# BIMM 143 Class 9 Structural Bioinformatics Pt1

## Daniel Gurholt (PID: A16767491)

The main database for structural data is called the PBD (Protein Data Bank). Let's see what it contains:

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

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

| | X.ray | EM | NMR | Multiple.methods | Neutron | Other |
|---|---|---|---|---|---|---|
| Protein (only) | 167,192 | 15,572 | 12,529 | 208 | 77 | 32 |
| Protein/Oligosaccharide | 9,639 | 2,635 | 34 | 8 | 2 | 0 |
| Protein/NA | 8,730 | 4,697 | 286 | 7 | 0 | 0 |
| Nucleic acid (only) | 2,869 | 137 | 1,507 | 14 | 3 | 1 |
| Other | 170 | 10 | 33 | 0 | 0 | 0 |
| Oligosaccharide (only) | 11 | 0 | 6 | 1 | 0 | 4 |

| | Total |
|---|---|
| Protein (only) | 195,610 |
| Protein/Oligosaccharide | 12,318 |
| Protein/NA | 13,720 |
| Nucleic acid (only) | 4,531 |
| Other | 213 |
| Oligosaccharide (only) | 22 |

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

```
pdbdb$Total
```

```
[1] "195,610" "12,318"  "13,720"  "4,531"   "213"       "22"
```

I need to remove the comma and convert to numeric to do math:

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

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

```
#as.numeric(pdbdb$Total)
```

I could turn this into a function to fix the whole table or any future table I read like this

```
x<- pdbdb$Total
as.numeric( sub(",", "", x))
```

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

```
comma2numeric<- function(x) {
  as.numeric( sub(",", "", x))
}
```

Test it

```
comma2numeric(pdbdb$X.ray)
```

```
[1] 167192   9639   8730   2869    170     11
```

```
apply(pdbdb,  2, comma2numeric)
```

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

##Or try a different read/import function

2

```
library(readr)
pdbdb<- 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.
```

```
pdbdb
```

```
# 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
```

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

```
sum(pdbdb$`X-ray`)/(sum(pdbdb$Total)) * 100
```

```
[1] 83.30359
```

From the calculations above, 83.3% of structures in the PDB are solved by X-Ray and Electron Microscopy.

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

```
pdbdb$Total[1]/(sum(pdbdb$Total)) * 100
```

```
[1] 86.39483
```

From the calculations above, 86.4% 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 4,563 protease structures in the current PDB that showed up in the results.

## Mol

Mol* (pronounced "molstar") is a new web-based molecule viewer that we will need to learn the basics of here.

https://molstar.org/viewer/
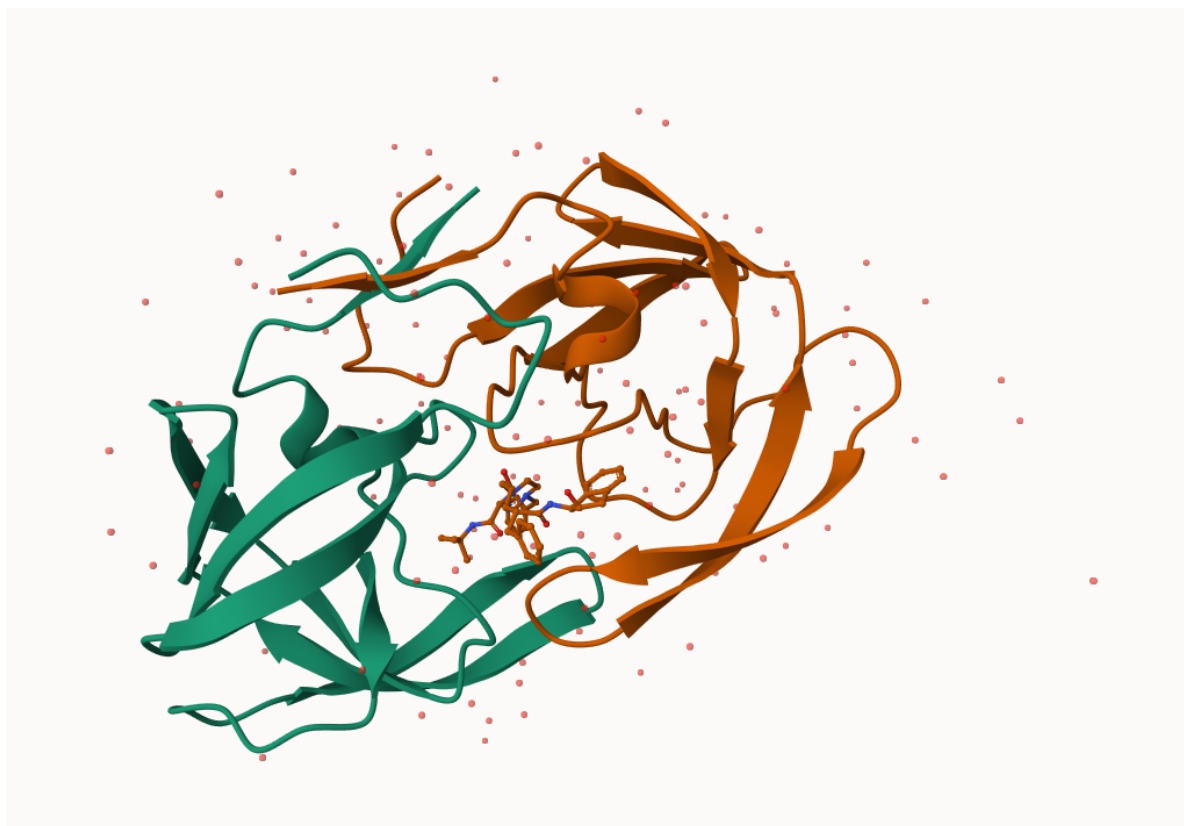
We will use PDB code: 1HSG



