

class09

Phoebe LI

2/13/2022

```
# import the csv file
db <- read.csv("Data Export Summary.csv", row.names = 1)
db
```

	X.ray	NMR	EM	Multiple.methods	Neutron	Other	Total
## Protein (only)	144433	11881	6732	182	70	32	163330
## Protein/Oligosaccharide	8543	31	1125	5	0	0	9704
## Protein/NA	7621	274	2165	3	0	0	10063
## Nucleic acid (only)	2396	1399	61	8	2	1	3867
## Other	150	31	3	0	0	0	184
## Oligosaccharide (only)	11	6	0	1	0	4	22

```
# Xray structure percent
xray.percent <- sum(db$X.ray)/sum(db$Total)*100
round(xray.percent, 2)
```

```
## [1] 87.17
```

```
# EM structure percent
EM.percent <- sum(db$EM)/sum(db$Total)*100
round(EM.percent, 2)
```

```
## [1] 5.39
```

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

x-ray:87.2% EM:5.39%

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

```
percent.protein <- (db$Total[1])/sum(db$Total)
round(percent.protein*100, 2)
```

```
## [1] 87.26
```

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?

4486 Structures

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

Because It does not show the H atom in this structure. The resolution is too low that they can not detect H atom.

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)?

OH308O

```
install.packages("bio3d", repos="http://cran.us.r-project.org")
```

Introduction to Bio3D in R

```
##  
## The downloaded binary packages are in  
## /var/folders/9d/xssg21015fq5rb8769f22wfw0000gn/T//RtmpOMLiWy/downloaded_packages
```

```
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:  
## PQITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWPKMIGGIGGFIKVRQYD  
## QILIEICGKAIGTVLVGPTPVIIGRNLLTQIGCTLNFPQITLWQRPLVTIKIGGQLKE
```

```
##      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>
```

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

```
PQITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWKPKMIGGIGGFIKVRQYD    QILIE-
ICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNFPQITLWQRPLVTIKIGGQLKE        ALLDTGAD-
DTVLEEMSLPGRWKPKMIGGIGGFIKVRQYDQILIEICGHKAIGTVLVGPTP VNIIGRNLLTQIGCTLNF
```

198 amino acid residues

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

HOH

Q9: How many protein chains are in this structure?

2 chains

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

Q13

```
options(repos = list(CRAN="http://cran.rstudio.com/"))
install.packages("bio3d")
```

```
##
## The downloaded binary packages are in
## /var/folders/9d/xssg21015fq5rb8769f22wfw0000gn/T//RtmpOMLiWy/downloaded_packages
```

```
install.packages("ggplot2")
```

```
##
## The downloaded binary packages are in
## /var/folders/9d/xssg21015fq5rb8769f22wfw0000gn/T//RtmpOMLiWy/downloaded_packages
```

```
install.packages("ggrepel")
```

```
##
## The downloaded binary packages are in
## /var/folders/9d/xssg21015fq5rb8769f22wfw0000gn/T//RtmpOMLiWy/downloaded_packages
```

```
install.packages("devtools")
```

```
##
## The downloaded binary packages are in
## /var/folders/9d/xssg21015fq5rb8769f22wfw0000gn/T//RtmpOMLiWy/downloaded_packages
```

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

```
##
## The downloaded binary packages are in
## /var/folders/9d/xssg21015fq5rb8769f22wfw0000gn/T//RtmpOMLiWy/downloaded_packages
```

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

```
## 'getOption("repos")' replaces Bioconductor standard repositories, see
## '?repositories' for details
##
## replacement repositories:
##   CRAN: http://cran.rstudio.com/
```

```
## Bioconductor version 3.14 (BiocManager 1.30.16), R 4.1.2 (2021-11-01)
```

```
## Warning: package(s) not installed when version(s) same as current; use 'force = TRUE' to
## re-install: 'msa'
```

```
devtools::install_bitbucket("Grantlab/bio3d-view")
```

```
## Skipping install of 'bio3d.view' from a bitbucket remote, the SHA1 (dd153987) has not changed since
## Use 'force = TRUE' to force installation
```

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

```
## Warning in get.seq("lake_A"): Removing existing file: seqs.fasta
```

```
## Fetching... Please wait. Done.
```

```
aa
```

```
##           1           .           .           .           .           .           60
## pdb|1AKE|A  MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMLRAAVKSGSELGKQAKDIMDAGKLV
##           1           .           .           .           .           .           60
##
##           61           .           .           .           .           .           120
## pdb|1AKE|A  DELVIALVKERIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFDVPDELIVDRI
##           61           .           .           .           .           .           120
##
##           121          .           .           .           .           .           180
## pdb|1AKE|A  VGRRVHAPSGRVYHVKNPPKVEGKDDVTGEELTRKDDQEETVRKRLVEYHQMTAPLIG
##           121          .           .           .           .           .           180
##
##           181          .           .           .           .           .           214
## pdb|1AKE|A  YYSKEAEAGNTKYAKVDGTPVAEVRADLEKILG
##           181          .           .           .           .           .           214
##
## Call:
## read.fasta(file = outfile)
##
## Class:
## fasta
##
## Alignment dimensions:
## 1 sequence rows; 214 position columns (214 non-gap, 0 gap)
##
## + attr: id, ali, call
```

