

Class 11: Structural Bioinformatics pt.2

Serena Quezada

AlphaFold DB

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

From last class 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
```

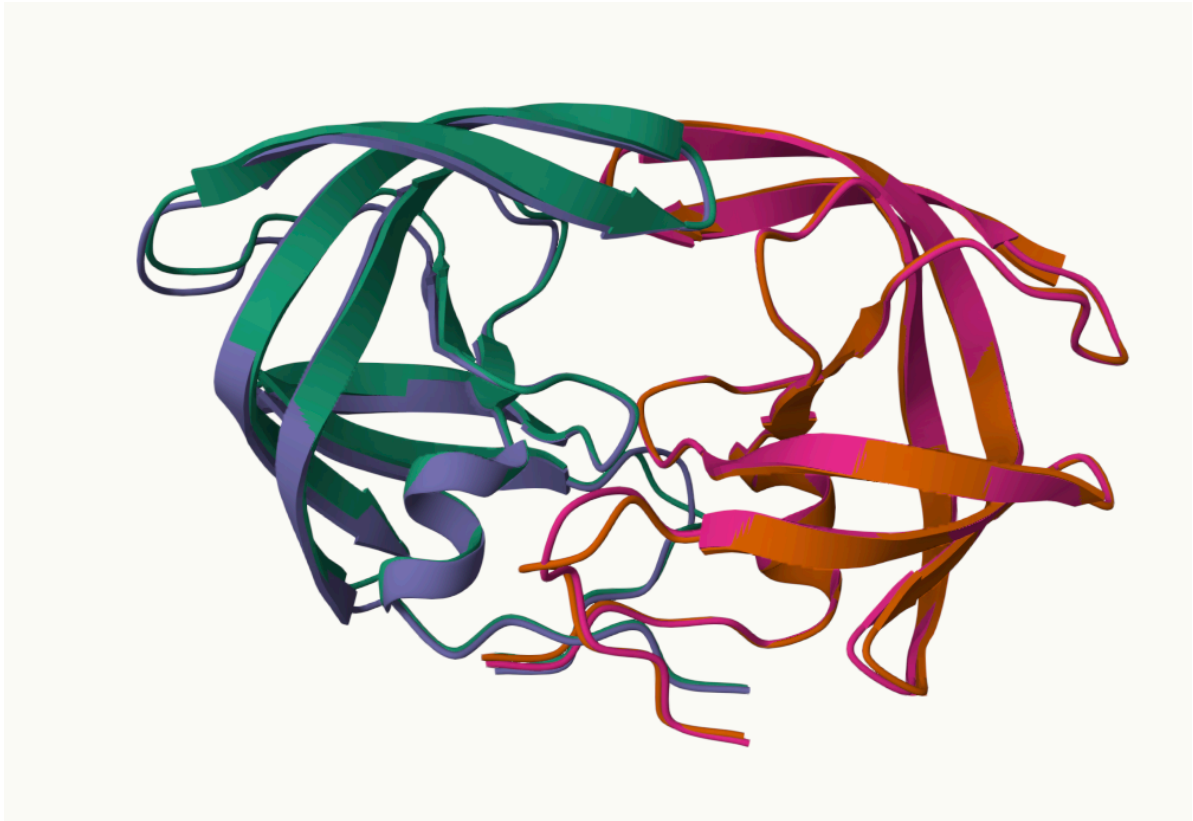
```
[1] 0.1224021
```

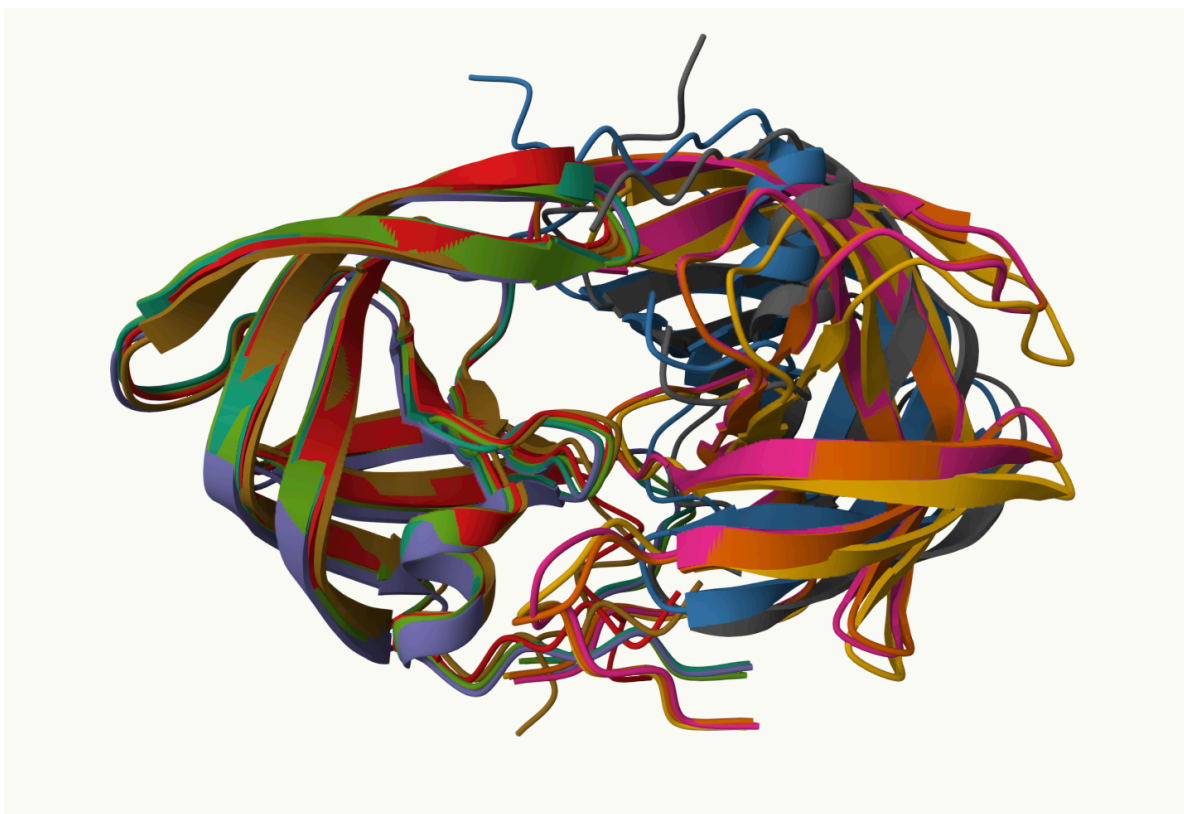
AFBD is attempting to address this gap...

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

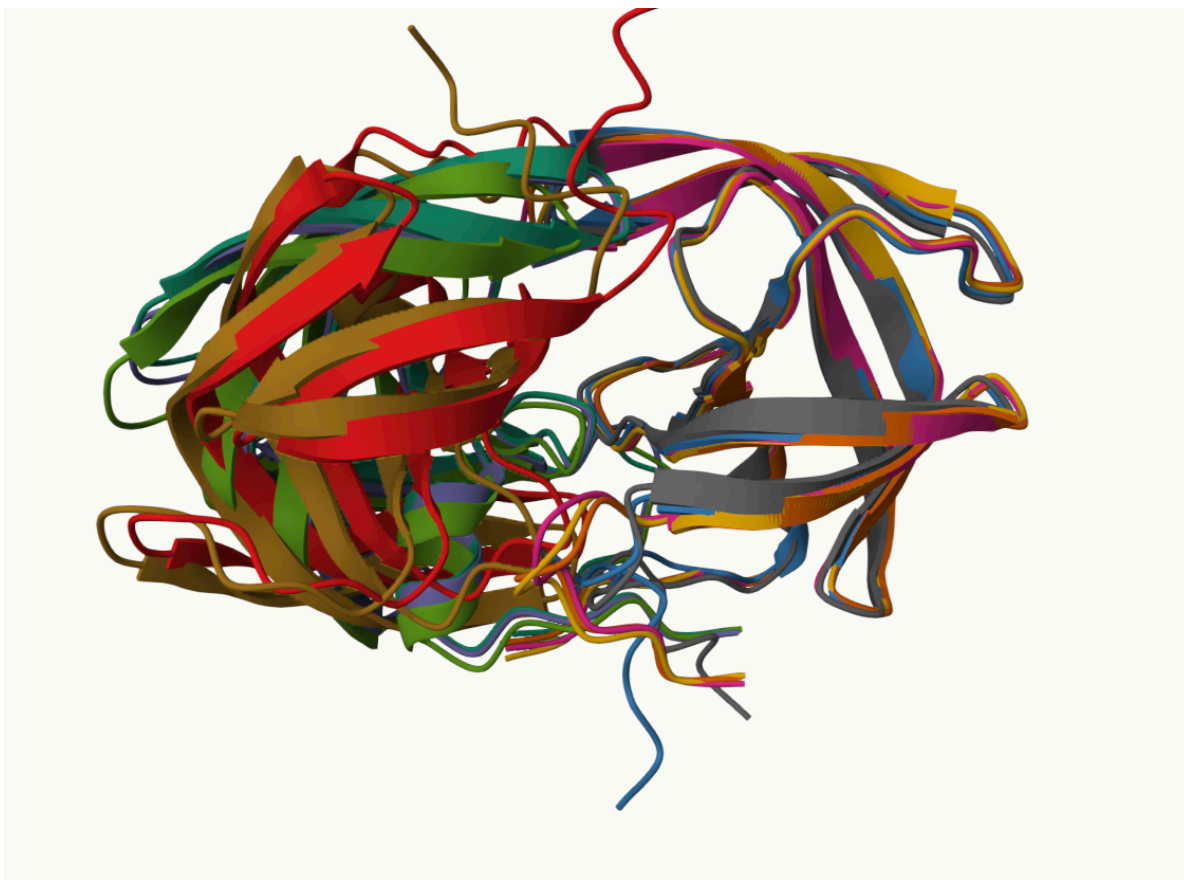
