PRACTICAL 6

**Question 1.** Using AL2CO server (http://prodata.swmed.edu/al2co/al2co.php), obtain the positional conservation scores from multiple sequence alignment (**MSA**) of given set of protein sequences (set1 and set2) using the methods given below:

(i) Unweighted frequency and entropy-based measure
(ii) Unweighted frequency and variance-based measure
(iii) Unweighted frequency and sum of pairs measure
(iv) Weighted frequency and variance-based measure
(v) Normalize the scores obtained with (i)

Sequences:
Set 1: P69905, P01946, P01942, P01966, P01958, P01959, P01965, P06635, P60529, P80043, P01980
Set 2: TPIS_HUMAN, TPIS_YEAST, TPIS_GRAGA, TPIS_TRYCR, TPIS_MAIZE, TPIS_MOUSE, TPIS_DROME, TPIS_RABIT, TPIS_CAEEL

**Solution.** Firstly, I need to compute the Multiple Sequence Alignment of the given sets of sequences. It is because it is the aligned sequences itself that is passed as input to the **AL2CO** server, for it to compute conservation scores. So, below is an image of the input given to the **Clustal Omega** to compute the Multiple Sequence Alignment of set 1 and set 2. The outputs of the Multiple Sequence Alignment have been enclosed ahead of these input images.



FIGURE 1. Input sequences to the CLUSTAL OMEGA from set 1



FIGURE 2. Input sequences to the CLUSTAL OMEGA from set 2

Figure 3. Output MSA of the CLUSTAL OMEGA from set 1

- The above image shows the multiple sequence alignment for the proteins enclosed in set 1. They are a total of 11 in number.
- The below image shows the multiple sequence alignment for the proteins enclosed in set 2. They are a total of 9 in number.
- These aligned sequences are those that will be passed to the AL2CO server tool to compute the conservation scores for these sets of proteins.

```
CLUSTAL O(1.2.4) multiple sequence alignment


TPIS_TRYCR      MASKPQPIAAANWKCNGSESLLVPLIETLNAATFDHD--VQCVVAPTFLHIPMTKARLTN      58
TPIS_GRAGA      -----------NWKCNLSKADIAELVSAFNAAPPIDAAHVQVVVAPPAVYLDSTRQAL-R      48
TPIS_YEAST      --MARTFFVGGNFKLNGSKQSIKEIVERLNTASIPEN--VEVVICPPATYLDYSVSLVKK      56
TPIS_MAIZE      --MGRKFFVGGNWKCNGTTDQVEKIVKTLNEGQVPPSDVVEVVVSPPYVFLPVVKSQL-R      57
TPIS_CAEEL      --MTRKFFVGGNWKMNGDYASVDGIVTFLNASADNSS--VDVVVAPPAPYLAYAKSKL-K      55
TPIS_DROME      --MSRKFCVGGNWKMNGDQKSIAEIAKTLSSAALDPN--TEVVIGCPAIYLMYARNLL-P      55
TPIS_HUMAN      MAPSRKFFVGGNWKMNGRKQSLGELIGTLNAAKVPAD--TEVVCAPPTAYIDFARQKL-D      57
TPIS_RABIT      MAPSRKFFVGGNWKMNGRKKNLGELITTLNAAKVPAD--TEVVCAPPTAYIDFARQKL-D      57
TPIS_MOUSE      MAPTRKFFVGGNWKMNGRKKCLGELICTLNAANVPAG--TEVVCAPPTAYIDFARQKL-D      57
                         *:*  *      :  :   :.. .       .: *       .:       :


TPIS_TRYCR      PKFQIAAQNAI-TRSGAFTGEVSLQILKDYGISWVVLGHSERRLY--YGETNEIVAEKVA     115
TPIS_GRAGA      ADFDTSAQNAWISKGGAFTGELDAAMVKDVGAEWVILGHSERRHIAQLKESDHTIAMKAA     108
TPIS_YEAST      PQVTVGAQNAYLKASGAFTGENSVDQIKDVGAKWVILGHSERRSY--FHEDDKFIADKTK     114
TPIS_MAIZE      QEFHVAAQNCWVKKGGAFTGEVSAEMLVNLGVPWVILGHSERRAL--LGESNEFVGDKVA     115
TPIS_CAEEL      AGVLVAAQNCYKVPKGAFTGEISPAMIKDLGLEWVILGHSERRHV--FGESDALIAEKTV     113
TPIS_DROME      CELGLAGQNAYKVAKGAFTGEISPAMLKDIGADWVILGHSERRAI--FGESDALIAEKAE     113
TPIS_HUMAN      PKIAVAAQNCYKVTNGAFTGEISPGMIKDCGATWVVLGHSERRHV--FGESDELIGQKVA     115
TPIS_RABIT      PKIAVAAQNCYKVTNGAFTGEISPGMIKDCGATWVVLGHSERRHV--FGESDELIGQKVA     115
TPIS_MOUSE      PKIAVAAQNCYKVTNGAFTGEISPGMIKDLGATWVVLGHSERRHV--FGESDELIGQKVS     115
                 .   ..**.      ******  .   :  :  *  **:*******    * :   :. *.


TPIS_TRYCR      QACA-AGFHVIVCVGETNEEREAGRTAAVVLTQLAAVAQKLSKEAWSRVVIAYEPVWAIG     174
TPIS_GRAGA      YALQHASLGVIYCIGELLEERESGQTIAVCERQLQALSDAI--SDWSDVVIAYEPVWAIG     166
TPIS_YEAST      FALG-QGVGVILCIGETLEEKKAGKTLDVVERQLNAVLEEV--KDWTNVVVAYEPVWAIG     171
TPIS_MAIZE      YALS-QGLKVIACVGETLEQREAGSTMDVVAAQTKAIAEKI--KDWSNVVVAYEPVWAIG     172
TPIS_CAEEL      HALE-AGIKVVFCIGEKLEEREAGHTKDVNFRQLQAIVDKG--VSWENIVIAYEPVWAIG     170
TPIS_DROME      HALA-EGLKVIACIGETLEEREAGKTNEVVARQMCAYAQKI--KDWKNVVVAYEPVWAIG     170
TPIS_HUMAN      HALA-EGLGVIACIGEKLDEREAGITEKVVFEQTKVIADNV--KDWSKVVLAYEPVWAIG     172
TPIS_RABIT      HALS-EGLGVIACIGEKLDEREAGITEKVVFEQTKVIADNV--KDWSKVVLAYEPVWAIG     172
TPIS_MOUSE      HALA-EGLGVIACIGEKLDEREAGITEKVVFEQTKVIADNV--KDWSKVVLAYEPVWAIG     172
                 *      .. *: *:**  ::::::* *  *    *  .  :       *  :*:********


TPIS_TRYCR      TGKVATPQQAQEVHELLRRWVRSKLGTDIAAQLRILYGGSVTAKNARTLYQMRDINGFLV     234
TPIS_GRAGA      TGKVATPEQAEQVHEAVRAWLANNVSPQVAASTRILYGGSVSPANCESLAKQPNIDGFLV     226
TPIS_YEAST      TGLAATPEDAQDIHASIRKFLASKLGDKAASELRILYGGSANGSNAVTFKDKADVDGFLV     231
TPIS_MAIZE      TGKVATPAQAQEVHASLRDWLKTNASPEVAESTRIIYGGSVTAANCKELAAQPDVDGFLV     232
TPIS_CAEEL      TGKTASGEQAQEVHEWIRAFLKEKVSPAVADATRIIYGGSVTADNAAELGKKPDIDGFLV     230
TPIS_DROME      TGQTATPDQAQEVHAFLRQWLSDNISKEVSASLRIQYGGSVTAANAKELAKKPDIDGFLV     230
TPIS_HUMAN      TGKTATPQQAQEVHEKLRGWLKSNVSDAVAQSTRIIYGGSVTGATCKELASQPDVDGFLV     232
TPIS_RABIT      TGKTATPQQAQEVHEKLRGWLKSNVSDAVAQSTRIIYGGSVTGATCKELASQPDVDGFLV     232
TPIS_MOUSE      TGKTATPQQAQEVHEKLRGWLKSNVNDGVAQSTRIIYGGSVTGATCKELASQPDVDGFLV     232
                **  .*:  :*:::*  :* ::  :  .   :    **  ****..  ..  :      :::****


TPIS_TRYCR      GGASLKPEFVEIIEATK-------         251
TPIS_GRAGA      GGASMKPTFLEIVDSYKATLAEAV         250
TPIS_YEAST      GGASLKPEFVDIINSRN-------         248
TPIS_MAIZE      GGASLKPEFIDIINAATVKSA---         253
TPIS_CAEEL      GGASLKPDFVKIINARS-------         247
TPIS_DROME      GGASLKPEFVDIINARQ-------         247
TPIS_HUMAN      GGASLKPEFVDIINAKQ-------         249
TPIS_RABIT      GGASLKPEFVDIINAKQ-------         249
TPIS_MOUSE      GGASLKPEFVDIINAKQ-------         249
                ****:** *:.*:::
```

