# Comparative Structure Analysis & Introduction to AlphaFold

Andy Hsu

## Setup

To begin with, we need to install some packages for this project. These include `bio3d`, `bio3d-view`, and `msa`. It's important to note that the `msa` package is managed by BioConductor, another package database with a focus on genomics work and adjacent fields. Similarly, the `bio3d-view` package is located on BitBucket, and can be accessed via the `devtools` package.

## Search and Retrieve Structures

Now, we can begin by accessing the sequence of our protein, Adenylate Kinase (AK).

```r
library(bio3d)
```

```
Warning: package 'bio3d' was built under R version 4.3.2
```

```r
aa <- get.seq("1ake_a")
```

```
Warning in get.seq("1ake_a"): Removing existing file: seqs.fasta
```

```
Fetching... Please wait. Done.
```

```r
aa
```

1

```
          1          .          .          .          .          .         60
pdb|1AKE|A    MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMLRAAVKSGSELGKQAKDIMDAGKLVT
          1          .          .          .          .          .         60


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


          121        .          .          .          .          .        180
pdb|1AKE|A    VGRRVHAPSGRVYHVKFNPPKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQMTAPLIG
          121        .          .          .          .          .        180


          181        .          .          .   214
pdb|1AKE|A    YYSKEAEAGNTKYAKVDGTKPVAEVRADLEKILG
          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
```

Next, we run a BLAST search of our sequence to find the corresponding protein.

```
  b <- blast.pdb(aa)
```

```
 Searching ... please wait (updates every 5 seconds) RID = MS46WVJB016
 ....
 Reporting 83 hits
```
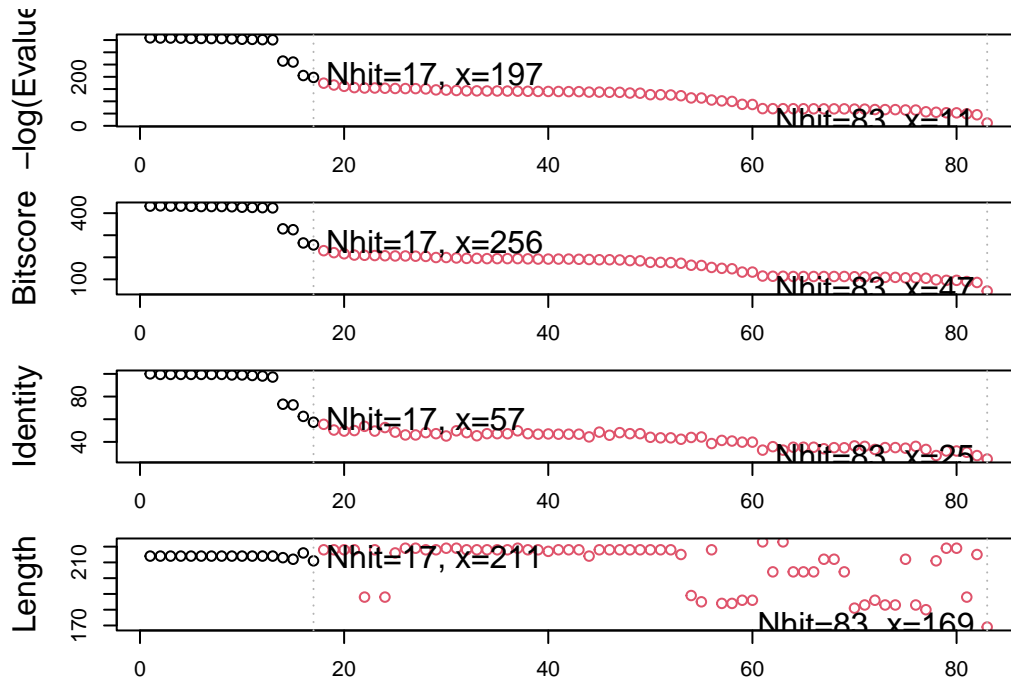
If we plot our results, we can see a summary of our BLAST results by alignment statistics.
We can also list the PDB IDs of some of the top results of our BLAST.

```
  hits <- plot(b)
```

```
 * Possible cutoff values:    197 11
          Yielding Nhits:    17 83

 * Chosen cutoff value of:    197
          Yielding Nhits:    17
```



```r
head(hits$pdb.id)
```

```
[1] "1AKE_A" "8BQF_A" "4X8M_A" "6S36_A" "6RZE_A" "4X8H_A"
```

Before we move on, let's annotate these top results for protein name, organism, method used, etc.

```r
an <- pdb.annotate(hits$pdb.id)
an
```

|        | structureId | chainId | macromoleculeType | chainLength | experimentalTechnique |
|--------|-------------|---------|-------------------|-------------|-----------------------|
| 1AKE_A | 1AKE        | A       | Protein           | 214         | X-ray                 |
| 8BQF_A | 8BQF        | A       | Protein           | 234         | X-ray                 |
| 4X8M_A | 4X8M        | A       | Protein           | 214         | X-ray                 |

