

# Class 10: Structural Bioinformatics (Pt. 1)

## PDB statistics

Download a CSV file from the PDB site (accessible from “Analyze” > “PDB Statistics” > “by Experimental Method and Molecular Type”

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

ANS: 93.26%

```
pdb.df <- read.csv("Data Export Summary.csv", row.names=1)
```

Creation of a function to read

```
sum_column <- function(data_frame, column_name) {  
  # Remove any non-numeric characters  
  cleaned_column <- gsub("[^0-9.]", "", data_frame[[column_name]])  
  
  # Convert the cleaned column to numeric  
  numeric_column <- as.numeric(cleaned_column)  
  
  # Sum the numeric values, handling any NA values that might arise during conversion  
  sum(numeric_column, na.rm = TRUE)  
}  
  
#Total xray  
total_xray<-sum_column(pdb.df, "X.ray")  
  
#Total EM  
total_EM<-sum_column(pdb.df, "EM")
```

```
total_Total <- sum_column(pdb.df, "Total")

Percentage<-sum(total_xray,total_EM)/total_Total*100
Percentage
```

[1] 93.26839

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

ANS: 86.65%

```
protein<-as.numeric(pdb.df[, "Total"])
```

Warning: NAs introduced by coercion

```
proportion <- protein/total_Total*100
proportion
```

[1] NA

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?

ANS: 23,471

## The PDB format

Now download the “PDB File” for the HIV-1 protease structure with the PDB identifier 1HSG. Using the terminal I typed the following command:

```
less 1hsg.pdb
```

## Visualizing the HIV-1 protease structure

### Using Mol to examine HIV-Pr

Working with Mol\* tool using 1HSG as the protein to visualize

Highlighting Asp 25 (D25) in both chains since they are critical for protease activity



Figure 1: A nice look of 1HSG.