FIGURE 4. Output MSA of the CLUSTAL OMEGA from set 2

The general format of providing the input to the **AL2CO** server is given below. It includes the syntax in which the input must be provided to the AL2CO. The aligned sequences should only be passed to the AL2Co server. The parameters involved in submitting the input have been enclosed later for each category mentioned above:



FIGURE 5. Input sequences to the AL2CO server from set 1



FIGURE 6. Input sequences to the AL2CO server from set 2

(i) Unweighted frequency and entropy-based measure

The parameters for this scenario are enclosed in the image below. The important parameters are as follows:

- sequence weighting scheme
- conservation calculation method
- scoring matrix (for sum of pairs method only)
- scoring matrix transformation (for sum of pairs method only)
- normalize conservation values



FIGURE 7. Parameters for the given scenario (i)

The AL2CO gives the list of positional conservation values and the alignment with integer conservation indices. The question asks to only calculate the positional conservation values.
The window of positional conservation values generates the following set of parameters which are the ones displayed in the above image.
It also displays some parameters taken into consideration to compute the desired positional conservation scores.

```
* gap fraction no less than  0.50; conservation set to M-S
  M: mean;  S: standard deviation

AL2CO parameters are:

Input alignment file: QUERY_qfIljn
Output conservation file: QUERY_qfIljn.csv.txt
Output alignment file with index: QUERY_qfIljn.csv.aln; Block size: 70
Weighting scheme: unweighted
Conservation calculation method: entropy-based
Window size: 1
Conservation not normalized
Gap fraction to suppress calculation:  0.50
```

FIGURE 8. Parameter output display in window

Figure 9. Positional conservation scores for set 1

**(A) 1-35 pos**

| Pos | Res | Score |
|---|---|---|
| 1 | M | 0.000 |
| 2 | V | -0.305 |
| 3 | L | 0.000 |
| 4 | S | 0.000 |
| 5 | P | -1.121 |
| 6 | A | -1.160 |
| 7 | D | 0.000 |
| 8 | K | 0.000 |
| 9 | T | -1.169 |
| 10 | N | -0.305 |
| 11 | V | -0.586 |
| 12 | K | -0.305 |
| 13 | A | -1.034 |
| 14 | A | -1.160 |
| 15 | W | -0.305 |
| 16 | G | -0.916 |
| 17 | K | 0.000 |
| 18 | V | -0.689 |
| 19 | G | -0.305 |
| 20 | A | -1.034 |
| 21 | H | -0.886 |
| 22 | A | -0.474 |
| 23 | G | -1.034 |
| 24 | E | -0.908 |
| 25 | Y | -0.886 |
| 26 | G | 0.000 |
| 27 | A | -0.886 |
| 28 | E | -0.305 |
| 29 | A | -0.305 |
| 30 | L | 0.000 |
| 31 | E | -0.886 |
| 32 | R | 0.000 |
| 33 | M | -0.305 |
| 34 | F | -0.305 |
| 35 | L | -1.295 |

**(B) 36-70 pos**

| Pos | Res | Score |
|---|---|---|
| 36 | S | -1.367 |
| 37 | F | -0.474 |
| 38 | P | 0.000 |
| 39 | T | -0.474 |
| 40 | T | 0.000 |
| 41 | K | 0.000 |
| 42 | T | 0.000 |
| 43 | Y | 0.000 |
| 44 | F | 0.000 |
| 45 | P | -0.474 |
| 46 | H | 0.000 |
| 47 | F | -0.600 |
| 48 | : | -0.857 * |
| 49 | D | -0.305 |
| 50 | L | -0.586 |
| 51 | S | -0.600 |
| 52 | H | -0.586 |
| 53 | G | 0.000 |
| 54 | S | 0.000 |
| 55 | A | -0.600 |
| 56 | Q | -0.305 |
| 57 | V | -0.305 |
| 58 | K | 0.000 |
| 59 | G | -0.860 |
| 60 | H | 0.000 |
| 61 | G | 0.000 |
| 62 | K | -0.600 |
| 63 | K | 0.000 |
| 64 | V | 0.000 |
| 65 | A | -0.760 |
| 66 | D | -0.760 |
| 67 | A | -0.305 |
| 68 | L | -0.305 |
| 69 | T | -0.860 |
| 70 | N | -1.673 |

**(C) 71-105 pos**

| Pos | Res | Score |
|---|---|---|
| 71 | A | 0.000 |
| 72 | V | -0.586 |
| 73 | A | -1.594 |
| 74 | H | -0.305 |
| 75 | V | -0.995 |
| 76 | D | -0.305 |
| 77 | D | 0.000 |
| 78 | M | -0.760 |
| 79 | P | -0.600 |
| 80 | N | -0.760 |
| 81 | A | -0.305 |
| 82 | L | 0.000 |
| 83 | S | -0.305 |
| 84 | A | -1.414 |
| 85 | L | 0.000 |
| 86 | S | 0.000 |
| 87 | D | -0.305 |
| 88 | L | -0.305 |
| 89 | H | 0.000 |
| 90 | A | 0.000 |
| 91 | H | -0.760 |
| 92 | K | 0.000 |
| 93 | L | 0.000 |
| 94 | R | 0.000 |
| 95 | V | 0.000 |
| 96 | D | 0.000 |
| 97 | P | 0.000 |
| 98 | V | -0.305 |
| 99 | N | 0.000 |
| 100 | F | 0.000 |
| 101 | K | 0.000 |
| 102 | L | -0.600 |
| 103 | L | 0.000 |
| 104 | S | -0.600 |
| 105 | H | -0.305 |

