

# Class09\_Structural1

AUTHOR

Zixuan Zeng (A16142927)

## PDB statistics

```
PDBstat <- read.csv("Data_Export_Summary.csv", row.names = 1)
PDBstat
```

|                         | X.ray   | EM     | NMR    | Integrative | Multiple.methods |
|-------------------------|---------|--------|--------|-------------|------------------|
| Protein (only)          | 176,204 | 20,299 | 12,708 | 342         | 218              |
| Protein/Oligosaccharide | 10,279  | 3,385  | 34     | 8           | 11               |
| Protein/NA              | 9,007   | 5,897  | 287    | 24          | 7                |
| Nucleic acid (only)     | 3,066   | 200    | 1,553  | 2           | 15               |
| Other                   | 173     | 13     | 33     | 3           | 0                |
| Oligosaccharide (only)  | 11      | 0      | 6      | 0           | 1                |

|                         | Neutron | Other | Total   |
|-------------------------|---------|-------|---------|
| Protein (only)          | 83      | 32    | 209,886 |
| Protein/Oligosaccharide | 1       | 0     | 13,718  |
| Protein/NA              | 0       | 0     | 15,222  |
| Nucleic acid (only)     | 3       | 1     | 4,840   |
| Other                   | 0       | 0     | 222     |
| Oligosaccharide (only)  | 0       | 4     | 22      |

```
# Convert numeric-looking columns from character to numeric
PDBstat[] <- lapply(PDBstat, function(x) as.numeric(gsub(",", "", x)))

# Check structure
str(PDBstat)
```

```
'data.frame': 6 obs. of 8 variables:
 $ X.ray      : num  176204 10279 9007 3066 173 ...
 $ EM         : num  20299 3385 5897 200 13 ...
 $ NMR        : num  12708 34 287 1553 33 ...
 $ Integrative : num  342 8 24 2 3 0
 $ Multiple.methods: num  218 11 7 15 0 1
 $ Neutron    : num  83 1 0 3 0 0
 $ Other      : num  32 0 0 1 0 4
 $ Total      : num  209886 13718 15222 4840 222 ...
```

```
(sum(PDBstat$X.ray) + sum(PDBstat$EM))/sum(PDBstat$Total) * 100
```

```
[1] 93.69604
```

Q1. 93.7% of the structures in the PDB were determined by X-ray and EM.

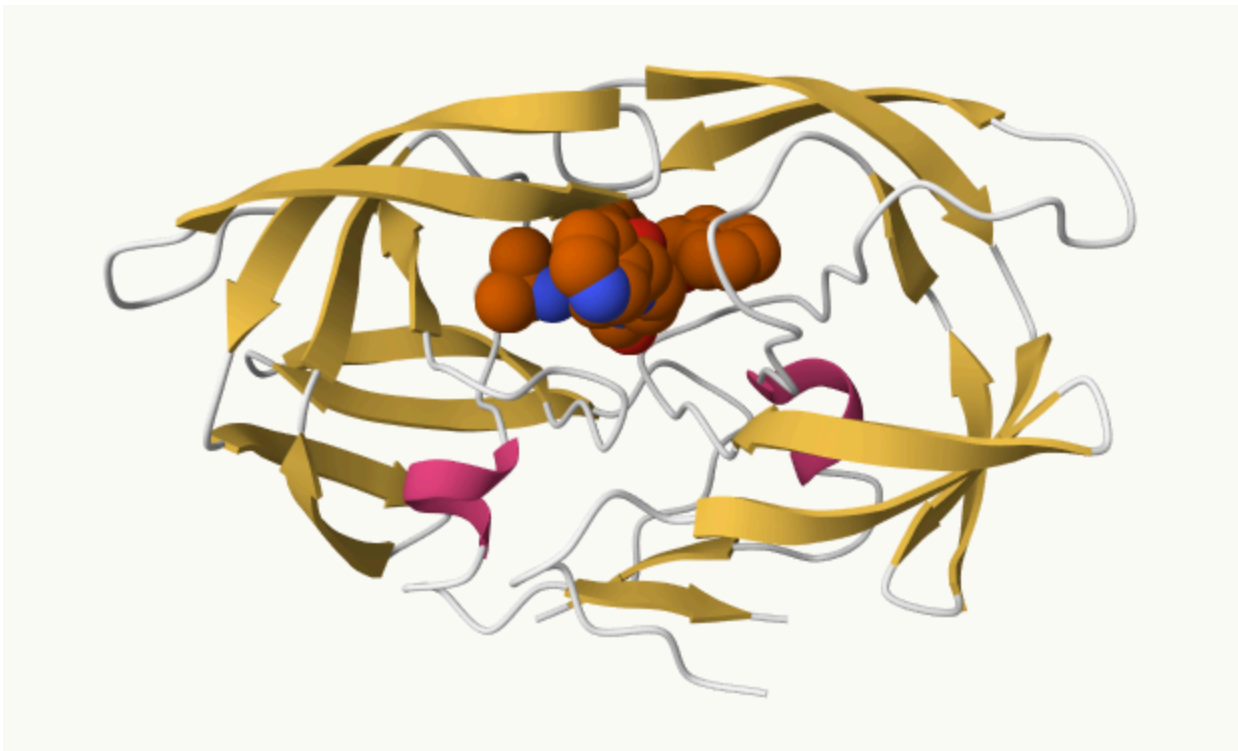
```
sum(PDBstat[c("Protein (only)", "Protein/Oligosaccharide", "Protein/NA"), "Total"])/sum(P
```