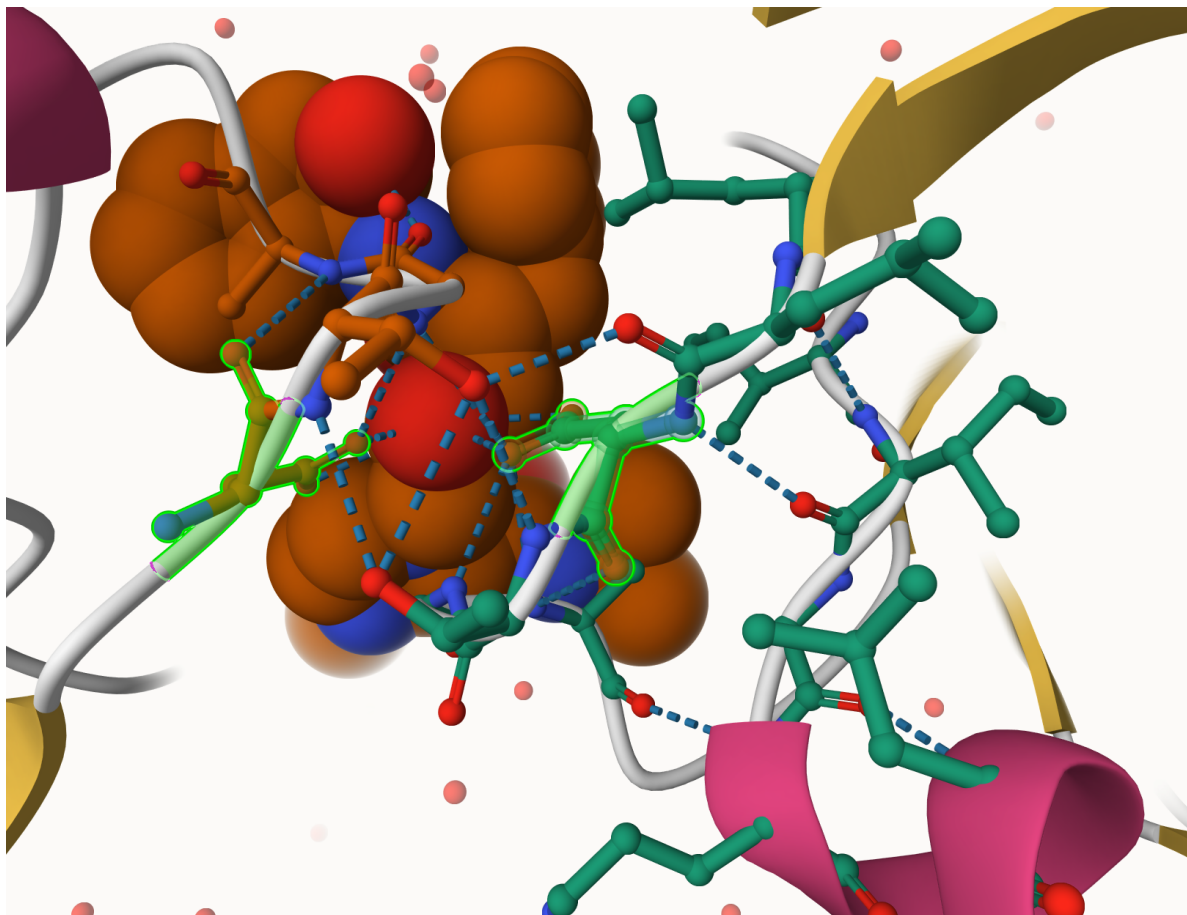


Figure 2: Asp 25 (D25) in both chains.

## The important role of water

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

The red spheres represent the Oxygen atoms. These visualization tools use just the oxygen atom to represent water molecule since it is larger and more significant than hydrogen atoms.

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

ANS: 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.