**(D) 106-143 pos**

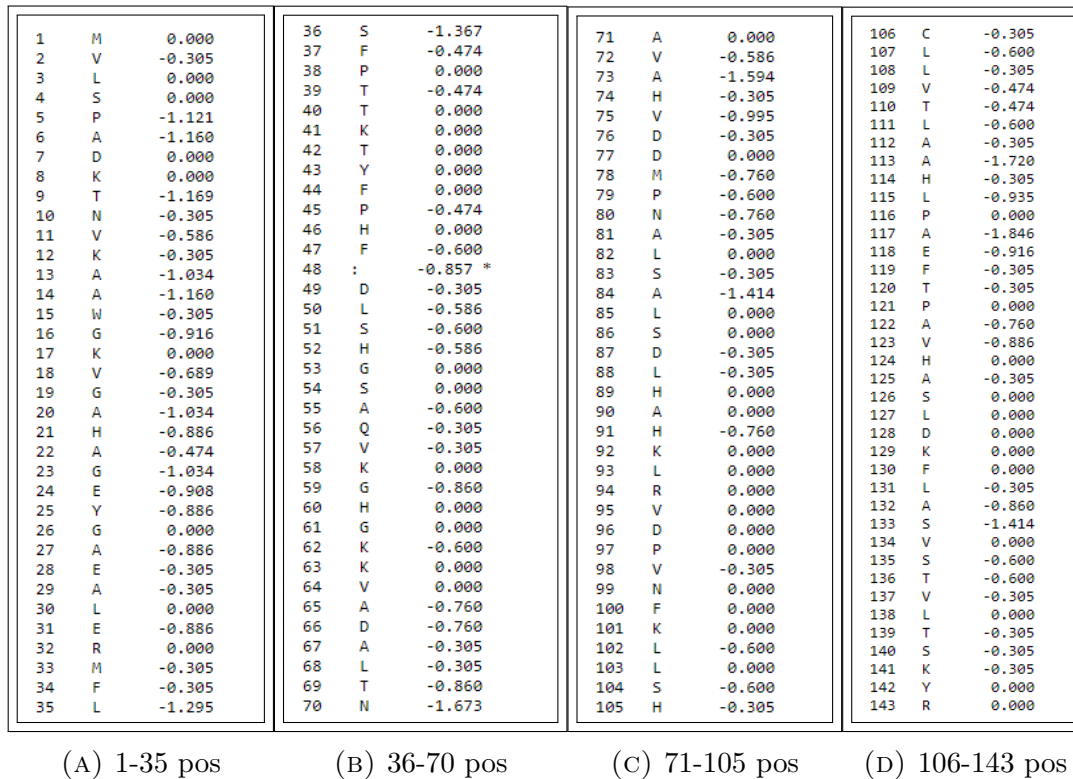| Pos | Res | Score |
|---|---|---|
| 106 | C | -0.305 |
| 107 | L | -0.600 |
| 108 | L | -0.305 |
| 109 | V | -0.474 |
| 110 | T | -0.474 |
| 111 | L | -0.600 |
| 112 | A | -0.305 |
| 113 | A | -1.720 |
| 114 | H | -0.305 |
| 115 | L | -0.935 |
| 116 | P | 0.000 |
| 117 | A | -1.846 |
| 118 | E | -0.916 |
| 119 | F | -0.305 |
| 120 | T | -0.305 |
| 121 | P | 0.000 |
| 122 | A | -0.760 |
| 123 | V | -0.886 |
| 124 | H | 0.000 |
| 125 | A | -0.305 |
| 126 | S | 0.000 |
| 127 | L | 0.000 |
| 128 | D | 0.000 |
| 129 | K | 0.000 |
| 130 | F | 0.000 |
| 131 | L | -0.305 |
| 132 | A | -0.860 |
| 133 | S | -1.414 |
| 134 | V | 0.000 |
| 135 | S | -0.600 |
| 136 | T | -0.600 |
| 137 | V | -0.305 |
| 138 | L | 0.000 |
| 139 | T | -0.305 |
| 140 | S | -0.305 |
| 141 | K | -0.305 |
| 142 | Y | 0.000 |
| 143 | R | 0.000 |

FIGURE 9. Positional conservation scores for set 1

Above is the image of the positional conservation scores for the set 1 proteins. Below is the image of the positional conservation scores for the set 2 proteins.

**(A) 1-54 pos**

| Pos | Res | Score |
|---|---|---|
| 1 | M | -1.216 * |
| 2 | A | -1.216 * |
| 3 | P | -0.974 |
| 4 | S | -1.494 |
| 5 | R | -0.377 |
| 6 | K | -0.736 |
| 7 | F | -0.377 |
| 8 | F | -0.736 |
| 9 | V | -0.377 |
| 10 | G | -0.377 |
| 11 | G | -0.377 |
| 12 | N | 0.000 |
| 13 | W | -0.349 |
| 14 | K | 0.000 |
| 15 | M | -0.937 |
| 16 | N | 0.000 |
| 17 | G | -0.349 |
| 18 | R | -1.311 |
| 19 | K | -1.303 |
| 20 | Q | -1.523 |
| 21 | S | -1.581 |
| 22 | L | -1.061 |
| 23 | G | -1.677 |
| 24 | E | -1.003 |
| 25 | L | -0.687 |
| 26 | I | -0.965 |
| 27 | G | -1.735 |
| 28 | T | -1.003 |
| 29 | L | -0.349 |
| 30 | N | -0.349 |
| 31 | A | -1.003 |
| 32 | A | -0.684 |
| 33 | K | -1.889 |
| 34 | V | -1.581 |
| 35 | P | -1.149 |
| 36 | A | -1.677 |
| 37 | D | -1.523 |
| 38 | : | -1.216 * |
| 39 | : | -1.216 * |
| 40 | T | -0.687 |
| 41 | E | -0.849 |
| 42 | V | -0.349 |
| 43 | V | 0.000 |
| 44 | C | -1.061 |
| 45 | A | -1.003 |
| 46 | P | -0.349 |
| 47 | P | -0.349 |
| 48 | T | -1.215 |
| 49 | A | -1.677 |
| 50 | Y | -0.684 |
| 51 | I | -0.687 |
| 52 | D | -1.149 |
| 53 | F | -1.465 |
| 54 | A | -1.149 |

**(B) 55-108 pos**

| Pos | Res | Score |
|---|---|---|
| 55 | R | -0.937 |
| 56 | Q | -1.215 |
| 57 | K | -1.427 |
| 58 | L | -0.349 |
| 59 | : | -1.216 * |
| 60 | D | -1.523 |
| 61 | P | -1.149 |
| 62 | K | -1.427 |
| 63 | I | -1.311 |
| 64 | A | -1.831 |
| 65 | V | -1.003 |
| 66 | A | -0.684 |
| 67 | A | -0.349 |
| 68 | Q | 0.000 |
| 69 | N | 0.000 |
| 70 | C | -0.687 |
| 71 | Y | -0.849 |
| 72 | K | -1.074 |
| 73 | V | -1.149 |
| 74 | T | -1.523 |
| 75 | N | -1.369 |
| 76 | G | 0.000 |
| 77 | A | 0.000 |
| 78 | F | 0.000 |
| 79 | T | 0.000 |
| 80 | G | 0.000 |
| 81 | E | 0.000 |
| 82 | I | -1.149 |
| 83 | S | -0.349 |
| 84 | P | -1.149 |
| 85 | G | -1.465 |
| 86 | M | -0.684 |
| 87 | I | -0.937 |
| 88 | K | -0.349 |
| 89 | D | -0.349 |
| 90 | C | -1.523 |
| 91 | G | 0.000 |
| 92 | A | -1.003 |
| 93 | T | -1.677 |
| 94 | W | 0.000 |
| 95 | V | 0.000 |
| 96 | V | -0.687 |
| 97 | L | 0.000 |
| 98 | G | 0.000 |
| 99 | H | 0.000 |
| 100 | S | 0.000 |
| 101 | E | 0.000 |
| 102 | R | 0.000 |
| 103 | R | 0.000 |
| 104 | H | -1.149 |
| 105 | V | -1.273 |
| 106 | : | -1.216 * |
| 107 | : | -1.216 * |
| 108 | F | -0.849 |

**(C) 109-162 pos**

| Pos | Res | Score |
|---|---|---|
| 109 | G | -0.684 |
| 110 | E | 0.000 |
| 111 | S | -0.684 |
| 112 | D | -0.530 |
| 113 | E | -1.149 |
| 114 | L | -1.149 |
| 115 | I | -0.530 |
| 116 | G | -0.687 |
| 117 | Q | -1.311 |
| 118 | K | 0.000 |
| 119 | V | -0.995 |
| 120 | A | -1.303 |
| 121 | H | -1.149 |
| 122 | A | 0.000 |
| 123 | L | -0.349 |
| 124 | A | -1.427 |
| 125 | : | -1.216 * |
| 126 | E | -1.061 |
| 127 | G | -0.349 |
| 128 | L | -1.003 |
| 129 | G | -0.937 |
| 130 | V | 0.000 |
| 131 | I | -0.349 |
| 132 | A | -1.303 |
| 133 | C | 0.000 |
| 134 | I | -0.530 |
| 135 | G | 0.000 |
| 136 | E | 0.000 |
| 137 | K | -0.965 |
| 138 | L | -0.349 |
| 139 | D | -0.637 |
| 140 | E | -0.349 |
| 141 | R | -0.349 |
| 142 | E | -0.349 |
| 143 | A | -0.349 |
| 144 | G | 0.000 |
| 145 | L | -1.677 |
| 146 | T | 0.000 |
| 147 | E | -1.831 |
| 148 | K | -1.311 |
| 149 | V | 0.000 |
| 150 | V | -0.684 |
| 151 | F | -1.273 |
| 152 | E | -1.215 |
| 153 | Q | 0.000 |
| 154 | T | -0.965 |
| 155 | K | -1.427 |
| 156 | V | -0.637 |
| 157 | I | -1.149 |
| 158 | A | -1.003 |
| 159 | D | -0.995 |
| 160 | N | -1.215 |
| 161 | V | -1.215 |
| 162 | : | -1.216 * |