[1] 97.91562

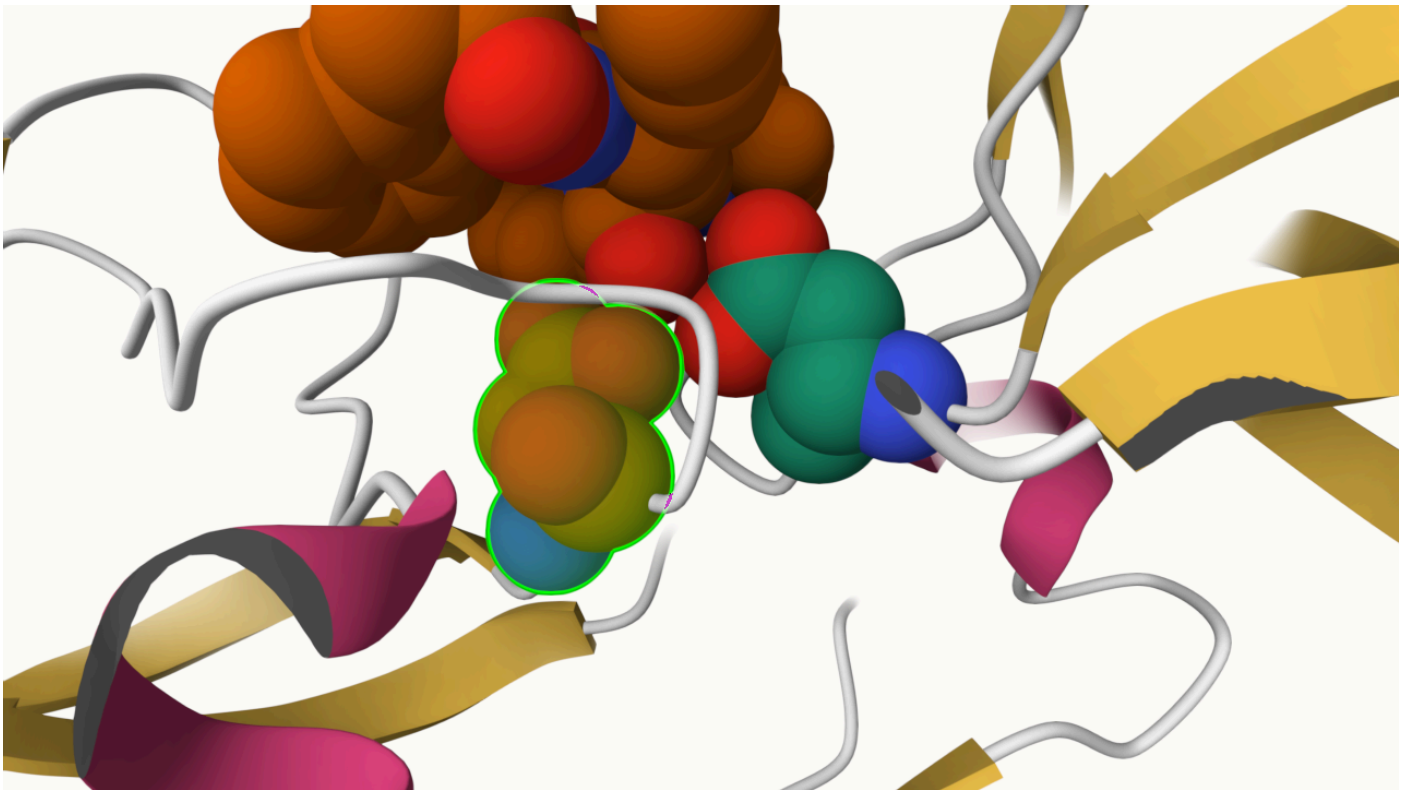
Q2. 97.9% of the structures in the PDB contain protein.

Q3. HIV currently has 4865 structure files in PDB, HIV-1 protease has 1150 structure files in PDB, which is about 23.6% of the total HIV structure files.

```
knitr::include_graphics("1HSG.png")
```



```
knitr::include_graphics("1HSG (1).png")
```