## Introduction to Bio3D in R

```
library(bio3d)
```

### Reading PDB file data into R

```
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)
```

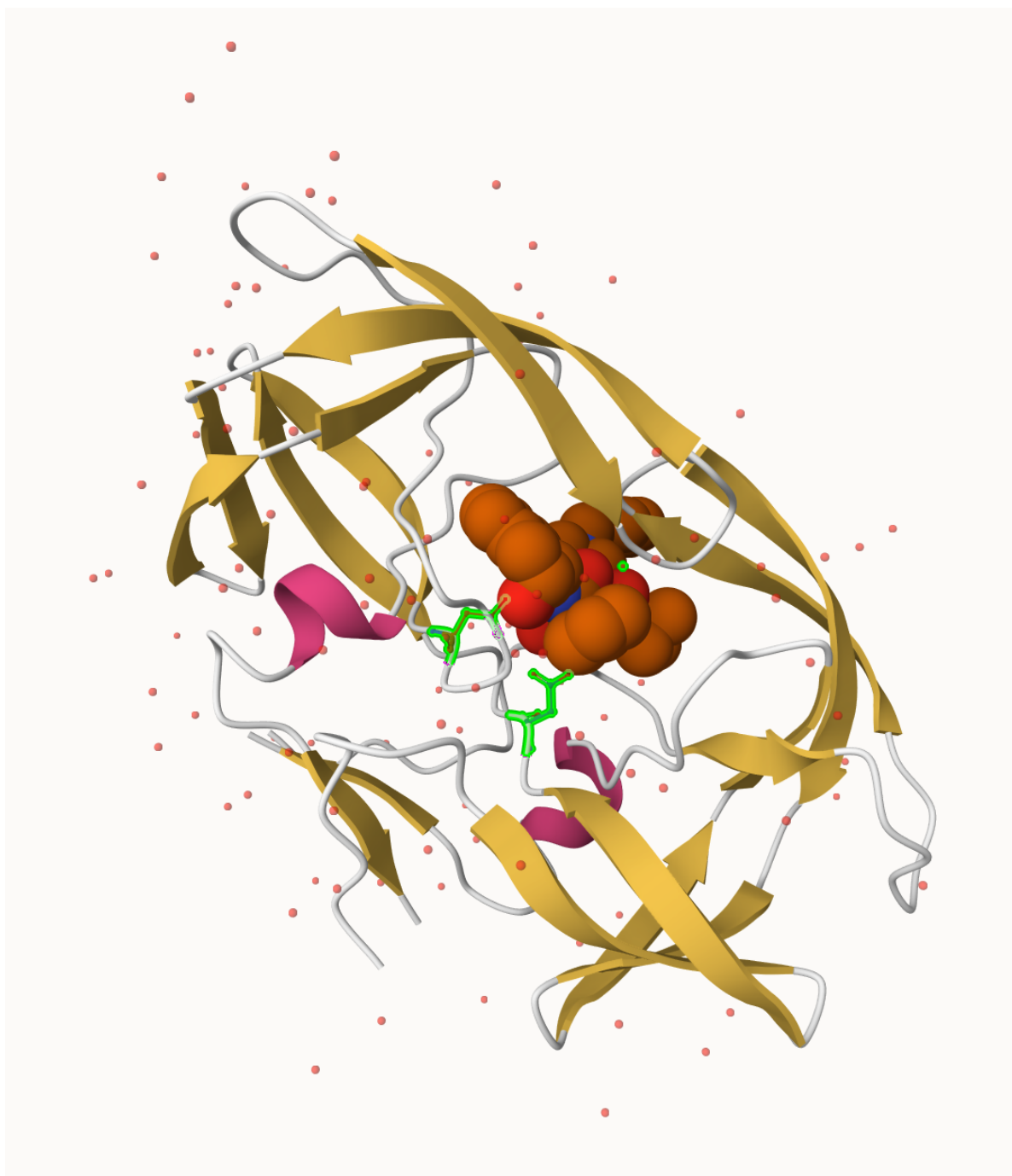


Figure 3: Asp 25 (D25) in both chains, critical water and two distinct chains of HIV-protease.

```
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?

ANS: 198

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

ANS: HOH

Q9: How many protein chains are in this structure?

ANS: 2

Note that the attributes (+ attr:) of this object are listed on the last couple of lines. To find the attributes of any such object you can use:

```
attributes(pdb)
```

```
$names
```

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

```
$class
```

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

To access these individual attributes we use the dollar-attribute name convention that is common with R list objects. For example, to access the atom attribute or component use `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>

## Predicting functional motions of a single structure

Let's read a new PDB structure of Adenylate Kinase and perform Normal mode analysis.

```
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



```

DELVIALVKERIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFDVPDELIVDKI
VGRRVHAPSGRVYHVKNPPKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQMTAPLIG
YYSKAEAEAGNTKYAKVDGTPVAEVRADLEKILG

```

```

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

```

Normal mode analysis (NMA) is a structural bioinformatics method to predict protein flexibility and potential functional motions (a.k.a. conformational changes).

```

# Perform flexibility prediction
m <- nma(adk)

```

```

Building Hessian...      Done in 0.03 seconds.
Diagonalizing Hessian... Done in 0.3 seconds.

```

```

m

```

Call:

```

nma.pdb(pdb = adk)

```

Class:

```

VibrationalModes (nma)

```

Number of modes:

```

642 (6 trivial)

```

Frequencies:

```

Mode 7:  0.005
Mode 8:  0.007
Mode 9:  0.009
Mode 10: 0.011
Mode 11: 0.013
Mode 12: 0.015

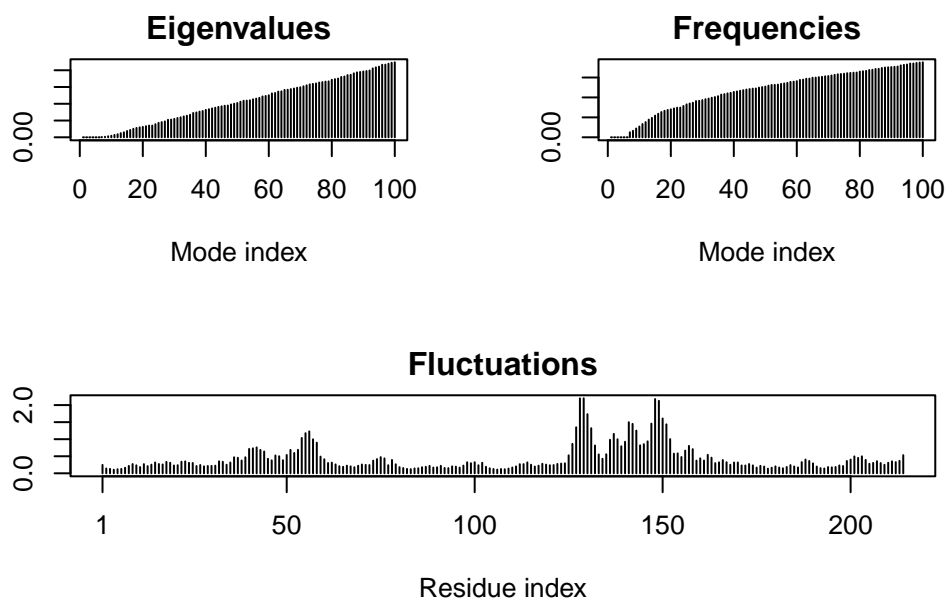
```

```

+ attr: modes, frequencies, force.constants, fluctuations,
      U, L, xyz, mass, temp, triv.modes, natoms, call

