

# classs09

Seong Tae Gwon

2/21/2022

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

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:

```
db<-read.csv("Data Export SUMmary.csv",row.names=1)
head(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

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

```
methodsums<-colSums(db)
round(methodsums/methodsums["Total"] *100,2)
```

##	X.ray	NMR	EM	Multiple.methods
##	87.17	7.28	5.39	0.11
##	Neutron	Other	Total	
##	0.04	0.02	100.00	

**Answer:** 87.17% of structures in the PDB are solved by X-Ray, and 5.39% are solved by Electron Microscopy.

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

```
typesums <- rowSums(db)
round( (db$Total/methodsums["Total"]) * 100,2)
```

```
## [1] 87.26  5.18  5.38  2.07  0.10  0.01
```

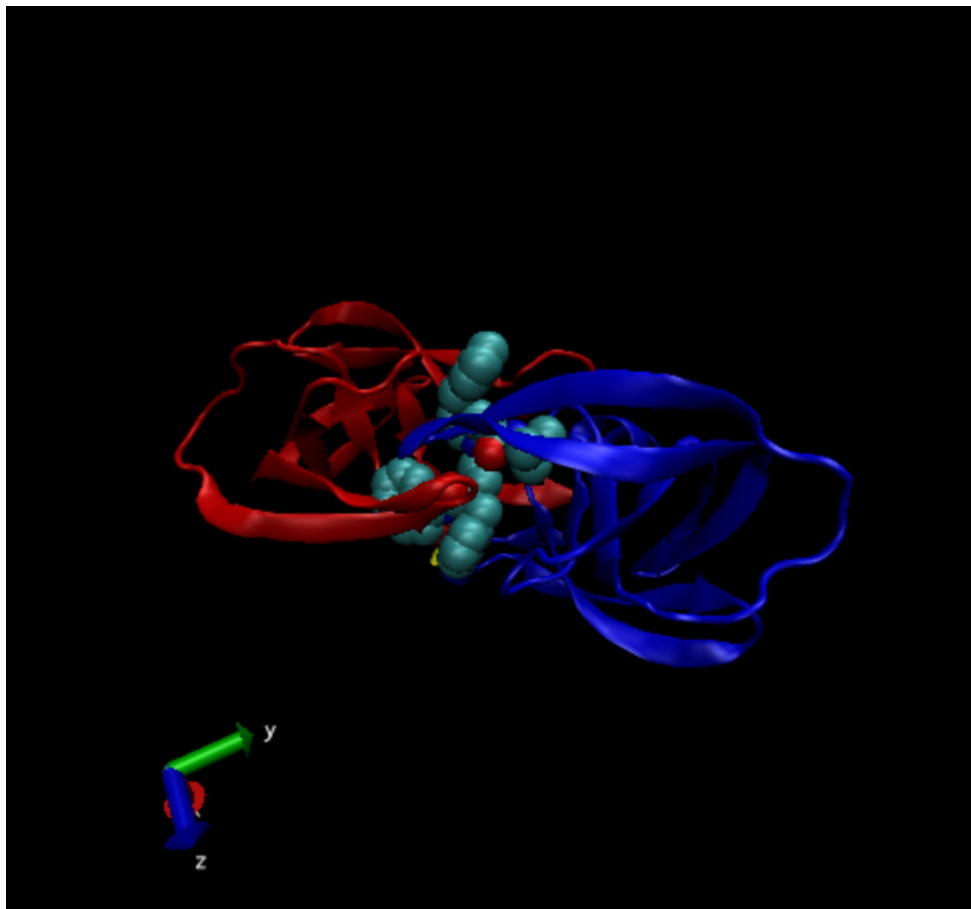
*# First element is protein (only)*

**Answer:** 87.26% 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?

**Answer:** There are 1828 HIV-1 protease structures in the current PDB.

## 2. Visualizing the HIV-1 protease structure

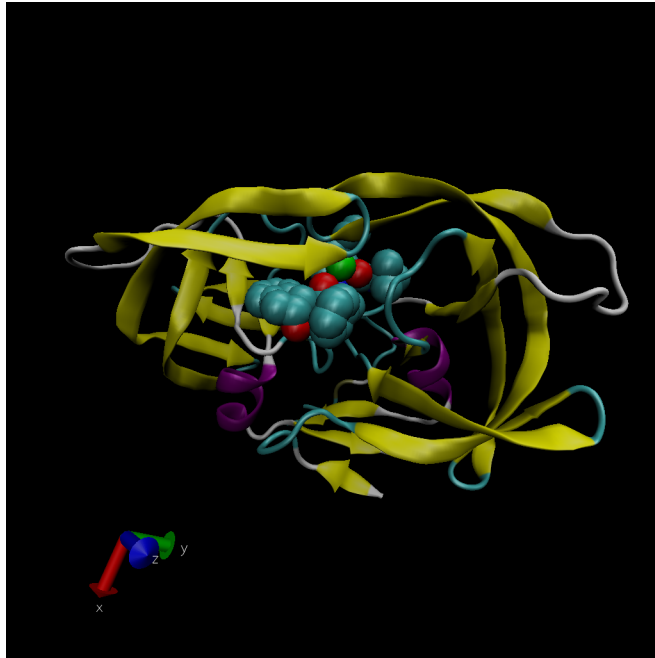


Q4. Water molecules normally have 3 atoms. Why do we see just one atom per water molecule in this structure?

**Answer:** We only see one atom per water molecule (H<sub>2</sub>O) in this structure because the resolution capacity is not high enough to display two hydrogen atoms that are significantly smaller than the oxygen atom. Hence, we only see the oxygen atom of the 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?

**Answer:** Conserved water molecule in the binding site has residue number 308.



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?

**Answer:** Yes. As shown on the generated VMD graphical representation, I observe multiple secondary structure elements in different colors other than sequences corresponding to B-factor and secondary structure elements. Noticeably, extended beta structure in yellow at chain B residue 3 and alpha helix structures in purple at residue 87-92 are likely to only form in the dimer rather than the monomer.

### 3. Introduction to Bio3D in R