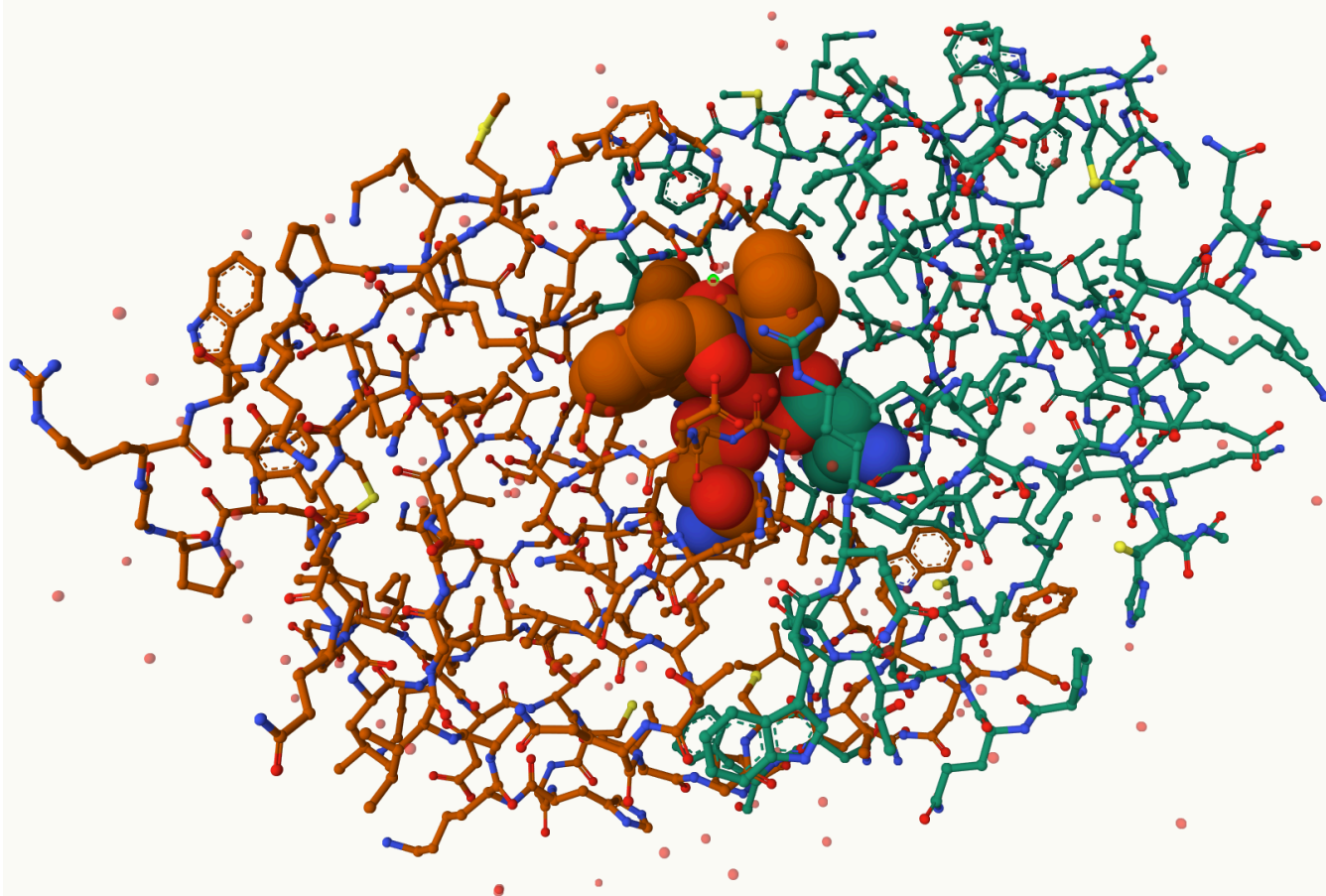


Q4. We can only see the oxygen atom from a water molecule because of the resolution of the structure. Hydrogen atom is not visible at this resolution.

Q5. HOH 308

Q6.

```
knitr::include_graphics("1HSG (2).png")
```



## Introduction to Bio3D in 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. 198 amino acid residues

Q8. HOH and MK1

Q9. 2 protein chains

```
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>   |

```
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  
TDELVIALVKERIAQEDCRNGFLLDGFPRITPQADAMKEAGINVDYVLEFDVPDELIVDKI  
VGRRVHAPSGRVYHVKFNPPKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQMTAPLIG  
YYSKEAEAGNTKYAKVDGTPVAEVRADLEKILG
```

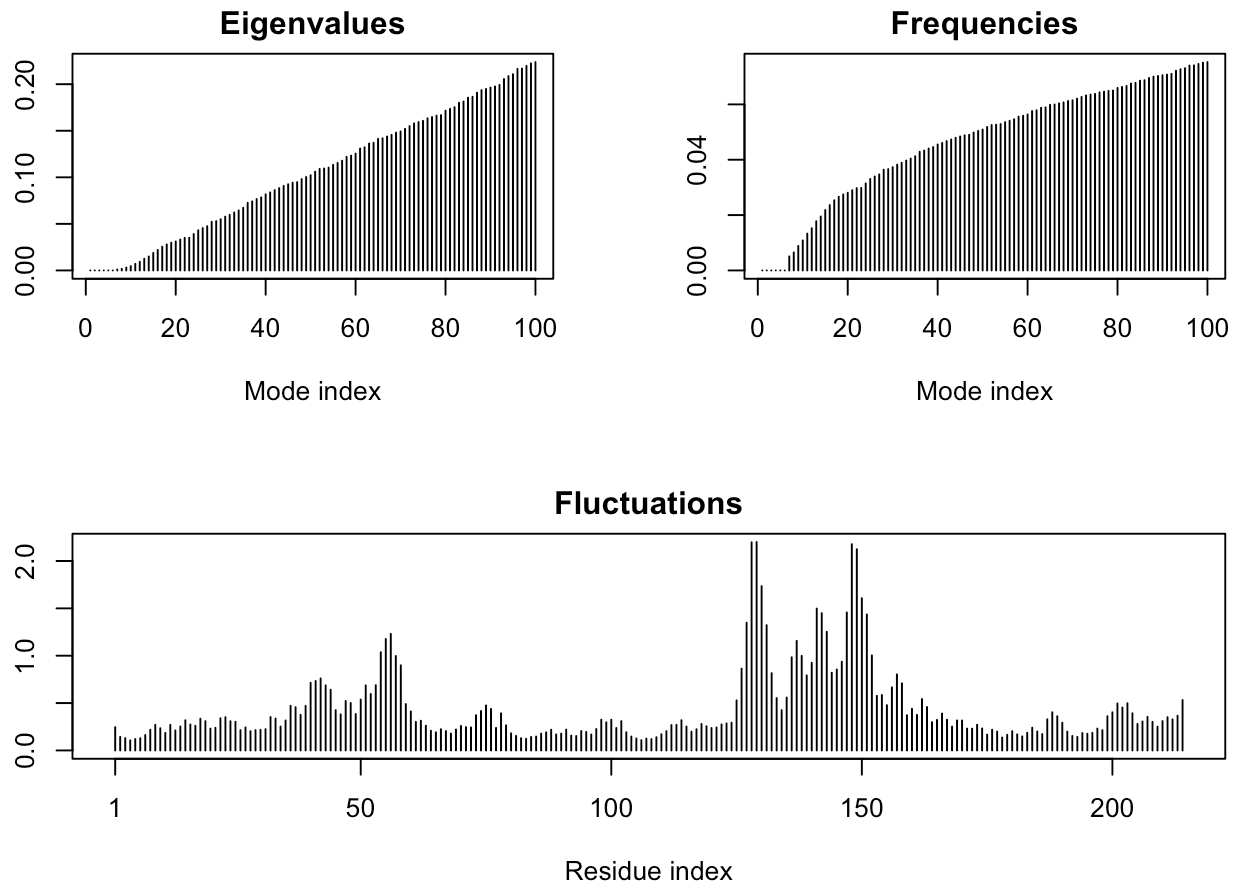
```
+ attr: atom, xyz, seqres, helix, sheet,  
      calpha, remark, call
```

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

```
Building Hessian... Done in 0.012 seconds.
```

```
Diagonalizing Hessian... Done in 0.263 seconds.
```

```
plot(m)
```



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

## Comparative struture analysis of Adenylate Kinase

Q10. msa

Q11. bio3d-view

Q12. True

```
library(bio3d)
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 MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMRLRAAVKSGSELGKQAKDIMDAGKLV
      1      .      .      .      .      .      60

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

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

Q13. 214 amino acid residues

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

Searching ... please wait (updates every 5 seconds) RID = FZ968N0M015

.....