Figure 1: A first image from molstar
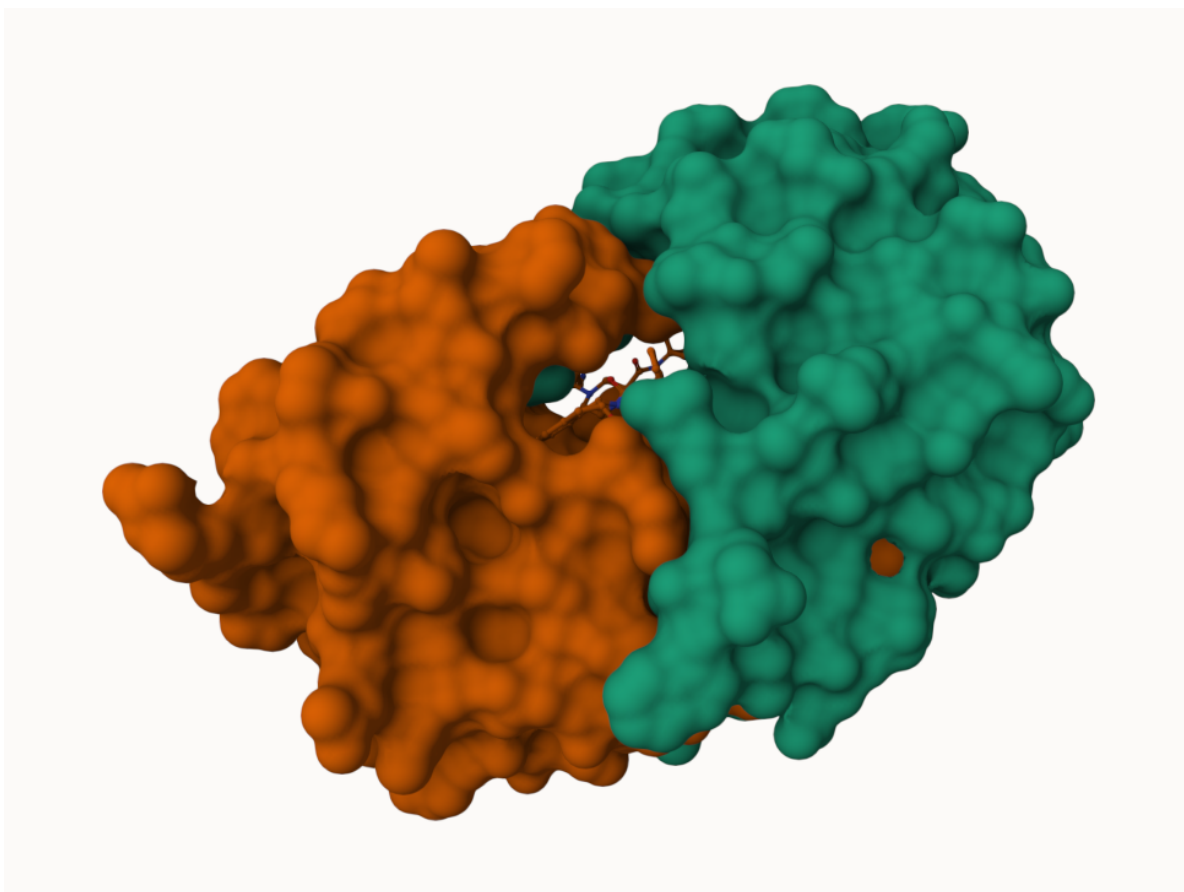
Figure 2: Modified 1HSG from molstar

Figure 3: Molecular Surface Pore 1HSG from molstar

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

We just see the oxygen atom of the water molecule in this structure because there is so much water around to stabilize the protein that adding the hydrogen atoms would make it hard to visualize and analyze. Additionally, all the atoms on the proteins do not show the hydrogen present as well so it reduces overall complexity while still showing bonds and interactions which mainly occurs on the oxygen atom anyways.

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

Yes I found the critical "conserved" water molecule in the binding site and it has water residue number 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.
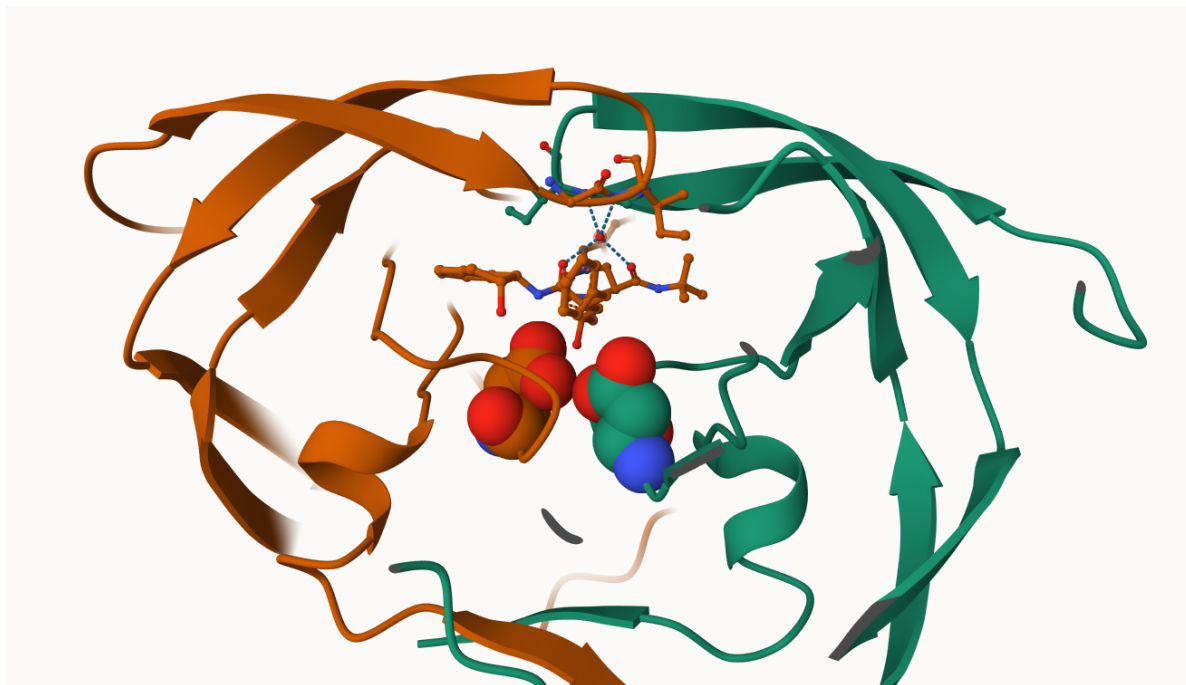