3

```
6S36_A        6S36       A        Protein        214            X-ray
6RZE_A        6RZE       A        Protein        214            X-ray
4X8H_A        4X8H       A        Protein        214            X-ray
3HPR_A        3HPR       A        Protein        214            X-ray
1E4V_A        1E4V       A        Protein        214            X-ray
5EJE_A        5EJE       A        Protein        214            X-ray
1E4Y_A        1E4Y       A        Protein        214            X-ray
3X2S_A        3X2S       A        Protein        214            X-ray
6HAP_A        6HAP       A        Protein        214            X-ray
6HAM_A        6HAM       A        Protein        214            X-ray
4K46_A        4K46       A        Protein        214            X-ray
4NP6_A        4NP6       A        Protein        217            X-ray
3GMT_A        3GMT       A        Protein        230            X-ray
4PZL_A        4PZL       A        Protein        242            X-ray
         resolution        scopDomain                                pfam
1AKE_A        2.000 Adenylate kinase Adenylate kinase, active site lid (ADK_lid)
8BQF_A        2.050           <NA> Adenylate kinase, active site lid (ADK_lid)
4X8M_A        2.600           <NA> Adenylate kinase, active site lid (ADK_lid)
6S36_A        1.600           <NA> Adenylate kinase, active site lid (ADK_lid)
6RZE_A        1.690           <NA> Adenylate kinase, active site lid (ADK_lid)
4X8H_A        2.500           <NA> Adenylate kinase, active site lid (ADK_lid)
3HPR_A        2.000           <NA> Adenylate kinase, active site lid (ADK_lid)
1E4V_A        1.850 Adenylate kinase Adenylate kinase, active site lid (ADK_lid)
5EJE_A        1.900           <NA> Adenylate kinase, active site lid (ADK_lid)
1E4Y_A        1.850 Adenylate kinase Adenylate kinase, active site lid (ADK_lid)
3X2S_A        2.800           <NA> Adenylate kinase, active site lid (ADK_lid)
6HAP_A        2.700           <NA> Adenylate kinase, active site lid (ADK_lid)
6HAM_A        2.550           <NA> Adenylate kinase, active site lid (ADK_lid)
4K46_A        2.010           <NA> Adenylate kinase, active site lid (ADK_lid)
4NP6_A        2.004           <NA> Adenylate kinase, active site lid (ADK_lid)
3GMT_A        2.100           <NA> Adenylate kinase, active site lid (ADK_lid)
4PZL_A        2.100           <NA> Adenylate kinase, active site lid (ADK_lid)
          ligandId
1AKE_A           AP5
8BQF_A           AP5
4X8M_A          <NA>
6S36_A CL (3),NA,MG (2)
6RZE_A    NA (3),CL (2)
4X8H_A          <NA>
3HPR_A           AP5
1E4V_A           AP5
5EJE_A        AP5,CO
1E4Y_A           AP5
```

```
3X2S_A    JPY (2),AP5,MG
6HAP_A            AP5
6HAM_A            AP5
4K46_A    ADP,AMP,PO4
4NP6_A           <NA>
3GMT_A        SO4 (2)
4PZL_A     CA,FMT,GOL
                                                                    ligandName
1AKE_A                                    BIS(ADENOSINE)-5'-PENTAPHOSPHATE
8BQF_A                                    BIS(ADENOSINE)-5'-PENTAPHOSPHATE
4X8M_A                                                                 <NA>
6S36_A                      CHLORIDE ION (3),SODIUM ION,MAGNESIUM ION (2)
6RZE_A                                  SODIUM ION (3),CHLORIDE ION (2)
4X8H_A                                                                 <NA>
3HPR_A                                    BIS(ADENOSINE)-5'-PENTAPHOSPHATE
1E4V_A                                    BIS(ADENOSINE)-5'-PENTAPHOSPHATE
5EJE_A                        BIS(ADENOSINE)-5'-PENTAPHOSPHATE,COBALT (II) ION
1E4Y_A                                    BIS(ADENOSINE)-5'-PENTAPHOSPHATE
3X2S_A N-(pyren-1-ylmethyl)acetamide (2),BIS(ADENOSINE)-5'-PENTAPHOSPHATE,MAGNESIUM ION
6HAP_A                                    BIS(ADENOSINE)-5'-PENTAPHOSPHATE
6HAM_A                                    BIS(ADENOSINE)-5'-PENTAPHOSPHATE
4K46_A          ADENOSINE-5'-DIPHOSPHATE,ADENOSINE MONOPHOSPHATE,PHOSPHATE ION
4NP6_A                                                                 <NA>
3GMT_A                                                       SULFATE ION (2)
4PZL_A                                    CALCIUM ION,FORMIC ACID,GLYCEROL
                                              source
1AKE_A                              Escherichia coli
8BQF_A                              Escherichia coli
4X8M_A                              Escherichia coli
6S36_A                              Escherichia coli
6RZE_A                              Escherichia coli
4X8H_A                              Escherichia coli
3HPR_A                         Escherichia coli K-12
1E4V_A                              Escherichia coli
5EJE_A          Escherichia coli O139:H28 str. E24377A
1E4Y_A                              Escherichia coli
3X2S_A        Escherichia coli str. K-12 substr. MDS42
6HAP_A          Escherichia coli O139:H28 str. E24377A
6HAM_A                         Escherichia coli K-12
4K46_A                       Photobacterium profundum
4NP6_A      Vibrio cholerae O1 biovar El Tor str. N16961
3GMT_A              Burkholderia pseudomallei 1710b
4PZL_A Francisella tularensis subsp. tularensis SCHU S4
```

| | title |
|---|---|
| 1AKE_A | STRUCTURE OF THE COMPLEX BETWEEN ADENYLATE KINASE FROM ESCHERICHIA COLI AND THE INHIB |
| 8BQF_A | |
| 4X8M_A | |
| 6S36_A | |
| 6RZE_A | |
| 4X8H_A | |
| 3HPR_A | |
| 1E4V_A | |
| 5EJE_A | Cryst |
| 1E4Y_A | |
| 3X2S_A | |
| 6HAP_A | |
| 6HAM_A | |
| 4K46_A | |
| 4NP6_A | |
| 3GMT_A | |
| 4PZL_A | The cryst |

| | citation | rObserved | rFree |
|---|---|---|---|
| 1AKE_A | Muller, C.W., et al. J Mol Biol (1992) | 0.19600 | NA |
| 8BQF_A | Scheerer, D., et al. Proc Natl Acad Sci U S A (2023) | 0.22073 | 0.25789 |
| 4X8M_A | Kovermann, M., et al. Nat Commun (2015) | 0.24910 | 0.30890 |
| 6S36_A | Rogne, P., et al. Biochemistry (2019) | 0.16320 | 0.23560 |
| 6RZE_A | Rogne, P., et al. Biochemistry (2019) | 0.18650 | 0.23500 |
| 4X8H_A | Kovermann, M., et al. Nat Commun (2015) | 0.19610 | 0.28950 |
| 3HPR_A | Schrank, T.P., et al. Proc Natl Acad Sci U S A (2009) | 0.21000 | 0.24320 |
| 1E4V_A | Muller, C.W., et al. Proteins (1993) | 0.19600 | NA |
| 5EJE_A | Kovermann, M., et al. Proc Natl Acad Sci U S A (2017) | 0.18890 | 0.23580 |
| 1E4Y_A | Muller, C.W., et al. Proteins (1993) | 0.17800 | NA |
| 3X2S_A | Fujii, A., et al. Bioconjug Chem (2015) | 0.20700 | 0.25600 |
| 6HAP_A | Kantaev, R., et al. J Phys Chem B (2018) | 0.22630 | 0.27760 |
| 6HAM_A | Kantaev, R., et al. J Phys Chem B (2018) | 0.20511 | 0.24325 |
| 4K46_A | Cho, Y.-J., et al. To be published | 0.17000 | 0.22290 |
| 4NP6_A | Kim, Y., et al. To be published | 0.18800 | 0.22200 |
| 3GMT_A | Buchko, G.W., et al. Biochem Biophys Res Commun (2010) | 0.23800 | 0.29500 |
| 4PZL_A | Tan, K., et al. To be published | 0.19360 | 0.23680 |