Generating your own structure predictions

picture

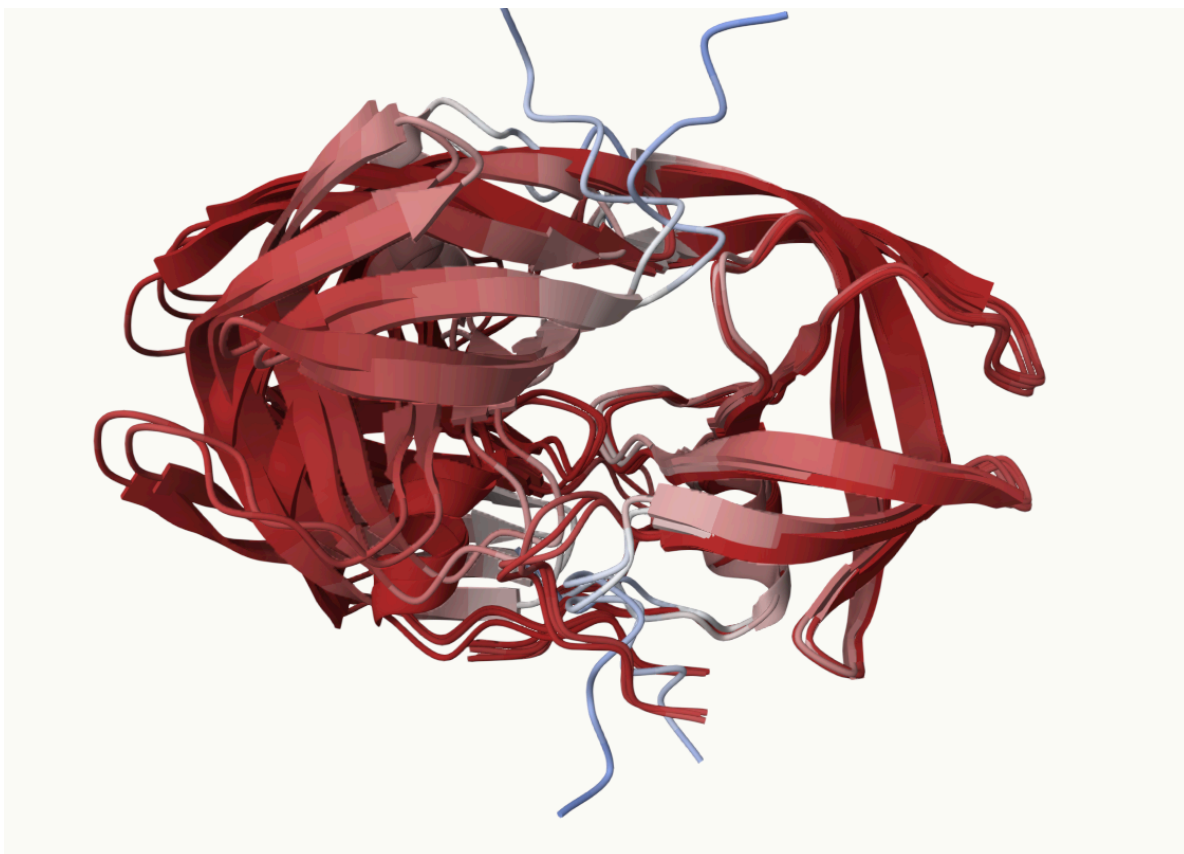




chain A is being superimposed



chain B is being superimposed



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"
```

