

# Class 11: Structural Bioinfo pt2

Fan Wu(PID: A15127541)

## Table of contents

AlphaFold DataBase . . . . .	1
Generating your own structure predictions . . . . .	2
<b>Custom analysis of resulting models in R</b>	<b>5</b>
Residue conservation from alignment file . . . . .	8

## AlphaFold DataBase

The EBI maintains the largest database of AlphaFold structure predictions models at:  
<https://alphafold.ebi.ac.uk/>

From last class, Class09, (before Halloween), we saw that the PDB had 244,290 protein structures(October 2025)

The total number of protein sequences in UniProtKB is 199,579,901

**Key Point:** This is a tiny fraction of sequence space that has structural coverage (0.12%)

```
244290 / 199579901 * 100
```

```
[1] 0.1224021
```

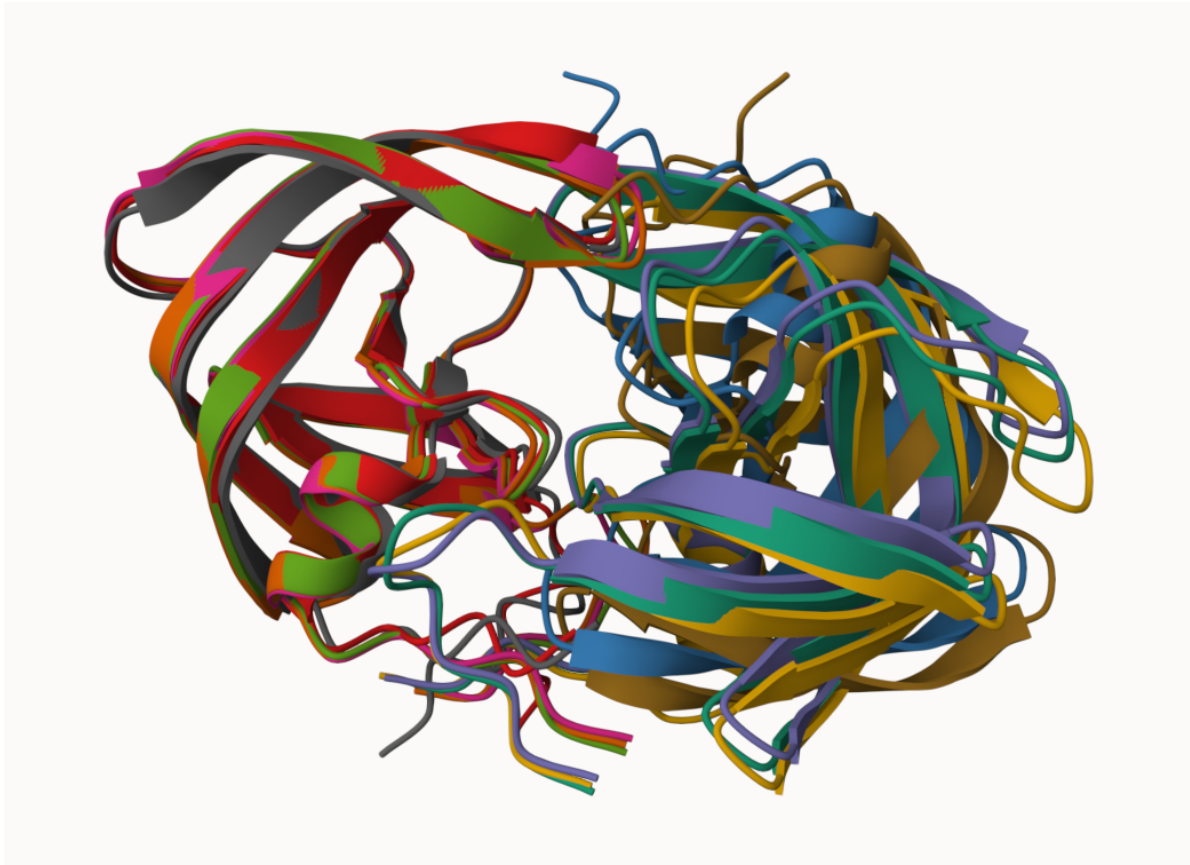
sequencing is way cheaper than experiment that cost 1M to get a sequence

