

# Lab 9

Gretel Warmuth (PID: A17595945)

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

Analyzing the Protein Data Bank (PDB):

from: <https://www.rcsb.org/stats/summary#>

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

|   | 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.

```
pdbdb$Total
```

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

Removing the commas to convert to numeric values:

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

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

Turning into a function to be able to use later:

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

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

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

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

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

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

Warning in FUN(newX[, i], ...): NAs introduced by coercion

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

Or use an import function:

```
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.

```
sum(pdbdb$Total)
```

```
[1] 226414
```

Percent structures by X-ray and electron microscopy:

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

```
[1] 83.30359
```

```
sum(pdbdb$EM)/sum(pdbdb$Total) * 100
```

```
[1] 10.18091
```

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

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

```
[1] 86.39483
```

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?

## Visualizing the HIV-1 protease structure

Mol(molstar) is a web-based molecular viewer that we will need to learn the basics of:

<https://molstar.org/viewer/>

Using PDB code 1HSG

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

The one “atom” is representing each water molecule interacting with the protein.

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

Asp25

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 1: 1HSG Protein



Figure 2: Aspartate Components

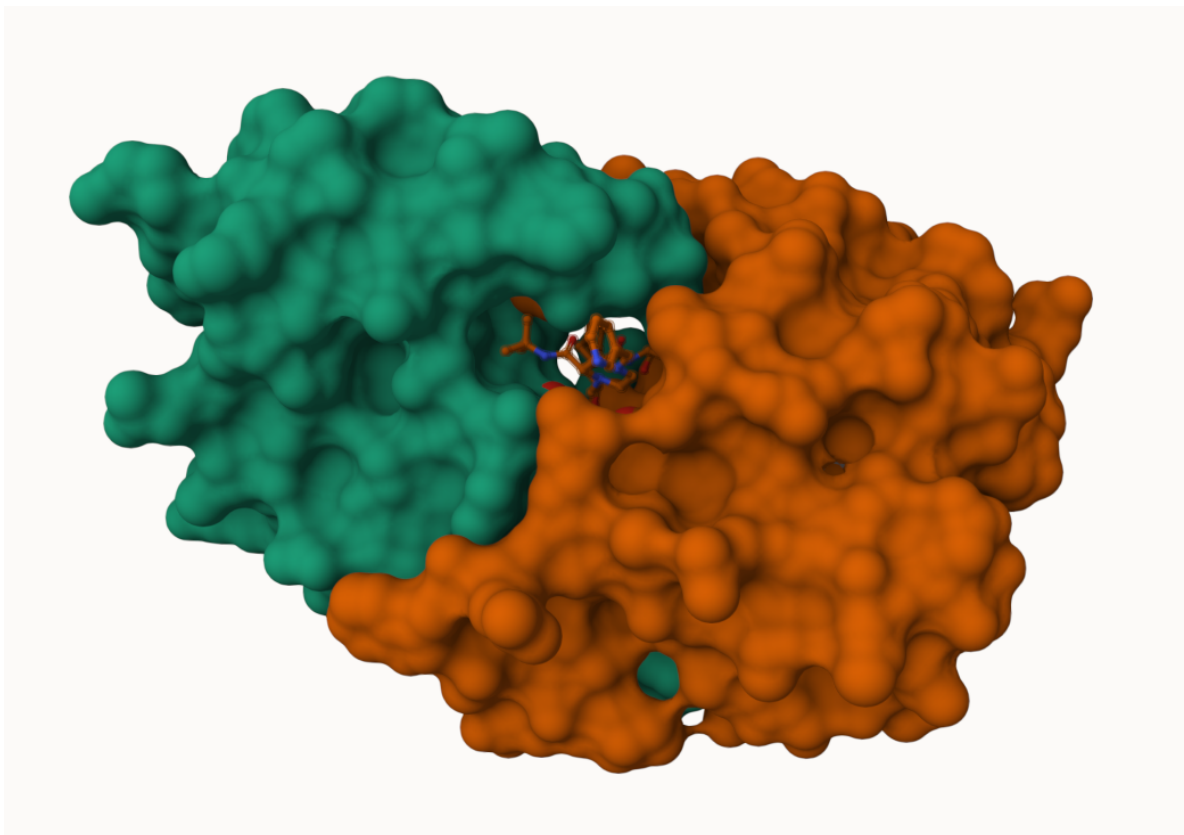


Figure 3: Surface of Protein

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

Ligands and substrates can enter the binding site when the protein is in the correct conformation and can allow binding.

## Introduction to Bio3D in R

bio3D allows for structural and bioinformatics work.

reading PDB files in bio3D:

```
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
```

## 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 | elemsy | 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
```

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

HOH and KM1

Q9: How many protein chains are in this structure?

```
2
```

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

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

Predicting functional motions of a single structure:

```
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  
YYSKEAEAGNTKYAKVDGTPVAEVRADLEKILG
```