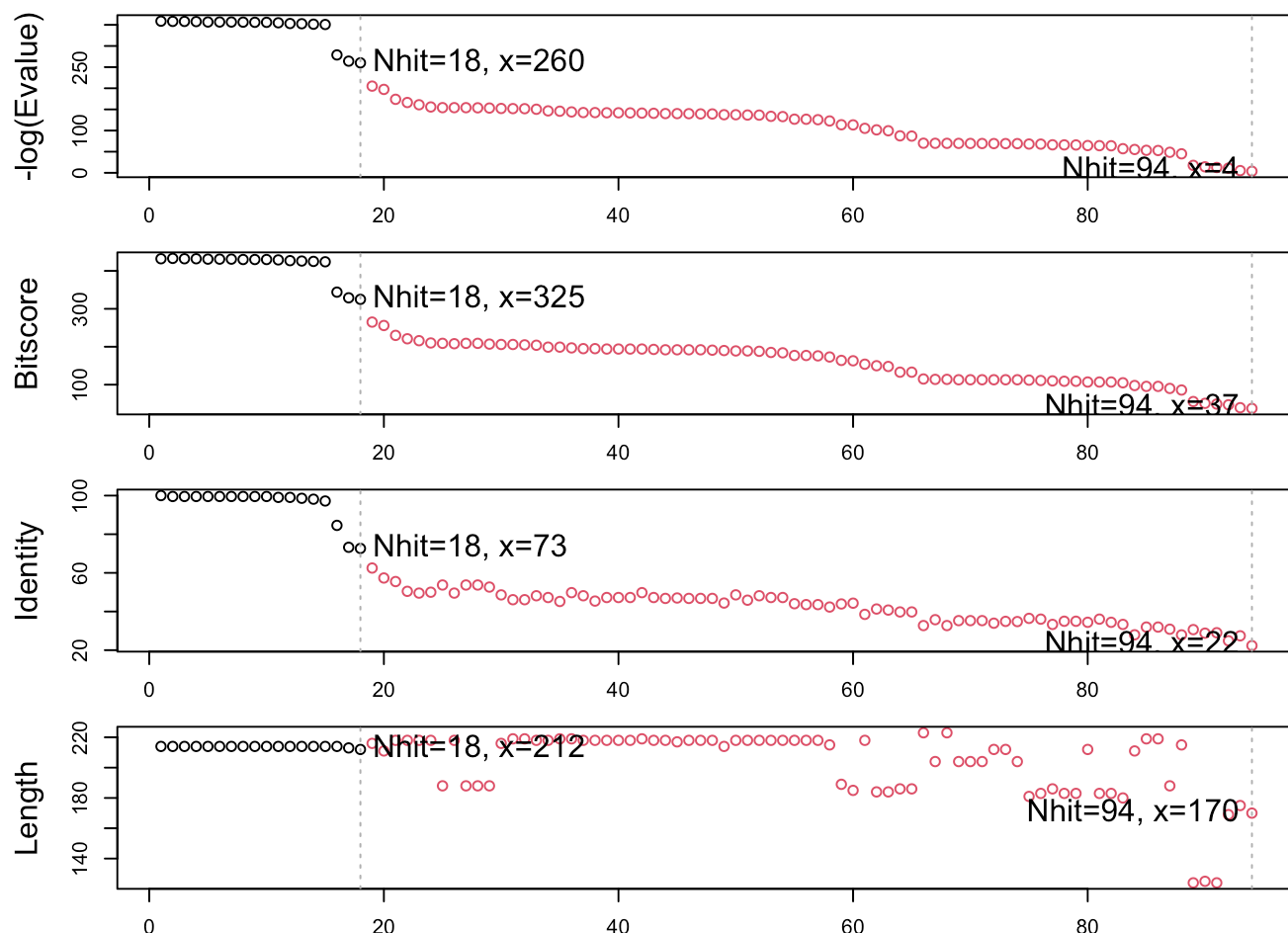
Reporting 94 hits

```
#Plot a summary of search results
hits <- plot(b)
```

```
* Possible cutoff values: 260 3
    Yielding Nhits:      18 94
```

```
* Chosen cutoff value of: 260
    Yielding Nhits:      18
```