| | rWork | spaceGroup |
|---|---|---|
| 1AKE_A | 0.19600 | P 21 2 21 |
| 8BQF_A | 0.21882 | P 2 21 21 |
| 4X8M_A | 0.24630 | C 1 2 1 |
| 6S36_A | 0.15940 | C 1 2 1 |
| 6RZE_A | 0.18190 | C 1 2 1 |
| 4X8H_A | 0.19140 | C 1 2 1 |

```
3HPR_A 0.20620  P 21 21 2
1E4V_A 0.19600  P 21 2 21
5EJE_A 0.18630  P 21 2 21
1E4Y_A 0.17800   P 1 21 1
3X2S_A 0.20700 P 21 21 21
6HAP_A 0.22370    I 2 2 2
6HAM_A 0.20311       P 43
4K46_A 0.16730 P 21 21 21
4NP6_A 0.18600       P 43
3GMT_A 0.23500   P 1 21 1
4PZL_A 0.19130       P 32
```

Finally for this step, we can fetch and store the structures of all these top results.

```r
files <- get.pdb(hits$pdb.id, path="pdbs", split=TRUE, gzip=TRUE)
```

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/1AKE.pdb exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/8BQF.pdb exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/4X8M.pdb exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/6S36.pdb exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/6RZE.pdb exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/4X8H.pdb exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/3HPR.pdb exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/1E4V.pdb exists. Skipping download

```
Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/5EJE.pdb exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/1E4Y.pdb exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/3X2S.pdb exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/6HAP.pdb exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/6HAM.pdb exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/4K46.pdb exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/4NP6.pdb exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/3GMT.pdb exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/4PZL.pdb exists. Skipping download


  |
  |                                                                    |   0%
  |
  |====                                                                |   6%
  |
  |=======                                                             |  12%
  |
  |===========                                                         |  18%
  |
  |================                                                    |  24%
  |
  |====================                                                |  29%
```
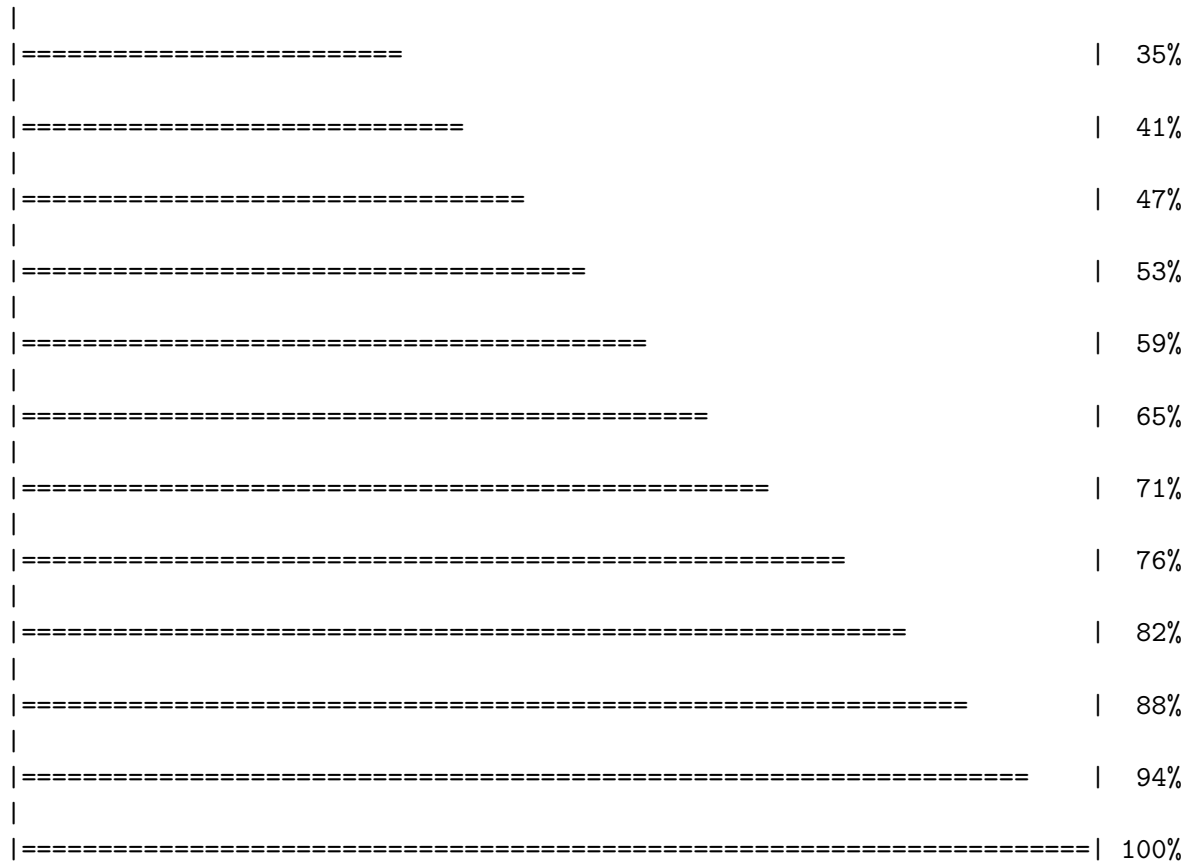
```
  |
  |========================                                            |   35%
  |
  |==========================                                          |   41%
  |
  |==============================                                      |   47%
  |
  |==================================                                  |   53%
  |
  |======================================                              |   59%
  |
  |==========================================                          |   65%
  |
  |==============================================                      |   71%
  |
  |==================================================                  |   76%
  |
  |======================================================              |   82%
  |
  |==========================================================          |   88%
  |
  |==============================================================      |   94%
  |
  |====================================================================|  100%
```