```
#install.packages("bio3d")
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
```

```
aa123(pdbseq(pdb))
```

```
##      [1] "PRO" "GLN" "ILE" "THR" "LEU" "TRP" "GLN" "ARG" "PRO" "LEU" "VAL" "THR"
##      [13] "ILE" "LYS" "ILE" "GLY" "GLY" "GLN" "LEU" "LYS" "GLU" "ALA" "LEU" "LEU"
##      [25] "ASP" "THR" "GLY" "ALA" "ASP" "ASP" "THR" "VAL" "LEU" "GLU" "GLU" "MET"
##      [37] "SER" "LEU" "PRO" "GLY" "ARG" "TRP" "LYS" "PRO" "LYS" "MET" "ILE" "GLY"
##      [49] "GLY" "ILE" "GLY" "GLY" "PHE" "ILE" "LYS" "VAL" "ARG" "GLN" "TYR" "ASP"
##      [61] "GLN" "ILE" "LEU" "ILE" "GLU" "ILE" "CYS" "GLY" "HIS" "LYS" "ALA" "ILE"
##      [73] "GLY" "THR" "VAL" "LEU" "VAL" "GLY" "PRO" "THR" "PRO" "VAL" "ASN" "ILE"
##      [85] "ILE" "GLY" "ARG" "ASN" "LEU" "LEU" "THR" "GLN" "ILE" "GLY" "CYS" "THR"
##      [97] "LEU" "ASN" "PHE" "PRO" "GLN" "ILE" "THR" "LEU" "TRP" "GLN" "ARG" "PRO"
##      [109] "LEU" "VAL" "THR" "ILE" "LYS" "ILE" "GLY" "GLY" "GLN" "LEU" "LYS" "GLU"
##      [121] "ALA" "LEU" "LEU" "ASP" "THR" "GLY" "ALA" "ASP" "ASP" "THR" "VAL" "LEU"
##      [133] "GLU" "GLU" "MET" "SER" "LEU" "PRO" "GLY" "ARG" "TRP" "LYS" "PRO" "LYS"
##      [145] "MET" "ILE" "GLY" "GLY" "ILE" "GLY" "GLY" "PHE" "ILE" "LYS" "VAL" "ARG"
##      [157] "GLN" "TYR" "ASP" "GLN" "ILE" "LEU" "ILE" "GLU" "ILE" "CYS" "GLY" "HIS"
##      [169] "LYS" "ALA" "ILE" "GLY" "THR" "VAL" "LEU" "VAL" "GLY" "PRO" "THR" "PRO"
##      [181] "VAL" "ASN" "ILE" "ILE" "GLY" "ARG" "ASN" "LEU" "LEU" "THR" "GLN" "ILE"
##      [193] "GLY" "CYS" "THR" "LEU" "ASN" "PHE"
```

```
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?

**Answer:** There are 198 amino acid residues in this pdb object.

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

**Answer:** One of the 2 non-protein residues is MK1 (residue # 1).

Q9. How many protein chains are in this structure?

