# Structural Bioinformatics (Pt. 1)

Vivian Cai

2/17/2022

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

**PDB statistics**

*Download a CSV file from the PDB site (accessible from "Analyze" > "PDB Statistics" > "by Experimental Method and Molecular Type". Move this CSV file into your RStudio project and use it to answer the following questions:*

**Q1: What percentage of structures in the PDB are solved by X-Ray and Electron Microscopy.** As shown below, about **87.2%** of structures in the PDB are solved by X-Ray, and **5.4%** by EM.

```
df <- read.csv("DataExportSummary.csv")
library(dplyr)
```

```
##
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':
##
##     filter, lag
```

```
## The following objects are masked from 'package:base':
##
##     intersect, setdiff, setequal, union
```

```
df %>%
  summarize(
    percent_xRay = sum(df$"X.ray")/sum(df$"Total"),
    percent_EM = sum(df$"EM")/sum(df$"Total")
    )
```

```
##   percent_xRay percent_EM
## 1    0.8716888 0.05388684
```

**Q2: What proportion of structures in the PDB are protein?** As shown from the calculations below, about 87.2% of all PDB structures are protein.

```
percent_prot = df[1,8]/sum(df$"Total")
percent_prot
```

```
## [1] 0.8726292
```

**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 **1225** structures of HIV-1 protease.

# 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?**   The PDB files downloaded don't have hydrogen atoms, so we only see the oxygen atom for every water molecule.

**Q5: There is a conserved water molecule in the binding site. Can you identify this water molecule? What residue number does this water molecule have (see note below)?**   The residue number for the water at the binding site is 308.

**Optional: 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 Licorice for these side-chains). Upload this figure to Piazza for some extra credit.**   (Posted on piazza)

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

**Sequence Viewer Extension [OPTIONAL]**

**Q6: 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 and the sequence viewer extension can you identify secondary structure elements that are likely to only form in the dimer rather than the monomer?**   The N-terminus of each chain intertwine with the other and forms a beta sheet structure. The same structure is unlikely to form if only one chain is present.

# 3. 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)
##
##        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**


```
# the attributes (+ attr:) of this object are listed
# on the last couple of lines.
attributes(pdb)
```

**Q9: How many protein chains are in this structure? 2**

```
## $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
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>
```

# 4. Comparative structure analysis of Adenylate Kinase

**Overview**

*Starting from only one Adk PDB identifier (PDB ID: 1AKE) we will search the entire PDB for related structures using BLAST, fetch, align and superpose the identified structures, perform PCA and finally calculate the normal modes of each individual structure in order to probe for potential differences in structural flexibility.*

**Setup**

```
# Install packages in the R console

# install.packages("bio3d")
# install.packages("ggplot2")
# install.packages("ggrepel")
# 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? TURE**

**Search and retrieve ADK structures**