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

```
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)
```

```
pdbbs <- 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_0
HIVPR_dimer_23119/HIVPR_dimer_23119_unrelaxed_rank_002_alphafold2_multimer_v3_model_1_seed_0
HIVPR_dimer_23119/HIVPR_dimer_23119_unrelaxed_rank_003_alphafold2_multimer_v3_model_5_seed_0
HIVPR_dimer_23119/HIVPR_dimer_23119_unrelaxed_rank_004_alphafold2_multimer_v3_model_2_seed_0
HIVPR_dimer_23119/HIVPR_dimer_23119_unrelaxed_rank_005_alphafold2_multimer_v3_model_3_seed_0
.....
```

Extracting sequences

```
pdb/seq: 1    name: HIVPR_dimer_23119/HIVPR_dimer_23119_unrelaxed_rank_001_alphafold2_multimer
pdb/seq: 2    name: HIVPR_dimer_23119/HIVPR_dimer_23119_unrelaxed_rank_002_alphafold2_multimer
pdb/seq: 3    name: HIVPR_dimer_23119/HIVPR_dimer_23119_unrelaxed_rank_003_alphafold2_multimer
pdb/seq: 4    name: HIVPR_dimer_23119/HIVPR_dimer_23119_unrelaxed_rank_004_alphafold2_multimer
pdb/seq: 5    name: HIVPR_dimer_23119/HIVPR_dimer_23119_unrelaxed_rank_005_alphafold2_multimer
```

```
pdbbs
```

```

1          .          .          .          .          50
[Truncated_Name:1]HIVPR_dime PQITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWKPKMIGGI
[Truncated_Name:2]HIVPR_dime PQITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWKPKMIGGI
[Truncated_Name:3]HIVPR_dime PQITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWKPKMIGGI
[Truncated_Name:4]HIVPR_dime PQITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWKPKMIGGI
[Truncated_Name:5]HIVPR_dime PQITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWKPKMIGGI
*****
1          .          .          .          .          50

51          .          .          .          .          100
[Truncated_Name:1]HIVPR_dime GGFIKVRQYDQILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNFP
[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

101         .          .          .          .          150
[Truncated_Name:1]HIVPR_dime  QITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWPKMIGGIG
[Truncated_Name:2]HIVPR_dime  QITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWPKMIGGIG
[Truncated_Name:3]HIVPR_dime  QITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWPKMIGGIG
[Truncated_Name:4]HIVPR_dime  QITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWPKMIGGIG
[Truncated_Name:5]HIVPR_dime  QITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWPKMIGGIG
*****
101         .          .          .          .          150

151         .          .          .          .          198
[Truncated_Name:1]HIVPR_dime  GFIKVRQYDQILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNF
[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:

```
pdbaln(files = pdb_files, fit = TRUE, exe_file = "msa")
```

Class:

```
pdb, fasta
```

Alignment dimensions:

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

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

