~* ** ***** * ^ *~ ~**~* *
	81 . . . 120
	121 . . . 160
[Truncated_Name:1] 1AKE_A.pdb	VPDELIVDRIVGRRVHAPSGRVYHV KFNPPKVEGKDDVTG
[Truncated_Name:2] 6S36_A.pdb	VPDELIVDKIVGRRVHAPSGRVYHV KFNPPKVEGKDDVTG
[Truncated_Name:3] 6RZE_A.pdb	VPDELIVDAIVGRRVHAPSGRVYHV KFNPPKVEGKDDVTG
[Truncated_Name:4] 3HPR_A.pdb	VPDELIVDRIVGRRVHAPSGRVYHV KFNPPKVEGKDDGTG
[Truncated_Name:5] 1E4V_A.pdb	VPDELIVDRIVGRRVHAPSGRVYHV KFNPPKVEGKDDVTG
[Truncated_Name:6] 5EJE_A.pdb	VPDELIVDRIVGRRVHAPSGRVYHV KFNPPKVEGKDDVTG
[Truncated_Name:7] 1E4Y_A.pdb	VPDELIVDRIVGRRVHAPSGRVYHV KFNPPKVEGKDDVTG
[Truncated_Name:8] 3X2S_A.pdb	VPDELIVDRIVGRRVHAPSGRVYHV KFNPPKVEGKDDVTG
[Truncated_Name:9] 6HAP_A.pdb	VPDELIVDRIVGRRVHAPSGRVYHV KFNPPKVEGKDDVTG
[Truncated_Name:10] 6HAM_A.pdb	VPDELIVDRIVGRRVHAPSGRVYHV KFNPPKVEGKDDVTG
[Truncated_Name:11] 4K46_A.pdb	VADSVIVERMAGRRAHLASGR TYHNVPKVEGKDDVTG
[Truncated_Name:12] 3GMT_A.pdb	VPFSEIIERMSGRRTHPASGR TYHV KFNPPKVEGKDDVTG
[Truncated_Name:13] 4PZL_A.pdb	VADNLLIERITGRRIHPASGR TYHTKFNPPKVADKDDVTG
	* ^^^ ^ *** * *** * ^***** *** *
	121 . . . 160
	161 . . . 200
[Truncated_Name:1] 1AKE_A.pdb	EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:2] 6S36_A.pdb	EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:3] 6RZE_A.pdb	EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:4] 3HPR_A.pdb	EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:5] 1E4V_A.pdb	EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:6] 5EJE_A.pdb	EELTTRKDDQEECVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:7] 1E4Y_A.pdb	EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:8] 3X2S_A.pdb	EELTTRKDDQEETVRKRLCEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:9] 6HAP_A.pdb	EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:10] 6HAM_A.pdb	EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:11] 4K46_A.pdb	EDLVIREDDKEETV LARLG VYHNQTAPLIAYYGKEAEAGN
[Truncated_Name:12] 3GMT_A.pdb	EPLVQRDDKEETVKKRLDVYEAQTKPLITYYGDWARRGA
[Truncated_Name:13] 4PZL_A.pdb	EPLITRTDDNEDTVKQRLSVYHAQTAKLIDFYRNFSSNT
	* * * * * ^ * * * * * ^ *

161 . . . 200

201 . 227

```
[Truncated_Name:1] 1AKE_A.pdb T--KYAKVDGTPVAEVRADLEKILG-
[Truncated_Name:2] 6S36_A.pdb T--KYAKVDGTPVAEVRADLEKILG-
[Truncated_Name:3] 6RZE_A.pdb T--KYAKVDGTPVAEVRADLEKILG-
[Truncated_Name:4] 3HPR_A.pdb T--KYAKVDGTPVAEVRADLEKILG-
[Truncated_Name:5] 1E4V_A.pdb T--KYAKVDGTPVAEVRADLEKILG-
[Truncated_Name:6] 5EJE_A.pdb T--KYAKVDGTPVAEVRADLEKILG-
[Truncated_Name:7] 1E4Y_A.pdb T--KYAKVDGTPVAEVRADLEKILG-
[Truncated_Name:8] 3X2S_A.pdb T--KYAKVDGTPVAEVRADLEKILG-
[Truncated_Name:9] 6HAP_A.pdb T--KYAKVDGTPVCEVRADLEKILG-
[Truncated_Name:10] 6HAM_A.pdb T--KYAKVDGTPVCEVRADLEKILG-
[Truncated_Name:11] 4K46_A.pdb T--QYLKFDGTKAFAEVSADLEKALA-
[Truncated_Name:12] 3GMT_A.pdb E-----NGLKAPA-----YRKISG-
[Truncated_Name:13] 4PZL_A.pdb KIPKYIKINGDQAVEKVSQDIFDQLNK
```

