# Class 9:Structural Bioinformatics (Pt. 1)

Xiaoyan Wang(A16454055)

## PDB statistics

The main database for structural data is called the PDB(Protein Data Bank). Lets see what it contains:

Data from: https://www.rcsb.org/stats/summary

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

```
          Molecular.Type   X.ray     EM    NMR Multiple.methods Neutron Other
1          Protein (only) 167,192 15,572 12,529              208      77    32
2 Protein/Oligosaccharide   9,639  2,635     34                8       2     0
3              Protein/NA   8,730  4,697    286                7       0     0
4      Nucleic acid (only)   2,869    137  1,507               14       3     1
5                   Other     170     10     33                0       0     0
6   Oligosaccharide (only)     11      0      6                1       0     4
    Total
1 195,610
2  12,318
3  13,720
4   4,531
5     213
6      22
```

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

```
as.numeric(sub(",", "" , pdbData$Total))
```

```
[1] 195610  12318  13720   4531    213     22
```

```r
# make this into a function
x<- pdbData$Total
RmvComma <- function(x) {
  sub(",", "" , x)
}
```

```r
#test
as.numeric(RmvComma(x))
```

```
[1] 195610  12318  13720   4531    213     22
```

```r
apply(pdbData, 2, RmvComma)
```

```
     Molecular.Type            X.ray     EM      NMR     Multiple.methods
[1,] "Protein (only)"          "167192" "15572" "12529" "208"
[2,] "Protein/Oligosaccharide" "9639"   "2635"  "34"    "  8"
[3,] "Protein/NA"              "8730"   "4697"  "286"   "  7"
[4,] "Nucleic acid (only)"     "2869"   "137"   "1507"  " 14"
[5,] "Other"                   "170"    "10"    "33"    "  0"
[6,] "Oligosaccharide (only)"  "11"     "0"     "6"     "  1"
     Neutron Other Total
[1,] "77"    "32"  "195610"
[2,] " 2"    " 0"  "12318"
[3,] " 0"    " 0"  "13720"
[4,] " 3"    " 1"  "4531"
[5,] " 0"    " 0"  "213"
[6,] " 0"    " 4"  "22"
```

```r
#Alternatively...
#install.packages("tidyverse")
#install.packages("readr")
library(readr)
```

```
Warning: package 'readr' was built under R version 4.3.3
```

```r
pdbData<- 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.
```

```
head(pdbData)
```

```
# 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)       167192 15572 12529                208      77    32 195610
2 Protein/Oligosacc~     9639  2635    34                  8       2     0  12318
3 Protein/NA             8730  4697   286                  7       0     0  13720
4 Nucleic acid (onl~     2869   137  1507                 14       3     1   4531
5 Other                   170    10    33                  0       0     0    213
6 Oligosaccharide (~       11     0     6                  1       0     4     22
```

```
xDivy <- function(pdbData) {
  sum(x)/sum(y)*100
}
```

```
x<-pdbData$'X-ray'
y<-pdbData$Total
xDivy()
```

```
[1] 83.30359
```

```
x<-pdbData$'EM'
y<-pdbData$Total
xDivy()
```

```
[1] 10.18091
```

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

```
# Extract rows where "Molecular Type" contains "protein"
protein_rows <- pdbData[grepl("Protein", pdbData$`Molecular Type`, ignore.case = TRUE), ]
x <- protein_rows$Total
y <- pdbData$Total
xDivy()
```

```
[1] 97.89501
```

**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?

This query matches 4,563 structures.