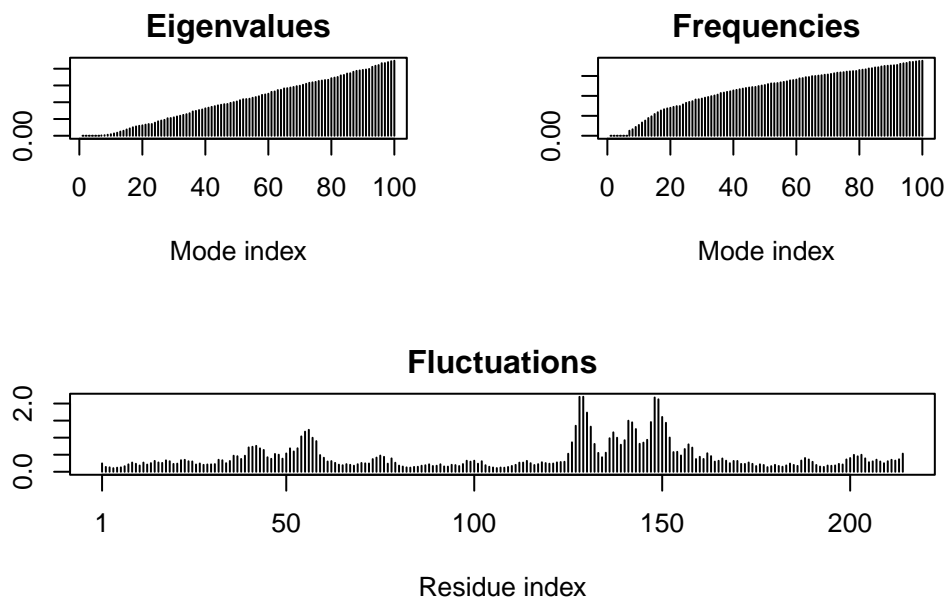
```
+ attr: atom, xyz, seqres, helix, sheet,  
      calpha, remark, call
```

```
# Perform flexibility prediction  
m <- nma(adk)
```

```
Building Hessian... Done in 0.06 seconds.
```

```
Diagonalizing Hessian... Done in 0.943 seconds.
```

```
plot(m)
```



Write out multi-model PDB file that can be used to make an animation of the predicted motions:

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

This file can be opened in Mol\*

## Comparative Analysis of Protein Structure

```
library(bio3d)
```

Analyzing all ADK structures in the PDB database, starting with "1ake\_A"

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

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

Fetching... Please wait. Done.

aa

```

      1      .      .      .      .      .      60
pdb|1AKE|A  MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMRLRAAVKSGSELGKQAKDIMDAGKLV
      1      .      .      .      .      .      60

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

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

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

The `msa` package

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

N/A

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

N/A

Packages installed:

```
install.packages("BiocManager") BiocManager::install("msa")
```

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

```
length(aa)
```

```
[1] 3
```

```
attributes(aa)
```

```
$names
```

```
[1] "id"   "ali"  "call"
```

```
$class
```

```
[1] "fasta"
```

```
ncol(aa$ali)
```

```
[1] 214
```

Finding related sequences:

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

```
Searching ... please wait (updates every 5 seconds) RID = JMYWKS42013
```

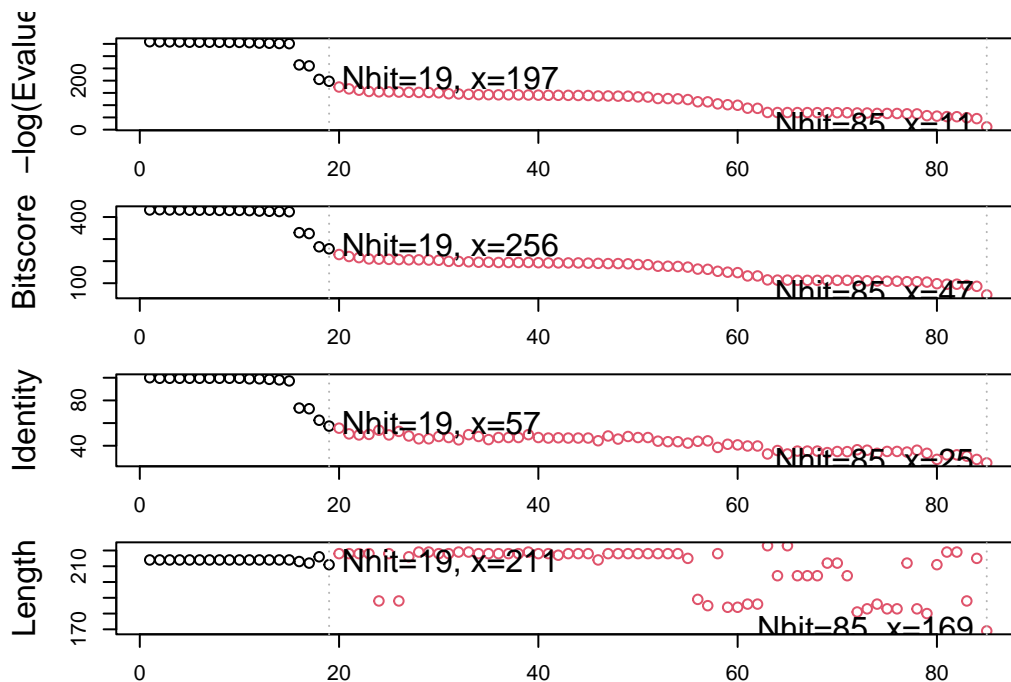
```
.
```

