

## Class 9: Structural Bioinformatics (Pt. 1)

Nataliana Hernandez (PID: A17096549)

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

Data from : <http://www.rcsb.org/>

Read this into R:

```
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					

and answer the following questions:

```
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
```

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)
```

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 try a different read/import function:**

```
#!/ message: false
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
```

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
```

```
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
```

## Mol\*

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

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

We will use PDB code: 1HSG

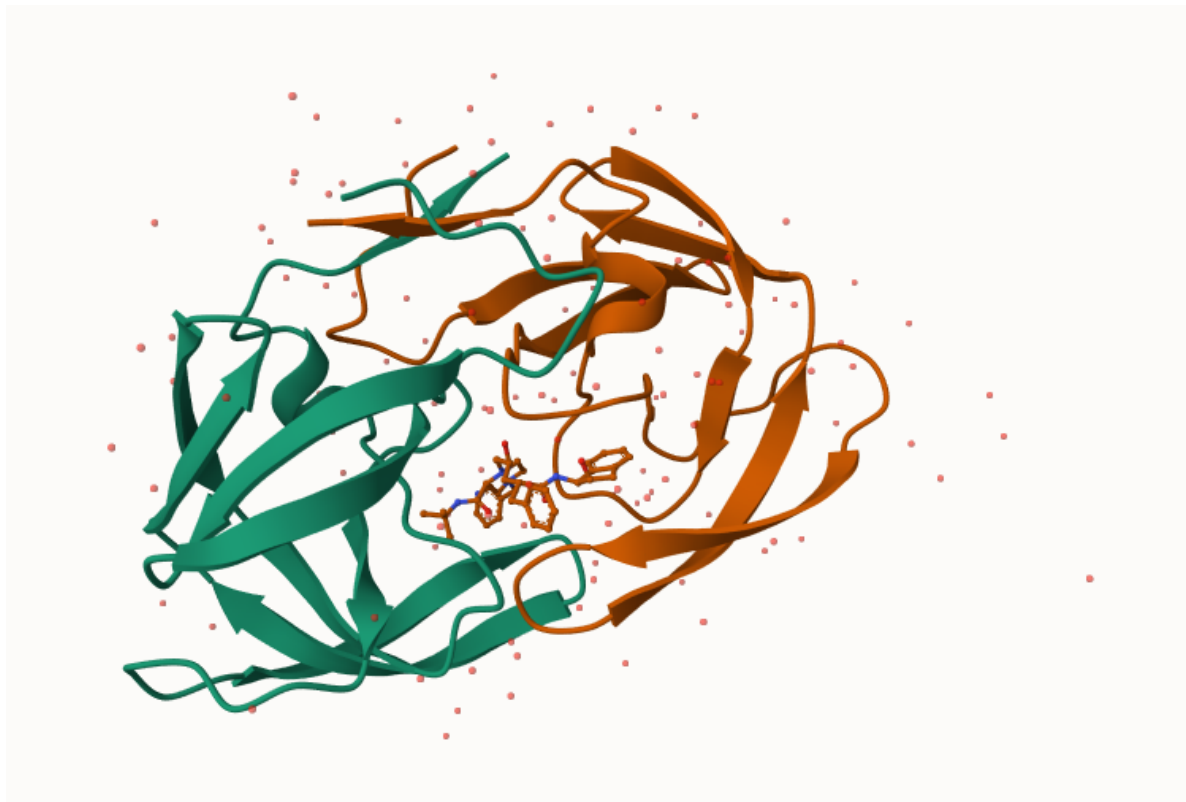


Figure 1: A first image from molstar

Some more custom images:

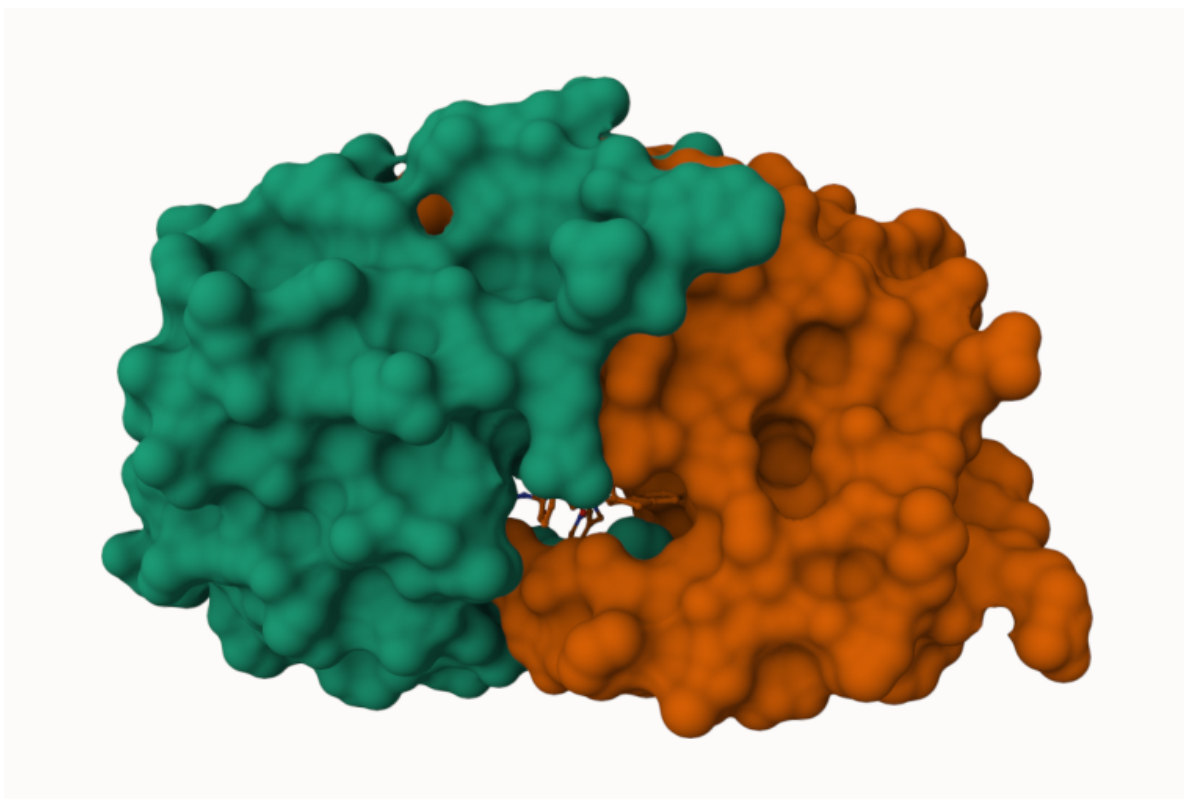
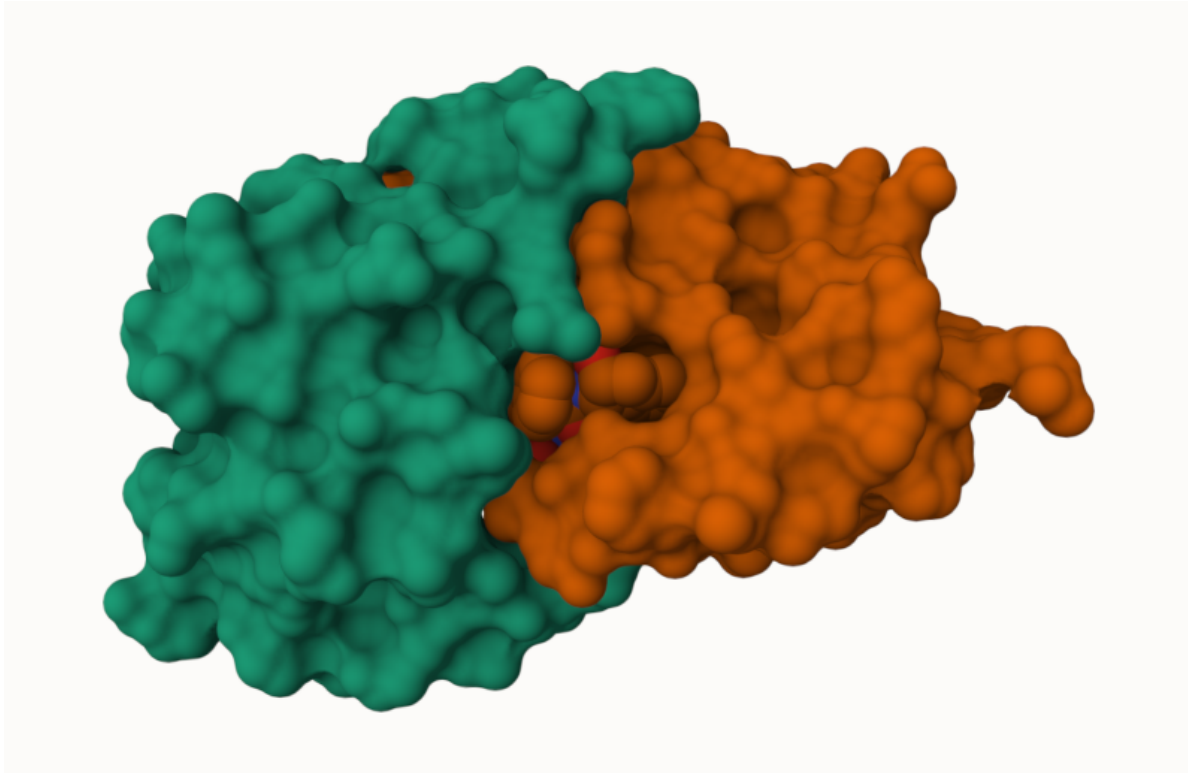


Figure 2: The all important catalytic ASP25 amino acids



### The Bio3D package