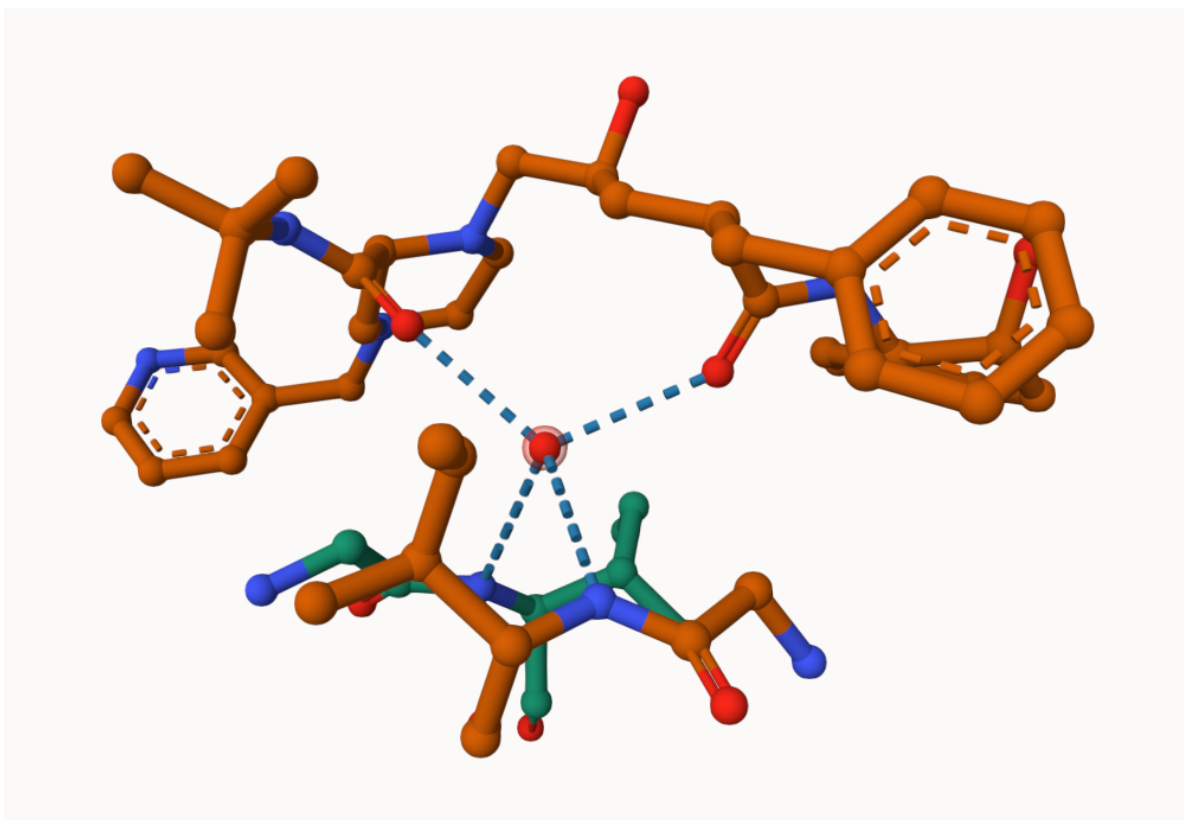Figure 4: water, chains, ASP 25 in 1HSG from molstar

Figure 5: Critical conserved water 308 in 1HSG from molstar

## The Bio3D package

The bio3d package allows us to do all sorts of structural bioinformatics work in R. Let's start with how it can read these PDB files:

```r
library (bio3d)

pdb<- read.pdb("1hsg")
```

```
Note: Accessing on-line PDB file
```

```r
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
```

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>
```

```
pdbseq(pdb)[25]
```

```
 25
"D"
```

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

```
sum(pdb$calpha)
```

```
[1] 198
```

There are 198 amino acid residues in this pdb object

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

HOH and MK1

Q9: How many protein chains are in this structure?

```
unique(pdb$atom$chain)
```

```
[1] "A" "B"
```

There are 2 unique protein chains are in this structure. chains A and B

##Predicting functional motions of a single structure

Let's do a bioinformatics prediction of functional motions - i.e the movements that one of these molecules needs to make to do its stuff.

```
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 a flexibility prediction

m<- nma(adk)
```

```
 Building Hessian...        Done in 0.08 seconds.
 Diagonalizing Hessian...   Done in 0.61 seconds.
```

```r
plot(m)
```

**Eigenvalues**

**Frequencies**

**Fluctuations**

Write out a multi_model PDB file that we can use to make an animation of the predicted motions.