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

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

```
files <- get.pdb(hits$pdb.id, path="pdbs", split=TRUE, gzip=TRUE)
```

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

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

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

```
## Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE): pdbs/  
## 3HPR.pdb.gz exists. Skipping download
```

```
## Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE): pdbs/  
## 1E4V.pdb.gz exists. Skipping download
```

```
## Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE): pdbs/  
## 5EJE.pdb.gz exists. Skipping download
```

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

```
## Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE): pdbs/  
## 3X2S.pdb.gz exists. Skipping download
```

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

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

```
## Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE): pdbs/  
## 4K46.pdb.gz exists. Skipping download
```

```
## 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
```

```
## |
```

```
pdbs <- pdbaln(files, fit = TRUE)
```

```

## 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

```

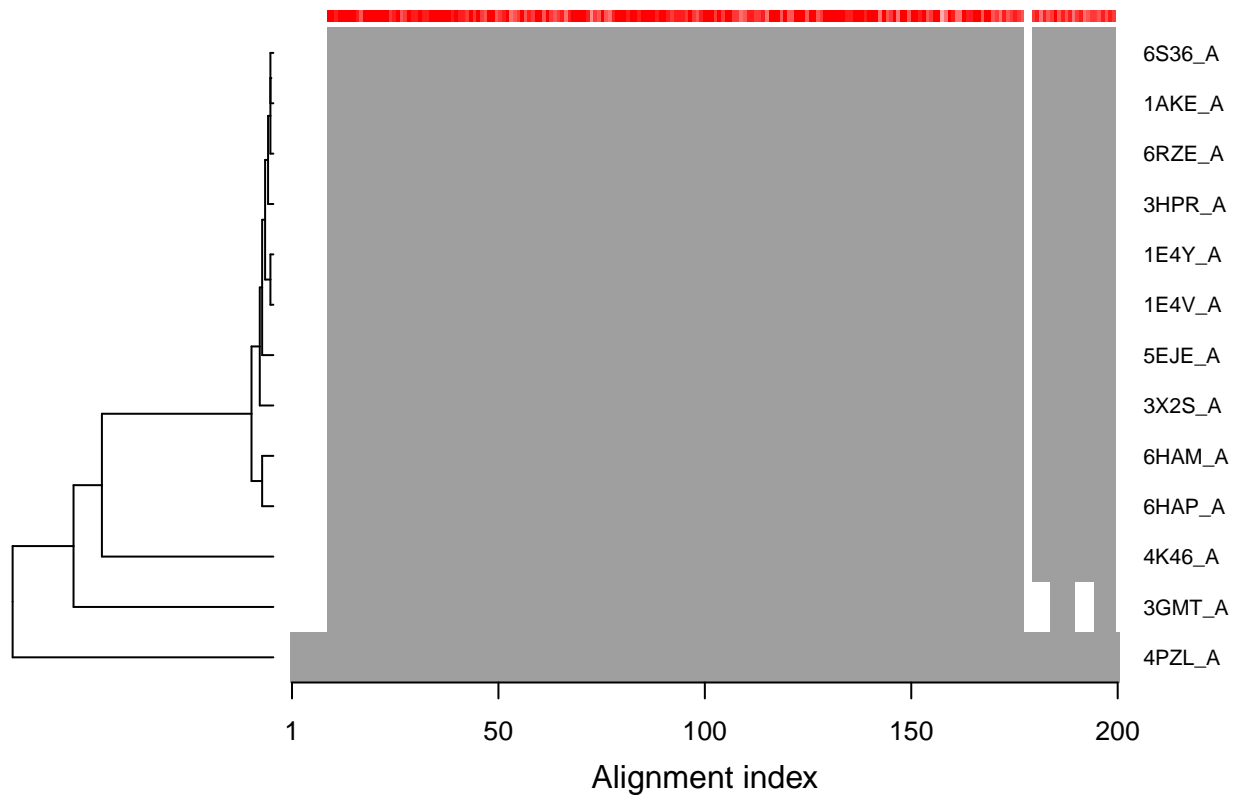
```

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

# Draw schematic alignment
plot(pdb, labels=ids)

```

Sequence Alignment Overview



Viewing our superposed structures

```
install.packages("rgl")
```

```
##
```

```
## The downloaded binary packages are in
```

```
## /var/folders/9d/xssg21015fq5rb8769f22wfw0000gn/T//RtmpOMLiWy/downloaded_packages
```

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

```
library(rgl)
```

```
view.pdbs(pdbs)
```