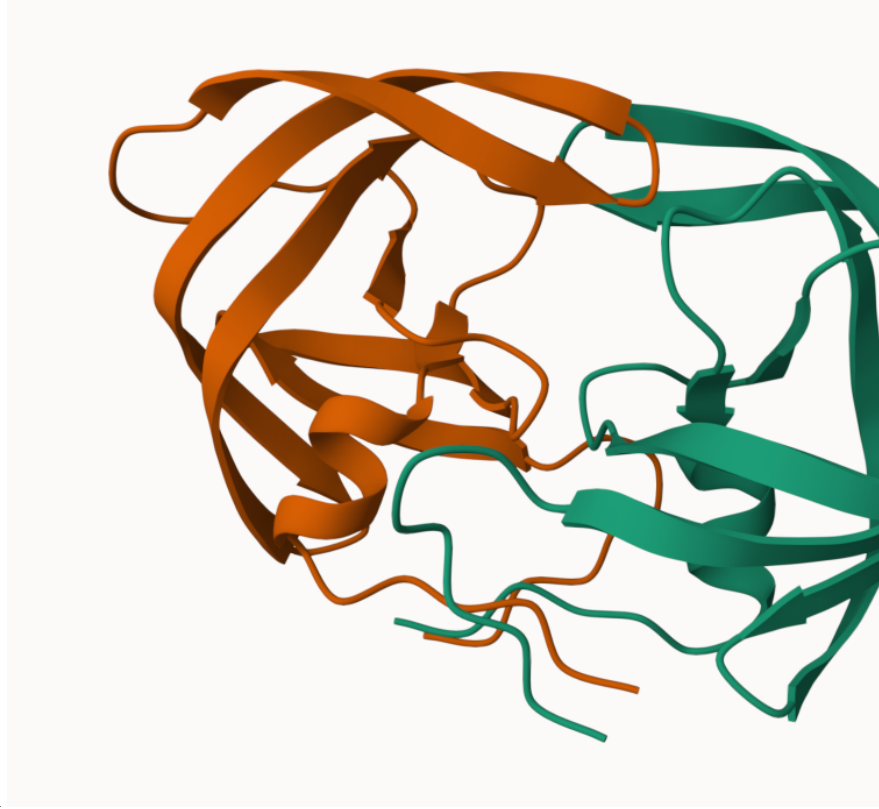
AFDB is attempting to address this gap...

There are two “Quality scores” from AlphaFold one for residues (i.e. each amino acid) called **pLDDT** score. The other **PAE** (Predicted Aligned Error) score measures the confidence in the relative position of two residues (i.e. score for every pair of residues)