```r
mktrj(m, file="adk.pdb")
```

I can open this in molstar to play the trajectory

##Comparative structure analysis of Adenylate Kinase

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

msa is found only on BioConductor and not CRAN

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

Bio3d-view 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, functions from the devtools package can be used to install packages from GitHub and BitBucket

## Comparative Analysis of protein structures ## Search and retrieve ADK structures

```
library(bio3d)

## Here we will find and analyze all ADK structures in the PBD database
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    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?

```
length(aa$ali)
```

```
[1] 214
```

There are 214 amino acids are in this sequence which means it is 214 amino acids in length just by looking at the sequencing results above.

```
#b <- blast.pdb(aa)
```

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

```
#head(hits$pdb.id)
```

##Pre calculated Results

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

```
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%
```

## Align and superpose structures

```
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
```

```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
```

pdbs

```
                                    1        .        .        .        40
[Truncated_Name:1]1AKE_A.pdb        ----------MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:2]6S36_A.pdb        ----------MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:3]6RZE_A.pdb        ----------MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:4]3HPR_A.pdb        ----------MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:5]1E4V_A.pdb        ----------MRIILLGAPVAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:6]5EJE_A.pdb        ----------MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:7]1E4Y_A.pdb        ----------MRIILLGALVAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:8]3X2S_A.pdb        ----------MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:9]6HAP_A.pdb        ----------MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:10]6HAM_A.pdb       ----------MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
```

```
[Truncated_Name:11]4K46_A.pdb    ----------MRIILLGAPGAGKGTQAQFIMAKFGIPQIS
[Truncated_Name:12]3GMT_A.pdb    ----------MRLILLGAPGAGKGTQANFIKEKFGIPQIS
[Truncated_Name:13]4PZL_A.pdb    TENLYFQSNAMRIILLGAPGAGKGTQAKIIEQKYNIAHIS
                                           **^***** *******  *  *^ *  **
                                 1         .         .         .        40


                                 41        .         .         .        80
[Truncated_Name:1]1AKE_A.pdb     TGDMLRAAVKSGSELGKQAKDIMDAGKLVTDELVIALVKE
[Truncated_Name:2]6S36_A.pdb     TGDMLRAAVKSGSELGKQAKDIMDAGKLVTDELVIALVKE
[Truncated_Name:3]6RZE_A.pdb     TGDMLRAAVKSGSELGKQAKDIMDAGKLVTDELVIALVKE
[Truncated_Name:4]3HPR_A.pdb     TGDMLRAAVKSGSELGKQAKDIMDAGKLVTDELVIALVKE
[Truncated_Name:5]1E4V_A.pdb     TGDMLRAAVKSGSELGKQAKDIMDAGKLVTDELVIALVKE
[Truncated_Name:6]5EJE_A.pdb     TGDMLRAAVKSGSELGKQAKDIMDACKLVTDELVIALVKE
[Truncated_Name:7]1E4Y_A.pdb     TGDMLRAAVKSGSELGKQAKDIMDAGKLVTDELVIALVKE
[Truncated_Name:8]3X2S_A.pdb     TGDMLRAAVKSGSELGKQAKDIMDCGKLVTDELVIALVKE
[Truncated_Name:9]6HAP_A.pdb     TGDMLRAAVKSGSELGKQAKDIMDAGKLVTDELVIALVRE
[Truncated_Name:10]6HAM_A.pdb    TGDMLRAAIKSGSELGKQAKDIMDAGKLVTDEIIIALVKE
[Truncated_Name:11]4K46_A.pdb    TGDMLRAAIKAGTELGKQAKSVIDAGQLVSDDIILGLVKE
[Truncated_Name:12]3GMT_A.pdb    TGDMLRAAVKAGTPLGVEAKTYMDEGKLVPDSLIIGLVKE
[Truncated_Name:13]4PZL_A.pdb    TGDMIRETIKSGSALGQELKKVLDAGELVSDEFIIKIVKD
                                 ****^*  ^* *^ **    *  ^*    ** *  ^^ ^*^^
                                 41        .         .         .        80


                                 81        .         .         .        120
[Truncated_Name:1]1AKE_A.pdb     RIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFD
[Truncated_Name:2]6S36_A.pdb     RIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFD
[Truncated_Name:3]6RZE_A.pdb     RIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFD
[Truncated_Name:4]3HPR_A.pdb     RIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFD
[Truncated_Name:5]1E4V_A.pdb     RIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFD
[Truncated_Name:6]5EJE_A.pdb     RIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFD
[Truncated_Name:7]1E4Y_A.pdb     RIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFD
[Truncated_Name:8]3X2S_A.pdb     RIAQEDSRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFD
[Truncated_Name:9]6HAP_A.pdb     RICQEDSRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFD
[Truncated_Name:10]6HAM_A.pdb    RICQEDSRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFD
[Truncated_Name:11]4K46_A.pdb    RIAQDDCAKGFLLDGFPRTIPQADGLKEVGVVVDYVIEFD
[Truncated_Name:12]3GMT_A.pdb    RLKEADCANGYLFDGFPRTIAQADAMKEAGVAIDYVLEID
[Truncated_Name:13]4PZL_A.pdb    RISKNDCNNGFLLDGVPRTIPQAQELDKLGVNIDYIVEVD
                                 *^    *   *^* ** **** **  ^    *^ ^**^^* *
                                 81        .         .         .        120


                                 121       .         .         .        160
[Truncated_Name:1]1AKE_A.pdb     VPDELIVDRIVGRRVHAPSGRVYHVKFNPPKVEGKDDVTG
[Truncated_Name:2]6S36_A.pdb     VPDELIVDKIVGRRVHAPSGRVYHVKFNPPKVEGKDDVTG
```