## Align and Superimpose Structures

Now, we can align our files using the `msa` package.

```
pdbs <- pdbaln(files, fit = TRUE, exefile="msa")
```

```
Reading PDB files:
pdbs/split_chain/1AKE_A.pdb
pdbs/split_chain/8BQF_A.pdb
pdbs/split_chain/4X8M_A.pdb
pdbs/split_chain/6S36_A.pdb
pdbs/split_chain/6RZE_A.pdb
pdbs/split_chain/4X8H_A.pdb
pdbs/split_chain/3HPR_A.pdb
pdbs/split_chain/1E4V_A.pdb
pdbs/split_chain/5EJE_A.pdb
```

```
pdbs/split_chain/1E4Y_A.pdb
pdbs/split_chain/3X2S_A.pdb
pdbs/split_chain/6HAP_A.pdb
pdbs/split_chain/6HAM_A.pdb
pdbs/split_chain/4K46_A.pdb
pdbs/split_chain/4NP6_A.pdb
pdbs/split_chain/3GMT_A.pdb
pdbs/split_chain/4PZL_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
.   PDB has ALT records, taking A only, rm.alt=TRUE
..   PDB has ALT records, taking A only, rm.alt=TRUE
.   PDB has ALT records, taking A only, rm.alt=TRUE
..   PDB has ALT records, taking A only, rm.alt=TRUE
..   PDB has ALT records, taking A only, rm.alt=TRUE
....   PDB has ALT records, taking A only, rm.alt=TRUE
.   PDB has ALT records, taking A only, rm.alt=TRUE
....

Extracting sequences

pdb/seq: 1   name: pdbs/split_chain/1AKE_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 2   name: pdbs/split_chain/8BQF_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 3   name: pdbs/split_chain/4X8M_A.pdb
pdb/seq: 4   name: pdbs/split_chain/6S36_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 5   name: pdbs/split_chain/6RZE_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 6   name: pdbs/split_chain/4X8H_A.pdb
pdb/seq: 7   name: pdbs/split_chain/3HPR_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 8   name: pdbs/split_chain/1E4V_A.pdb
pdb/seq: 9   name: pdbs/split_chain/5EJE_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 10   name: pdbs/split_chain/1E4Y_A.pdb
pdb/seq: 11   name: pdbs/split_chain/3X2S_A.pdb
pdb/seq: 12   name: pdbs/split_chain/6HAP_A.pdb
pdb/seq: 13   name: pdbs/split_chain/6HAM_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 14   name: pdbs/split_chain/4K46_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 15   name: pdbs/split_chain/4NP6_A.pdb
```

```
pdb/seq: 16    name: pdbs/split_chain/3GMT_A.pdb
pdb/seq: 17    name: pdbs/split_chain/4PZL_A.pdb
```

