

Class 11: Structural Bioinformatics pt2

Maria Tavares

##Background

We saw last day that PDB has 209,886 entries (Oct/Nov 2025). UniProtKB (i.e. protein sequence database) has 199,579,901 entries.

```
209886/199579901 * 100
```

```
[1] 0.1051639
```

So the PDB has only 0.1% coverage of the main sequence database.

Enter AlphaFold data base (AFDB) <<https://alphafold.ebi.ac.uk>> that attempts to provide computed models for all sequences in UniProt.

“AlphaFold DB provides open access to over 200 million protein structure predictions to accelerate scientific research”

AlphaFold

AlphaFold has 3 main outputs

- the predicted coordinates (PDB files)
- A local quality score called **plDDT** (one for each amino acid)
- A second quality score **PAE** Predicted Alignes Error (for each pair of amino-acid)

We can run AlphaFold ourselves if we not happy with AFDB (i.e. no coverage or poor model)

Interpreting/analyzing AF results in 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"
```

```
library(bio3d)
```

```
# Read all data from Models
# and superpose/fit coords
pdbs <- pdbaln(pdb_files, fit=TRUE, exefile="msa")
```

Reading PDB files:

```
HIVPR_dimer_23119//HIVPR_dimer_23119_unrelaxed_rank_001_alphaFold2_multimer_v3_model_4_seed_000
HIVPR_dimer_23119//HIVPR_dimer_23119_unrelaxed_rank_002_alphaFold2_multimer_v3_model_1_seed_000
HIVPR_dimer_23119//HIVPR_dimer_23119_unrelaxed_rank_003_alphaFold2_multimer_v3_model_5_seed_000
HIVPR_dimer_23119//HIVPR_dimer_23119_unrelaxed_rank_004_alphaFold2_multimer_v3_model_2_seed_000
HIVPR_dimer_23119//HIVPR_dimer_23119_unrelaxed_rank_005_alphaFold2_multimer_v3_model_3_seed_000
....
```

Extracting sequences

```
pdb/seq: 1 name: HIVPR_dimer_23119//HIVPR_dimer_23119_unrelaxed_rank_001_alphaFold2_multimer_v3_model_4_seed_000
pdb/seq: 2 name: HIVPR_dimer_23119//HIVPR_dimer_23119_unrelaxed_rank_002_alphaFold2_multimer_v3_model_1_seed_000
pdb/seq: 3 name: HIVPR_dimer_23119//HIVPR_dimer_23119_unrelaxed_rank_003_alphaFold2_multimer_v3_model_5_seed_000
pdb/seq: 4 name: HIVPR_dimer_23119//HIVPR_dimer_23119_unrelaxed_rank_004_alphaFold2_multimer_v3_model_2_seed_000
pdb/seq: 5 name: HIVPR_dimer_23119//HIVPR_dimer_23119_unrelaxed_rank_005_alphaFold2_multimer_v3_model_3_seed_000
```

```
pdbs
```

```
1 . . . .
[Truncated_Name:1] HIVPR_dime PQITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWPKPMIGGI 50
```

[Truncated_Name:2] HIVPR_dime	PQITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWPKPMIGGI	
[Truncated_Name:3] HIVPR_dime	PQITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWPKPMIGGI	
[Truncated_Name:4] HIVPR_dime	PQITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWPKPMIGGI	
[Truncated_Name:5] HIVPR_dime	PQITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWPKPMIGGI	

	1	50
[Truncated_Name:1] HIVPR_dime	GGFIKVRQYDQILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNFP	51 100
[Truncated_Name:2] HIVPR_dime	GGFIKVRQYDQILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNFP	
[Truncated_Name:3] HIVPR_dime	GGFIKVRQYDQILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNFP	
[Truncated_Name:4] HIVPR_dime	GGFIKVRQYDQILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNFP	
[Truncated_Name:5] HIVPR_dime	GGFIKVRQYDQILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNFP	

	51	100
[Truncated_Name:1] HIVPR_dime	QITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWPKPMIGGIG	101 150
[Truncated_Name:2] HIVPR_dime	QITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWPKPMIGGIG	
[Truncated_Name:3] HIVPR_dime	QITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWPKPMIGGIG	
[Truncated_Name:4] HIVPR_dime	QITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWPKPMIGGIG	
[Truncated_Name:5] HIVPR_dime	QITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWPKPMIGGIG	

	101	150
[Truncated_Name:1] HIVPR_dime	GFIKVRQYDQILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNF	151 198
[Truncated_Name:2] HIVPR_dime	GFIKVRQYDQILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNF	
[Truncated_Name:3] HIVPR_dime	GFIKVRQYDQILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNF	
[Truncated_Name:4] HIVPR_dime	GFIKVRQYDQILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNF	
[Truncated_Name:5] HIVPR_dime	GFIKVRQYDQILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNF	

	151	198

Call:

```
pdabaln(files = pdb_files, fit = TRUE, exefile = "msa")
```

Class:

pdbs, fasta

Alignment dimensions:

5 sequence rows; 198 position columns (198 non-gap, 0 gap)

```
+ attr: xyz, resno, b, chain, id, ali, resid, sse, call
```