```
library(bio3d)
aa <- get.seq("1ake_A") # chain A only
```

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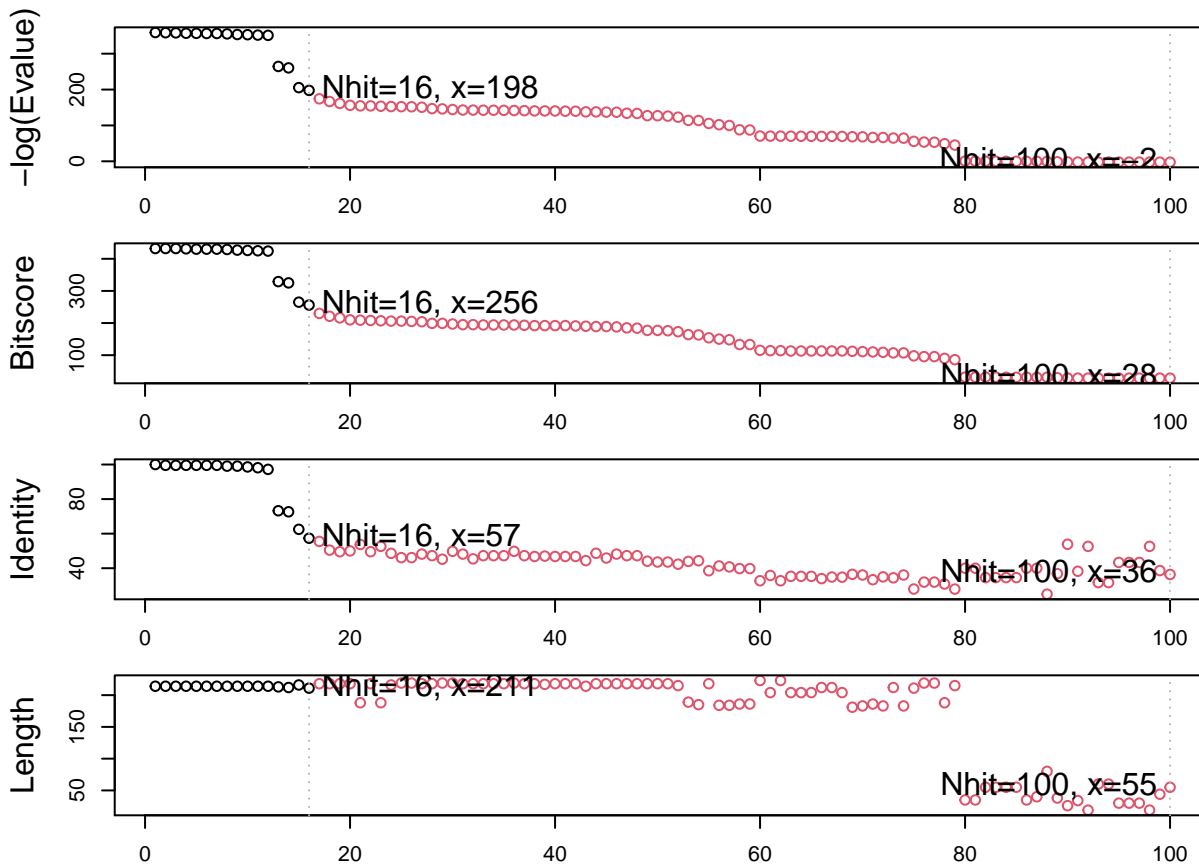
```
# Blast or hmmer search
b <- blast.pdb(aa)
```

**Q13. How many amino acids are in this sequence, i.e. how long is this sequence? 214**

```
##  Searching ... please wait (updates every 5 seconds) RID = 13HC0SFA013
##  .
##  Reporting 100 hits
```

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

```
##    * Possible cutoff values:    197 -3
##              Yielding Nhits:    16 100
##
##    * Chosen cutoff value of:    197
##              Yielding Nhits:    16
```

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

```
## [1] "1AKE_A" "4X8M_A" "6S36_A" "6RZE_A" "4X8H_A" "3HPR_A"
```

```
# in case blast doesn't run
# hits <- NULL
# hits$pdb.id <- c('1AKE_A','6S36_A','6RZE_A','3HPR_A','1E4V_A','5EJE_A','1E4Y_A','3X2S_A','6HAP_A','6H
```

```
# 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.gz exists. Skipping download
```

```
## Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE): pdbs/
## 4X8M.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/
## 4X8H.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/
## 4NP6.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 releated PDBs
pdbs <- pdbaln(files, fit = TRUE)#, exefile="msa")
```