**(D) 163-216 pos**

| Pos | Res | Score |
|---|---|---|
| 163 | : | -1.216 * |
| 164 | K | -1.003 |
| 165 | D | -0.684 |
| 166 | W | 0.000 |
| 167 | S | -1.003 |
| 168 | K | -1.215 |
| 169 | V | -0.349 |
| 170 | V | 0.000 |
| 171 | L | -1.099 |
| 172 | A | 0.000 |
| 173 | Y | 0.000 |
| 174 | E | 0.000 |
| 175 | P | 0.000 |
| 176 | V | 0.000 |
| 177 | W | 0.000 |
| 178 | A | 0.000 |
| 179 | I | 0.000 |
| 180 | G | 0.000 |
| 181 | T | 0.000 |
| 182 | G | 0.000 |
| 183 | K | -0.684 |
| 184 | T | -0.937 |
| 185 | A | 0.000 |
| 186 | T | -0.349 |
| 187 | P | -0.349 |
| 188 | Q | -1.215 |
| 189 | Q | -0.349 |
| 190 | A | 0.000 |
| 191 | Q | -0.349 |
| 192 | E | -0.684 |
| 193 | V | -0.349 |
| 194 | H | 0.000 |
| 195 | E | -0.637 |
| 196 | K | -1.677 |
| 197 | L | -0.849 |
| 198 | R | 0.000 |
| 199 | G | -1.677 |
| 200 | W | -0.530 |
| 201 | L | -0.349 |
| 202 | K | -1.149 |
| 203 | S | -1.303 |
| 204 | N | -0.637 |
| 205 | V | -1.149 |
| 206 | S | -0.849 |
| 207 | D | -1.215 |
| 208 | A | -1.677 |
| 209 | V | -0.684 |
| 210 | A | -0.349 |
| 211 | Q | -1.465 |
| 212 | S | -1.003 |
| 213 | T | -0.637 |
| 214 | R | 0.000 |
| 215 | I | 0.000 |
| 216 | I | -0.937 |

**(E) 216-264 pos**

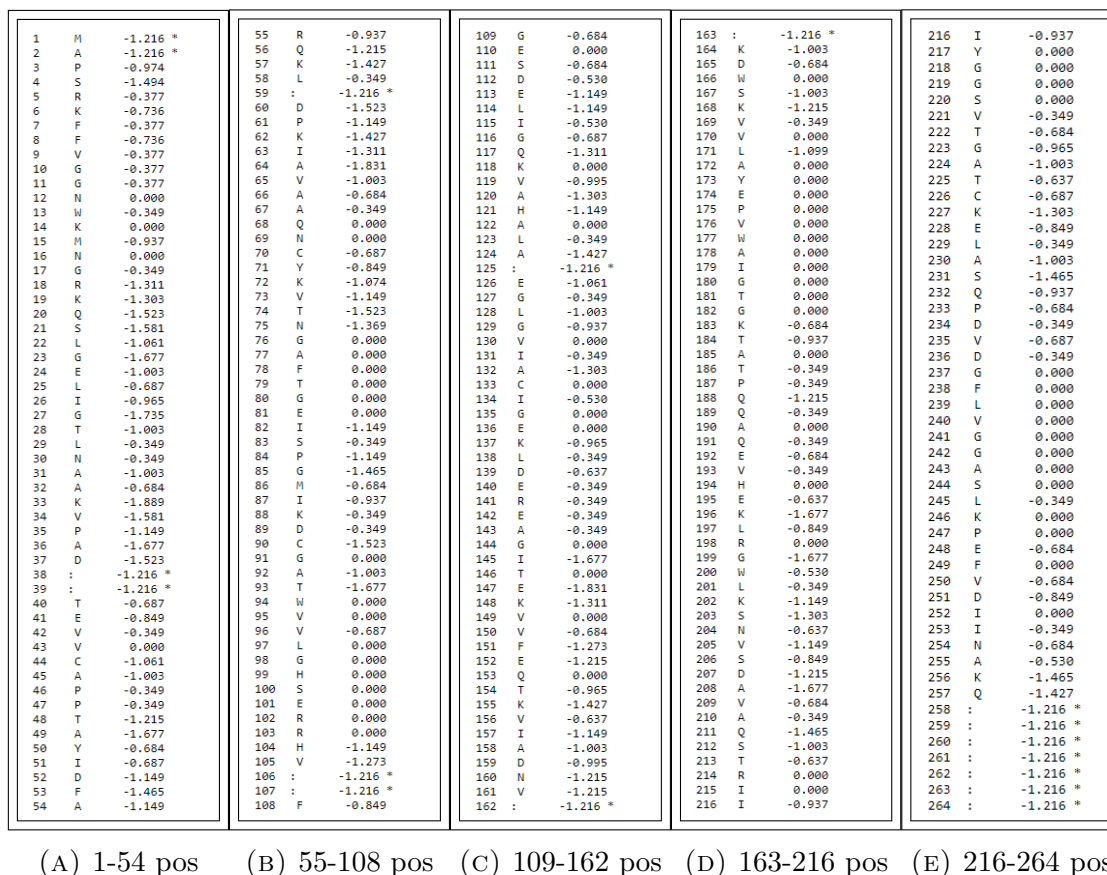| Pos | Res | Score |
|---|---|---|
| 216 | I | -0.937 |
| 217 | Y | 0.000 |
| 218 | G | 0.000 |
| 219 | G | 0.000 |
| 220 | S | 0.000 |
| 221 | V | -0.349 |
| 222 | T | -0.684 |
| 223 | G | -0.965 |
| 224 | A | -1.003 |
| 225 | T | -0.637 |
| 226 | C | -0.687 |
| 227 | K | -1.303 |
| 228 | E | -0.849 |
| 229 | L | -0.349 |
| 230 | A | -1.003 |
| 231 | S | -1.465 |
| 232 | Q | -0.937 |
| 233 | P | -0.684 |
| 234 | D | -0.349 |
| 235 | V | -0.687 |
| 236 | D | -0.349 |
| 237 | G | 0.000 |
| 238 | F | 0.000 |
| 239 | L | 0.000 |
| 240 | V | 0.000 |
| 241 | G | 0.000 |
| 242 | G | 0.000 |
| 243 | A | 0.000 |
| 244 | S | 0.000 |
| 245 | L | -0.349 |
| 246 | K | 0.000 |
| 247 | P | 0.000 |
| 248 | E | -0.684 |
| 249 | F | 0.000 |
| 250 | V | -0.684 |
| 251 | D | -0.849 |
| 252 | I | 0.000 |
| 253 | I | -0.349 |
| 254 | N | -0.684 |
| 255 | A | -0.530 |
| 256 | K | -1.465 |
| 257 | Q | -1.427 |
| 258 | : | -1.216 * |
| 259 | : | -1.216 * |
| 260 | : | -1.216 * |
| 261 | : | -1.216 * |
| 262 | : | -1.216 * |
| 263 | : | -1.216 * |
| 264 | : | -1.216 * |

FIGURE 10. Positional conservation scores for set 2

(ii) Unweighted frequency and variance-based measure

The parameters for this scenario are enclosed in the image below. The important parameters are as follows:

- sequence weighting scheme
- conservation calculation method
- scoring matrix (for sum of pairs method only)
- scoring matrix transformation (for sum of pairs method only)
- normalize conservation values



FIGURE 11. Parameters for the given scenario (ii)

The AL2CO gives the list of positional conservation values and the alignment with integer conservation indices. The question asks to only calculate the positional conservation values.
The window of positional conservation values generates the following set of parameters which are the ones displayed in the above image.
It also displays some parameters taken into consideration to compute the desired positional conservation scores.