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

```
[1] "1AKE_A" "8BQF_A" "4X8M_A" "6S36_A" "8Q2B_A" "8RJ9_A"
```

```
# 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

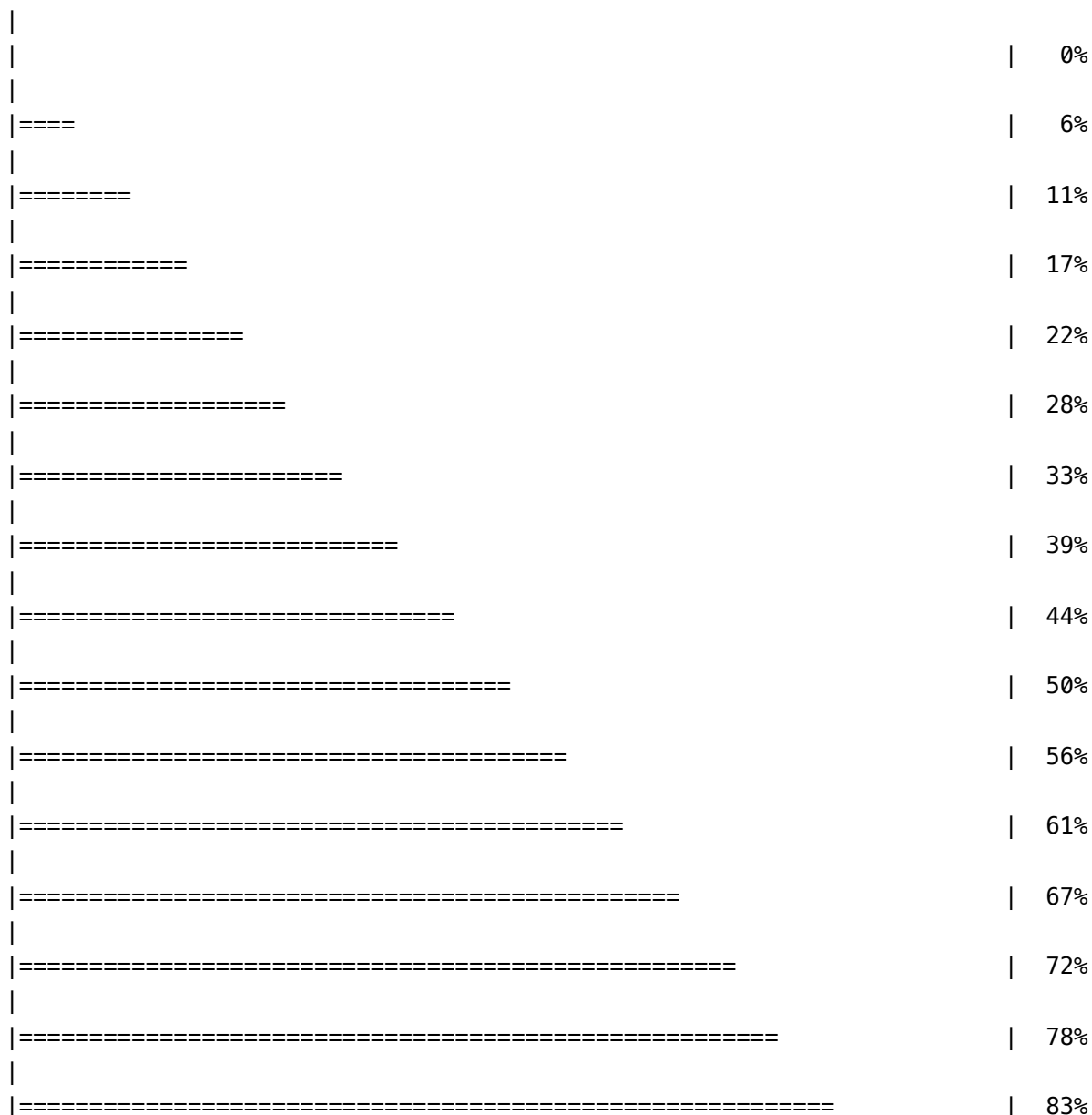
Warning in get.pdb(hits\$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):  
pdbs/1E4Y.pdb.gz exists. Skipping download

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



```

|
|=====| 89%
|
|=====| 94%
|
|=====| 100%

```

```

# Align related PDBs
pdbs <- pdbaln(files, fit = TRUE, exefile="msa")

```

Reading PDB files:

```

pdbs/split_chain/1AKE_A.pdb
pdbs/split_chain/8BQF_A.pdb
pdbs/split_chain/4X8M_A.pdb
pdbs/split_chain/6S36_A.pdb
pdbs/split_chain/8Q2B_A.pdb
pdbs/split_chain/8RJ9_A.pdb
pdbs/split_chain/6RZE_A.pdb
pdbs/split_chain/4X8H_A.pdb
pdbs/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/8PVW_A.pdb
pdbs/split_chain/4K46_A.pdb
pdbs/split_chain/4NP6_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
..  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/8BQF_A.pdb
  PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 3  name: pdbs/split_chain/4X8M_A.pdb
pdb/seq: 4  name: pdbs/split_chain/6S36_A.pdb