*

201 . 227

Call:

```
pdbaln(files = files, fit = TRUE, exefile = "msa")
```

Class:

```
pdbs, fasta
```

Alignment dimensions:

```
13 sequence rows; 227 position columns (204 non-gap, 23 gap)
```

```
+ attr: xyz, resno, b, chain, id, ali, resid, sse, call
```

```
# Vector containing PDB codes for figure axis
ids <- basename.pdb(pdb$ids)
```

```
# Draw schematic alignment
#plot(pdb, labels=ids)
#par(mar=c(1,1,1,1))
```

##Annotate collected PDB structures

```
anno <- pdb.annotate(ids)
unique(anno$source)
```

```

[1] "Escherichia coli"
[2] "Escherichia coli K-12"
[3] "Escherichia coli O139:H28 str. E24377A"
[4] "Escherichia coli str. K-12 substr. MDS42"
[5] "Photobacterium profundum"
[6] "Burkholderia pseudomallei 1710b"
[7] "Francisella tularensis subsp. tularensis SCHU S4"

```

```
head(anno)
```

	structureId	chainId	macromoleculeType	chainLength	experimentalTechnique			
	1AKE_A	1AKE	A	Protein	214	X-ray		
	6S36_A	6S36	A	Protein	214	X-ray		
	6RZE_A	6RZE	A	Protein	214	X-ray		
	3HPR_A	3HPR	A	Protein	214	X-ray		
	1E4V_A	1E4V	A	Protein	214	X-ray		
	5EJE_A	5EJE	A	Protein	214	X-ray		
	resolution		scopDomain		pfam			
	1AKE_A	2.00	Adenylate kinase	Adenylate kinase, active site lid (ADK_lid)				
	6S36_A	1.60	<NA>	Adenylate kinase, active site lid (ADK_lid)				
	6RZE_A	1.69	<NA>	Adenylate kinase, active site lid (ADK_lid)				
	3HPR_A	2.00	<NA>	Adenylate kinase, active site lid (ADK_lid)				
	1E4V_A	1.85	Adenylate kinase	Adenylate kinase, active site lid (ADK_lid)				
	5EJE_A	1.90	<NA>	Adenylate kinase, active site lid (ADK_lid)				
	ligandId		ligandName					
	1AKE_A	AP5	BIS(ADENOSINE)-5'-PENTAPHOSPHATE					
	6S36_A	CL (3),NA,MG (2)	CHLORIDE ION (3),SODIUM ION,MAGNESIUM ION (2)					
	6RZE_A	NA (3),CL (2)	SODIUM ION (3),CHLORIDE ION (2)					
	3HPR_A	AP5	BIS(ADENOSINE)-5'-PENTAPHOSPHATE					
	1E4V_A	AP5	BIS(ADENOSINE)-5'-PENTAPHOSPHATE					
	5EJE_A	AP5,CO	BIS(ADENOSINE)-5'-PENTAPHOSPHATE,COBALT (II) ION					
	source							
	1AKE_A	Escherichia coli						
	6S36_A	Escherichia coli						
	6RZE_A	Escherichia coli						
	3HPR_A	Escherichia coli K-12						
	1E4V_A	Escherichia coli						
	5EJE_A	Escherichia coli O139:H28 str. E24377A						

```

1AKE_A STRUCTURE OF THE COMPLEX BETWEEN ADENYLATE KINASE FROM ESCHERICHIA COLI AND THE INHIB.
6S36_A