```

```
plot(m)
```



To view a “movie” of these predicted motions we can generate a molecular “trajectory” with the `mktrj()` function. We can now visualize the animation in Mol\*

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

## Comparative structure analysis of Adenylate Kinase

We will begin by first installing the packages we need for today’s session.

### Install packages in the R console NOT your Rmd/Quarto file

```
install.packages("bio3d") install.packages("devtools") install.packages("BiocManager")  
BiocManager::install("msa") devtools::install_bitbucket("Grantlab/bio3d-view")
```

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

ANS: msa

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

ANS: bio3d-view

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

ANS: 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  MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMRLAAVKSSELGKQAKDIMDAGKLV
      1      .      .      .      .      .      .      60

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

      121      .      .      .      .      .      .      180
pdb|1AKE|A  VGRRVHAPSGRVYHVKNPPKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQMTAPLIG
      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?

ANS: 214

Now we can use this sequence as a query to BLAST search the PDB to find similar sequences and structures.

```
# Blast or hmmer search. It is not working when I do render
#b <- blast.pdb(aa)
```

```
#hits <- plot.blast(b)
```

```
# List out some 'top hits'
#head(hits$ pdb.id)
```

Blast did not work.

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

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

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

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

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

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

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

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

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

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

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

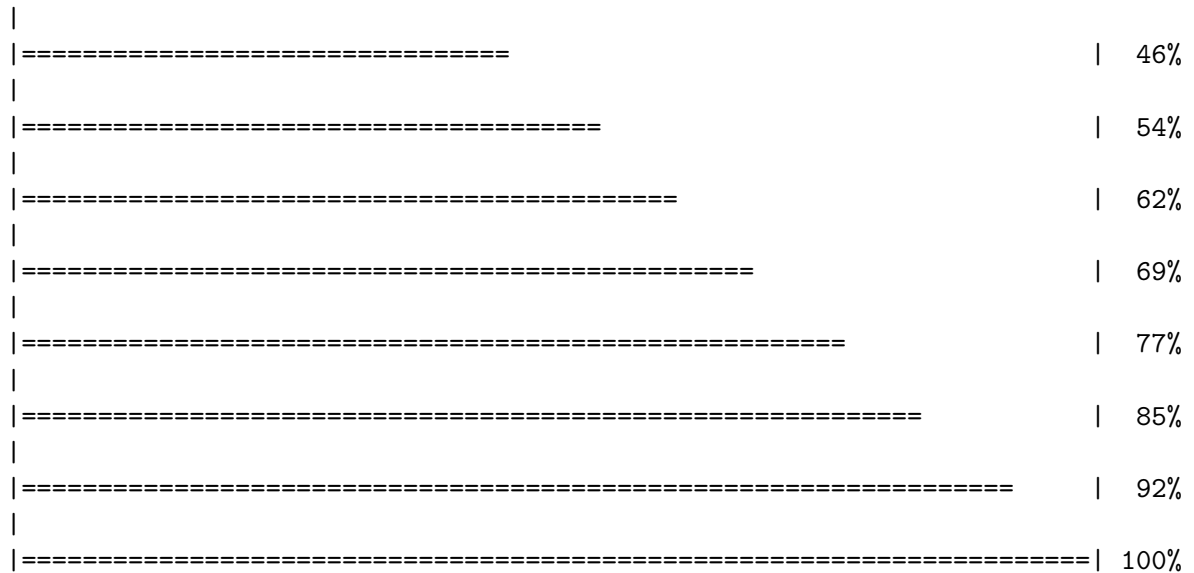
Warning in get.pdb(hits\$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):  
pdbs/6HAM.pdb exists. Skipping download