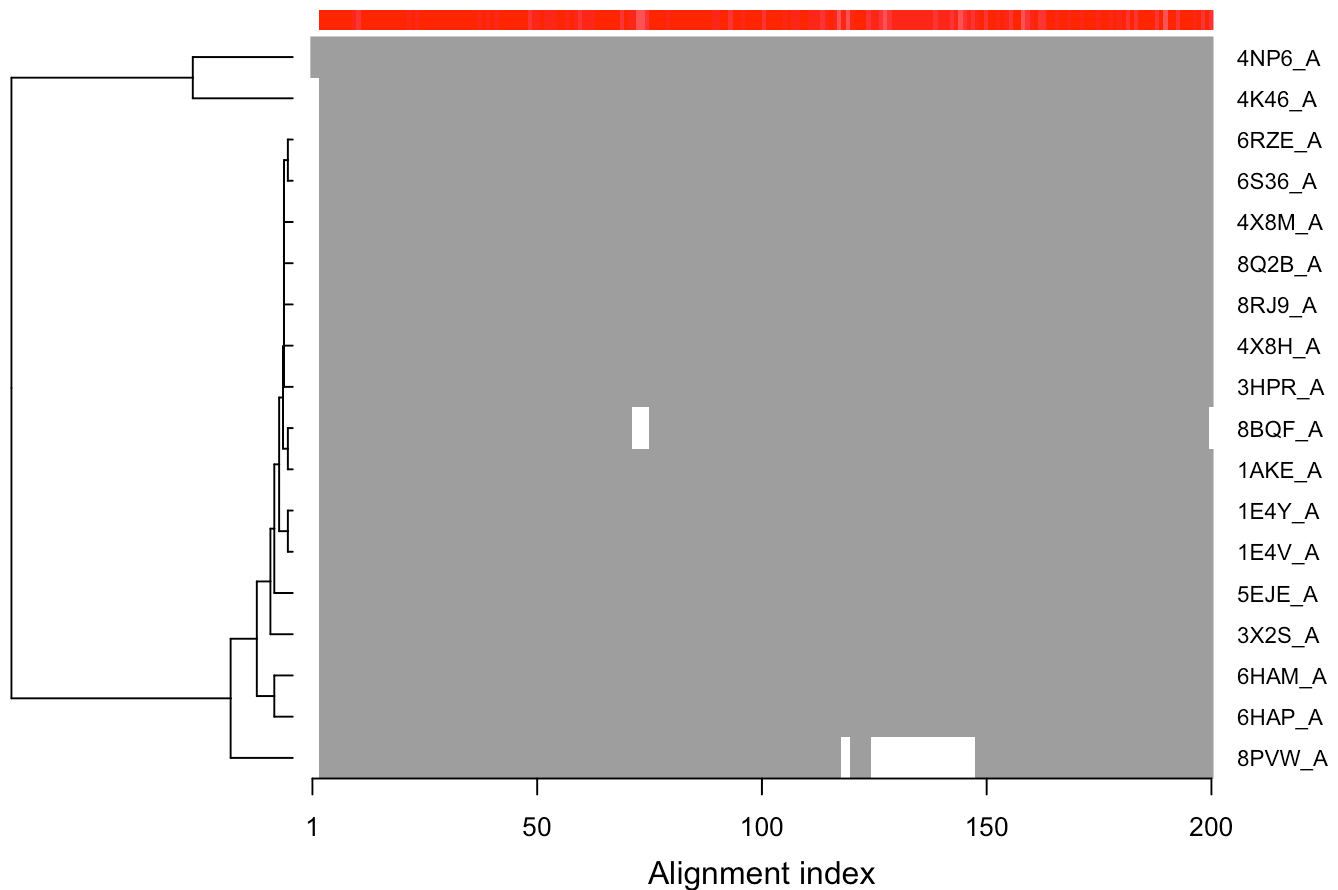
```

```
PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 5   name: pdbc/split_chain/8Q2B_A.pdb
PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 6   name: pdbc/split_chain/8RJ9_A.pdb
PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 7   name: pdbc/split_chain/6RZE_A.pdb
PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 8   name: pdbc/split_chain/4X8H_A.pdb
pdb/seq: 9   name: pdbc/split_chain/3HPR_A.pdb
PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 10  name: pdbc/split_chain/1E4V_A.pdb
pdb/seq: 11  name: pdbc/split_chain/5EJE_A.pdb
PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 12  name: pdbc/split_chain/1E4Y_A.pdb
pdb/seq: 13  name: pdbc/split_chain/3X2S_A.pdb
pdb/seq: 14  name: pdbc/split_chain/6HAP_A.pdb
pdb/seq: 15  name: pdbc/split_chain/6HAM_A.pdb
PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 16  name: pdbc/split_chain/8PVW_A.pdb
PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 17  name: pdbc/split_chain/4K46_A.pdb
PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 18  name: pdbc/split_chain/4NP6_A.pdb
```

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

# Draw schematic alignment
plot(pdbc, labels=ids)
```

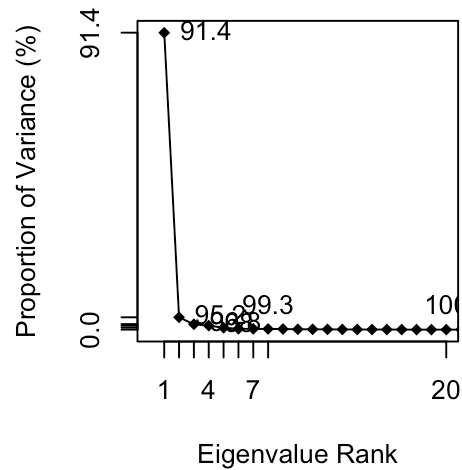
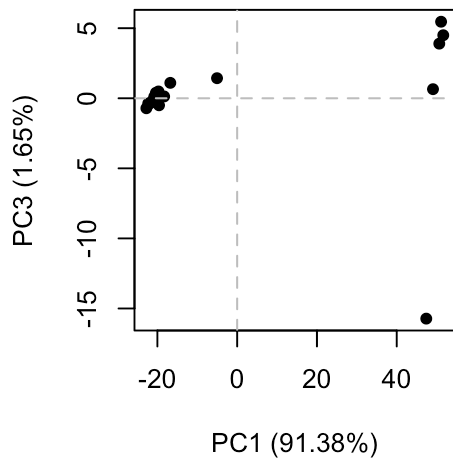
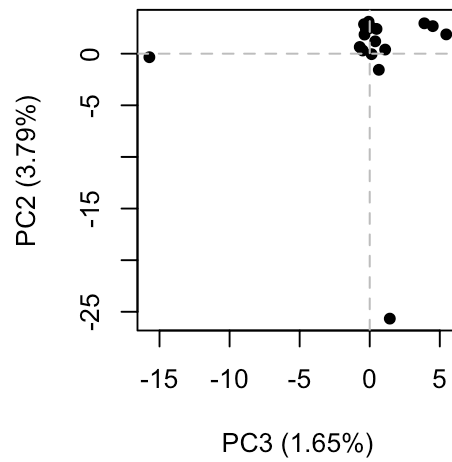
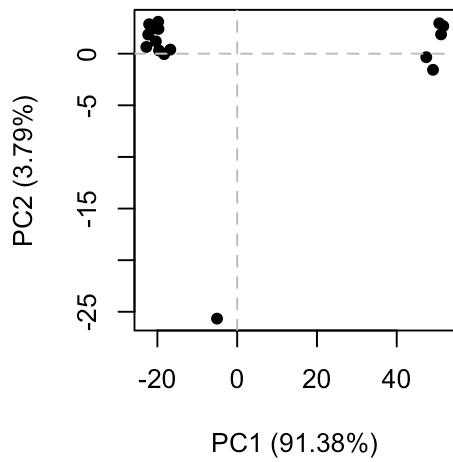
## Sequence Alignment Overview



