# Class 11: Structural Bioinformatics pt 2

Nicole (PID: A18116280)

**AlphaFold Data Base (AFDB)**

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

From last class (before Halloween) we saw that the PDB had 244,290 (Oct 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
```

AFDB is attempting to address this gap...

There are two "Quality Scores" from Alphafold. Due for residues (i.e. each amino acid) called **pLDDT** score. The other **PAE** score measures the confidence in the relative position of two residues (i.e. a score for every pair of residues).

**Generating your own structure predictions**
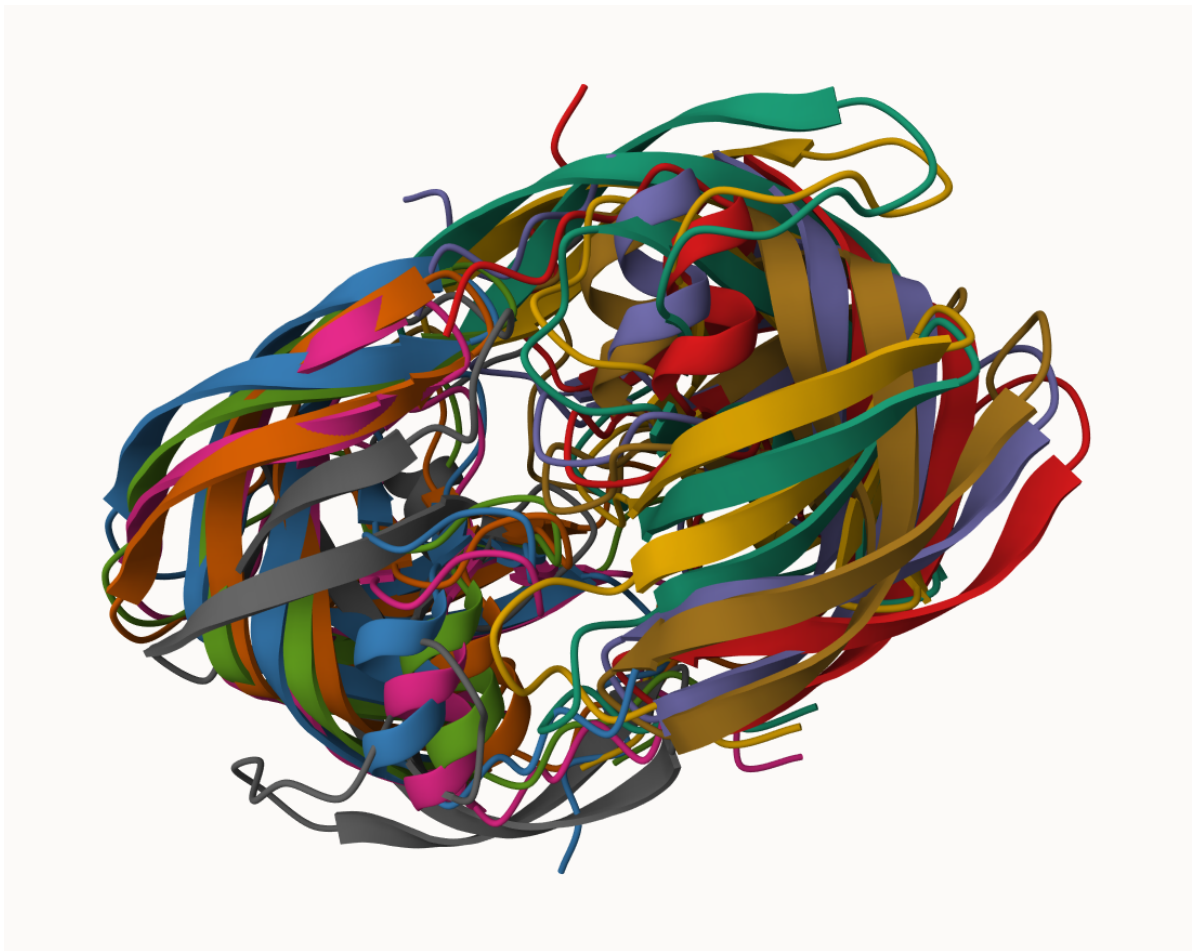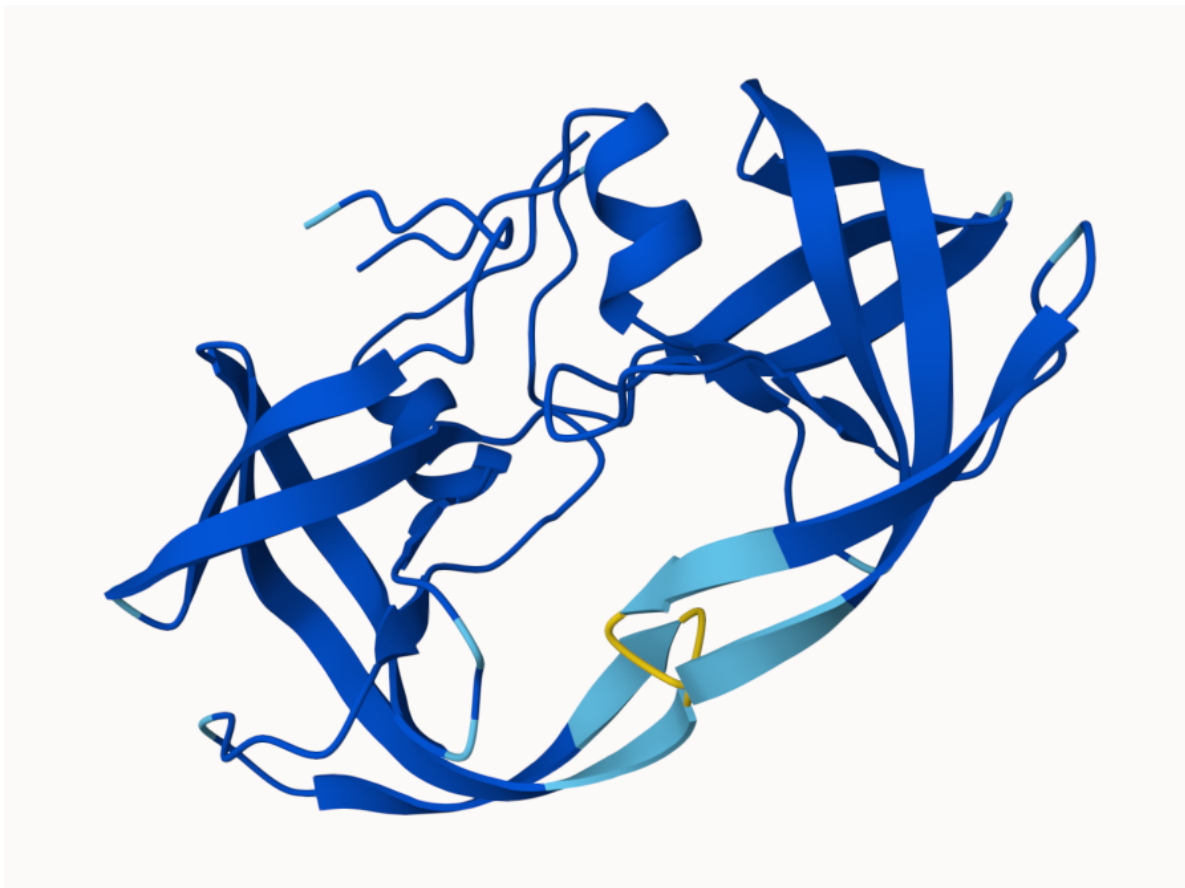
Image of all 5 models

Image of the 1st models

pLDDT score of 1st model

pLDDT score of 4th model

## Custom analysis of resulting models in R

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