```
Reporting 85 hits
```

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

```
* Possible cutoff values:    197 11
      Yielding Nhits:        19 85
```

```
* Chosen cutoff value of:    197
      Yielding Nhits:        19
```



```
hits$pdb.id
```

```
[1] "1AKE_A" "8BQF_A" "4X8M_A" "6S36_A" "8Q2B_A" "8RJ9_A" "6RZE_A" "4X8H_A"
[9] "3HPR_A" "1E4V_A" "5EJE_A" "1E4Y_A" "3X2S_A" "6HAP_A" "6HAM_A" "4K46_A"
[17] "4NP6_A" "3GMT_A" "4PZL_A"
```

Downloading files:

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

```
Warning in get.pdb(hits$pdb.id, path = "pdb", split = TRUE, gzip = TRUE):
pdb/8BQF.pdb.gz exists. Skipping download
```

```
Warning in get.pdb(hits$pdb.id, path = "pdb", split = TRUE, gzip = TRUE):
pdb/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/8Q2B.pdb.gz exists. Skipping download

Warning in get.pdb(hits\$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):  
pdbs/8RJ9.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

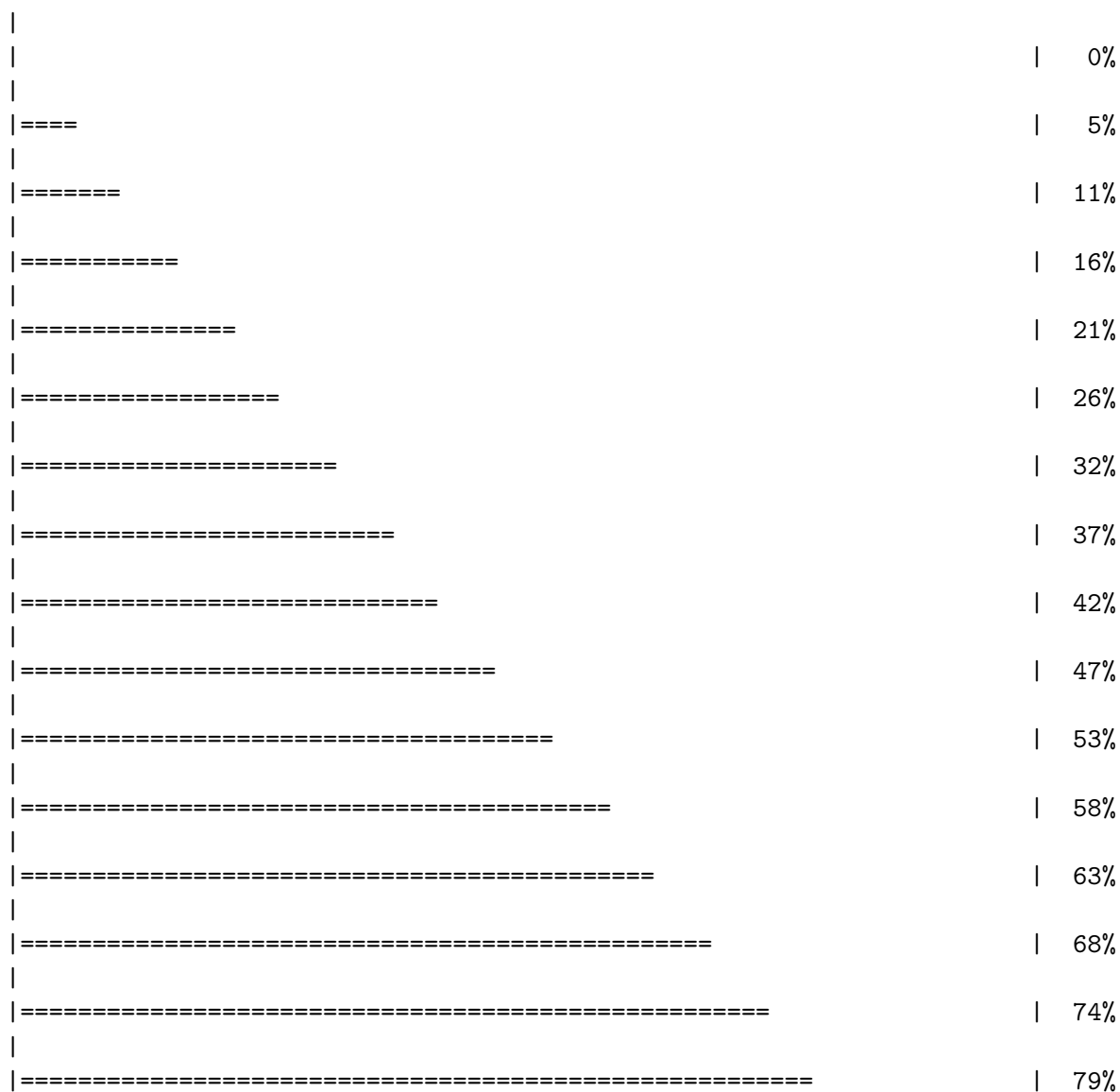
Warning in get.pdb(hits\$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):  
pdbs/6HAM.pdb.gz exists. Skipping download