```

6RZE_A
3HPR_A
1E4V_A
5EJE_A

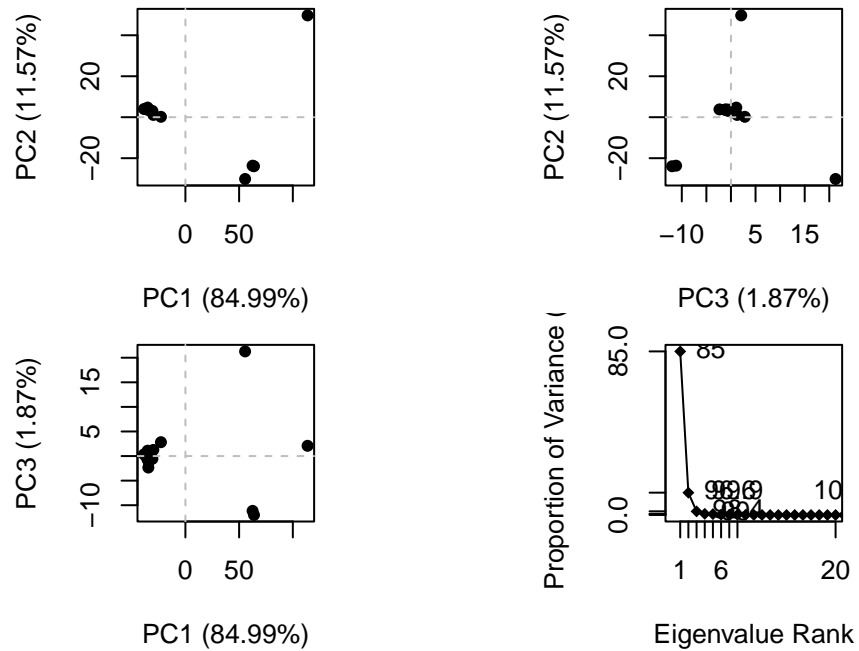
Cryst

	citation	rObserved	rFree
1AKE_A	Muller, C.W., et al. J Mol Biol (1992)	0.1960	NA
6S36_A	Rogne, P., et al. Biochemistry (2019)	0.1632	0.2356
6RZE_A	Rogne, P., et al. Biochemistry (2019)	0.1865	0.2350
3HPR_A	Schrank, T.P., et al. Proc Natl Acad Sci U S A (2009)	0.2100	0.2432
1E4V_A	Muller, C.W., et al. Proteins (1993)	0.1960	NA
5EJE_A	Kovermann, M., et al. Proc Natl Acad Sci U S A (2017)	0.1889	0.2358

	rWork	spaceGroup
1AKE_A	0.1960	P 21 2 21
6S36_A	0.1594	C 1 2 1
6RZE_A	0.1819	C 1 2 1
3HPR_A	0.2062	P 21 21 2
1E4V_A	0.1960	P 21 2 21
5EJE_A	0.1863	P 21 2 21

Principal component analysis (PCA)

```
# Perform PCA
pc.xray <- pca(pdbx)
plot(pc.xray)
```

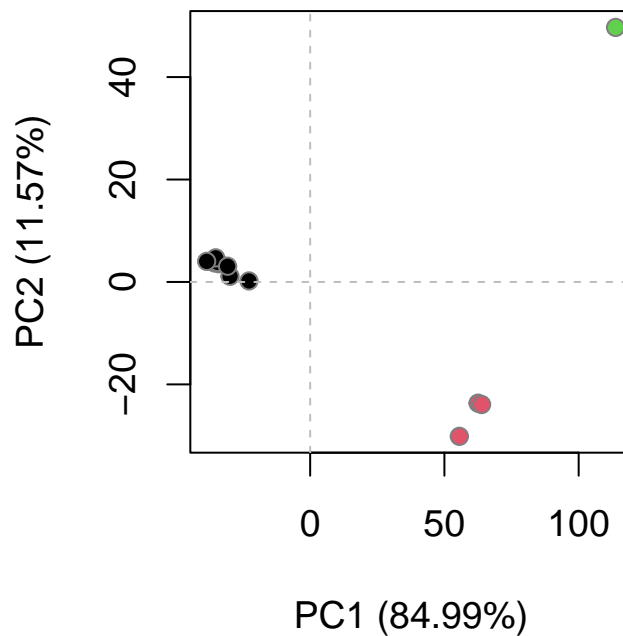


```
# Calculate RMSD
rd <- rmsd(pdbbs)
```

Warning in rmsd(pdbbs): No indices provided, using the 204 non NA positions

```
# Structure-based clustering
hc.rd <- hclust(dist(rd))
grps.rd <- cutree(hc.rd, k=3)
```

```
plot(pc.xray, 1:2, col="grey50", bg=grps.rd, pch=21, cex=1)
```



5. Optional further visualization

To visualize the major structural variations in the ensemble the function `mktrj()` can be used to generate a trajectory PDB file by interpolating along to give a PC (eigenvector):

```
# Visualize first principal component
pc1 <- mktrj(pc.xray, pc=1, file="pc_1.pdb")
```