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

```
[1] "Escherichia coli"
[2] "Escherichia coli K-12"
[3] "Escherichia coli 0139:H28 str. E24377A"
[4] "Escherichia coli str. K-12 substr. MDS42"
[5] "Photobacterium profundum"
[6] "Vibrio cholerae 01 biovar El Tor str. N16961"
```

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

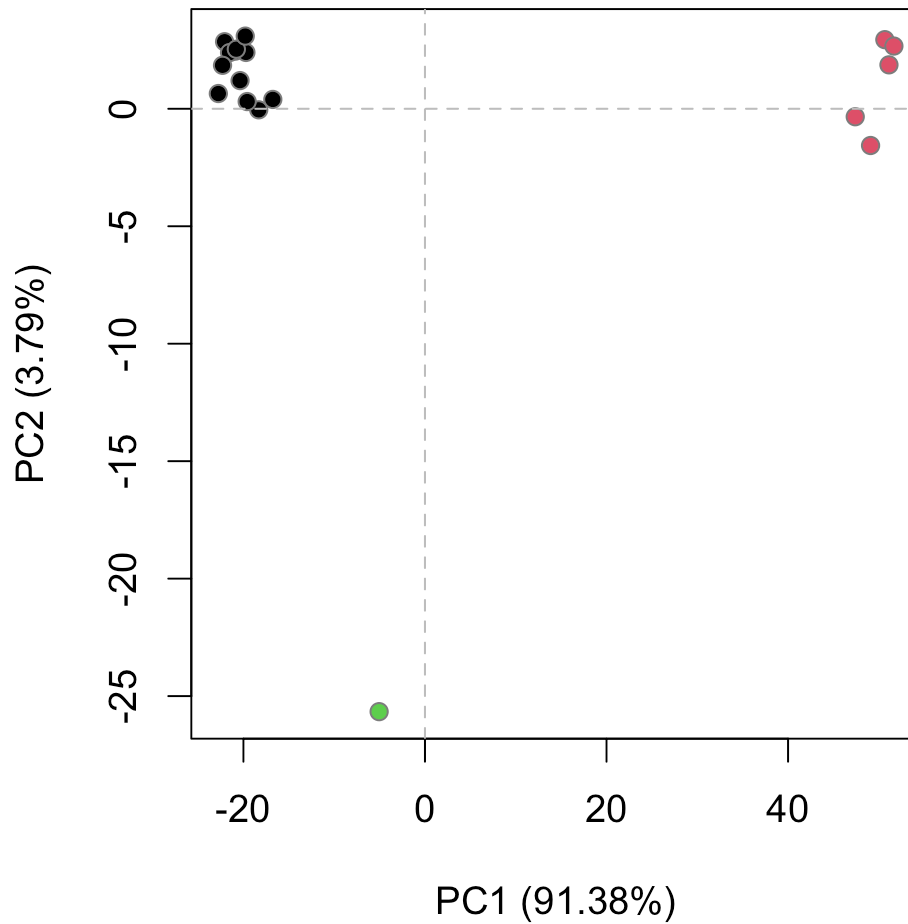


```
# Calculate RMSD
rd <- rmsd(pdb)
```

Warning in rmsd(pdb): No indices provided, using the 182 non NA positions

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

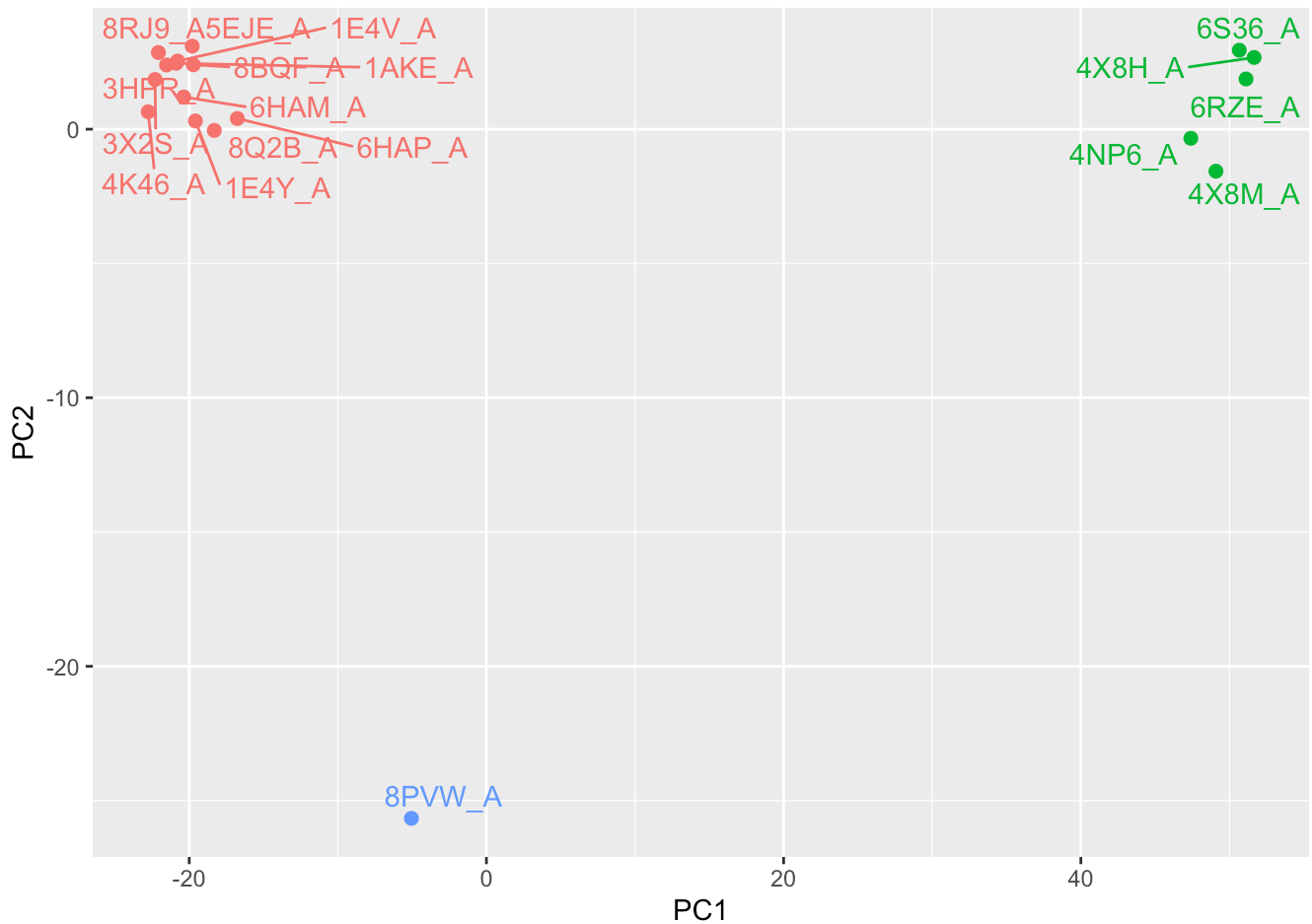
plot(pc.xray, 1:2, col="grey50", bg=grps.rd, pch=21, cex=1)
```