FIGURE 12. Parameter output display in window

FIGURE 13. Positional conservation scores for set 1

(A) 1-35 pos

| Pos | Res | Score |
|---|---|---|
| 1 | M | 1.021 |
| 2 | V | 0.865 |
| 3 | L | 0.909 |
| 4 | S | 0.953 |
| 5 | P | 0.527 |
| 6 | A | 0.561 |
| 7 | D | 0.973 |
| 8 | K | 0.956 |
| 9 | T | 0.541 |
| 10 | N | 0.917 |
| 11 | V | 0.744 |
| 12 | K | 0.871 |
| 13 | A | 0.567 |
| 14 | A | 0.559 |
| 15 | W | 0.943 |
| 16 | G | 0.612 |
| 17 | K | 0.956 |
| 18 | V | 0.691 |
| 19 | G | 0.886 |
| 20 | A | 0.613 |
| 21 | H | 0.713 |
| 22 | A | 0.740 |
| 23 | G | 0.624 |
| 24 | E | 0.665 |
| 25 | Y | 0.757 |
| 26 | G | 0.977 |
| 27 | A | 0.650 |
| 28 | E | 0.922 |
| 29 | A | 0.817 |
| 30 | L | 0.909 |
| 31 | E | 0.749 |
| 32 | R | 1.015 |
| 33 | M | 0.933 |
| 34 | F | 0.904 |
| 35 | L | 0.470 |

(B) 36-70 pos

| Pos | Res | Score |
|---|---|---|
| 36 | S | 0.472 |
| 37 | F | 0.830 |
| 38 | P | 0.993 |
| 39 | T | 0.821 |
| 40 | T | 0.977 |
| 41 | K | 0.956 |
| 42 | T | 0.977 |
| 43 | Y | 1.016 |
| 44 | F | 0.988 |
| 45 | P | 0.822 |
| 46 | H | 0.969 |
| 47 | F | 0.823 |
| 48 | : | 0.657 * |
| 49 | D | 0.888 |
| 50 | L | 0.674 |
| 51 | S | 0.776 |
| 52 | H | 0.745 |
| 53 | G | 0.977 |
| 54 | S | 0.953 |
| 55 | A | 0.730 |
| 56 | Q | 0.936 |
| 57 | V | 0.872 |
| 58 | K | 0.956 |
| 59 | G | 0.596 |
| 60 | H | 0.969 |
| 61 | G | 0.977 |
| 62 | K | 0.777 |
| 63 | K | 0.956 |
| 64 | V | 0.956 |
| 65 | A | 0.660 |
| 66 | D | 0.704 |
| 67 | A | 0.817 |
| 68 | L | 0.825 |
| 69 | T | 0.643 |
| 70 | N | 0.360 |

(C) 71-105 pos

| Pos | Res | Score |
|---|---|---|
| 71 | A | 0.906 |
| 72 | V | 0.701 |
| 73 | A | 0.406 |
| 74 | H | 0.879 |
| 75 | V | 0.543 |
| 76 | D | 0.888 |
| 77 | D | 0.973 |
| 78 | M | 0.675 |
| 79 | P | 0.807 |
| 80 | N | 0.733 |
| 81 | A | 0.817 |
| 82 | L | 0.909 |
| 83 | S | 0.869 |
| 84 | A | 0.451 |
| 85 | L | 0.909 |
| 86 | S | 0.953 |
| 87 | D | 0.888 |
| 88 | L | 0.826 |
| 89 | H | 0.969 |
| 90 | A | 0.906 |
| 91 | H | 0.731 |
| 92 | K | 0.956 |
| 93 | L | 0.909 |
| 94 | R | 1.015 |
| 95 | V | 0.956 |
| 96 | D | 0.973 |
| 97 | P | 0.993 |
| 98 | V | 0.860 |
| 99 | N | 1.012 |
| 100 | F | 0.988 |
| 101 | K | 0.956 |
| 102 | L | 0.739 |
| 103 | L | 0.909 |
| 104 | S | 0.781 |
| 105 | H | 0.886 |

(D) 106-143 pos

| Pos | Res | Score |
|---|---|---|
| 106 | C | 0.937 |
| 107 | L | 0.739 |
| 108 | L | 0.825 |
| 109 | V | 0.784 |
| 110 | T | 0.806 |
| 111 | L | 0.735 |
| 112 | A | 0.814 |
| 113 | A | 0.345 |
| 114 | H | 0.885 |
| 115 | L | 0.559 |
| 116 | P | 0.993 |
| 117 | A | 0.302 |
| 118 | E | 0.633 |
| 119 | F | 0.893 |
| 120 | T | 0.891 |
| 121 | P | 0.993 |
| 122 | A | 0.660 |
| 123 | V | 0.695 |
| 124 | H | 0.969 |
| 125 | A | 0.814 |
| 126 | S | 0.953 |
| 127 | L | 0.909 |
| 128 | D | 0.973 |
| 129 | K | 0.956 |
| 130 | F | 0.988 |
| 131 | L | 0.822 |
| 132 | A | 0.595 |
| 133 | S | 0.459 |
| 134 | V | 0.956 |
| 135 | S | 0.769 |
| 136 | T | 0.796 |
| 137 | V | 0.860 |
| 138 | L | 0.909 |
| 139 | T | 0.881 |
| 140 | S | 0.868 |
| 141 | K | 0.871 |
| 142 | Y | 1.016 |
| 143 | R | 1.015 |

Above is the image of the positional conservation scores for the set 1 proteins. Below is the image of the positional conservation scores for the set 2 proteins.

(A) 1-54 pos

| Pos | Res | Score |
|---|---|---|
| 1 | M | 0.529 * |
| 2 | A | 0.529 * |
| 3 | P | 0.648 |
| 4 | S | 0.426 |
| 5 | R | 0.882 |
| 6 | K | 0.726 |
| 7 | F | 0.884 |
| 8 | F | 0.769 |
| 9 | V | 0.801 |
| 10 | G | 0.812 |
| 11 | G | 0.812 |
| 12 | N | 0.995 |
| 13 | W | 0.909 |
| 14 | K | 0.962 |
| 15 | M | 0.679 |
| 16 | N | 0.995 |
| 17 | G | 0.835 |
| 18 | R | 0.510 |
| 19 | K | 0.553 |
| 20 | Q | 0.401 |
| 21 | S | 0.481 |
| 22 | L | 0.526 |
| 23 | G | 0.317 |
| 24 | E | 0.640 |
| 25 | L | 0.664 |
| 26 | I | 0.550 |
| 27 | G | 0.362 |
| 28 | T | 0.663 |
| 29 | L | 0.859 |
| 30 | N | 0.888 |
| 31 | A | 0.591 |
| 32 | A | 0.691 |
| 33 | K | 0.305 |
| 34 | V | 0.426 |
| 35 | P | 0.604 |
| 36 | A | 0.372 |
| 37 | D | 0.429 |
| 38 | : | 0.529 * |
| 39 | : | 0.529 * |
| 40 | T | 0.645 |
| 41 | E | 0.671 |
| 42 | V | 0.829 |
| 43 | V | 0.930 |
| 44 | C | 0.545 |
| 45 | A | 0.594 |
| 46 | P | 0.893 |
| 47 | P | 0.889 |
| 48 | T | 0.498 |
| 49 | A | 0.321 |
| 50 | Y | 0.808 |
| 51 | I | 0.664 |
| 52 | D | 0.595 |
| 53 | F | 0.504 |
| 54 | A | 0.509 |