****Note:** The animation would not format properly for the pdf, so it is not included here.

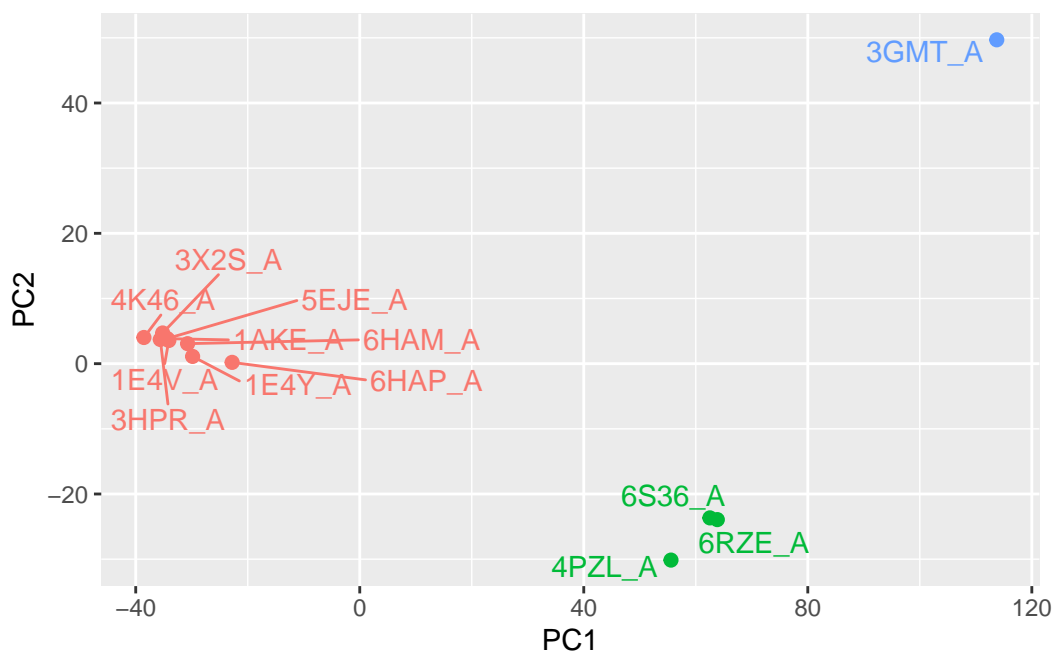
```
#Plotting results with ggplot2
library(ggplot2)
library(ggrepel)

df <- data.frame(PC1=pc.xray$z[,1],
                 PC2=pc.xray$z[,2],
                 col=as.factor(grps.rd),
                 ids=ids)
```



```
p <- ggplot(df) +
  aes(PC1, PC2, col=col, label=ids) +
  geom_point(size=2) +
  geom_text_repel(max.overlaps = 20) +
  theme(legend.position = "none")
```