19

```
[Truncated_Name:3]6RZE_A.pdb     VPDELIVDAIVGRRVHAPSGRVYHVKFNPPKVEGKDDVTG
[Truncated_Name:4]3HPR_A.pdb     VPDELIVDRIVGRRVHAPSGRVYHVKFNPPKVEGKDDGTG
[Truncated_Name:5]1E4V_A.pdb     VPDELIVDRIVGRRVHAPSGRVYHVKFNPPKVEGKDDVTG
[Truncated_Name:6]5EJE_A.pdb     VPDELIVDRIVGRRVHAPSGRVYHVKFNPPKVEGKDDVTG
[Truncated_Name:7]1E4Y_A.pdb     VPDELIVDRIVGRRVHAPSGRVYHVKFNPPKVEGKDDVTG
[Truncated_Name:8]3X2S_A.pdb     VPDELIVDRIVGRRVHAPSGRVYHVKFNPPKVEGKDDVTG
[Truncated_Name:9]6HAP_A.pdb     VPDELIVDRIVGRRVHAPSGRVYHVKFNPPKVEGKDDVTG
[Truncated_Name:10]6HAM_A.pdb    VPDELIVDRIVGRRVHAPSGRVYHVKFNPPKVEGKDDVTG
[Truncated_Name:11]4K46_A.pdb    VADSVIVERMAGRRAHLASGRTYHNVYNPPKVEGKDDVTG
[Truncated_Name:12]3GMT_A.pdb    VPFSEIIERMSGRRTHPASGRTYHVKFNPPKVEGKDDVTG
[Truncated_Name:13]4PZL_A.pdb    VADNLLIERITGRRIHPASGRTYHTKFNPPKVADKDDVTG
                                 *      ^^^ ^ *** *  *** **  ^*****  *** **
                         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    EDLVIREDDKEETVLARLGVYHNQTAPLIAYYGKEAEAGN
[Truncated_Name:12]3GMT_A.pdb    EPLVQRDDDKEETVKKRLDVYEAQTKPLITYYGDWARRGA
[Truncated_Name:13]4PZL_A.pdb    EPLITRTDDNEDTVKQRLSVYHAQTAKLIDFYRNFSSTNT
                                 * *   * ** *^ *  **   *    *  ** ^*
                         161           .         .         .          200


                         201           .         .      227
[Truncated_Name:1]1AKE_A.pdb     T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:2]6S36_A.pdb     T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:3]6RZE_A.pdb     T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:4]3HPR_A.pdb     T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:5]1E4V_A.pdb     T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:6]5EJE_A.pdb     T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:7]1E4Y_A.pdb     T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:8]3X2S_A.pdb     T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:9]6HAP_A.pdb     T--KYAKVDGTKPVCEVRADLEKILG-
[Truncated_Name:10]6HAM_A.pdb    T--KYAKVDGTKPVCEVRADLEKILG-
[Truncated_Name:11]4K46_A.pdb    T--QYLKFDGTKAVAEVSAELEKALA-
```

```
[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
```

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

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

```r
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"
```

```r
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
1E4Y_A        1E4Y        A           Protein        214              X-ray
3X2S_A        3X2S        A           Protein        214              X-ray
6HAP_A        6HAP        A           Protein        214              X-ray
6HAM_A        6HAM        A           Protein        214              X-ray
4K46_A        4K46        A           Protein        214              X-ray
3GMT_A        3GMT        A           Protein        230              X-ray
4PZL_A        4PZL        A           Protein        242              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)
1E4Y_A        1.85 Adenylate kinase Adenylate kinase, active site lid (ADK_lid)
3X2S_A        2.80              <NA>                Adenylate kinase (ADK)
6HAP_A        2.70              <NA> Adenylate kinase, active site lid (ADK_lid)
6HAM_A        2.55              <NA>                Adenylate kinase (ADK)
4K46_A        2.01              <NA> Adenylate kinase, active site lid (ADK_lid)
3GMT_A        2.10              <NA>                Adenylate kinase (ADK)
4PZL_A        2.10              <NA>                Adenylate kinase (ADK)
                 ligandId
1AKE_A                AP5
6S36_A CL (3),NA,MG (2)
6RZE_A    NA (3),CL (2)
3HPR_A                AP5
1E4V_A                AP5
5EJE_A             AP5,CO
1E4Y_A                AP5
3X2S_A    AP5,JPY (2),MG
6HAP_A                AP5
6HAM_A                AP5
4K46_A      ADP,AMP,PO4
3GMT_A           SO4 (2)
4PZL_A        CA,FMT,GOL
                                                             ligandName
1AKE_A                              BIS(ADENOSINE)-5'-PENTAPHOSPHATE
6S36_A                 CHLORIDE ION (3),SODIUM ION,MAGNESIUM ION (2)
6RZE_A                            SODIUM ION (3),CHLORIDE ION (2)
3HPR_A                              BIS(ADENOSINE)-5'-PENTAPHOSPHATE
1E4V_A                              BIS(ADENOSINE)-5'-PENTAPHOSPHATE
5EJE_A              BIS(ADENOSINE)-5'-PENTAPHOSPHATE,COBALT (II) ION
```

| | |
|---|---|
| 1E4Y_A | BIS(ADENOSINE)-5'-PENTAPHOSPHATE |
| 3X2S_A | BIS(ADENOSINE)-5'-PENTAPHOSPHATE,N-(pyren-1-ylmethyl)acetamide (2),MAGNESIUM ION |
| 6HAP_A | BIS(ADENOSINE)-5'-PENTAPHOSPHATE |
| 6HAM_A | BIS(ADENOSINE)-5'-PENTAPHOSPHATE |
| 4K46_A | ADENOSINE-5'-DIPHOSPHATE,ADENOSINE MONOPHOSPHATE,PHOSPHATE ION |
| 3GMT_A | SULFATE ION (2) |
| 4PZL_A | CALCIUM ION,FORMIC ACID,GLYCEROL |