```
#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
```



## Normal mode analysis

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

Warning in nma.pdb(pdb): 8BQF\_A.pdb might have missing residue(s) in structure:  
Fluctuations at neighboring positions may be affected.

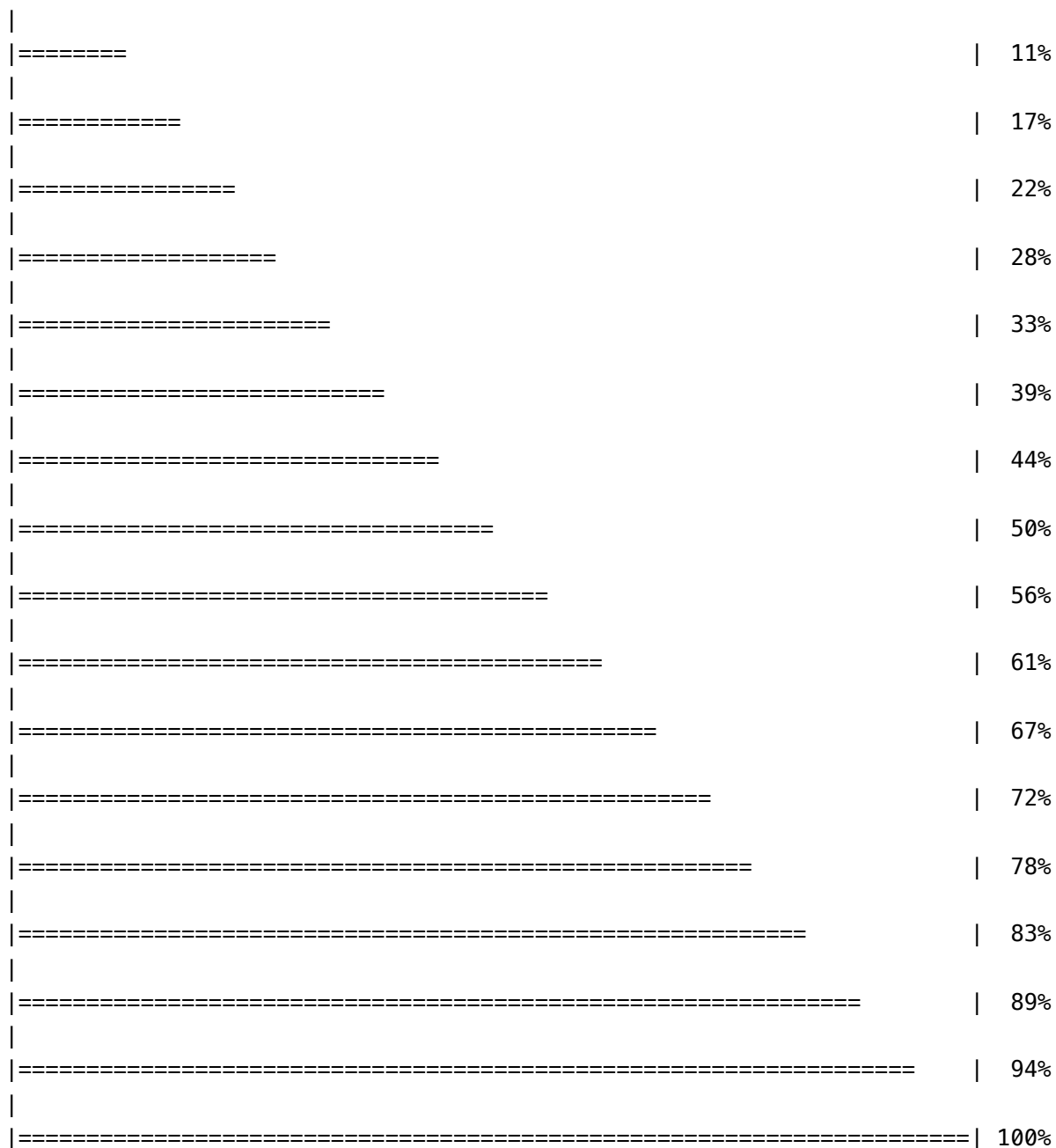
### Details of Scheduled Calculation:

```
... 18 input structures
... storing 540 eigenvectors for each structure
... dimension of x$U.subspace: ( 546x540x18 )
... coordinate superposition prior to NM calculation
... aligned eigenvectors (gap containing positions removed)
... estimated memory usage of final 'eNMA' object: 40.6 Mb
```

```
|
|
|
|====
```

```
| 0%
| 6%
```





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

Extracting SSE from pdbc\$sse attribute