pdbs

```
                              1        .         .         .        40
[Truncated_Name:1]1AKE_A.pdb  ----------MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:2]8BQF_A.pdb  ----------MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:3]4X8M_A.pdb  ----------MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:4]6S36_A.pdb  ----------MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:5]6RZE_A.pdb  ----------MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:6]4X8H_A.pdb  ----------MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:7]3HPR_A.pdb  ----------MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:8]1E4V_A.pdb  ----------MRIILLGAPVAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:9]5EJE_A.pdb  ----------MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:10]1E4Y_A.pdb ----------MRIILLGALVAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:11]3X2S_A.pdb ----------MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:12]6HAP_A.pdb ----------MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:13]6HAM_A.pdb ----------MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:14]4K46_A.pdb ----------MRIILLGAPGAGKGTQAQFIMAKFGIPQIS
[Truncated_Name:15]4NP6_A.pdb --------NAMRIILLGAPGAGKGTQAQFIMEKFGIPQIS
[Truncated_Name:16]3GMT_A.pdb ----------MRLILLGAPGAGKGTQANFIKEKFGIPQIS
[Truncated_Name:17]4PZL_A.pdb TENLYFQSNAMRIILLGAPGAGKGTQAKIIEQKYNIAHIS
                                        **^*****  *******  *  *^ *   **
                              1        .         .         .        40

                              41       .         .         .        80
[Truncated_Name:1]1AKE_A.pdb  TGDMLRAAVKSGSELGKQAKDIMDAGKLVTDELVIALVKE
[Truncated_Name:2]8BQF_A.pdb  TGDMLRAAVKSGSELGKQAKDIMDAGKLVTDELVIALVKE
[Truncated_Name:3]4X8M_A.pdb  TGDMLRAAVKSGSELGKQAKDIMDAGKLVTDELVIALVKE
[Truncated_Name:4]6S36_A.pdb  TGDMLRAAVKSGSELGKQAKDIMDAGKLVTDELVIALVKE
[Truncated_Name:5]6RZE_A.pdb  TGDMLRAAVKSGSELGKQAKDIMDAGKLVTDELVIALVKE
[Truncated_Name:6]4X8H_A.pdb  TGDMLRAAVKSGSELGKQAKDIMDAGKLVTDELVIALVKE
[Truncated_Name:7]3HPR_A.pdb  TGDMLRAAVKSGSELGKQAKDIMDAGKLVTDELVIALVKE
[Truncated_Name:8]1E4V_A.pdb  TGDMLRAAVKSGSELGKQAKDIMDAGKLVTDELVIALVKE
[Truncated_Name:9]5EJE_A.pdb  TGDMLRAAVKSGSELGKQAKDIMDACKLVTDELVIALVKE
[Truncated_Name:10]1E4Y_A.pdb TGDMLRAAVKSGSELGKQAKDIMDAGKLVTDELVIALVKE
[Truncated_Name:11]3X2S_A.pdb TGDMLRAAVKSGSELGKQAKDIMDCGKLVTDELVIALVKE
[Truncated_Name:12]6HAP_A.pdb TGDMLRAAVKSGSELGKQAKDIMDAGKLVTDELVIALVRE
[Truncated_Name:13]6HAM_A.pdb TGDMLRAAIKSGSELGKQAKDIMDAGKLVTDEIIIALVKE
[Truncated_Name:14]4K46_A.pdb TGDMLRAAIKAGTELGKQAKSVIDAGQLVSDDIILGLVKE
```

```
[Truncated_Name:15]4NP6_A.pdb      TGDMLRAAIKAGTELGKQAKAVIDAGQLVSDDIILGLIKE
[Truncated_Name:16]3GMT_A.pdb      TGDMLRAAVKAGTPLGVEAKTYMDEGKLVPDSLIIGLVKE
[Truncated_Name:17]4PZL_A.pdb      TGDMIRETIKSGSALGQELKKVLDAGELVSDEFIIKIVKD
                                   ****^*  ^* *^ **    *  ^*    ** *  ^^ ^^^^
```

```
[Truncated_Name:1]1AKE_A.pdb       RIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFD
[Truncated_Name:2]8BQF_A.pdb       RIAQE----GFLLDGFPRTIPQADAMKEAGINVDYVIEFD
[Truncated_Name:3]4X8M_A.pdb       RIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFD
[Truncated_Name:4]6S36_A.pdb       RIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFD
[Truncated_Name:5]6RZE_A.pdb       RIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFD
[Truncated_Name:6]4X8H_A.pdb       RIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFD
[Truncated_Name:7]3HPR_A.pdb       RIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFD
[Truncated_Name:8]1E4V_A.pdb       RIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFD
[Truncated_Name:9]5EJE_A.pdb       RIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFD
[Truncated_Name:10]1E4Y_A.pdb      RIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFD
[Truncated_Name:11]3X2S_A.pdb      RIAQEDSRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFD
[Truncated_Name:12]6HAP_A.pdb      RICQEDSRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFD
[Truncated_Name:13]6HAM_A.pdb      RICQEDSRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFD
[Truncated_Name:14]4K46_A.pdb      RIAQDDCAKGFLLDGFPRTIPQADGLKEVGVVVDYVIEFD
[Truncated_Name:15]4NP6_A.pdb      RIAQADCEKGFLLDGFPRTIPQADGLKEMGINVDYVIEFD
[Truncated_Name:16]3GMT_A.pdb      RLKEADCANGYLFDGFPRTIAQADAMKEAGVAIDYVLEID
[Truncated_Name:17]4PZL_A.pdb      RISKNDCNNGFLLDGVPRTIPQAQELDKLGVNIDYIVEVD
                                   *^         *^* ** **** **   ^     *^ ^**^^* *
```

```
[Truncated_Name:1]1AKE_A.pdb       VPDELIVDRIVGRRVHAPSGRVYHVKFNPPKVEGKDDVTG
[Truncated_Name:2]8BQF_A.pdb       VPDELIVDRIVGRRVHAPSGRVYHVKFNPPKVEGKDDVTG
[Truncated_Name:3]4X8M_A.pdb       VPDELIVDRIVGRRVHAPSGRVYHVKFNPPKVEGKDDVTG
[Truncated_Name:4]6S36_A.pdb       VPDELIVDKIVGRRVHAPSGRVYHVKFNPPKVEGKDDVTG
[Truncated_Name:5]6RZE_A.pdb       VPDELIVDAIVGRRVHAPSGRVYHVKFNPPKVEGKDDVTG
[Truncated_Name:6]4X8H_A.pdb       VPDELIVDRIVGRRVHAPSGRVYHVKFNPPKVEGKDDVTG
[Truncated_Name:7]3HPR_A.pdb       VPDELIVDRIVGRRVHAPSGRVYHVKFNPPKVEGKDDGTG
[Truncated_Name:8]1E4V_A.pdb       VPDELIVDRIVGRRVHAPSGRVYHVKFNPPKVEGKDDVTG
[Truncated_Name:9]5EJE_A.pdb       VPDELIVDRIVGRRVHAPSGRVYHVKFNPPKVEGKDDVTG
[Truncated_Name:10]1E4Y_A.pdb      VPDELIVDRIVGRRVHAPSGRVYHVKFNPPKVEGKDDVTG
[Truncated_Name:11]3X2S_A.pdb      VPDELIVDRIVGRRVHAPSGRVYHVKFNPPKVEGKDDVTG
[Truncated_Name:12]6HAP_A.pdb      VPDELIVDRIVGRRVHAPSGRVYHVKFNPPKVEGKDDVTG
[Truncated_Name:13]6HAM_A.pdb      VPDELIVDRIVGRRVHAPSGRVYHVKFNPPKVEGKDDVTG
[Truncated_Name:14]4K46_A.pdb      VADSVIVERMAGRRAHLASGRTYHNVYNPPKVEGKDDVTG
[Truncated_Name:15]4NP6_A.pdb      VADDVIVERMAGRRAHLPSGRTYHVVYNPPKVEGKDDVTG
```

```
[Truncated_Name:16]3GMT_A.pdb    VPFSEIIERMSGRRTHPASGRTYHVKFNPPKVEGKDDVTG
[Truncated_Name:17]4PZL_A.pdb    VADNLLIERITGRRIHPASGRTYHTKFNPPKVADKDDVTG
                                 *     ^^^ ^ *** *  *** **  ^*****  *** **
                              121        .       .        .       160


                              161        .       .        .       200
[Truncated_Name:1]1AKE_A.pdb     EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:2]8BQF_A.pdb     EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:3]4X8M_A.pdb     EELTTRKDDQEETVRKRLVEWHQMTAPLIGYYSKEAEAGN
[Truncated_Name:4]6S36_A.pdb     EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:5]6RZE_A.pdb     EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:6]4X8H_A.pdb     EELTTRKDDQEETVRKRLVEYHQMTAALIGYYSKEAEAGN
[Truncated_Name:7]3HPR_A.pdb     EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:8]1E4V_A.pdb     EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:9]5EJE_A.pdb     EELTTRKDDQEECVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:10]1E4Y_A.pdb    EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:11]3X2S_A.pdb    EELTTRKDDQEETVRKRLCEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:12]6HAP_A.pdb    EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:13]6HAM_A.pdb    EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:14]4K46_A.pdb    EDLVIREDDKEETVLARLGVYHNQTAPLIAYYGKEAEAGN
[Truncated_Name:15]4NP6_A.pdb    EDLVIREDDKEETVRARLNVYHTQTAPLIEYYGKEAAAGK
[Truncated_Name:16]3GMT_A.pdb    EPLVQRDDDKEETVKKRLDVYEAQTKPLITYYGDWARRGA
[Truncated_Name:17]4PZL_A.pdb    EPLITRTDDNEDTVKQRLSVYHAQTAKLIDFYRNFSSTNT
                                 * *   * ** *^ *   **   ^    *   ** ^*
                              161        .       .        .       200


                              201        .       .    227
[Truncated_Name:1]1AKE_A.pdb     T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:2]8BQF_A.pdb     T--KYAKVDGTKPVAEVRADLEKIL--
[Truncated_Name:3]4X8M_A.pdb     T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:4]6S36_A.pdb     T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:5]6RZE_A.pdb     T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:6]4X8H_A.pdb     T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:7]3HPR_A.pdb     T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:8]1E4V_A.pdb     T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:9]5EJE_A.pdb     T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:10]1E4Y_A.pdb    T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:11]3X2S_A.pdb    T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:12]6HAP_A.pdb    T--KYAKVDGTKPVCEVRADLEKILG-
[Truncated_Name:13]6HAM_A.pdb    T--KYAKVDGTKPVCEVRADLEKILG-
[Truncated_Name:14]4K46_A.pdb    T--QYLKFDGTKAVAEVSAELEKALA-
[Truncated_Name:15]4NP6_A.pdb    T--QYLKFDGTKQVSEVSADIAKALA-
[Truncated_Name:16]3GMT_A.pdb    E-------NGLKAPA-----YRKISG-
```

13

```
[Truncated_Name:17]4PZL_A.pdb   KIPKYIKINGDQAVEKVSQDIFDQLNK
                                          *
                           201         .        .       227
```

```
Call:
  pdbaln(files = files, fit = TRUE, exefile = "msa")

Class:
  pdbs, fasta

Alignment dimensions:
  17 sequence rows; 227 position columns (199 non-gap, 28 gap)

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

## PCA

Next, we will perform PCA on the alignment to find the relationships between the structures.