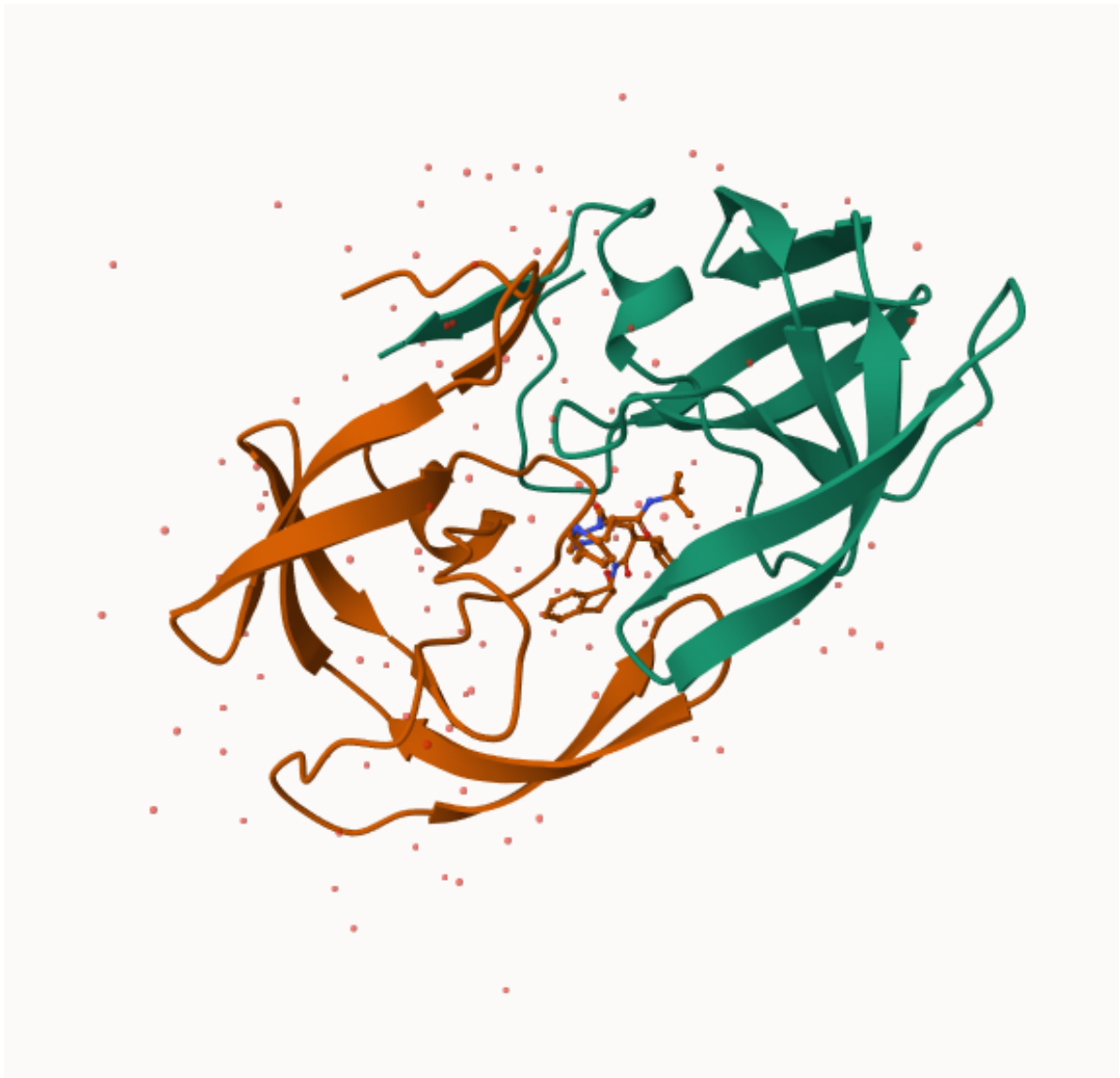## The PDB format

##Mol*

We will use the PDB code: 1HSG

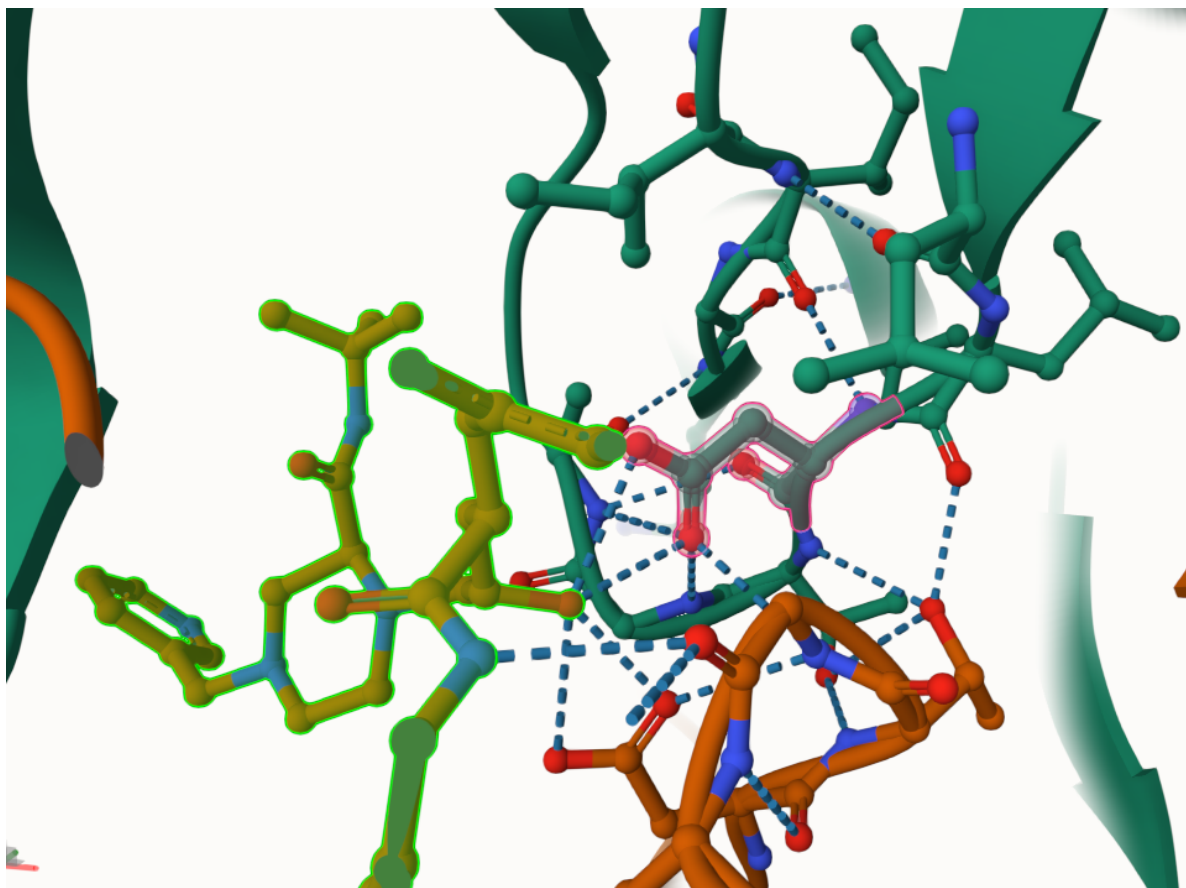Figure 1. A overview of 1HSG

Figure 2. Ligand shown in spacefill

Figure 3. Focus of D25 position(shown in pink)

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

The one "atom" here represents a residue, which is HOH for water. In other words, the 3 atoms of water is considered as one residue and the 2 hydrogen were hidden to not be presented.