p



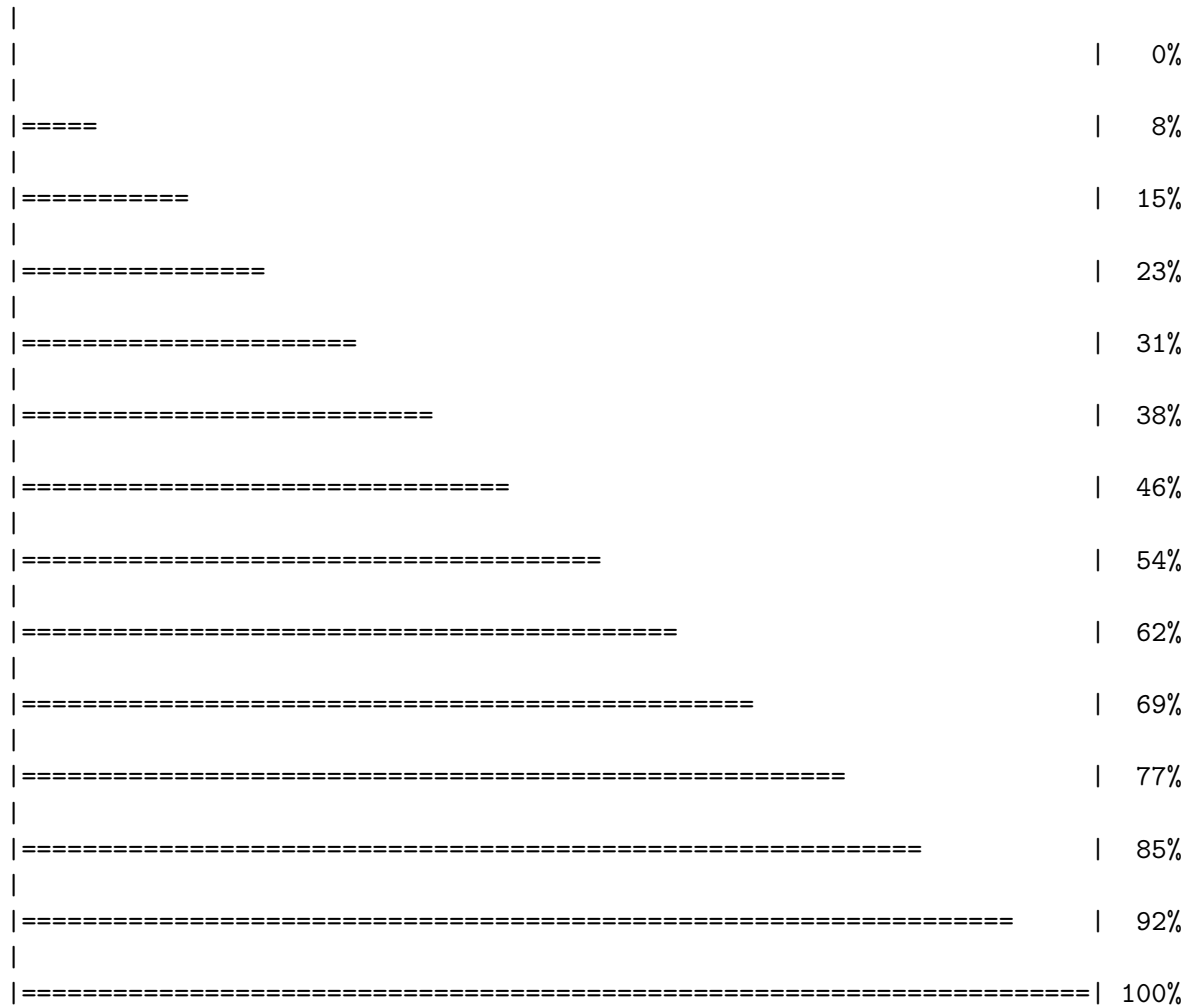
6. Normal mode analysis [optional]

```
# NMA of all structures
modes <- nma(pdbbs)
```

Details of Scheduled Calculation:

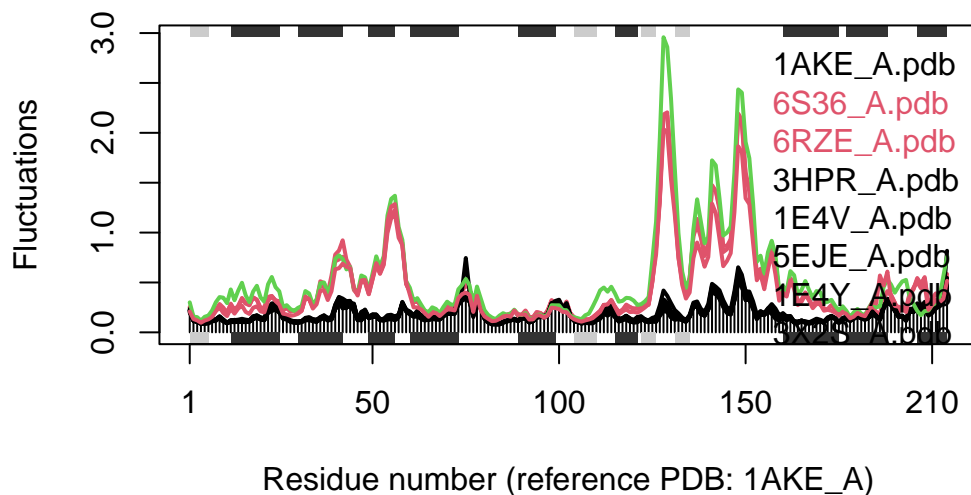
```
... 13 input structures
... storing 606 eigenvectors for each structure
... dimension of x$U.subspace: ( 612x606x13 )
... coordinate superposition prior to NM calculation
... aligned eigenvectors (gap containing positions removed)
```

... estimated memory usage of final 'eNMA' object: 36.9 Mb



```
plot(modes, pdba, col=grps.rd)
```

Extracting SSE from pdba\$sse attribute



Q14. What do you note about this plot? Are the black and colored lines similar or different? Where do you think they differ most and why?

In this plot, the black and colored lines tend to follow a similar pattern, however it seems that the colored lines generally have much higher fluctuations than the black line. They differ the most around residue number 150. This could be due to the fact that the colored lines represent proteins that have binding sites at this location, and the black line does not.