```
pc.xray <- pca(pdbs)
plot(pc.xray)
```

And that's it for this analysis of a couple of homologous structures.

## Analysis of AlphaFold Predictions

This next section will focus on analyzing structure predictions of a specific dimer found by AlphaFold. The results have already been loaded into the project folder. The following code will store the names of PDB files in the results as a vector.

```
results_dir <- "HIVPrDi_23119.result/HIVPrDi_23119"
pdb_files <- list.files(path=results_dir,
                        pattern="*.pdb",
                        full.names = TRUE)
pdb_files
```

```
[1] "HIVPrDi_23119.result/HIVPrDi_23119/HIVPrDi_23119_unrelaxed_rank_001_alphafold2_multimer_
[2] "HIVPrDi_23119.result/HIVPrDi_23119/HIVPrDi_23119_unrelaxed_rank_002_alphafold2_multimer_
[3] "HIVPrDi_23119.result/HIVPrDi_23119/HIVPrDi_23119_unrelaxed_rank_003_alphafold2_multimer_
[4] "HIVPrDi_23119.result/HIVPrDi_23119/HIVPrDi_23119_unrelaxed_rank_004_alphafold2_multimer_
[5] "HIVPrDi_23119.result/HIVPrDi_23119/HIVPrDi_23119_unrelaxed_rank_005_alphafold2_multimer_
```

Next, we use Bio3D to align the sequences. We can view the resulting alignment to check that everything is in order.

```
library(bio3d)
pdbs <- pdbaln(pdb_files, fit=TRUE, exefile="msa")
```

```
Reading PDB files:
HIVPrDi_23119.result/HIVPrDi_23119/HIVPrDi_23119_unrelaxed_rank_001_alphafold2_multimer_v3_mo
HIVPrDi_23119.result/HIVPrDi_23119/HIVPrDi_23119_unrelaxed_rank_002_alphafold2_multimer_v3_mo
HIVPrDi_23119.result/HIVPrDi_23119/HIVPrDi_23119_unrelaxed_rank_003_alphafold2_multimer_v3_mo
HIVPrDi_23119.result/HIVPrDi_23119/HIVPrDi_23119_unrelaxed_rank_004_alphafold2_multimer_v3_mo
HIVPrDi_23119.result/HIVPrDi_23119/HIVPrDi_23119_unrelaxed_rank_005_alphafold2_multimer_v3_mo
.....

Extracting sequences

pdb/seq: 1   name: HIVPrDi_23119.result/HIVPrDi_23119/HIVPrDi_23119_unrelaxed_rank_001_alpha
pdb/seq: 2   name: HIVPrDi_23119.result/HIVPrDi_23119/HIVPrDi_23119_unrelaxed_rank_002_alpha
pdb/seq: 3   name: HIVPrDi_23119.result/HIVPrDi_23119/HIVPrDi_23119_unrelaxed_rank_003_alpha
pdb/seq: 4   name: HIVPrDi_23119.result/HIVPrDi_23119/HIVPrDi_23119_unrelaxed_rank_004_alpha
pdb/seq: 5   name: HIVPrDi_23119.result/HIVPrDi_23119/HIVPrDi_23119_unrelaxed_rank_005_alpha
```

15

```
pdbs
```

```
                                          1         .         .         .         .         50
[Truncated_Name:1]HIVPrDi_23    PQITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWKPKMIGGI
[Truncated_Name:2]HIVPrDi_23    PQITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWKPKMIGGI
[Truncated_Name:3]HIVPrDi_23    PQITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWKPKMIGGI
[Truncated_Name:4]HIVPrDi_23    PQITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWKPKMIGGI
[Truncated_Name:5]HIVPrDi_23    PQITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWKPKMIGGI
                                **************************************************
                                          1         .         .         .         .         50


                                          51        .         .         .         .         100
[Truncated_Name:1]HIVPrDi_23     GGFIKVRQYDQILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNFP
[Truncated_Name:2]HIVPrDi_23     GGFIKVRQYDQILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNFP
[Truncated_Name:3]HIVPrDi_23     GGFIKVRQYDQILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNFP
[Truncated_Name:4]HIVPrDi_23     GGFIKVRQYDQILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNFP
[Truncated_Name:5]HIVPrDi_23     GGFIKVRQYDQILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNFP
                                 **************************************************
                                          51        .         .         .         .         100


                                          101       .         .         .         .         150
[Truncated_Name:1]HIVPrDi_23     QITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWKPKMIGGIG
[Truncated_Name:2]HIVPrDi_23     QITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWKPKMIGGIG
[Truncated_Name:3]HIVPrDi_23     QITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWKPKMIGGIG
[Truncated_Name:4]HIVPrDi_23     QITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWKPKMIGGIG
[Truncated_Name:5]HIVPrDi_23     QITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWKPKMIGGIG
                                 **************************************************
                                          101       .         .         .         .         150


                                          151       .         .         .         .         198
[Truncated_Name:1]HIVPrDi_23     GFIKVRQYDQILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNF
[Truncated_Name:2]HIVPrDi_23     GFIKVRQYDQILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNF
[Truncated_Name:3]HIVPrDi_23     GFIKVRQYDQILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNF
[Truncated_Name:4]HIVPrDi_23     GFIKVRQYDQILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNF
[Truncated_Name:5]HIVPrDi_23     GFIKVRQYDQILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNF
                                 ***********************************************
                                          151       .         .         .         .         198
```

```
Call:
  pdbaln(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
```

We can also calculate the RMSD to find relative distance between the structures.