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





```

|
|=====| 84%
|
|=====| 89%
|
|=====| 95%
|
|=====| 100%

```

Using the `pdbaln` function to align the found sequences:

```
pdbbs <- pdbaln(files, fit = TRUE, exefile="msa")
```

Reading PDB files:

```

pdbbs/split_chain/1AKE_A.pdb
pdbbs/split_chain/8BQF_A.pdb
pdbbs/split_chain/4X8M_A.pdb
pdbbs/split_chain/6S36_A.pdb
pdbbs/split_chain/8Q2B_A.pdb
pdbbs/split_chain/8RJ9_A.pdb
pdbbs/split_chain/6RZE_A.pdb
pdbbs/split_chain/4X8H_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/4NP6_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
..   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: pdbc/split_chain/1AKE_A.pdb
            PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 2   name: pdbc/split_chain/8BQF_A.pdb
            PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 3   name: pdbc/split_chain/4X8M_A.pdb
pdb/seq: 4   name: pdbc/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/4K46_A.pdb
            PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 17  name: pdbc/split_chain/4NP6_A.pdb
pdb/seq: 18  name: pdbc/split_chain/3GMT_A.pdb
pdb/seq: 19  name: pdbc/split_chain/4PZL_A.pdb
```

## pdbc

|                               |       |                                |   |   |    |
|-------------------------------|-------|--------------------------------|---|---|----|
|                               | 1     | .                              | . | . | 40 |
| [Truncated_Name:1] 1AKE_A.pdb | ----- | MRIILLGAPGAGKGTQAQFIMEKYGIPQIS |   |   |    |
| [Truncated_Name:2] 8BQF_A.pdb | ----- | MRIILLGAPGAGKGTQAQFIMEKYGIPQIS |   |   |    |
| [Truncated_Name:3] 4X8M_A.pdb | ----- | MRIILLGAPGAGKGTQAQFIMEKYGIPQIS |   |   |    |
| [Truncated_Name:4] 6S36_A.pdb | ----- | MRIILLGAPGAGKGTQAQFIMEKYGIPQIS |   |   |    |

|                                |  |
|--------------------------------|--|
| [Truncated_Name:5] 8Q2B_A.pdb  | -----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS      |
| [Truncated_Name:6] 8RJ9_A.pdb  | -----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS      |
| [Truncated_Name:7] 6RZE_A.pdb  | -----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS      |
| [Truncated_Name:8] 4X8H_A.pdb  | -----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS      |
| [Truncated_Name:9] 3HPR_A.pdb  | -----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS      |
| [Truncated_Name:10] 1E4V_A.pdb | -----MRIILLGAPVAGKGTQAQFIMEKYGIPQIS      |
| [Truncated_Name:11] 5EJE_A.pdb | -----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS      |
| [Truncated_Name:12] 1E4Y_A.pdb | -----MRIILLGALVAGKGTQAQFIMEKYGIPQIS      |
| [Truncated_Name:13] 3X2S_A.pdb | -----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS      |
| [Truncated_Name:14] 6HAP_A.pdb | -----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS      |
| [Truncated_Name:15] 6HAM_A.pdb | -----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS      |
| [Truncated_Name:16] 4K46_A.pdb | -----MRIILLGAPGAGKGTQAQFIMAKFGIPQIS      |
| [Truncated_Name:17] 4NP6_A.pdb | -----NAMRIILLGAPGAGKGTQAQFIMEKFGIPQIS    |
| [Truncated_Name:18] 3GMT_A.pdb | -----MRLILLGAPGAGKGTQANFIKEKFGIPQIS      |
| [Truncated_Name:19] 4PZL_A.pdb | TENLYFQSNAMRIILLGAPGAGKGTQAKIIEQYINIAHIS |