| | 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 |
| 1E4Y_A | Escherichia coli |
| 3X2S_A | Escherichia coli str. K-12 substr. MDS42 |
| 6HAP_A | Escherichia coli O139:H28 str. E24377A |
| 6HAM_A | Escherichia coli K-12 |
| 4K46_A | Photobacterium profundum |
| 3GMT_A | Burkholderia pseudomallei 1710b |
| 4PZL_A | Francisella tularensis subsp. tularensis SCHU S4 |

| | |
|---|---|
| 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 |
| 1E4Y_A | |
| 3X2S_A | |
| 6HAP_A | |
| 6HAM_A | |
| 4K46_A | |
| 3GMT_A | |
| 4PZL_A | The cryst |

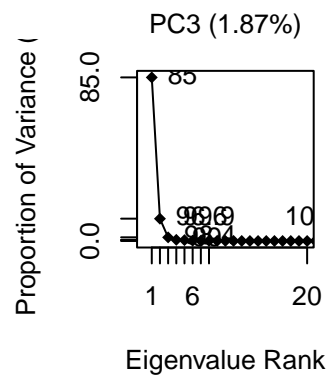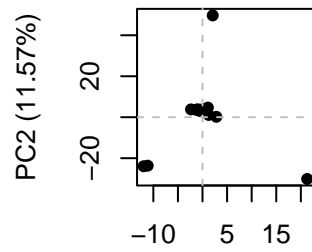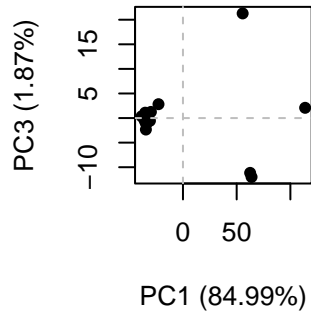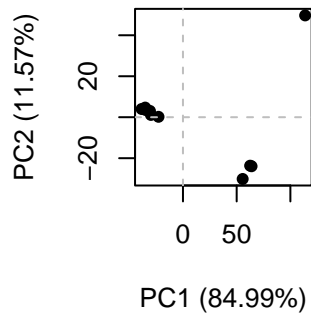| | citation | rObserved | rFree |
|---|---|---|---|
| 1AKE_A | Muller, C.W., et al. J Mol Biol (1992) | 0.19600 | NA |
| 6S36_A | Rogne, P., et al. Biochemistry (2019) | 0.16320 | 0.23560 |
| 6RZE_A | Rogne, P., et al. Biochemistry (2019) | 0.18650 | 0.23500 |
| 3HPR_A | Schrank, T.P., et al. Proc Natl Acad Sci U S A (2009) | 0.21000 | 0.24320 |
| 1E4V_A | Muller, C.W., et al. Proteins (1993) | 0.19600 | NA |
| 5EJE_A | Kovermann, M., et al. Proc Natl Acad Sci U S A (2017) | 0.18890 | 0.23580 |
| 1E4Y_A | Muller, C.W., et al. Proteins (1993) | 0.17800 | NA |

```
3X2S_A                  Fujii, A., et al. Bioconjug Chem (2015)   0.20700 0.25600
6HAP_A                 Kantaev, R., et al. J Phys Chem B (2018)   0.22630 0.27760
6HAM_A                 Kantaev, R., et al. J Phys Chem B (2018)   0.20511 0.24325
4K46_A                    Cho, Y.-J., et al. To be published   0.17000 0.22290
3GMT_A Buchko, G.W., et al. Biochem Biophys Res Commun (2010)   0.23800 0.29500
4PZL_A                        Tan, K., et al. To be published   0.19360 0.23680
         rWork spaceGroup
1AKE_A 0.19600  P 21 2 21
6S36_A 0.15940    C 1 2 1
6RZE_A 0.18190    C 1 2 1
3HPR_A 0.20620  P 21 21 2
1E4V_A 0.19600  P 21 2 21
5EJE_A 0.18630  P 21 2 21
1E4Y_A 0.17800   P 1 21 1
3X2S_A 0.20700 P 21 21 21
6HAP_A 0.22370    I 2 2 2
6HAM_A 0.20311      P 43
4K46_A 0.16730 P 21 21 21
3GMT_A 0.23500   P 1 21 1
4PZL_A 0.19130      P 32
```
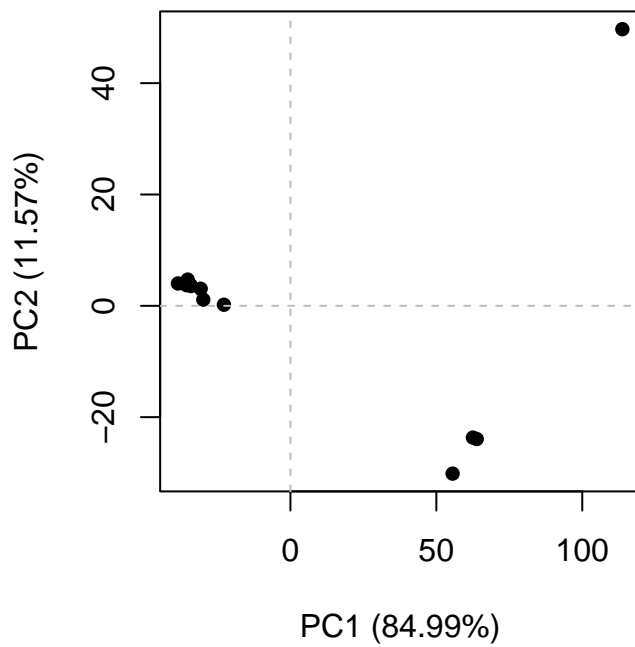
## Principal component analysis

```
pc.xray <- pca(pdbs)
plot(pc.xray)
```

```
# Visualize first principal component
pc1 <- mktrj(pc.xray, pc=1, file="pc_1.pdb")
```
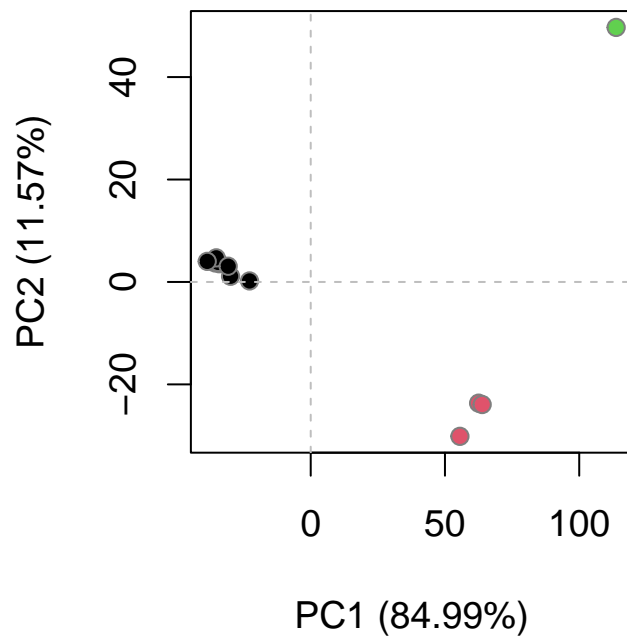
```
plot(pc.xray, pc.axes = c(1,2))
```