```r
rd <- rmsd(pdbs)
```

```
Warning in rmsd(pdbs): No indices provided, using the 198 non NA positions
```
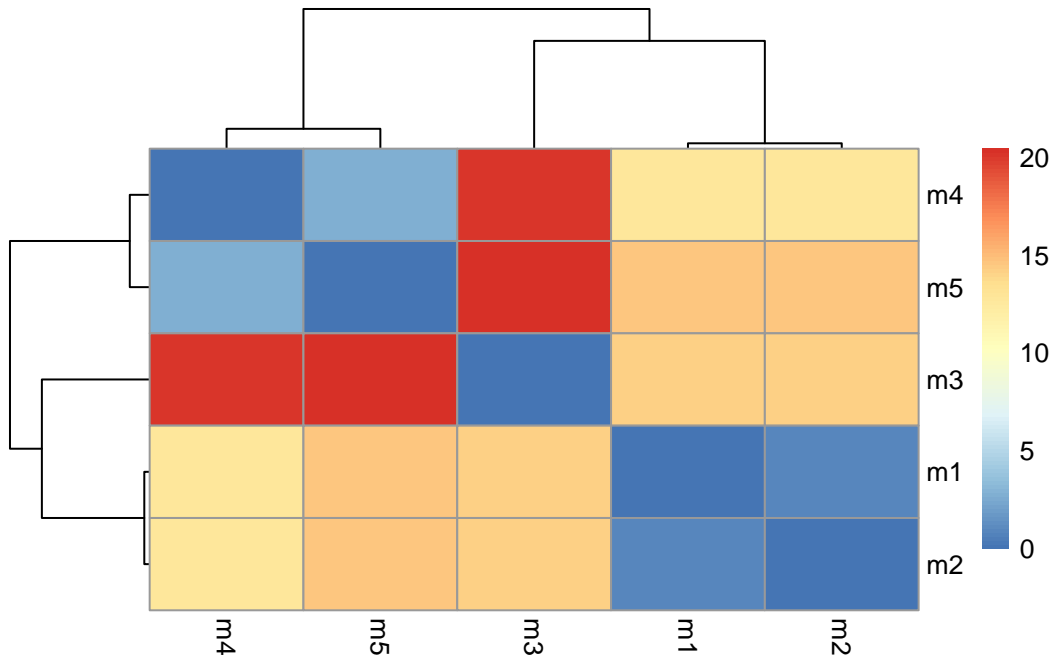
```r
range(rd)
```

```
[1]  0.000 20.431
```

Now, we can use the following code to plot a heat map of our values.

```r
library(pheatmap)
```

```
Warning: package 'pheatmap' was built under R version 4.3.2
```

```r
colnames(rd) <- paste0("m",1:5)
rownames(rd) <- paste0("m",1:5)
pheatmap(rd)
```

A plot of pLDDT values across all models is also easily created.

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

```
Note: Accessing on-line PDB file
```

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

To improve our superpositions, we can employ the `core.find()` function as follows.

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

```
core size 197 of 198  vol = 6154.839
core size 196 of 198  vol = 5399.676
core size 195 of 198  vol = 5074.795
core size 194 of 198  vol = 4802.518
core size 193 of 198  vol = 4520.256
core size 192 of 198  vol = 4305.362
core size 191 of 198  vol = 4089.792
core size 190 of 198  vol = 3886.145
core size 189 of 198  vol = 3758.321
core size 188 of 198  vol = 3620.18
core size 187 of 198  vol = 3496.698
core size 186 of 198  vol = 3389.985
core size 185 of 198  vol = 3320.114
core size 184 of 198  vol = 3258.683
core size 183 of 198  vol = 3208.591
core size 182 of 198  vol = 3156.736
core size 181 of 198  vol = 3141.668
core size 180 of 198  vol = 3136.574
```

```
core size 179 of 198   vol = 3155.52
core size 178 of 198   vol = 3185.362
core size 177 of 198   vol = 3204.487
core size 176 of 198   vol = 3211.978
core size 175 of 198   vol = 3234.993
core size 174 of 198   vol = 3244.062
core size 173 of 198   vol = 3237.845
core size 172 of 198   vol = 3218.77
core size 171 of 198   vol = 3180.743
core size 170 of 198   vol = 3130.369
core size 169 of 198   vol = 3067.881
core size 168 of 198   vol = 2989.546
core size 167 of 198   vol = 2928.272
core size 166 of 198   vol = 2851.193
core size 165 of 198   vol = 2780.877
core size 164 of 198   vol = 2708.433
core size 163 of 198   vol = 2636.516
core size 162 of 198   vol = 2563.25
core size 161 of 198   vol = 2478.024
core size 160 of 198   vol = 2404.793
core size 159 of 198   vol = 2330.997
core size 158 of 198   vol = 2250.477
core size 157 of 198   vol = 2159.432
core size 156 of 198   vol = 2070.759
core size 155 of 198   vol = 1983.579
core size 154 of 198   vol = 1917.913
core size 153 of 198   vol = 1842.556
core size 152 of 198   vol = 1775.398
core size 151 of 198   vol = 1695.133
core size 150 of 198   vol = 1632.173
core size 149 of 198   vol = 1570.391
core size 148 of 198   vol = 1497.238
core size 147 of 198   vol = 1434.802
core size 146 of 198   vol = 1367.706
core size 145 of 198   vol = 1302.596
core size 144 of 198   vol = 1251.985
core size 143 of 198   vol = 1207.976
core size 142 of 198   vol = 1167.112
core size 141 of 198   vol = 1118.27
core size 140 of 198   vol = 1081.664
core size 139 of 198   vol = 1029.75
core size 138 of 198   vol = 981.766
core size 137 of 198   vol = 944.446
```

```
core size 136 of 198   vol = 899.224
core size 135 of 198   vol = 859.402
core size 134 of 198   vol = 814.694
core size 133 of 198   vol = 771.862
core size 132 of 198   vol = 733.807
core size 131 of 198   vol = 702.053
core size 130 of 198   vol = 658.757
core size 129 of 198   vol = 622.574
core size 128 of 198   vol = 578.29
core size 127 of 198   vol = 543.07
core size 126 of 198   vol = 510.934
core size 125 of 198   vol = 481.595
core size 124 of 198   vol = 464.672
core size 123 of 198   vol = 451.721
core size 122 of 198   vol = 430.417
core size 121 of 198   vol = 409.141
core size 120 of 198   vol = 378.942
core size 119 of 198   vol = 348.325
core size 118 of 198   vol = 324.738
core size 117 of 198   vol = 312.394
core size 116 of 198   vol = 300.89
core size 115 of 198   vol = 279.976
core size 114 of 198   vol = 263.434
core size 113 of 198   vol = 250.263
core size 112 of 198   vol = 229.592
core size 111 of 198   vol = 209.929
core size 110 of 198   vol = 196.379
core size 109 of 198   vol = 180.628
core size 108 of 198   vol = 167.088
core size 107 of 198   vol = 155.875
core size 106 of 198   vol = 142.595
core size 105 of 198   vol = 128.924
core size 104 of 198   vol = 114.054
core size 103 of 198   vol = 100.936
core size 102 of 198   vol = 90.431
core size 101 of 198   vol = 81.972
core size 100 of 198   vol = 74.017
core size 99 of 198   vol = 66.855
core size 98 of 198   vol = 59.525
core size 97 of 198   vol = 52.263
core size 96 of 198   vol = 43.699
core size 95 of 198   vol = 35.813
core size 94 of 198   vol = 28.888
```