```
## Reading PDB files:
## pdbs/split_chain/1AKE_A.pdb
## pdbs/split_chain/4X8M_A.pdb
## pdbs/split_chain/6S36_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/4K46_A.pdb
## pdbs/split_chain/4NP6_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/4X8M_A.pdb
## pdb/seq: 3   name: pdbs/split_chain/6S36_A.pdb
##     PDB has ALT records, taking A only, rm.alt=TRUE
## pdb/seq: 4   name: pdbs/split_chain/6RZE_A.pdb
##     PDB has ALT records, taking A only, rm.alt=TRUE
## pdb/seq: 5   name: pdbs/split_chain/4X8H_A.pdb
## pdb/seq: 6   name: pdbs/split_chain/3HPR_A.pdb
##     PDB has ALT records, taking A only, rm.alt=TRUE
## pdb/seq: 7   name: pdbs/split_chain/1E4V_A.pdb
## pdb/seq: 8   name: pdbs/split_chain/5EJE_A.pdb
##     PDB has ALT records, taking A only, rm.alt=TRUE
## pdb/seq: 9   name: pdbs/split_chain/1E4Y_A.pdb
## pdb/seq: 10   name: pdbs/split_chain/3X2S_A.pdb
## pdb/seq: 11   name: pdbs/split_chain/6HAP_A.pdb
## pdb/seq: 12   name: pdbs/split_chain/6HAM_A.pdb
##     PDB has ALT records, taking A only, rm.alt=TRUE
## pdb/seq: 13   name: pdbs/split_chain/4K46_A.pdb
##     PDB has ALT records, taking A only, rm.alt=TRUE
## pdb/seq: 14   name: pdbs/split_chain/4NP6_A.pdb
## pdb/seq: 15   name: pdbs/split_chain/3GMT_A.pdb
## pdb/seq: 16   name: pdbs/split_chain/4PZL_A.pdb
```
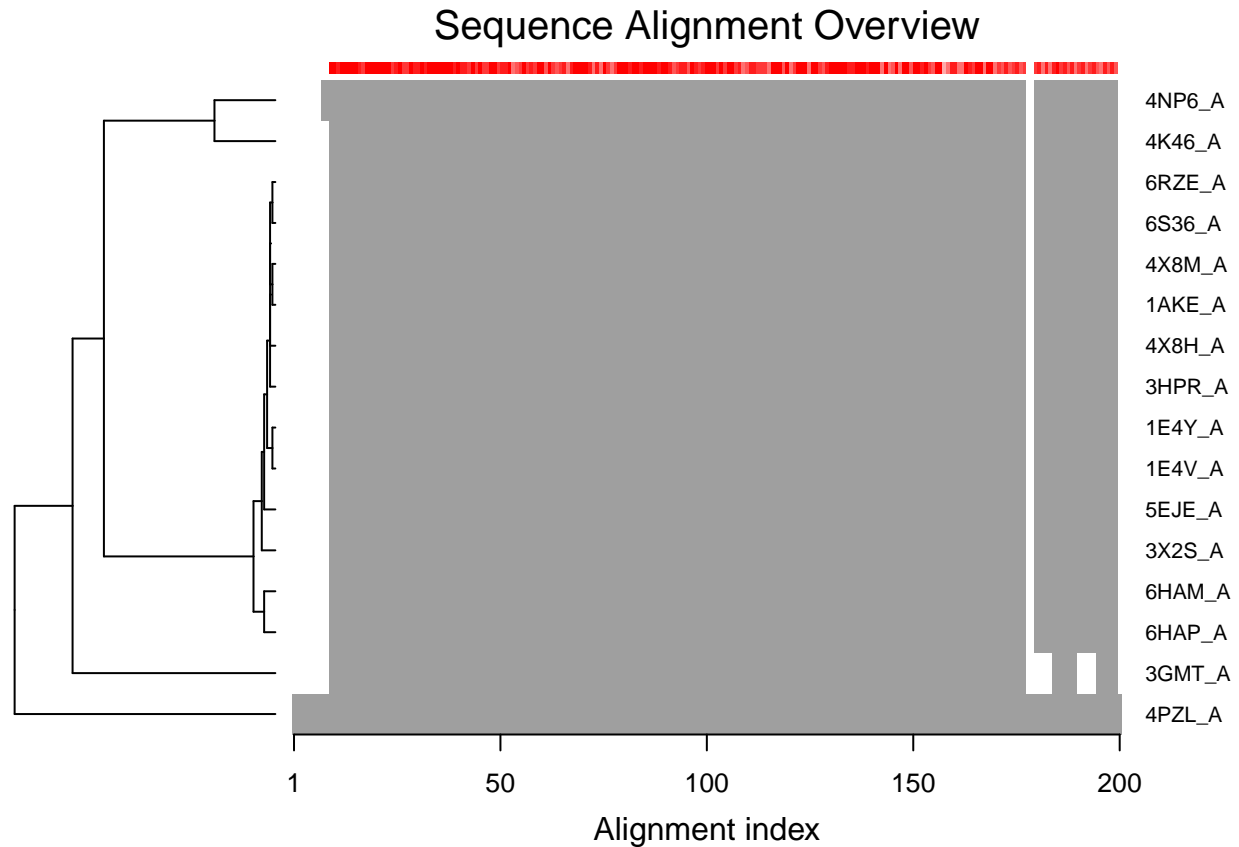
```r
# Vector containing PDB codes for figure axis
ids <- basename.pdb(pdbs$id)