**Answer:** There are 2 protein chains in this structure.

## 4. Comparative structure analysis of Adenylate Kinase

```
# Setup: Install packages in the R console not your Rmd

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?

**Answer:** “msa” package is found only on BioConductor and not CRAN.

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

**Answer:** “bio3d-view” package 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?

**Answer:** True

```
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  MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMMLRAAVKSGSELGKQAKDIMDAGKLV
##           1           .           .           .           .           .           60
##
##           61           .           .           .           .           .           120
## pdb|1AKE|A  DELVIALVKERIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFDVPDELIVDRI
##           61           .           .           .           .           .           120
##
##           121          .           .           .           .           .           180
## pdb|1AKE|A  VGRRVHAPSGRVYHVKFNPVKVEGKDDVTGEELTRKDDQEETVRKRLVEYHQMTAPLIG
##           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?

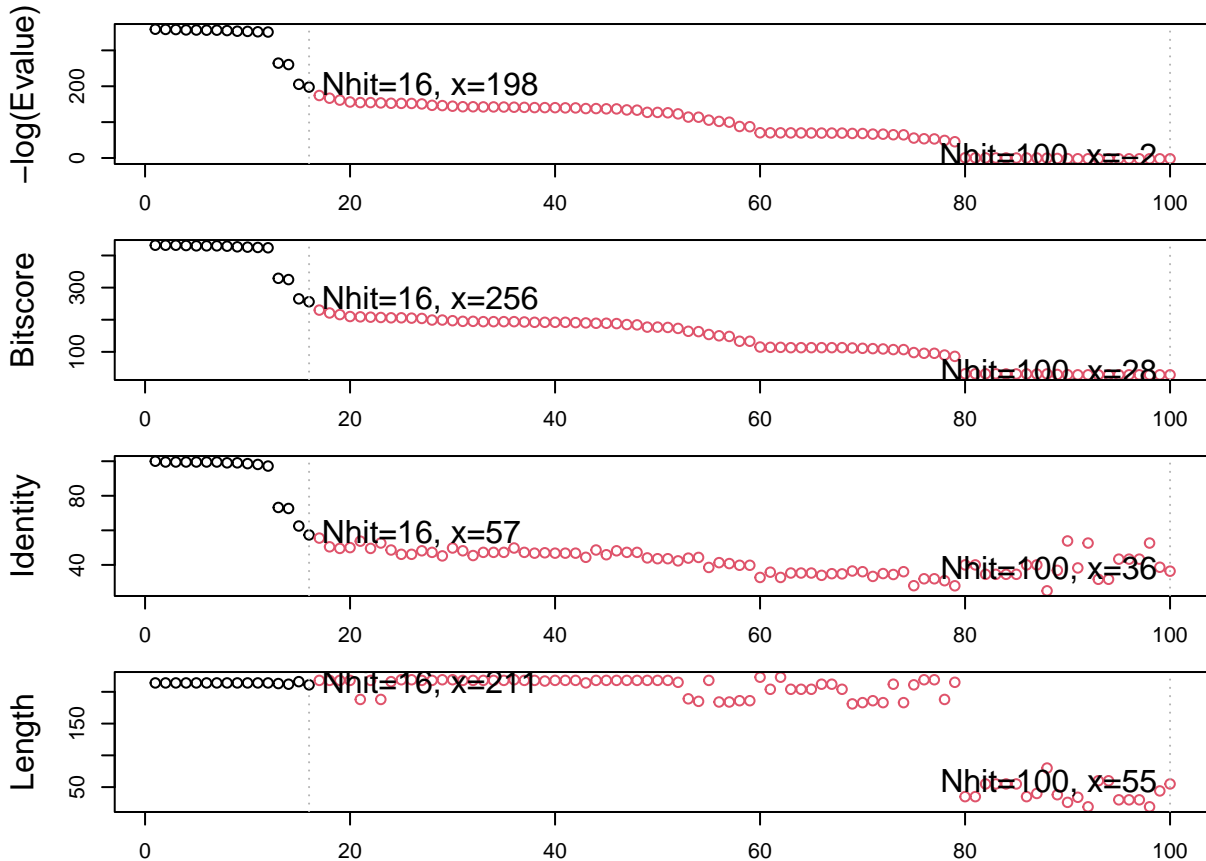
**Answer:** There are 214 amino acids in this sequence.

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

```
## Searching ... please wait (updates every 5 seconds) RID = 1KN89YTX013
## .
## 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"
```

```
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')
```

```
# 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/
## 6S36.pdb.gz exists. Skipping download
```

```
## Warning in get.pdb(hits$ pdb.id, path = "pdb", split = TRUE, gzip = TRUE): pdb/
## 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

##      |

```

\*\* Stop here: what we covered in class \*\*