```
rd <- rmsd(pdb, fit=T)
```

Warning in rmsd(pdb, fit = T): No indices provided, using the 198 non NA positions

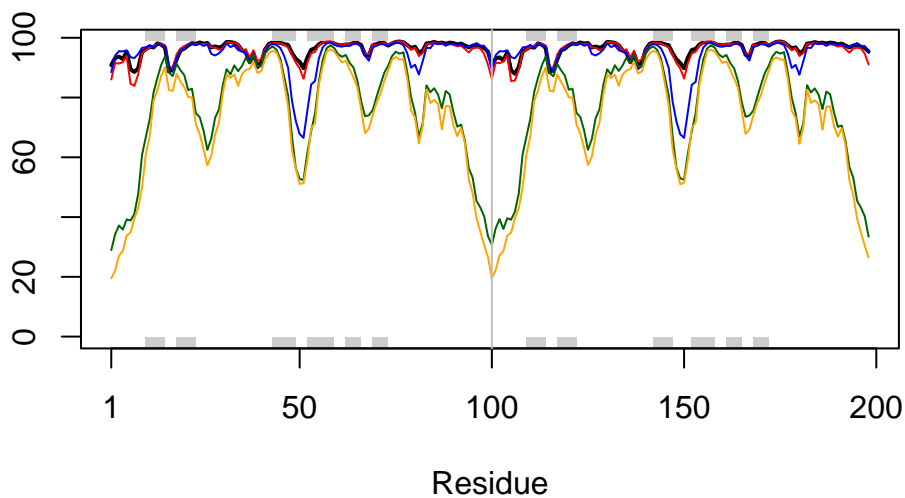
```
range(rd)
```

```
[1] 0.000 14.754
```

```
pdb <- read.pdb("1hsg")
```

Note: Accessing on-line PDB file

```
plotb3(pdb$b[1,], typ="l", lwd=2, sse=pdb)
points(pdb$b[2,], typ="l", col="red")
points(pdb$b[3,], typ="l", col="blue")
points(pdb$b[4,], typ="l", col="darkgreen")
points(pdb$b[5,], typ="l", col="orange")
abline(v=100, col="gray")
```



```
core <- core.find(pdb)
```

```
core size 197 of 198 vol = 9885.515
core size 196 of 198 vol = 6904.596
core size 195 of 198 vol = 1338.769
core size 194 of 198 vol = 1040.692
core size 193 of 198 vol = 951.871
core size 192 of 198 vol = 899.094
core size 191 of 198 vol = 834.741
core size 190 of 198 vol = 771.349
```


core size 189 of 198	vol = 733.076
core size 188 of 198	vol = 697.291
core size 187 of 198	vol = 659.754
core size 186 of 198	vol = 625.285
core size 185 of 198	vol = 589.554
core size 184 of 198	vol = 568.266
core size 183 of 198	vol = 545.027
core size 182 of 198	vol = 512.901
core size 181 of 198	vol = 490.735
core size 180 of 198	vol = 470.278
core size 179 of 198	vol = 450.742
core size 178 of 198	vol = 434.747
core size 177 of 198	vol = 420.349
core size 176 of 198	vol = 406.67
core size 175 of 198	vol = 393.345
core size 174 of 198	vol = 382.406
core size 173 of 198	vol = 372.869
core size 172 of 198	vol = 357.005
core size 171 of 198	vol = 346.579
core size 170 of 198	vol = 337.458
core size 169 of 198	vol = 326.671
core size 168 of 198	vol = 314.962
core size 167 of 198	vol = 304.139
core size 166 of 198	vol = 294.564
core size 165 of 198	vol = 285.661
core size 164 of 198	vol = 278.896
core size 163 of 198	vol = 266.777
core size 162 of 198	vol = 259.006
core size 161 of 198	vol = 247.734
core size 160 of 198	vol = 239.852
core size 159 of 198	vol = 234.975
core size 158 of 198	vol = 230.074
core size 157 of 198	vol = 221.997
core size 156 of 198	vol = 215.632
core size 155 of 198	vol = 206.801
core size 154 of 198	vol = 196.992
core size 153 of 198	vol = 188.547
core size 152 of 198	vol = 182.27
core size 151 of 198	vol = 176.961
core size 150 of 198	vol = 170.72
core size 149 of 198	vol = 166.128
core size 148 of 198	vol = 159.805
core size 147 of 198	vol = 153.775

core size 146 of 198	vol = 149.101
core size 145 of 198	vol = 143.664
core size 144 of 198	vol = 137.145
core size 143 of 198	vol = 132.523
core size 142 of 198	vol = 127.237
core size 141 of 198	vol = 121.579
core size 140 of 198	vol = 116.78
core size 139 of 198	vol = 112.575
core size 138 of 198	vol = 108.175
core size 137 of 198	vol = 105.137
core size 136 of 198	vol = 101.254
core size 135 of 198	vol = 97.379
core size 134 of 198	vol = 92.978
core size 133 of 198	vol = 88.188
core size 132 of 198	vol = 84.032
core size 131 of 198	vol = 81.902
core size 130 of 198	vol = 78.023
core size 129 of 198	vol = 75.276
core size 128 of 198	vol = 73.057
core size 127 of 198	vol = 70.699
core size 126 of 198	vol = 68.976
core size 125 of 198	vol = 66.707
core size 124 of 198	vol = 64.376
core size 123 of 198	vol = 61.145
core size 122 of 198	vol = 59.029
core size 121 of 198	vol = 56.625
core size 120 of 198	vol = 54.369
core size 119 of 198	vol = 51.826
core size 118 of 198	vol = 49.651
core size 117 of 198	vol = 48.19
core size 116 of 198	vol = 46.644
core size 115 of 198	vol = 44.748
core size 114 of 198	vol = 43.288
core size 113 of 198	vol = 41.089
core size 112 of 198	vol = 39.143
core size 111 of 198	vol = 36.468
core size 110 of 198	vol = 34.114
core size 109 of 198	vol = 31.467
core size 108 of 198	vol = 29.445
core size 107 of 198	vol = 27.323
core size 106 of 198	vol = 25.82
core size 105 of 198	vol = 24.149
core size 104 of 198	vol = 22.647