## Generating your own structure predictions

Figure of 5 generated HIV-PR models

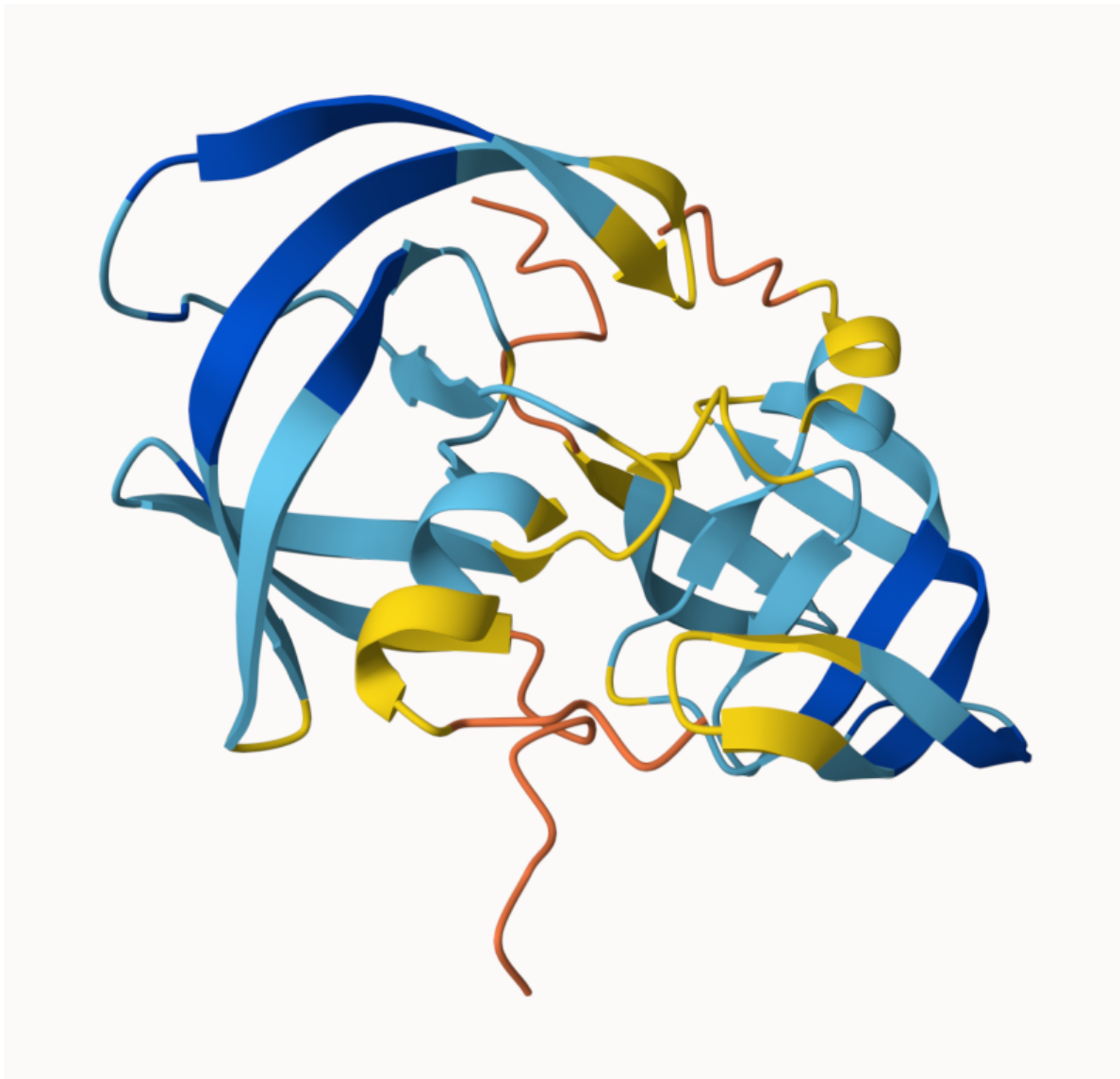




And the top ranked model colored by chain  
pLDTT score for model 1



and model 5



## Custom analysis of resulting models in R

Read key result files into R. The first thing I need to know is what my results directory/folder is called (i.e. its name is different for every AlphaFold run/job)

```
results_dir <- "HIVPR_dimer_23119/"  
  
# File names for all PDB models
```

```

pdb_files <- list.files(path=results_dir,
                        pattern="*.pdb",
                        full.names = TRUE)

```

```

# Print our PDB file names
basename(pdb_files)

```

```

[1] "HIVPR_dimer_23119_unrelaxed_rank_001_alphafold2_multimer_v3_model_4_seed_000.pdb"
[2] "HIVPR_dimer_23119_unrelaxed_rank_002_alphafold2_multimer_v3_model_1_seed_000.pdb"
[3] "HIVPR_dimer_23119_unrelaxed_rank_003_alphafold2_multimer_v3_model_5_seed_000.pdb"
[4] "HIVPR_dimer_23119_unrelaxed_rank_004_alphafold2_multimer_v3_model_2_seed_000.pdb"
[5] "HIVPR_dimer_23119_unrelaxed_rank_005_alphafold2_multimer_v3_model_3_seed_000.pdb"

```

```

library(bio3d)

```

```

m1 <- read.pdb(pdb_files[2])
m1

```

Call: read.pdb(file = pdb\_files[2])

Total Models#: 1

Total Atoms#: 1514, XYZs#: 4542 Chains#: 2 (values: A B)