```
#readLines("hsg.pdb")
```

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?

- 5 structures

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

- The water molecules are only represented by the oxygen atom in the entire molecule, rather than all three atoms present based on x-ray crystallography.

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

- The conserved water molecule is the one in the binding site that stabilizes the interaction. HOH 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.

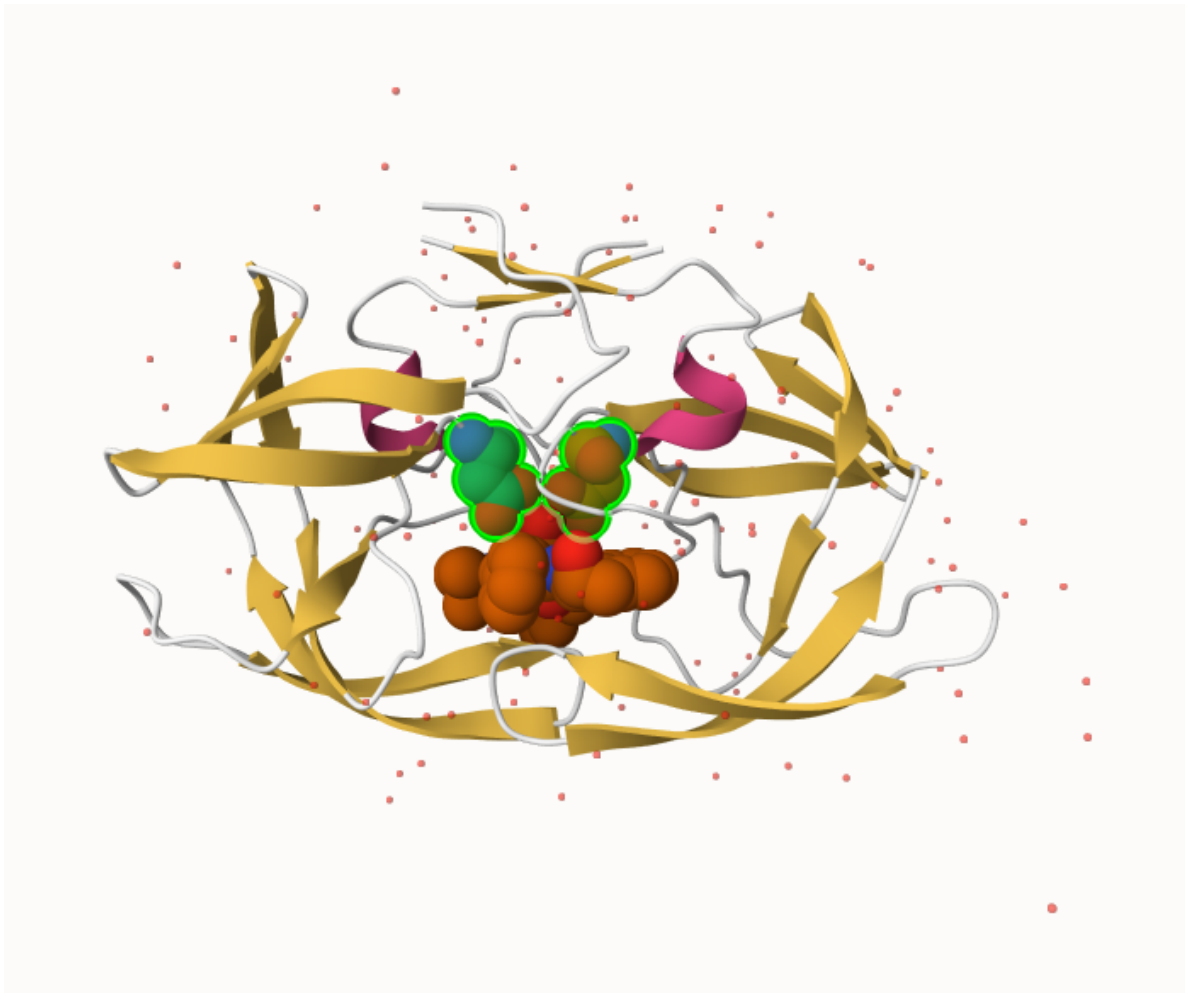


Figure 3: Both ASP 25

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