**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

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.
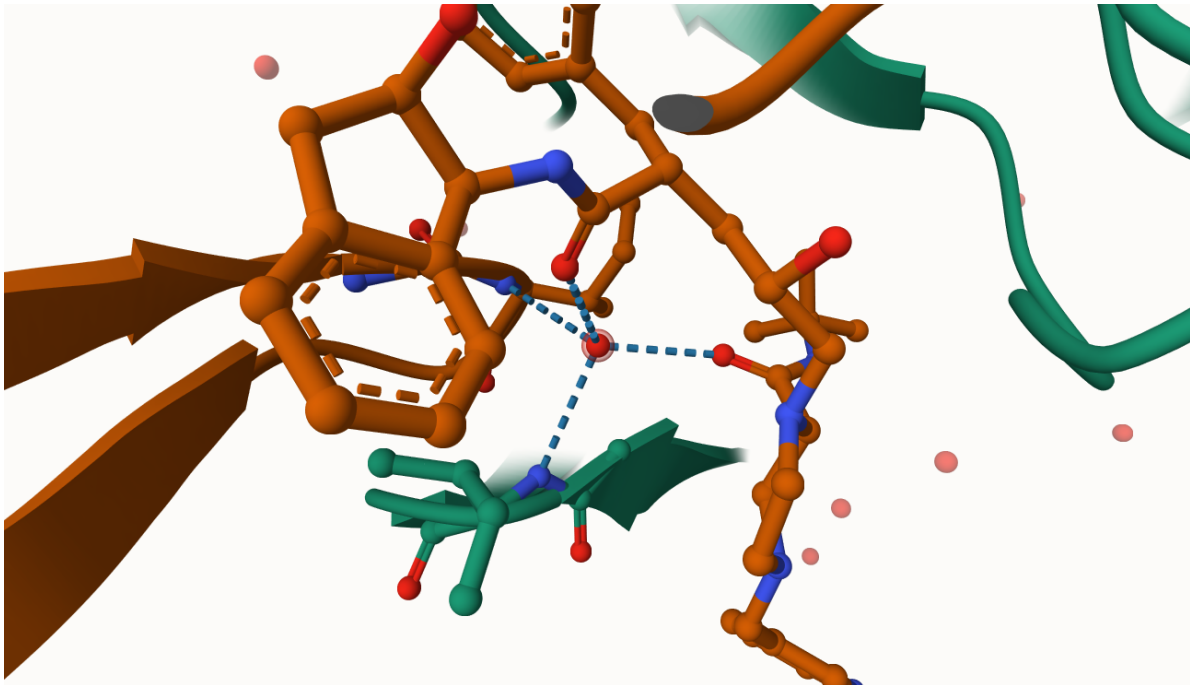
Figure 4. A representation of HOH 308

## Introduction to Bio3D in R

```
library(bio3d)
```

Warning: package 'bio3d' was built under R version 4.3.3

```
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

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

MK1

**Q9:** How many protein chains are in this structure?

2

```
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:
      MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMLRAAVKSGSELGKQAKDIMDAGKLVT
      DELVIALVKERIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFDVPDELIVDKI
      VGRRVHAPSGRVYHVKFNPPKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQMTAPLIG
      YYSKEAEAGNTKYAKVDGTKPVAEVRADLEKILG

+ attr: atom, xyz, seqres, helix, sheet,
        calpha, remark, call
```
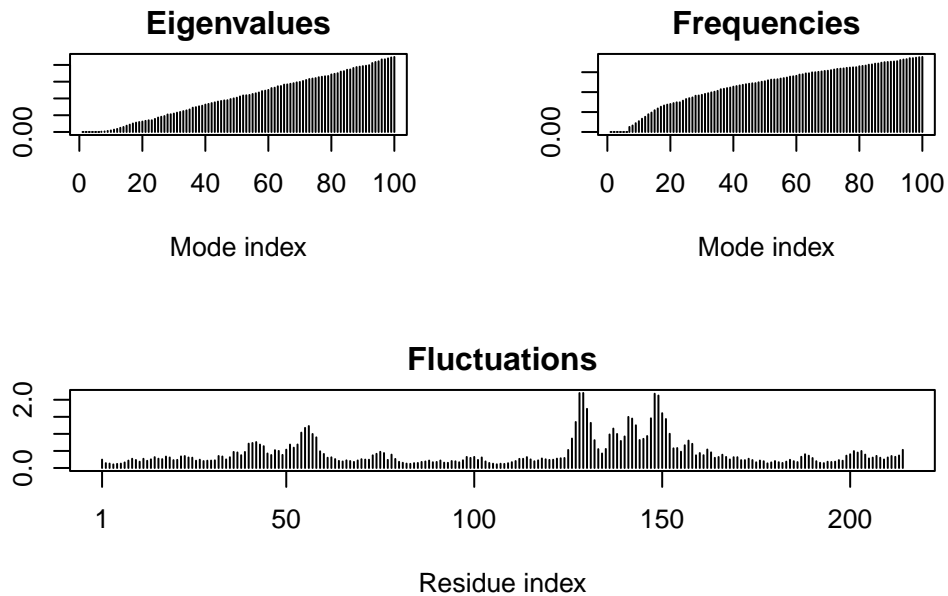
```
# Perform flexiblity prediction
m <- nma(adk)
```