```

core size 103 of 198  vol = 21.068
core size 102 of 198  vol = 19.953
core size 101 of 198  vol = 18.3
core size 100 of 198  vol = 15.723
core size 99 of 198   vol = 14.841
core size 98 of 198   vol = 11.646
core size 97 of 198   vol = 9.435
core size 96 of 198   vol = 7.354
core size 95 of 198   vol = 6.181
core size 94 of 198   vol = 5.667
core size 93 of 198   vol = 4.706
core size 92 of 198   vol = 3.664
core size 91 of 198   vol = 2.77
core size 90 of 198   vol = 2.151
core size 89 of 198   vol = 1.715
core size 88 of 198   vol = 1.15
core size 87 of 198   vol = 0.874
core size 86 of 198   vol = 0.685
core size 85 of 198   vol = 0.528
core size 84 of 198   vol = 0.37
FINISHED: Min vol ( 0.5 ) reached

```

```
core.inds <- print(core, vol = 0.5)
```

```

# 85 positions (cumulative volume <= 0.5 Angstrom^3)
  start end length
1      9  49     41
2     52  95     44

```

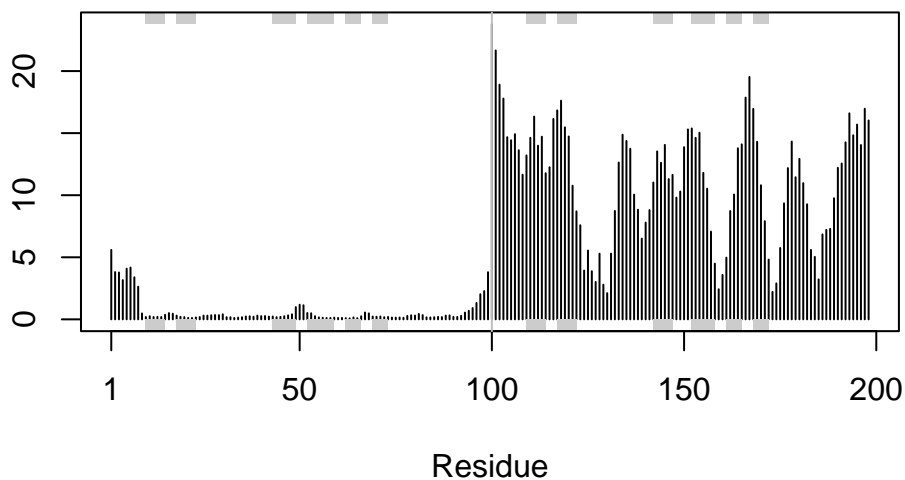
```
xyz <- pdbfit(pdb, core.inds, outpath="corefit_structures")
```

```
rf <- rmsf(xyz)
```

```

plotb3(rf, sse=pdb)
abline(v=100, col="gray", ylab="RMSF")

```



Residue conservation from alignment file

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

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

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

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

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

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

```
[1] 5397 132
```

```
sim <- conserv(aln)

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