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

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

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

	0%
=====	8%
=====	15%
=====	23%
=====	31%
=====	38%



## Align and superpose structures

Next we will use the `pdbaln()` function to align and also optionally fit (i.e. superpose) the identified PDB structures.

```
# Align related 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
pdbbs/split_chain/1E4V_A.pdb
pdbbs/split_chain/5EJE_A.pdb
pdbbs/split_chain/1E4Y_A.pdb
pdbbs/split_chain/3X2S_A.pdb
pdbbs/split_chain/6HAP_A.pdb
pdbbs/split_chain/6HAM_A.pdb
pdbbs/split_chain/4K46_A.pdb
pdbbs/split_chain/3GMT_A.pdb
pdbbs/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

```

```

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

# Draw schematic alignment
plot(pdb, labels=ids)

```

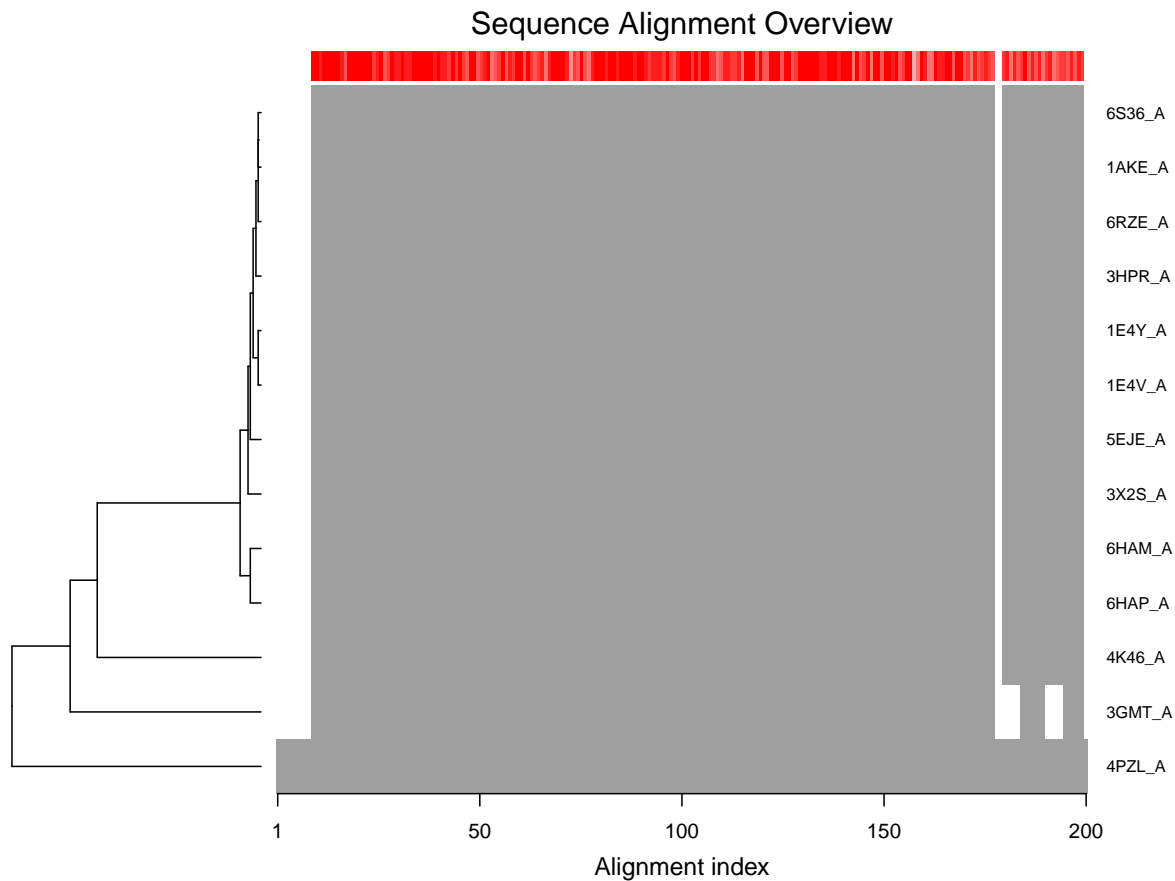


Figure 7: Schematic representation of alignment. Grey regions depict aligned residues, while white depict gap regions. The red bar at the top depict sequence conservation.