```
core size 93 of 198  vol = 20.692
core size 92 of 198  vol = 14.975
core size 91 of 198  vol = 9.146
core size 90 of 198  vol = 5.232
core size 89 of 198  vol = 3.53
core size 88 of 198  vol = 2.657
core size 87 of 198  vol = 1.998
core size 86 of 198  vol = 1.333
core size 85 of 198  vol = 1.141
core size 84 of 198  vol = 1.012
core size 83 of 198  vol = 0.891
core size 82 of 198  vol = 0.749
core size 81 of 198  vol = 0.618
core size 80 of 198  vol = 0.538
core size 79 of 198  vol = 0.479
FINISHED: Min vol ( 0.5 ) reached
```

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

```
# 80 positions (cumulative volume <= 0.5 Angstrom^3)
  start end length
1    10  25     16
2    27  48     22
3    53  94     42
```

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

This code generates a collection of PDB files at a directory in the project folder with the improved superpositions, which can be viewed in Mol*. Our updated RMSD heatmap is displayed below.

```r
rd <- rmsd(xyz)
```

```
Warning in rmsd(xyz): No indices provided, using the 198 non NA positions
```

```r
colnames(rd) <- paste0("m",1:5)
rownames(rd) <- paste0("m",1:5)
pheatmap(rd)
```

An RMSF plot can also be created to compare differences in the chains.

```
rf <- rmsf(xyz)

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

## Visualizing Predicted Alignment Error

AlphaFold also provides files documenting the Predicted Alignment Error, located in JSON files that we can access via the `jsonlite` package.

```r
library(jsonlite)
```

```
Warning: package 'jsonlite' was built under R version 4.3.2
```
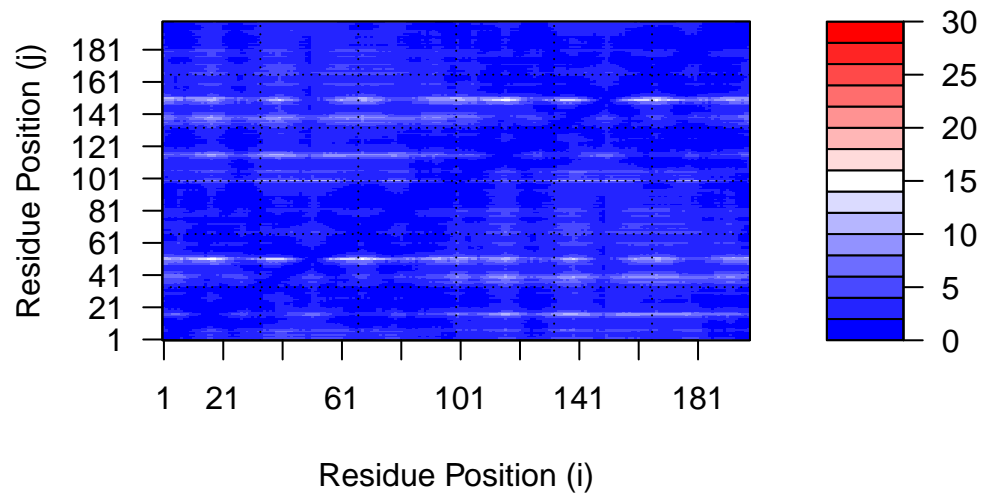
```r
pae_files <- list.files(path=results_dir,
                        pattern=".*model.*\\.json",
                        full.names = TRUE)
pae1 <- read_json(pae_files[1],simplifyVector = TRUE)
pae5 <- read_json(pae_files[5],simplifyVector = TRUE)
```

We can plot these PAE values using the Bio3D package.

```r
plot.dmat(pae1$pae,
          xlab="Residue Position (i)",
          ylab="Residue Position (j)",
          grid.col = "black",
```

```
          zlim=c(0,30))
```



```
plot.dmat(pae5$pae,
          xlab="Residue Position (i)",
          ylab="Residue Position (j)",
          grid.col = "black",
          zlim=c(0,30))
```

## Measuring Residue Conservation

Another thing AlphaFold allows us to do is a measure of residue conservation, derived from the sequences stored in a .a3m file.
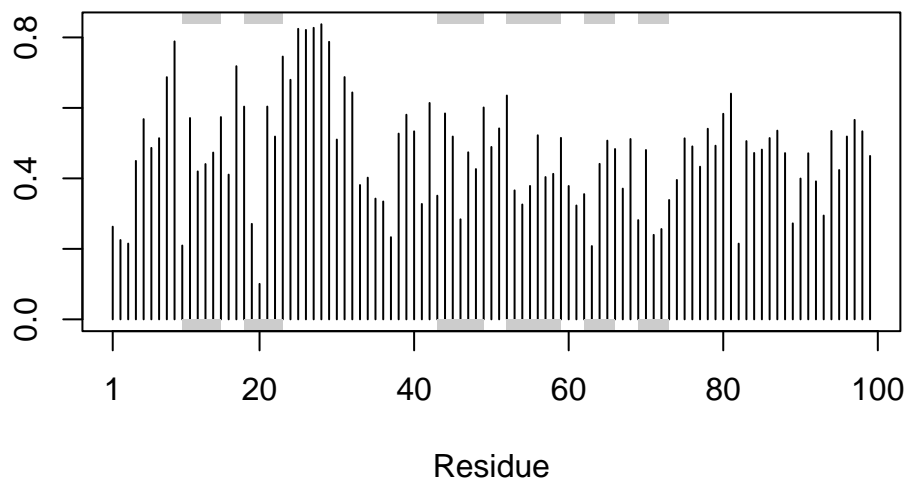
```
aln_file <- list.files(path=results_dir,
                       pattern=".a3m$",
                       full.names = TRUE)
aln <- read.fasta(aln_file[1], to.upper = TRUE)
```

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

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

We can plot the resulting residue conservations to visualize them.

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

Finally, we can create a pdb file to view these results in Mol*.

```
m1.pdb <- read.pdb(pdb_files[1])
occ <- vec2resno(c(sim[1:99], sim[1:99]), m1.pdb$atom$resno)
write.pdb(m1.pdb, o=occ, file="m1_conserv.pdb")
```

And that's all for the structure prediction of a protein from the sequence using AlphaFold.