```r
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"
```

```r
library(bio3d)

m1 <- read.pdb(pdb_files[1])
m1
```

```
 Call:  read.pdb(file = pdb_files[1])

   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
```
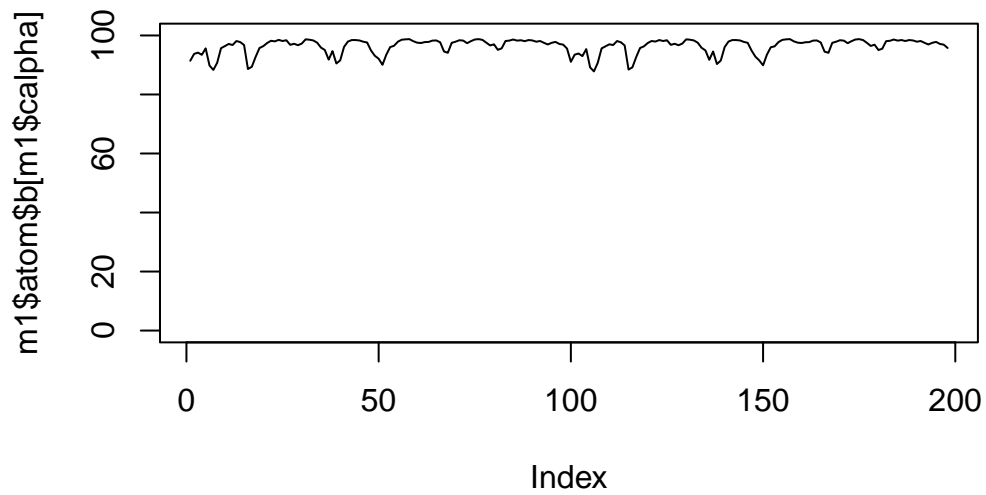
```r
head(m1$atom)
```