                  \*\*~\*\*\*\*\*   \*\*\*\*\*   \*   \*~\*   \*\*

1                   .                   .                   .                   40

|                                |   |   |
|--------------------------------|---|---|
| [Truncated_Name:1] 1AKE_A.pdb  | 41                   .                   .                   .                   80 | TGDMRLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE  |
| [Truncated_Name:2] 8BQF_A.pdb  |   | TGDMRLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE  |
| [Truncated_Name:3] 4X8M_A.pdb  |   | TGDMRLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE  |
| [Truncated_Name:4] 6S36_A.pdb  |   | TGDMRLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE  |
| [Truncated_Name:5] 8Q2B_A.pdb  |   | TGDMRLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE  |
| [Truncated_Name:6] 8RJ9_A.pdb  |   | TGDMRLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE  |
| [Truncated_Name:7] 6RZE_A.pdb  |   | TGDMRLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE  |
| [Truncated_Name:8] 4X8H_A.pdb  |   | TGDMRLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE  |
| [Truncated_Name:9] 3HPR_A.pdb  |   | TGDMRLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE  |
| [Truncated_Name:10] 1E4V_A.pdb |   | TGDMRLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE  |
| [Truncated_Name:11] 5EJE_A.pdb |   | TGDMRLRAAVKSGSELGKQAKDIMDACKLVDELVIALVKE  |
| [Truncated_Name:12] 1E4Y_A.pdb |   | TGDMRLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE  |
| [Truncated_Name:13] 3X2S_A.pdb |   | TGDMRLRAAVKSGSELGKQAKDIMDCGKLVDELVIALVKE  |
| [Truncated_Name:14] 6HAP_A.pdb |   | TGDMRLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVRE  |
| [Truncated_Name:15] 6HAM_A.pdb |   | TGDMRLRAAIKSGSELGKQAKDIMDAGKLVDEIIIALVKE  |
| [Truncated_Name:16] 4K46_A.pdb |   | TGDMRLRAAIKAGTELGKQAKSVIDAGQLVSDDIILGLVKE |
| [Truncated_Name:17] 4NP6_A.pdb |   | TGDMRLRAAIKAGTELGKQAKAVIDAGQLVSDDIILGLIKE |
| [Truncated_Name:18] 3GMT_A.pdb |   | TGDMRLRAAVKAGTPLGVEAKTYMDEGKLPVDSLIIGLVKE |
| [Truncated_Name:19] 4PZL_A.pdb |   | TGDMIRETIKSGSALGQELKKVLDAGELVSDEFIIVKIVKD |

\*\*\*\*~\*   ~\* \*~\*\*   \*   ~\*   \*\* \*   ^^ ~~~~

41                   .                   .                   .                   80

|                               |  |   |
|-------------------------------|--|---|
| [Truncated_Name:1] 1AKE_A.pdb | 81                   .                   .                   .                   120 | RIAQEDCRNGFLLDGFPRTPQADAMKEAGINVDYVLEFD |
|-------------------------------|--|---|