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

## Sequence Alignment Overview



**Viewing our superposed structures**

```
library(bio3d.view)
library(rgl)

view.pdbs(pdbs)
```
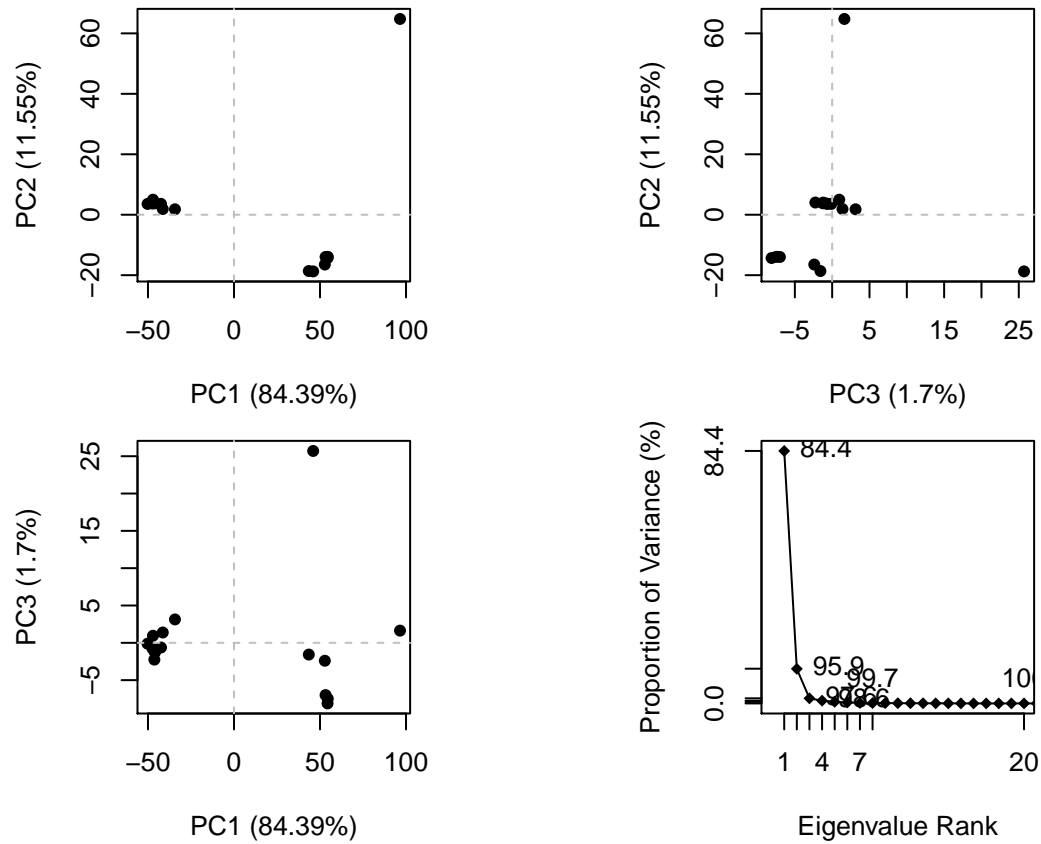
**Annotate collected PDB structures**

```
# I couldn't execute this, I kept getting:
# Error in split.default(X, group) : first argument must be a vector

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

**Principal component analysis**

```
# Perform PCA
pc.xray <- pca(pdbs)
plot(pc.xray)
```
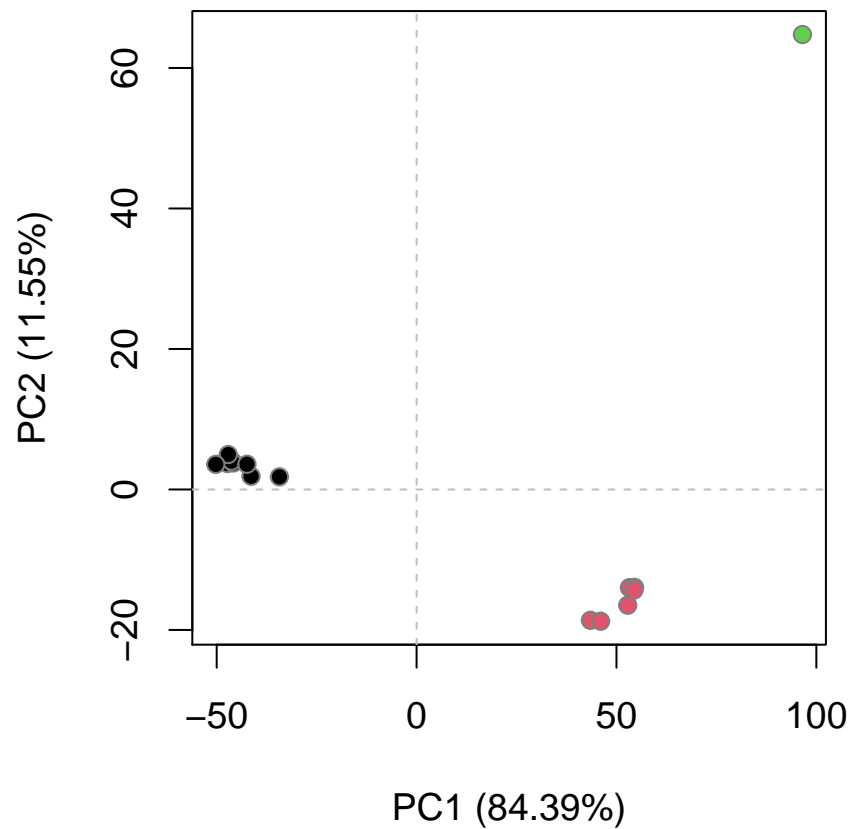
```
# Calculate RMSD
rd <- rmsd(pdbs)
```