25

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

Warning in rmsd(pdbs): No indices provided, using the 204 non NA positions
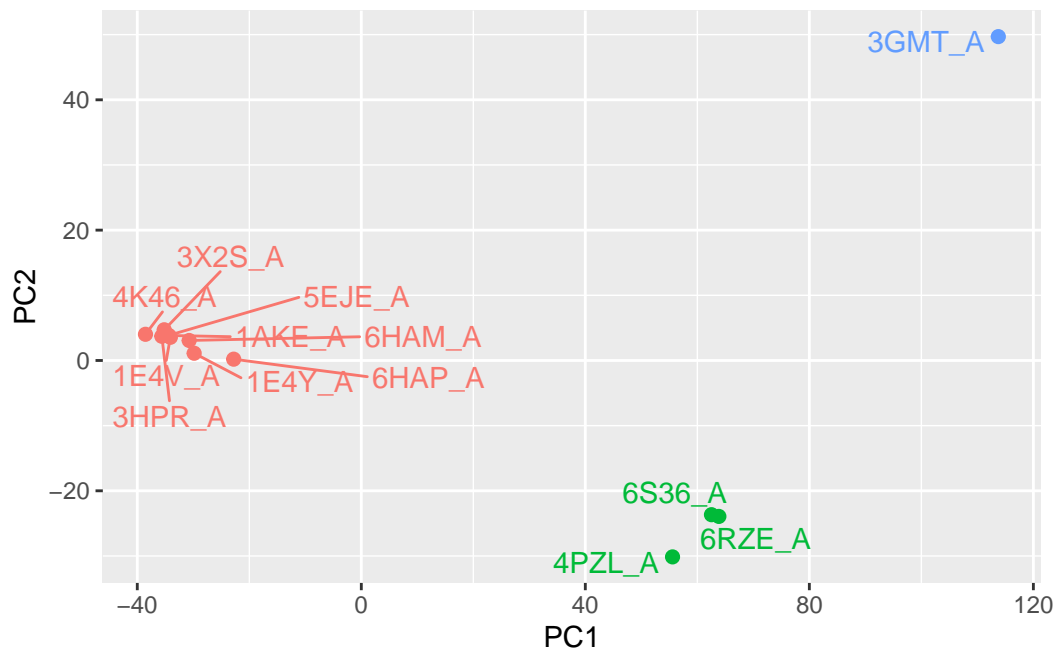
```
# Structure-based clustering
hc.rd <- hclust(dist(rd))
grps.rd <- cutree(hc.rd, k=3)

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