|                               |   |
|-------------------------------|---|
| [Truncated_Name:2]8BQF_A.pdb  | RIAQE----GFLLDGFPRTPQADAMKEAGINVDYVIEFD       |
| [Truncated_Name:3]4X8M_A.pdb  | RIAQEDCRNGFLLDGFPRTPQADAMKEAGINVDYVLEFD       |
| [Truncated_Name:4]6S36_A.pdb  | RIAQEDCRNGFLLDGFPRTPQADAMKEAGINVDYVLEFD       |
| [Truncated_Name:5]8Q2B_A.pdb  | RIAQEDCRNGFLLDGFPRTPQADAMKEAGINVDYVLEFD       |
| [Truncated_Name:6]8RJ9_A.pdb  | RIAQEDCRNGFLLAGFPRTIPQADAMKEAGINVDYVLEFD      |
| [Truncated_Name:7]6RZE_A.pdb  | RIAQEDCRNGFLLDGFPRTPQADAMKEAGINVDYVLEFD       |
| [Truncated_Name:8]4X8H_A.pdb  | RIAQEDCRNGFLLDGFPRTPQADAMKEAGINVDYVLEFD       |
| [Truncated_Name:9]3HPR_A.pdb  | RIAQEDCRNGFLLDGFPRTPQADAMKEAGINVDYVLEFD       |
| [Truncated_Name:10]1E4V_A.pdb | RIAQEDCRNGFLLDGFPRTPQADAMKEAGINVDYVLEFD       |
| [Truncated_Name:11]5EJE_A.pdb | RIAQEDCRNGFLLDGFPRTPQADAMKEAGINVDYVLEFD       |
| [Truncated_Name:12]1E4Y_A.pdb | RIAQEDCRNGFLLDGFPRTPQADAMKEAGINVDYVLEFD       |
| [Truncated_Name:13]3X2S_A.pdb | RIAQEDSRNGFLLDGFPRTPQADAMKEAGINVDYVLEFD       |
| [Truncated_Name:14]6HAP_A.pdb | RICQEDSRNGFLLDGFPRTPQADAMKEAGINVDYVLEFD       |
| [Truncated_Name:15]6HAM_A.pdb | RICQEDSRNGFLLDGFPRTPQADAMKEAGINVDYVLEFD       |
| [Truncated_Name:16]4K46_A.pdb | RIAQDDCAKGFLLDGFPRTPQADGLKEVGVVVDYVIEFD       |
| [Truncated_Name:17]4NP6_A.pdb | RIAQADCEKGFLLDGFPRTPQADGLKEMGINVDYVIEFD       |
| [Truncated_Name:18]3GMT_A.pdb | RLKEADCANGYLFDGFPRTIAQADAMKEAGVAIDYVLEID      |
| [Truncated_Name:19]4PZL_A.pdb | RISKNCNNGFLLDGVPRTIPQAQELDKLGVNIDYIVEVD       |
|                               | *~      *~*   * ***** **    ^     *~ ^**^^* * |
| 81                            | . . .   |
|                               |   |
| 121                           | . . .   |
|                               |   |
| [Truncated_Name:1]1AKE_A.pdb  | VPDELIVDRIVGRRVHAPSGRVYHVKFNNPPKVEGKD DVTG    |
| [Truncated_Name:2]8BQF_A.pdb  | VPDELIVDRIVGRRVHAPSGRVYHVKFNNPPKVEGKD DVTG    |
| [Truncated_Name:3]4X8M_A.pdb  | VPDELIVDRIVGRRVHAPSGRVYHVKFNNPPKVEGKD DVTG    |
| [Truncated_Name:4]6S36_A.pdb  | VPDELIVDKIVGRRVHAPSGRVYHVKFNNPPKVEGKD DVTG    |
| [Truncated_Name:5]8Q2B_A.pdb  | VPDELIVDRIVGRRVHAPSGRVYHVKFNNPPKVEGKD DVTG    |
| [Truncated_Name:6]8RJ9_A.pdb  | VPDELIVDRIVGRRVHAPSGRVYHVKFNNPPKVEGKD DVTG    |
| [Truncated_Name:7]6RZE_A.pdb  | VPDELIVDAIVGRRVHAPSGRVYHVKFNNPPKVEGKD DVTG    |
| [Truncated_Name:8]4X8H_A.pdb  | VPDELIVDRIVGRRVHAPSGRVYHVKFNNPPKVEGKD DVTG    |
| [Truncated_Name:9]3HPR_A.pdb  | VPDELIVDRIVGRRVHAPSGRVYHVKFNNPPKVEGKD DGTG    |
| [Truncated_Name:10]1E4V_A.pdb | VPDELIVDRIVGRRVHAPSGRVYHVKFNNPPKVEGKD DVTG    |
| [Truncated_Name:11]5EJE_A.pdb | VPDELIVDRIVGRRVHAPSGRVYHVKFNNPPKVEGKD DVTG    |
| [Truncated_Name:12]1E4Y_A.pdb | VPDELIVDRIVGRRVHAPSGRVYHVKFNNPPKVEGKD DVTG    |
| [Truncated_Name:13]3X2S_A.pdb | VPDELIVDRIVGRRVHAPSGRVYHVKFNNPPKVEGKD DVTG    |
| [Truncated_Name:14]6HAP_A.pdb | VPDELIVDRIVGRRVHAPSGRVYHVKFNNPPKVEGKD DVTG    |
| [Truncated_Name:15]6HAM_A.pdb | VADS VIVERMAGRRAHLASGRTYHN VYNPPKVEGKD DVTG   |
| [Truncated_Name:16]4K46_A.pdb | VADD VIVERMAGRRAHLPSGRTYH VYNPPKVEGKD DVTG    |
| [Truncated_Name:17]4NP6_A.pdb | VPFSEI IERMSGRRTHPASGRTYH VKFNNPPKVEGKD DVTG  |
| [Truncated_Name:18]3GMT_A.pdb | VADNLLIERITGRRIH PASGRTYHT KFNPPKVADKD DVTG   |
| [Truncated_Name:19]4PZL_A.pdb |   |
|                               | * ~~~ ^ *** * *** ** ^***** *** **            |
| 121                           | . . .   |
|                               |   |
|                               |   |