```
  type eleno elety  alt resid chain resno insert       x     y     z o     b
1 ATOM     1     N <NA>   PRO     A     1   <NA> -16.656 5.527 4.973 1 91.44
2 ATOM     2    CA <NA>   PRO     A     1   <NA> -17.000 4.836 3.725 1 91.44
3 ATOM     3     C <NA>   PRO     A     1   <NA> -16.375 3.453 3.623 1 91.44
4 ATOM     4    CB <NA>   PRO     A     1   <NA> -16.453 5.773 2.643 1 91.44
5 ATOM     5     O <NA>   PRO     A     1   <NA> -15.445 3.137 4.375 1 91.44
6 ATOM     6    CG <NA>   PRO     A     1   <NA> -15.336 6.512 3.307 1 91.44
  segid elesy charge
1  <NA>     N   <NA>
```
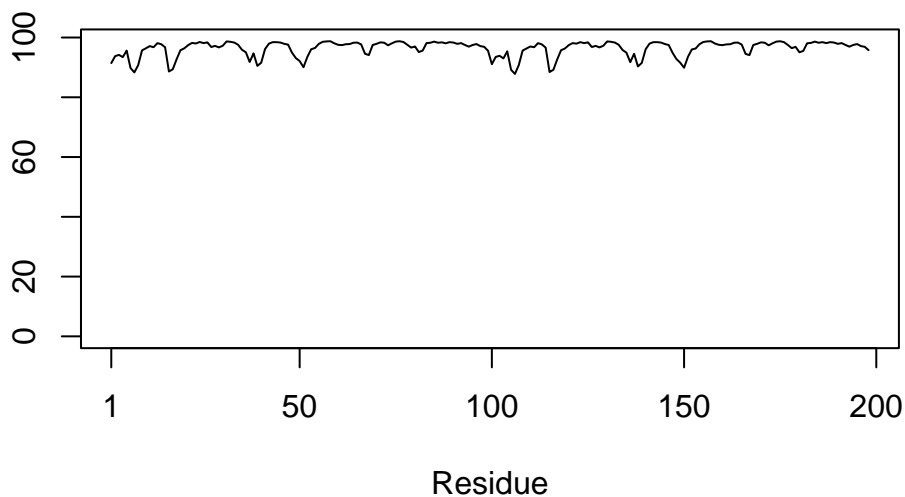
```
2   <NA>      C   <NA>
3   <NA>      C   <NA>
4   <NA>      C   <NA>
5   <NA>      O   <NA>
6   <NA>      C   <NA>
```

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



```r
plot.bio3d(m1$atom$b[m1$calpha], type="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 **"
```
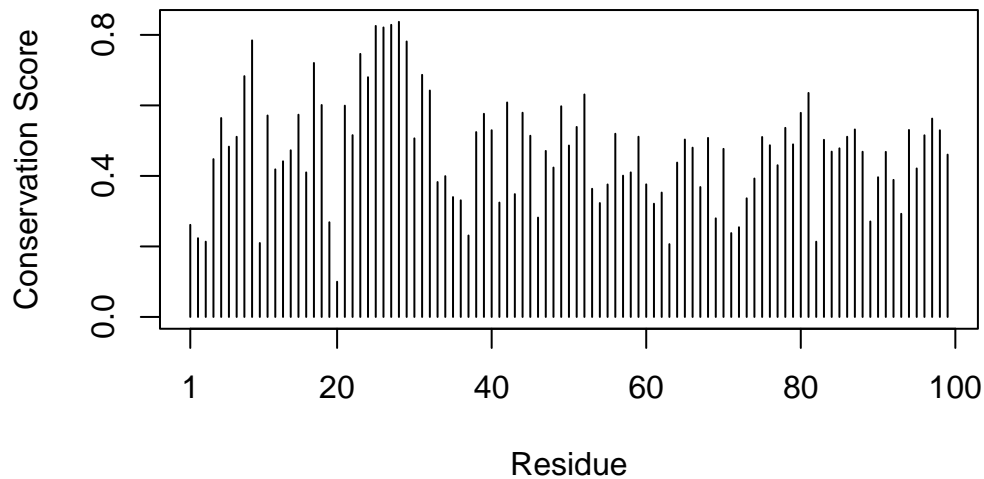
How many sequences are in this alignment

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

```
[1] 5397  132
```

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

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

```
plotb3(sim[1:99], ylab="Conservation Score")
```



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

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