```
## 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)
```

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

# 5. Optional further visualization

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

```
view.xyz(pc1)
```

```
## Potential all C-alpha atom structure(s) detected: Using calpha.connectivity()
```

(Figure 12: Visualization of PC-1 trajectory generated using mktrj().)

```
view.xyz(pc1, col=vec2color( rmsf(pc1) ))
```

```
## Potential all C-alpha atom structure(s) detected: Using calpha.connectivity()
```

```
#Plotting results with ggplot2
library(ggplot2)
library(ggrepel)
```
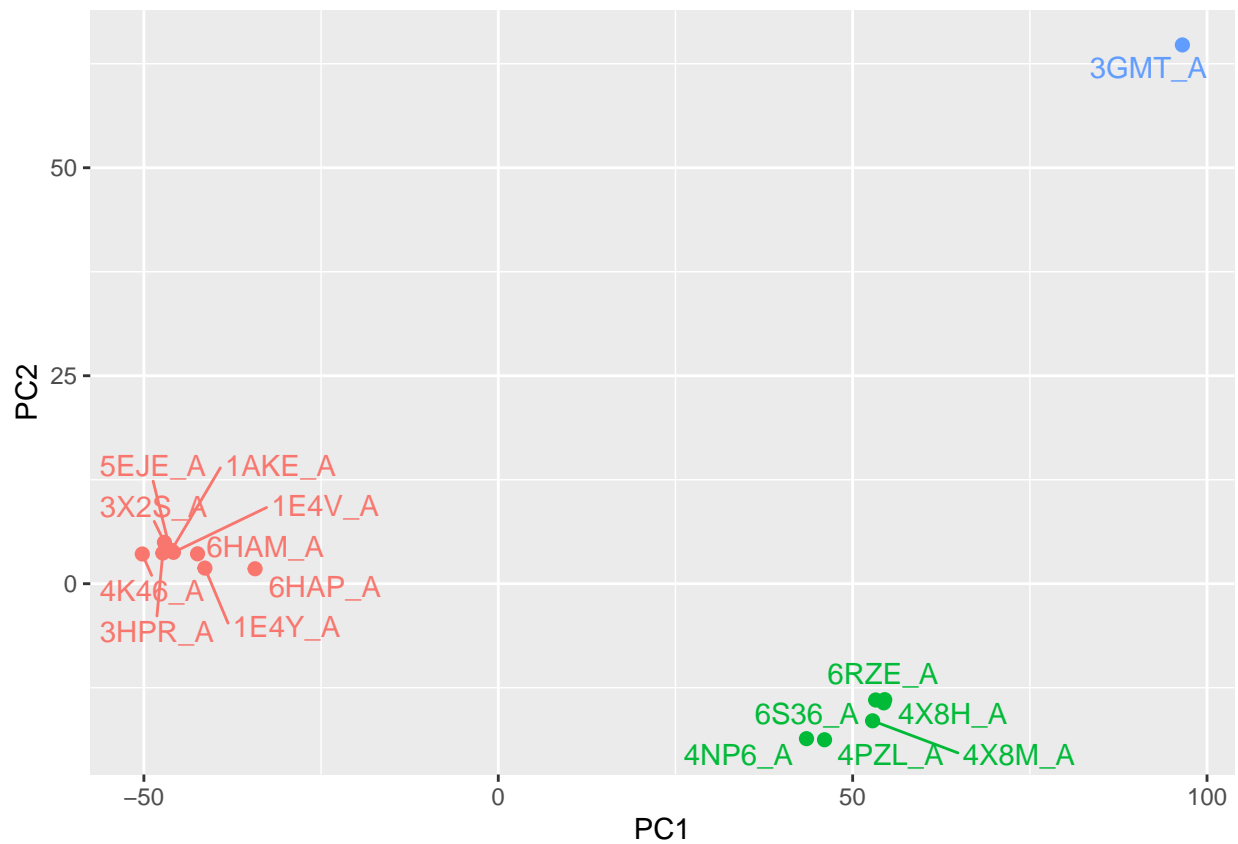
```
df <- data.frame(PC1=pc.xray$z[,1],
```