|                                |  |   |   |   |     |
|--------------------------------|--|---|---|---|-----|
|                                | 161                                      | . | . | . | 200 |
| [Truncated_Name:1] 1AKE_A.pdb  | EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN |   |   |   |     |
| [Truncated_Name:2] 8BQF_A.pdb  | EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN |   |   |   |     |
| [Truncated_Name:3] 4X8M_A.pdb  | EELTTRKDDQEETVRKRLVEWHQMTAPLIGYYSKEAEAGN |   |   |   |     |
| [Truncated_Name:4] 6S36_A.pdb  | EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN |   |   |   |     |
| [Truncated_Name:5] 8Q2B_A.pdb  | EELTTRKADQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN |   |   |   |     |
| [Truncated_Name:6] 8RJ9_A.pdb  | EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN |   |   |   |     |
| [Truncated_Name:7] 6RZE_A.pdb  | EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN |   |   |   |     |
| [Truncated_Name:8] 4X8H_A.pdb  | EELTTRKDDQEETVRKRLVEYHQMTAALIGYYSKEAEAGN |   |   |   |     |
| [Truncated_Name:9] 3HPR_A.pdb  | EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN |   |   |   |     |
| [Truncated_Name:10] 1E4V_A.pdb | EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN |   |   |   |     |
| [Truncated_Name:11] 5EJE_A.pdb | EELTTRKDDQEECVRKRLVEYHQMTAPLIGYYSKEAEAGN |   |   |   |     |
| [Truncated_Name:12] 1E4Y_A.pdb | EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN |   |   |   |     |
| [Truncated_Name:13] 3X2S_A.pdb | EELTTRKDDQEETVRKRLCEYHQMTAPLIGYYSKEAEAGN |   |   |   |     |
| [Truncated_Name:14] 6HAP_A.pdb | EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN |   |   |   |     |
| [Truncated_Name:15] 6HAM_A.pdb | EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN |   |   |   |     |
| [Truncated_Name:16] 4K46_A.pdb | EDLVIREDDKEETVLARLGVYHNQTAPLIAYYGKEAEAGN |   |   |   |     |
| [Truncated_Name:17] 4NP6_A.pdb | EDLVIREDDKEETVRARLNVYHTQTAPLIEYYGKEAAAGK |   |   |   |     |
| [Truncated_Name:18] 3GMT_A.pdb | EPLVQRDDKEETVKKRLDVYEAQTKPLITYYGDWARRGA  |   |   |   |     |
| [Truncated_Name:19] 4PZL_A.pdb | EPLITRTDDNEDTVKQRLSVYHAQTAKLIDFYRNFSSNT  |   |   |   |     |
|                                | * * * * ~ * ** ^ * ** ^ *                |   |   |   |     |

|     |   |   |   |     |
|-----|---|---|---|-----|
| 161 | . | . | . | 200 |
|-----|---|---|---|-----|

|                                |                             |   |   |     |
|--------------------------------|-----------------------------|---|---|-----|
|                                | 201                         | . | . | 227 |
| [Truncated_Name:1] 1AKE_A.pdb  | T--KYAKVDGTPVAEVRADLEKILG-  |   |   |     |
| [Truncated_Name:2] 8BQF_A.pdb  | T--KYAKVDGTPVAEVRADLEKIL--  |   |   |     |
| [Truncated_Name:3] 4X8M_A.pdb  | T--KYAKVDGTPVAEVRADLEKILG-  |   |   |     |
| [Truncated_Name:4] 6S36_A.pdb  | T--KYAKVDGTPVAEVRADLEKILG-  |   |   |     |
| [Truncated_Name:5] 8Q2B_A.pdb  | T--KYAKVDGTPVAEVRADLEKILG-  |   |   |     |
| [Truncated_Name:6] 8RJ9_A.pdb  | T--KYAKVDGTPVAEVRADLEKILG-  |   |   |     |
| [Truncated_Name:7] 6RZE_A.pdb  | T--KYAKVDGTPVAEVRADLEKILG-  |   |   |     |
| [Truncated_Name:8] 4X8H_A.pdb  | T--KYAKVDGTPVAEVRADLEKILG-  |   |   |     |
| [Truncated_Name:9] 3HPR_A.pdb  | T--KYAKVDGTPVAEVRADLEKILG-  |   |   |     |
| [Truncated_Name:10] 1E4V_A.pdb | T--KYAKVDGTPVAEVRADLEKILG-  |   |   |     |
| [Truncated_Name:11] 5EJE_A.pdb | T--KYAKVDGTPVAEVRADLEKILG-  |   |   |     |
| [Truncated_Name:12] 1E4Y_A.pdb | T--KYAKVDGTPVAEVRADLEKILG-  |   |   |     |
| [Truncated_Name:13] 3X2S_A.pdb | T--KYAKVDGTPVAEVRADLEKILG-  |   |   |     |
| [Truncated_Name:14] 6HAP_A.pdb | T--KYAKVDGTPVCEVRADLEKILG-  |   |   |     |
| [Truncated_Name:15] 6HAM_A.pdb | T--KYAKVDGTPVCEVRADLEKILG-  |   |   |     |
| [Truncated_Name:16] 4K46_A.pdb | T--QYLKFDGTKAVAEVSAELEKALA- |   |   |     |
| [Truncated_Name:17] 4NP6_A.pdb | T--QYLKFDGTKQVSEVSADIAKALA- |   |   |     |
| [Truncated_Name:18] 3GMT_A.pdb | E-----NGLKAPA-----YRKISG-   |   |   |     |

```
[Truncated_Name:19]4PZL_A.pdb      KIPKYIKINGDQAVEKVSQDIFDQLNK
                                *
                                .          .          227
                                201
```