Protein Atoms#: 1514 (residues/Calpha atoms#: 198)

Nucleic acid Atoms#: 0 (residues/phosphate atoms#: 0)

Non-protein/nucleic Atoms#: 0 (residues: 0)

Non-protein/nucleic resid values: [ none ]

Protein sequence:

```

PQITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWKPKMIGGIGGFIKVRQYD
QILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNFPQITLWQRPLVTIKIGGQLKE
ALLDTGADDTVLEEMSLPGRWKPKMIGGIGGFIKVRQYDQILIEICGHKAIGTVLVGPTP
VNIIGRNLLTQIGCTLNF

```

+ attr: atom, xyz, calpha, call

```

head(m1$atom)

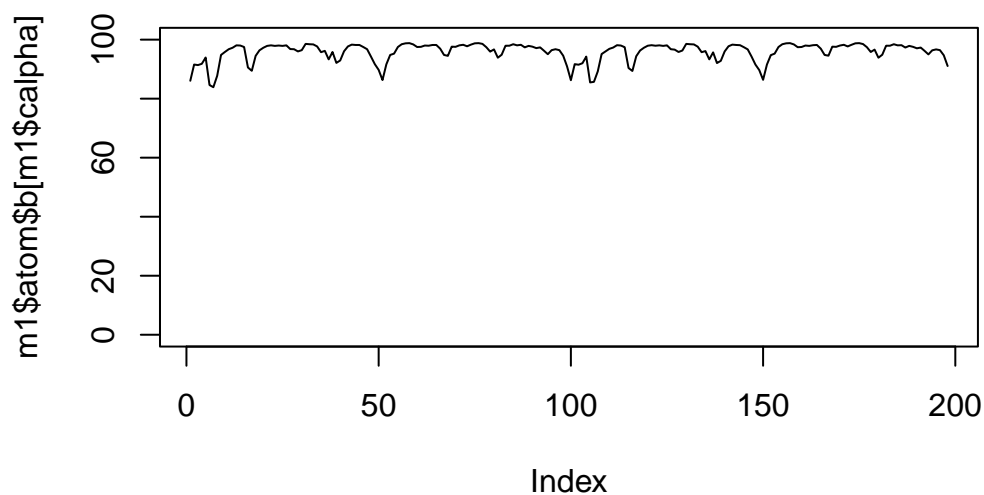
```

	type	eleno	elety	alt	resid	chain	resno	insert	x	y	z	o	b
1	ATOM	1	N	<NA>	PRO	A	1	<NA>	16.531	-4.008	-8.734	1	86.06
2	ATOM	2	CA	<NA>	PRO	A	1	<NA>	16.391	-2.600	-9.125	1	86.06
3	ATOM	3	C	<NA>	PRO	A	1	<NA>	15.875	-1.722	-7.984	1	86.06
4	ATOM	4	CB	<NA>	PRO	A	1	<NA>	15.391	-2.646	-10.281	1	86.06
5	ATOM	5	O	<NA>	PRO	A	1	<NA>	15.359	-2.234	-6.988	1	86.06
6	ATOM	6	CG	<NA>	PRO	A	1	<NA>	14.656	-3.936	-10.094	1	86.06