(B) 55-108 pos

| Pos | Res | Score |
|---|---|---|
| 55 | R | 0.625 |
| 56 | Q | 0.541 |
| 57 | K | 0.462 |
| 58 | L | 0.850 |
| 59 | : | 0.529 * |
| 60 | D | 0.455 |
| 61 | P | 0.584 |
| 62 | K | 0.467 |
| 63 | I | 0.477 |
| 64 | A | 0.321 |
| 65 | V | 0.611 |
| 66 | A | 0.691 |
| 67 | A | 0.798 |
| 68 | Q | 0.994 |
| 69 | N | 0.995 |
| 70 | C | 0.669 |
| 71 | Y | 0.723 |
| 72 | K | 0.596 |
| 73 | V | 0.532 |
| 74 | T | 0.408 |
| 75 | N | 0.453 |
| 76 | G | 0.944 |
| 77 | A | 0.909 |
| 78 | F | 1.001 |
| 79 | T | 0.983 |
| 80 | G | 0.944 |
| 81 | E | 0.962 |
| 82 | I | 0.559 |
| 83 | S | 0.873 |
| 84 | P | 0.562 |
| 85 | G | 0.391 |
| 86 | M | 0.811 |
| 87 | I | 0.601 |
| 88 | K | 0.849 |
| 89 | D | 0.883 |
| 90 | C | 0.430 |
| 91 | G | 0.944 |
| 92 | A | 0.580 |
| 93 | T | 0.406 |
| 94 | W | 1.013 |
| 95 | V | 0.930 |
| 96 | V | 0.644 |
| 97 | L | 0.963 |
| 98 | G | 0.944 |
| 99 | H | 1.017 |
| 100 | S | 0.979 |
| 101 | E | 0.962 |
| 102 | R | 0.999 |
| 103 | R | 0.999 |
| 104 | H | 0.592 |
| 105 | V | 0.484 |
| 106 | : | 0.529 * |
| 107 | : | 0.529 * |
| 108 | F | 0.703 |

(C) 109-162 pos

| Pos | Res | Score |
|---|---|---|
| 109 | G | 0.735 |
| 110 | E | 0.962 |
| 111 | S | 0.768 |
| 112 | D | 0.796 |
| 113 | E | 0.546 |
| 114 | L | 0.577 |
| 115 | I | 0.762 |
| 116 | G | 0.601 |
| 117 | Q | 0.505 |
| 118 | K | 0.962 |
| 119 | V | 0.535 |
| 120 | A | 0.475 |
| 121 | H | 0.638 |
| 122 | A | 0.909 |
| 123 | L | 0.861 |
| 124 | A | 0.419 |
| 125 | : | 0.529 * |
| 126 | E | 0.514 |
| 127 | G | 0.837 |
| 128 | L | 0.642 |
| 129 | G | 0.592 |
| 130 | V | 0.930 |
| 131 | I | 0.858 |
| 132 | A | 0.492 |
| 133 | C | 1.017 |
| 134 | I | 0.762 |
| 135 | G | 0.944 |
| 136 | E | 0.962 |
| 137 | K | 0.592 |
| 138 | L | 0.859 |
| 139 | D | 0.706 |
| 140 | E | 0.857 |
| 141 | R | 0.890 |
| 142 | E | 0.853 |
| 143 | A | 0.802 |
| 144 | G | 0.944 |
| 145 | I | 0.414 |
| 146 | T | 0.983 |
| 147 | E | 0.350 |
| 148 | K | 0.452 |
| 149 | V | 0.930 |
| 150 | V | 0.726 |
| 151 | F | 0.497 |
| 152 | E | 0.533 |
| 153 | Q | 0.994 |
| 154 | T | 0.604 |
| 155 | K | 0.481 |
| 156 | V | 0.629 |
| 157 | I | 0.562 |
| 158 | A | 0.581 |
| 159 | D | 0.612 |
| 160 | N | 0.519 |
| 161 | V | 0.484 |
| 162 | : | 0.529 * |

(D) 163-216 pos

| Pos | Res | Score |
|---|---|---|
| 163 | : | 0.529 * |
| 164 | K | 0.635 |
| 165 | D | 0.766 |
| 166 | W | 1.013 |
| 167 | S | 0.658 |
| 168 | K | 0.549 |
| 169 | V | 0.822 |
| 170 | V | 0.930 |
| 171 | L | 0.495 |
| 172 | A | 0.909 |
| 173 | Y | 1.013 |
| 174 | E | 0.962 |
| 175 | P | 0.995 |
| 176 | V | 0.930 |
| 177 | W | 1.013 |
| 178 | A | 0.909 |
| 179 | I | 0.971 |
| 180 | G | 0.944 |
| 181 | T | 0.983 |
| 182 | G | 0.944 |
| 183 | K | 0.750 |
| 184 | T | 0.591 |
| 185 | A | 0.909 |
| 186 | T | 0.876 |
| 187 | P | 0.884 |
| 188 | Q | 0.530 |
| 189 | Q | 0.888 |
| 190 | A | 0.909 |
| 191 | Q | 0.885 |
| 192 | E | 0.753 |
| 193 | V | 0.822 |
| 194 | H | 1.017 |
| 195 | E | 0.669 |
| 196 | K | 0.389 |
| 197 | L | 0.656 |
| 198 | R | 0.999 |
| 199 | G | 0.350 |
| 200 | W | 0.822 |
| 201 | L | 0.850 |
| 202 | K | 0.545 |
| 203 | S | 0.564 |
| 204 | N | 0.724 |
| 205 | V | 0.516 |
| 206 | S | 0.673 |
| 207 | D | 0.555 |
| 208 | A | 0.333 |
| 209 | V | 0.707 |
| 210 | A | 0.802 |
| 211 | Q | 0.426 |
| 212 | S | 0.651 |
| 213 | T | 0.713 |
| 214 | R | 0.999 |
| 215 | I | 0.971 |
| 216 | I | 0.612 |

(E) 216-264 pos

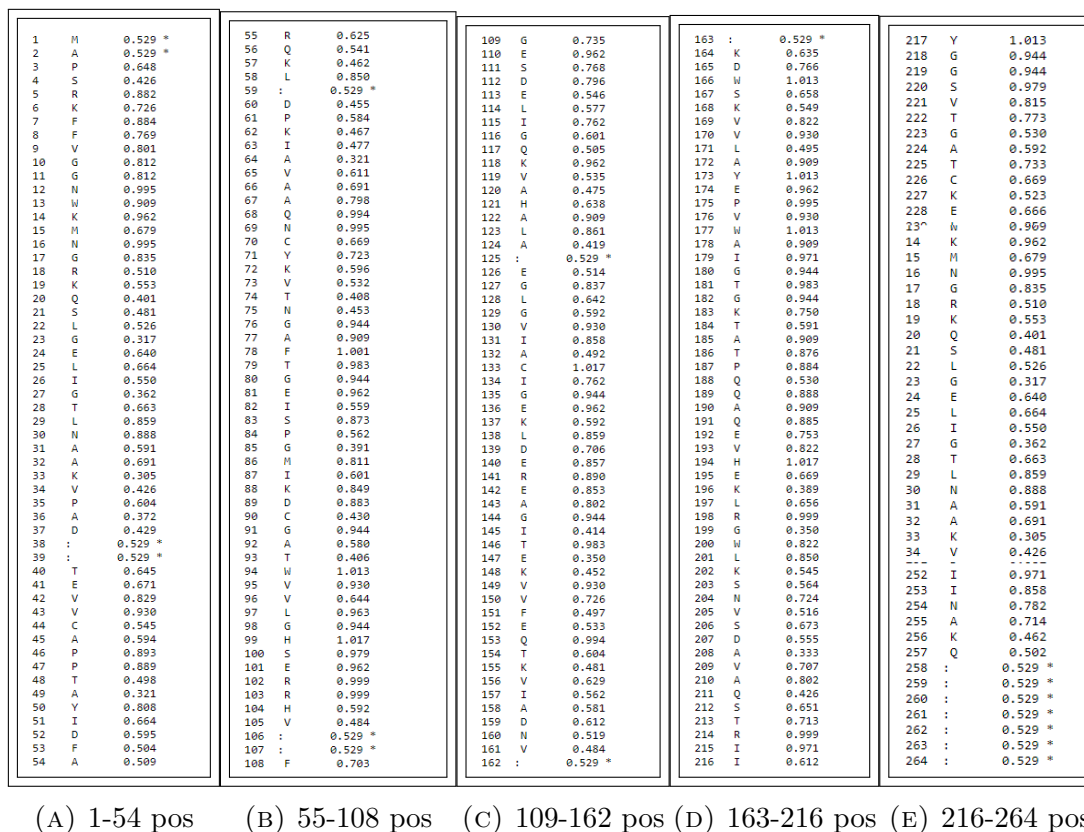| Pos | Res | Score |
|---|---|---|
| 217 | Y | 1.013 |
| 218 | G | 0.944 |
| 219 | G | 0.944 |
| 220 | S | 0.979 |
| 221 | V | 0.815 |
| 222 | T | 0.773 |
| 223 | G | 0.530 |
| 224 | A | 0.592 |
| 225 | T | 0.733 |
| 226 | C | 0.669 |
| 227 | K | 0.523 |
| 228 | E | 0.666 |
| 23^ | N | 0.909 |
| 14 | K | 0.962 |
| 15 | M | 0.679 |
| 16 | N | 0.995 |
| 17 | G | 0.835 |
| 18 | R | 0.510 |
| 19 | K | 0.553 |
| 20 | Q | 0.401 |
| 21 | S | 0.481 |
| 22 | L | 0.526 |
| 23 | G | 0.317 |
| 24 | E | 0.640 |
| 25 | L | 0.664 |
| 26 | I | 0.550 |
| 27 | G | 0.362 |
| 28 | T | 0.663 |
| 29 | L | 0.859 |
| 30 | N | 0.888 |
| 31 | A | 0.591 |
| 32 | A | 0.691 |
| 33 | K | 0.305 |
| 34 | V | 0.426 |
| 252 | I | 0.971 |
| 253 | I | 0.858 |
| 254 | N | 0.782 |
| 255 | A | 0.714 |
| 256 | K | 0.462 |
| 257 | Q | 0.502 |
| 258 | : | 0.529 * |
| 259 | : | 0.529 * |
| 260 | : | 0.529 * |
| 261 | : | 0.529 * |
| 262 | : | 0.529 * |
| 263 | : | 0.529 * |
| 264 | : | 0.529 * |