```
                    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

```
# NMA of all structures
modes <- nma(pdbs)
```
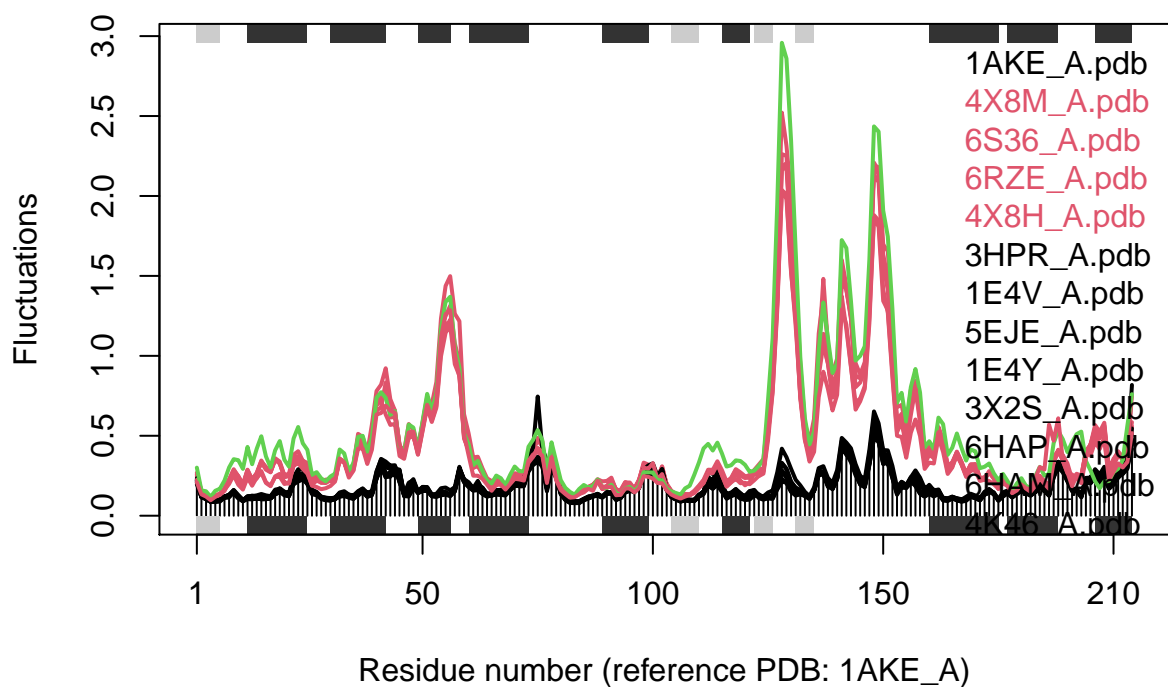
```
##
## Details of Scheduled Calculation:
##    ... 16 input structures
##    ... storing 606 eigenvectors for each structure
##    ... dimension of x$U.subspace: ( 612x606x16 )
##    ... coordinate superposition prior to NM calculation
```

```
##     ... aligned eigenvectors (gap containing positions removed)
##     ... estimated memory usage of final 'eNMA' object: 45.4 Mb
##
##     |                                                                |
```

```
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?** The most flexible regions are the loopy regions (with no secondary structure prediction). The difference between the black and colored lines shows two different conformations of the protein. After nucleotide binding, the loopy regions becomes a lot more dynamic.