## Principal component analysis

PCA can be performed on the structural ensemble (stored in the `pdb`s object) with the function `pca.xyz()`, or more simply `pca()`.

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



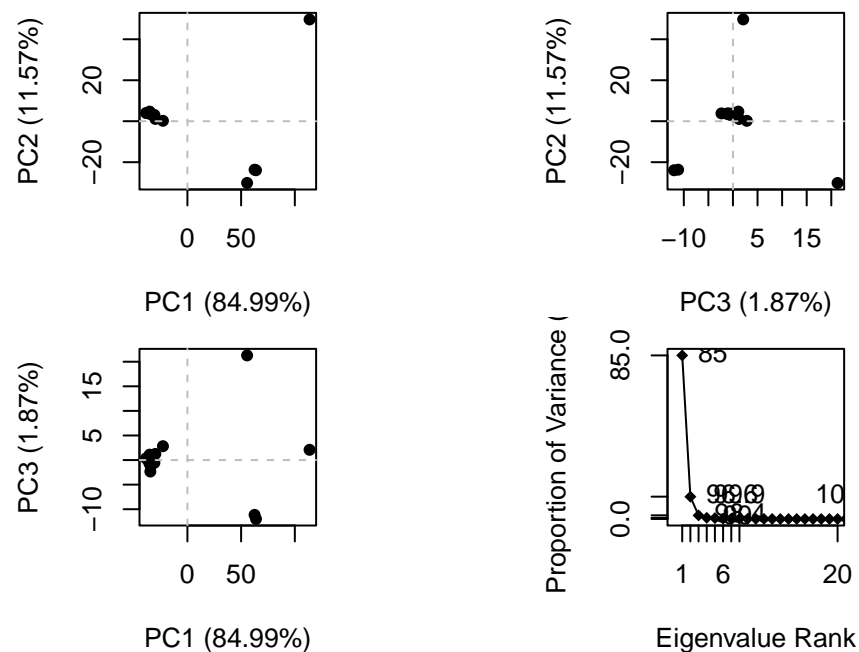


Figure 9: Results of PCA on Adenylate kinase X-ray structures. Each dot represents one PDB structure.

Function `rmsd()` will calculate all pairwise RMSD values of the structural ensemble. This facilitates clustering analysis based on the pairwise structural deviation

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

Warning in `rmsd(pdb)`: 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)
```

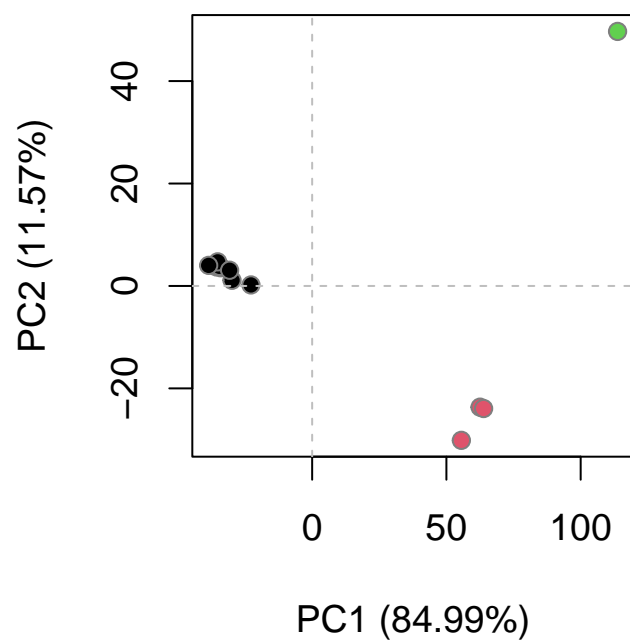


Figure 10: Projection of Adenylate kinase X-ray structures. Each dot represents one PDB structure.