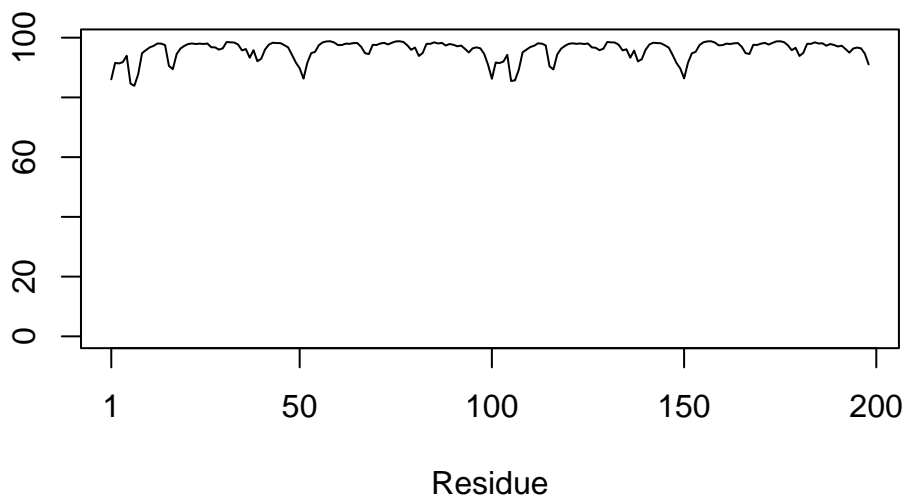
  

	segid	elesy	charge
1	<NA>	N	<NA>
2	<NA>	C	<NA>
3	<NA>	C	<NA>
4	<NA>	C	<NA>
5	<NA>	O	<NA>
6	<NA>	C	<NA>

```
plot( m1$atom$b[m1$alpha], typ = "l", ylim = c(0,100))
```



```
plot.bio3d(m1$atom$b[m1$alpha], typ = "l")
```



### Residue conservation from alignment file

Find the large AlphaFold alignment file

```
aln_file <- list.files(path=results_dir,  
                      pattern=".a3m$",  
                      full.names = TRUE)  
aln_file
```

```
[1] "HIVPR_dimer_23119//HIVPR_dimer_23119.a3m"
```

Read this into R

```
aln <- read.fasta(aln_file[1], to.upper = TRUE)
```

```
[1] " ** Duplicated sequence id's: 101 **"  
[2] " ** Duplicated sequence id's: 101 **"
```

How many sequences are in this alignment



```
dim(aln$ali)
```

```
[1] 5397 132
```

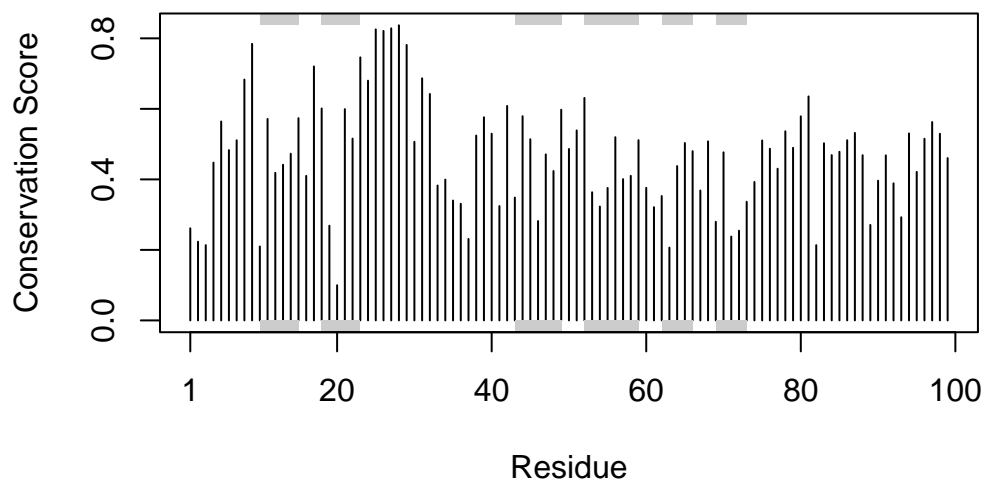
We can score residue conservation in the alignment with the `conserv()` function.

```
sim <- conserv(aln)
```

```
# Read a reference PDB structure  
pdb <- read.pdb("1hsg")
```

Note: Accessing on-line PDB file

```
plotb3(sim[1:99], sse=trim.pdb(pdb, chain="A"),  
        ylab="Conservation Score")
```



```
con <- consensus(aln, cutoff = 0.9)  
con$seq
```

```

[1] "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-"
[19] "-" "-" "-" "-" "-" "-" "D" "T" "G" "A" "-" "-" "-" "-" "-" "-" "-" "-"
[37] "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-"
[55] "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-"
[73] "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-"
[91] "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-"
[109] "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-" "-"
[127] "-" "-" "-" "-" "-" "-"

```