```
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 [optional]

```
modes <- nma(pdbs)
```
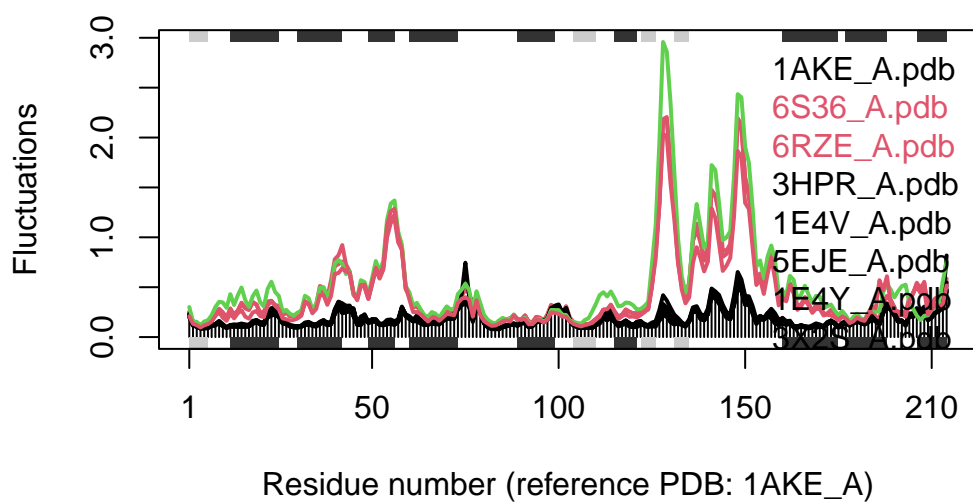
```
Details of Scheduled Calculation:
  ... 13 input structures
  ... storing 606 eigenvectors for each structure
  ... dimension of x$U.subspace: ( 612x606x13 )
  ... coordinate superposition prior to NM calculation
  ... aligned eigenvectors (gap containing positions removed)
  ... estimated memory usage of final 'eNMA' object: 36.9 Mb


  |
  |                                                              |   0%
  |
  |=====                                                         |   8%
  |
  |==========                                                    |  15%
  |
  |===============                                               |  23%
  |
```

```
|====================                                                 |  31%
|
|=========================                                            |  38%
|
|==============================                                       |  46%
|
|=====================================                                |  54%
|
|==========================================                           |  62%
|
|===============================================                      |  69%
|
|====================================================                 |  77%
|
|==========================================================           |  85%
|
|=================================================================    |  92%
|
|=====================================================================| 100%
```

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

Extracting SSE from pdbs$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?

When looking at this plot, I notice there are two main areas of high peaks or fluctuations where the black line is, for the most part, always below the colored lines and do not peak that much. However the two colored lines peak quite dramatically at certain residues of the protein which could possibly point to certain areas of the reference protein where there is a lot of conformation possibilities of similar proteins when the protein folds which dictates its function.

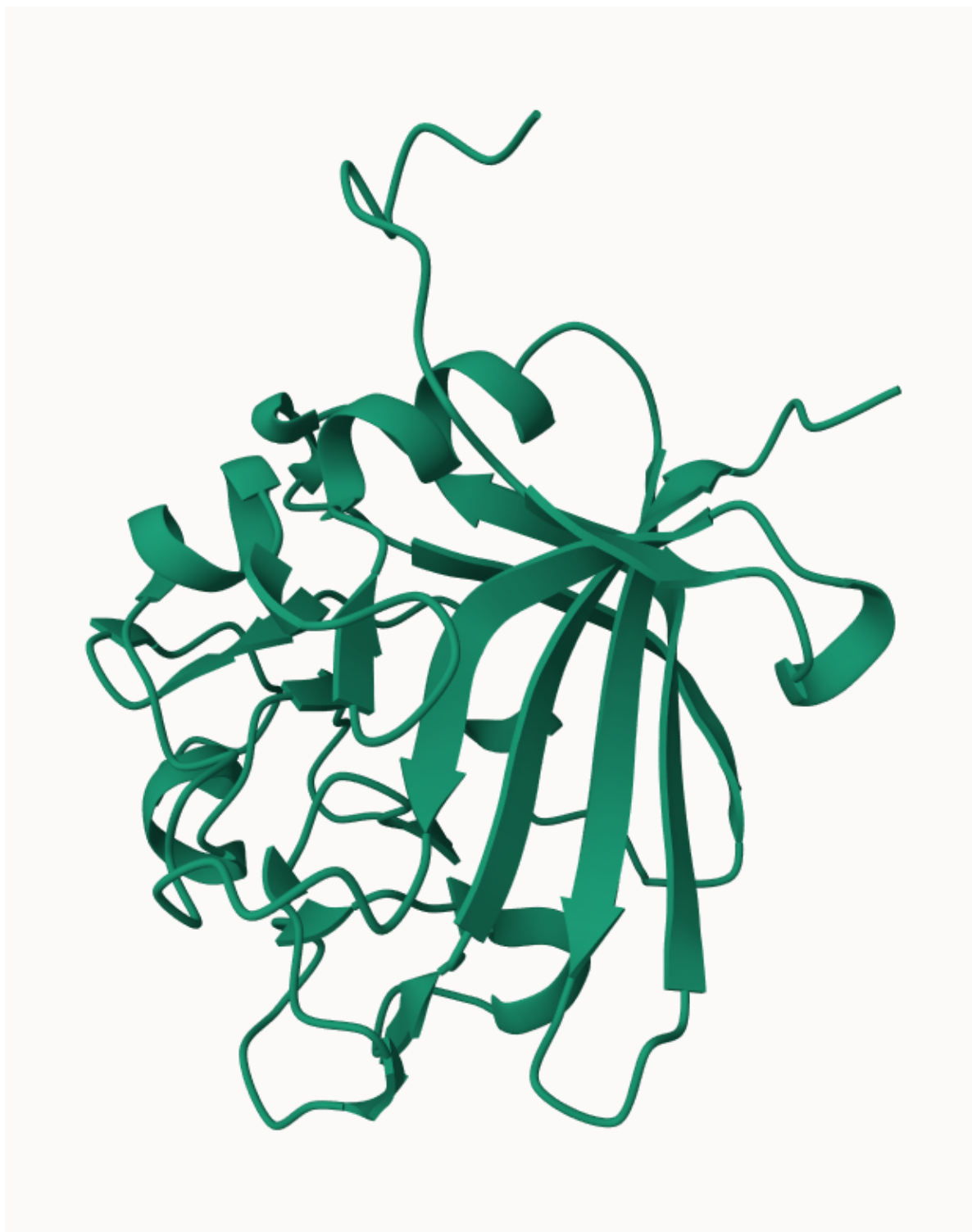##Using ALpha Fold to predict protein structure to use in Molstar

Figure 6: Molecular Surface Pore 1HSG from molstar