Call:

```
pdbaln(files = files, fit = TRUE, exefile = "msa")
```

Class:

```
pdb, fasta
```

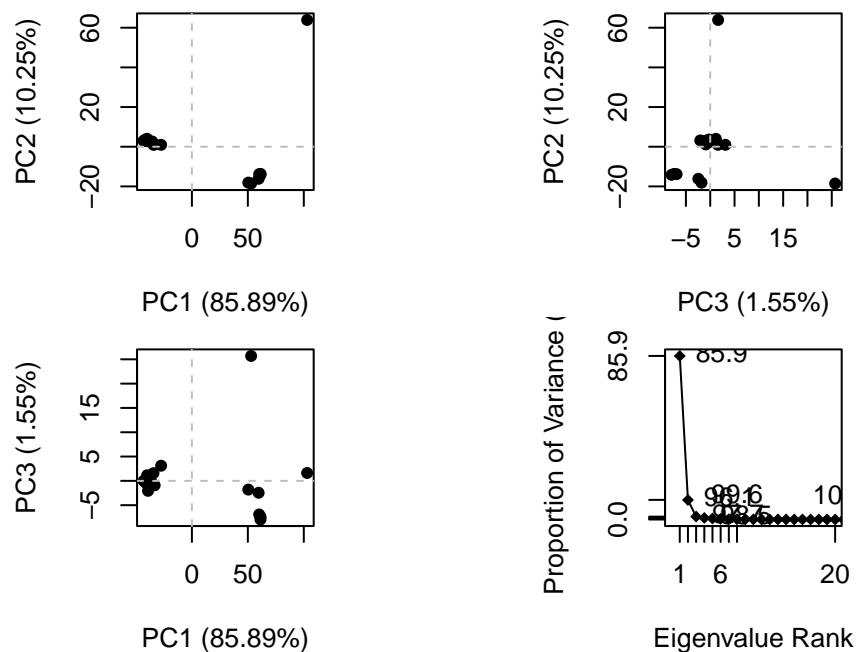
Alignment dimensions:

```
19 sequence rows; 227 position columns (199 non-gap, 28 gap)
```

```
+ attr: xyz, resno, b, chain, id, ali, resid, sse, call
```

## Principal Component Analysis

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



Making an animatino of these results:

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

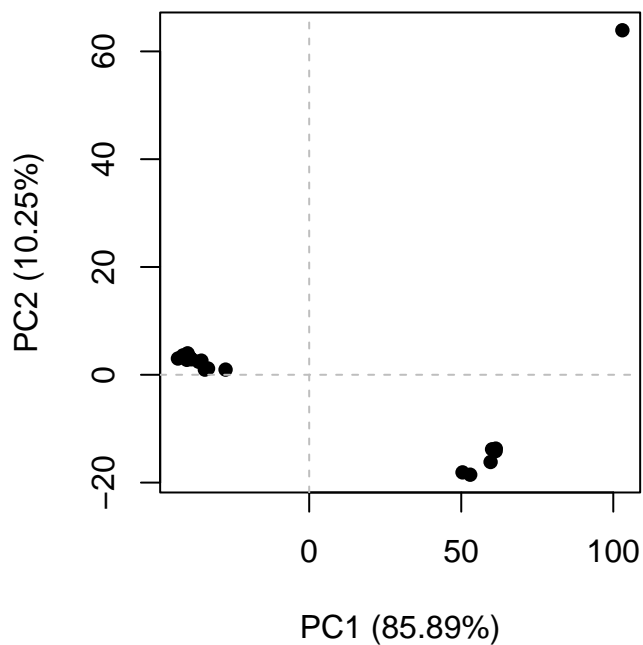
```
Total Frames#: 34
Total XYZs#: 597, (Atoms#: 199)
```

```
[1] 26.963 52.028 40.591 <...> 17.149 50.681 40.366 [20298]
```

```
+ attr: Matrix DIM = 34 x 597
```

```
pc.xray <- pca(pdb)
```

```
plot(pc.xray, pc.axes = c(1,2))
```



Percent of protein pdb sequences compared to uniprot sequences:

```
uniprot <- 24883887
pdb <- 195610
pdb/uniprot * 100
```

[1] 0.0786091