```
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>

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

- 198

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

- HOH (127) and MK1 (1)

Q9: How many protein chains are in this structure?

- 2

## 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  
DELVIALVKERIAQEDCRNGFLDGFRTIPQADAMKEAGINVDYVLEFDVPDELIVDKI  
VGRRVHAPSGRVYHVKFNPPKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQMTAPLIG  
YYSKEAEAGNTKYAKVDGTPVAEVRADLEKILG

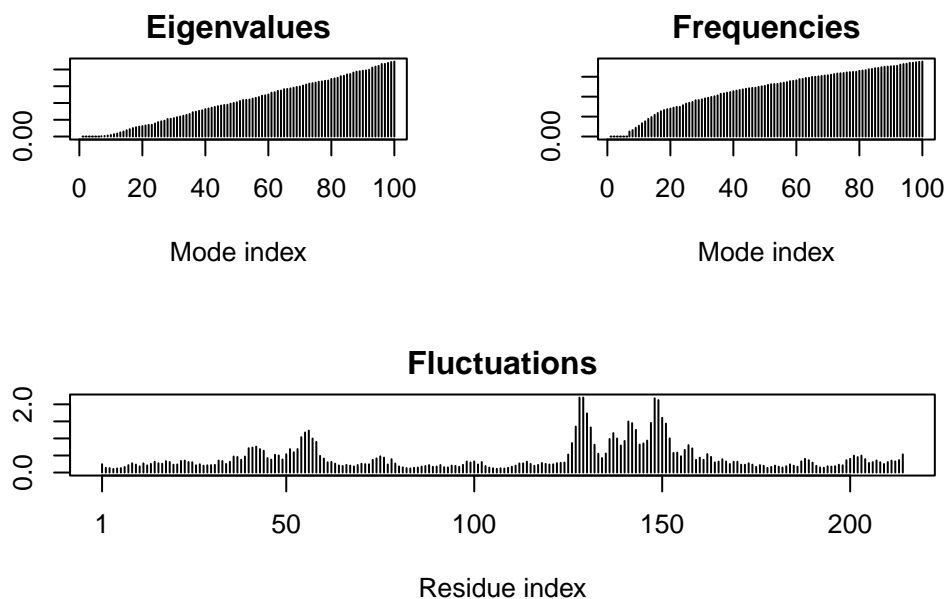
+ attr: atom, xyz, seqres, helix, sheet,  
calpha, remark, call

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

Building Hessian... Done in 0.03 seconds.

Diagonalizing Hessian... Done in 0.19 seconds.

```
plot(m)
```



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

```
# Install packages in the R console NOT your Rmd/Quarto file

#install.packages("bio3d")
#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?

- True

```
library(bio3d)
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  MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMLRAAVKSGSELGKQAKDIMDAGKLV
      1      .      .      .      .      .      60

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

     121      .      .      .      .      .      180
```

```

pdb|1AKE|A    VGRRVHAPSGRVYHVKNPPKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQM TAPLIG
              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