FIGURE 14. Positional conservation scores for set 2

(iii) Unweighted frequency and sum of pairs measure

The parameters for this scenario are enclosed in the image below. The important parameters are as follows:

- sequence weighting scheme
- conservation calculation method
- scoring matrix (for sum of pairs method only)
- scoring matrix transformation (for sum of pairs method only)
- normalize conservation values

PARAMETERS

- sequence weighting scheme:  ○ henikoff-henikoff ○ independent count ● unweighted
- conservation calculation method:  ○ entropy ○ variance ● sum-of-pairs
- For sum-of-pairs method only:
  scoring matrix:
    ● BLOSUM62 matrix ○ identity matrix
  scoring matrix transformation :
    ● no transformation ○ normalization ○ adjustment
- normalize conservation values:  ○ True ● False
- window size used for averaging conservation (for smoothing purpose):  1
- gap fraction above which conservation calculation is not performed:  0.5
- exclude the first sequence from calculation:  ○ True ● False
- output alignment block size:  70
- pdb file for which b-factor field is replaced with conservation (optional):
  Choose File   No file chosen

FIGURE 15. Parameters for the given scenario (iii)

The AL2CO gives the list of positional conservation values and the alignment with integer conservation indices. The question asks to only calculate the positional conservation values.
The window of positional conservation values generates the following set of parameters which are the ones displayed in the above image.
It also displays some parameters taken into consideration to compute the desired positional conservation scores.

```
* gap fraction no less than  0.50; conservation set to M-S
  M: mean;  S: standard deviation

AL2CO parameters are:

Input alignment file: QUERY_phGSQi
Output conservation file: QUERY_phGSQi.csv.txt
Output alignment file with index: QUERY_phGSQi.csv.aln; Block size: 70
Input matrix file: BLOSUM62
Matrix transformation:  no transformation
Weighting scheme: unweighted
Conservation calculation method: sum-of-pairs measure
Window size: 1
Conservation not normalized
Gap fraction to suppress calculation:  0.50
```

FIGURE 16. Parameter output display in window

(A) 1-35 pos  (B) 36-70 pos  (C) 71-105 pos  (D) 106-143 pos

FIGURE 17. Positional conservation scores for set 1

Above is the image of the positional conservation scores for the set 1 proteins. Below is the image of the positional conservation scores for the set 2 proteins.



(A) 1-54 pos  (B) 55-108 pos  (C) 109-162 pos  (D) 163-216 pos  (E) 216-264 pos

FIGURE 18. Positional conservation scores for set 2

(iv) Weighted frequency and variance-based measure

The parameters for this scenario are enclosed in the image below. The important parameters are as follows:

- sequence weighting scheme
- conservation calculation method
- scoring matrix (for sum of pairs method only)
- scoring matrix transformation (for sum of pairs method only)
- normalize conservation values



FIGURE 19. Parameters for the given scenario (iv)

The AL2CO gives the list of positional conservation values and the alignment with integer conservation indices. The question asks to only calculate the positional conservation values.
The window of positional conservation values generates the following set of parameters which are the ones displayed in the above image.
It also displays some parameters taken into consideration to compute the desired positional conservation scores.



FIGURE 20. Parameter output display in window

Figure 21 — Positional conservation scores for set 1

**(A) 1-35 pos**

| Pos | AA | Score |
|---|---|---|
| 1 | M | 1.019 |
| 2 | V | 0.729 |
| 3 | L | 0.916 |
| 4 | S | 0.951 |
| 5 | P | 0.492 |
| 6 | A | 0.439 |
| 7 | D | 0.972 |
| 8 | K | 0.950 |
| 9 | T | 0.465 |
| 10 | N | 0.770 |
| 11 | V | 0.763 |
| 12 | K | 0.744 |
| 13 | A | 0.566 |
| 14 | A | 0.415 |
| 15 | W | 0.888 |
| 16 | G | 0.604 |
| 17 | K | 0.950 |
| 18 | V | 0.719 |
| 19 | G | 0.831 |
| 20 | A | 0.502 |
| 21 | H | 0.627 |
| 22 | A | 0.767 |
| 23 | G | 0.517 |
| 24 | E | 0.593 |
| 25 | Y | 0.679 |
| 26 | G | 0.977 |
| 27 | A | 0.561 |
| 28 | E | 0.784 |
| 29 | A | 0.768 |
| 30 | L | 0.916 |
| 31 | E | 0.628 |
| 32 | R | 1.011 |
| 33 | M | 0.944 |
| 34 | F | 0.780 |
| 35 | L | 0.373 |

**(B) 36-70 pos**

| Pos | AA | Score |
|---|---|---|
| 36 | S | 0.386 |
| 37 | F | 0.718 |
| 38 | P | 0.993 |
| 39 | T | 0.715 |
| 40 | T | 0.979 |
| 41 | K | 0.950 |
| 42 | T | 0.979 |
| 43 | Y | 1.012 |
| 44 | F | 0.989 |
| 45 | P | 0.712 |
| 46 | H | 0.972 |
| 47 | F | 0.711 |
| 48 | : | 0.582 * |
| 49 | D | 0.901 |
| 50 | L | 0.621 |
| 51 | S | 0.605 |
| 52 | H | 0.689 |
| 53 | G | 0.977 |
| 54 | S | 0.951 |
| 55 | A | 0.628 |
| 56 | Q | 0.800 |
| 57 | V | 0.748 |
| 58 | K | 0.950 |
| 59 | G | 0.687 |
| 60 | H | 0.972 |
| 61 | G | 0.977 |
| 62 | K | 0.812 |
| 63 | K | 0.950 |
| 64 | V | 0.956 |
| 65 | A | 0.612 |
| 66 | D | 0.561 |
| 67 | A | 0.692 |
| 68 | L | 0.710 |
| 69 | T | 0.545 |
| 70 | N | 0.382 |