```
Building Hessian...        Done in 0.02 seconds.
Diagonalizing Hessian...   Done in 0.14 seconds.
```

```
plot(m)
```

**Eigenvalues**

**Frequencies**

**Fluctuations**

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

**Comparative structure analysis of Adenylate Kinase**

```
# 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?

  msa

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

  bio3d-view

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

  TRUE

```
library(bio3d)
aa <- "1ake_A"
get.seq(aa)
```

```
Warning in get.seq(aa): Removing existing file: seqs.fasta

Fetching... Please wait. Done.

            1        .         .         .         .         .          60
pdb|1AKE|A    MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMLRAAVKSGSELGKQAKDIMDAGKLVT
            1        .         .         .         .         .          60

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

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

            181      .         .         .   214
pdb|1AKE|A    YYSKEAEAGNTKYAKVDGTKPVAEVRADLEKILG
            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?

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

```r
# Download releated 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 exists. Skipping download
```

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

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

```
Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/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%
  |
  |==============================                                    |  46%
  |
  |===================================                               |  54%
  |
  |========================================                          |  62%
  |
  |=============================================                     |  69%
  |
  |==================================================                |  77%
  |
  |=======================================================           |  85%
  |
  |============================================================      |  92%
  |
  |==================================================================| 100%
```

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

```
Reading PDB files:
pdbs/split_chain/1AKE_A.pdb
pdbs/split_chain/6S36_A.pdb
pdbs/split_chain/6RZE_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/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
```

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

# Draw schematic alignment
#plot(pdbs, labels=ids)
# Error: Figure margins too large
```

```
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 (ADK)
6RZE_A        1.69              <NA>                   Adenylate kinase (ADK)
3HPR_A        2.00              <NA>                   Adenylate kinase (ADK)
1E4V_A        1.85 Adenylate kinase                   Adenylate kinase (ADK)
5EJE_A        1.90              <NA>                   Adenylate kinase (ADK)
            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
```
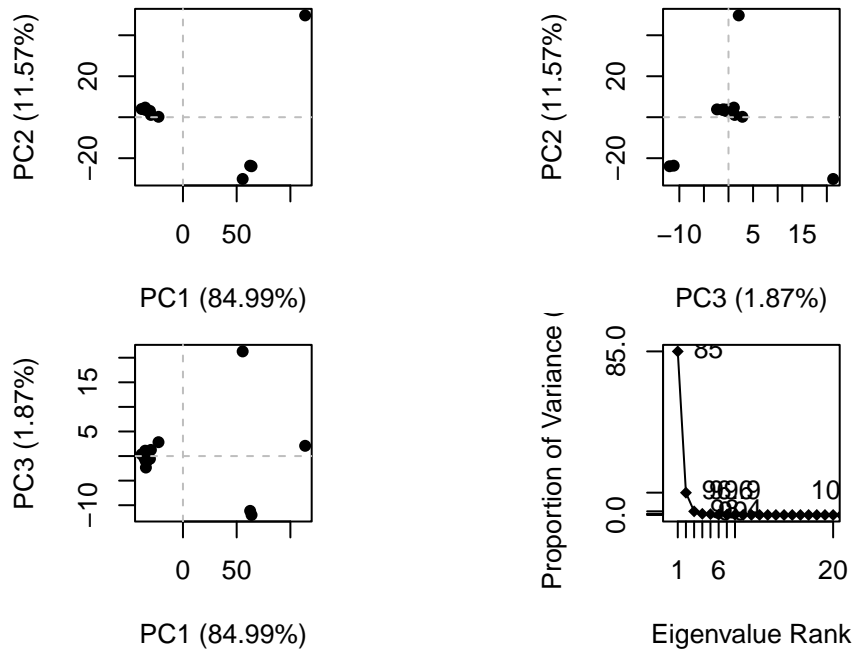
```r
# Perform PCA
pc.xray <- pca(pdbs)
plot(pc.xray)
```

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

```
Warning in rmsd(pdbs): 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)
```