**(C) 71-105 pos**

| Pos | AA | Score |
|---|---|---|
| 71 | A | 0.909 |
| 72 | V | 0.682 |
| 73 | A | 0.359 |
| 74 | H | 0.745 |
| 75 | V | 0.534 |
| 76 | D | 0.888 |
| 77 | D | 0.972 |
| 78 | M | 0.683 |
| 79 | P | 0.618 |
| 80 | N | 0.664 |
| 81 | A | 0.692 |
| 82 | L | 0.916 |
| 83 | S | 0.749 |
| 84 | A | 0.428 |
| 85 | L | 0.916 |
| 86 | S | 0.951 |
| 87 | D | 0.763 |
| 88 | L | 0.715 |
| 89 | H | 0.972 |
| 90 | A | 0.909 |
| 91 | H | 0.618 |
| 92 | K | 0.950 |
| 93 | L | 0.916 |
| 94 | R | 1.011 |
| 95 | V | 0.956 |
| 96 | D | 0.972 |
| 97 | P | 0.993 |
| 98 | V | 0.715 |
| 99 | N | 1.011 |
| 100 | F | 0.989 |
| 101 | K | 0.950 |
| 102 | L | 0.632 |
| 103 | L | 0.916 |
| 104 | S | 0.619 |
| 105 | H | 0.839 |

**(D) 106-143 pos**

| Pos | AA | Score |
|---|---|---|
| 106 | C | 0.962 |
| 107 | L | 0.588 |
| 108 | L | 0.783 |
| 109 | V | 0.844 |
| 110 | T | 0.677 |
| 111 | L | 0.580 |
| 112 | A | 0.683 |
| 113 | A | 0.355 |
| 114 | H | 0.769 |
| 115 | L | 0.529 |
| 116 | P | 0.993 |
| 117 | A | 0.299 |
| 118 | E | 0.584 |
| 119 | F | 0.837 |
| 120 | T | 0.907 |
| 121 | P | 0.993 |
| 122 | A | 0.589 |
| 123 | V | 0.519 |
| 124 | H | 0.972 |
| 125 | A | 0.685 |
| 126 | S | 0.951 |
| 127 | L | 0.916 |
| 128 | D | 0.972 |
| 129 | K | 0.950 |
| 130 | F | 0.989 |
| 131 | L | 0.839 |
| 132 | A | 0.528 |
| 133 | S | 0.400 |
| 134 | V | 0.956 |
| 135 | S | 0.577 |
| 136 | T | 0.615 |
| 137 | V | 0.715 |
| 138 | L | 0.916 |
| 139 | T | 0.738 |
| 140 | S | 0.743 |
| 141 | K | 0.744 |
| 142 | Y | 1.012 |
| 143 | R | 1.011 |

FIGURE 21. Positional conservation scores for set 1

Above is the image of the positional conservation scores for the set 1 proteins. Below is the image of the positional conservation scores for the set 2 proteins.



**(A) 1-54 pos**

| Pos | AA | Score |
|---|---|---|
| 1 | M | 0.494 * |
| 2 | A | 0.494 * |
| 3 | P | 0.666 |
| 4 | S | 0.365 |
| 5 | R | 0.827 |
| 6 | K | 0.641 |
| 7 | F | 0.830 |
| 8 | F | 0.713 |
| 9 | V | 0.739 |
| 10 | G | 0.754 |
| 11 | V | 0.754 |
| 12 | N | 0.994 |
| 13 | W | 0.876 |
| 14 | K | 0.966 |
| 15 | M | 0.639 |
| 16 | N | 0.994 |
| 17 | G | 0.810 |
| 18 | R | 0.530 |
| 19 | K | 0.515 |
| 20 | Q | 0.372 |
| 21 | S | 0.482 |
| 22 | L | 0.514 |
| 23 | G | 0.270 |
| 24 | E | 0.597 |
| 25 | L | 0.656 |
| 26 | I | 0.557 |
| 27 | G | 0.393 |
| 28 | T | 0.595 |
| 29 | L | 0.826 |
| 30 | N | 0.885 |
| 31 | A | 0.549 |
| 32 | A | 0.676 |
| 33 | K | 0.302 |
| 34 | V | 0.364 |
| 35 | P | 0.553 |
| 36 | A | 0.348 |
| 37 | D | 0.421 |
| 38 | : | 0.494 * |
| 39 | : | 0.494 * |
| 40 | T | 0.684 |
| 41 | E | 0.623 |
| 42 | V | 0.795 |
| 43 | V | 0.930 |
| 44 | C | 0.563 |
| 45 | A | 0.557 |
| 46 | P | 0.894 |
| 47 | P | 0.855 |
| 48 | T | 0.510 |
| 49 | A | 0.294 |
| 50 | Y | 0.771 |
| 51 | I | 0.686 |
| 52 | D | 0.564 |
| 53 | F | 0.492 |
| 54 | A | 0.450 |

**(B) 55-108 pos**

| Pos | AA | Score |
|---|---|---|
| 55 | R | 0.582 |
| 56 | Q | 0.505 |
| 57 | K | 0.399 |
| 58 | L | 0.809 |
| 59 | : | 0.494 * |
| 60 | D | 0.440 |
| 61 | P | 0.551 |
| 62 | K | 0.421 |
| 63 | I | 0.497 |
| 64 | A | 0.292 |
| 65 | V | 0.549 |
| 66 | A | 0.623 |
| 67 | A | 0.793 |
| 68 | Q | 0.993 |
| 69 | N | 0.994 |
| 70 | C | 0.649 |
| 71 | Y | 0.666 |
| 72 | K | 0.506 |
| 73 | V | 0.465 |
| 74 | T | 0.380 |
| 75 | N | 0.446 |
| 76 | G | 0.951 |
| 77 | A | 0.905 |
| 78 | F | 1.006 |
| 79 | T | 0.984 |
| 80 | G | 0.951 |
| 81 | E | 0.960 |
| 82 | I | 0.478 |
| 83 | S | 0.840 |
| 84 | P | 0.469 |
| 85 | G | 0.375 |
| 86 | M | 0.744 |
| 87 | I | 0.557 |
| 88 | K | 0.848 |
| 89 | D | 0.878 |
| 90 | C | 0.415 |
| 91 | G | 0.951 |
| 92 | A | 0.529 |
| 93 | T | 0.368 |
| 94 | W | 1.013 |
| 95 | V | 0.930 |
| 96 | V | 0.676 |
| 97 | L | 0.959 |
| 98 | G | 0.951 |
| 99 | H | 1.017 |
| 100 | S | 0.977 |
| 101 | E | 0.960 |
| 102 | R | 1.000 |
| 103 | R | 1.000 |
| 104 | H | 0.522 |
| 105 | V | 0.463 |
| 106 | : | 0.494 * |
| 107 | : | 0.494 * |
| 108 | F | 0.650 |

**(C) 109-162 pos**

| Pos | AA | Score |
|---|---|---|
| 109 | G | 0.682 |
| 110 | E | 0.960 |
| 111 | S | 0.700 |
| 112 | D | 0.765 |
| 113 | E | 0.484 |
| 114 | L | 0.508 |
| 115 | V | 0.726 |
| 116 | G | 0.643 |
| 117 | Q | 0.499 |
| 118 | K | 0.966 |
| 119 | V | 0.492 |
| 120 | A | 0.465 |
| 121 | H | 0.569 |
| 122 | A | 0.905 |
| 123 | L | 0.824 |
| 124 | A | 0.381 |
| 125 | : | 0.494 * |
| 126 | E | 0.491 |
| 127 | G | 0.813 |
| 128 | L | 0.566 |
| 129 | G | 0.573 |
| 130 | V | 0.930 |
| 131 | I | 0.849 |
| 132 | A | 0.401 |
| 133 | C | 1.018 |
| 134 | I | 0.726 |
| 135 | G | 0.951 |
| 136 | E | 0.960 |
| 137 | K | 0.598 |
| 138 | L | 0.819 |
| 139 | D | 0.782 |
| 140 | E | 0.850 |
| 141 | R | 0.856 |
| 142 | E | 0.816 |
| 143 | A | 0.766 |
| 144 | G | 0.951 |
| 145 | I | 0.387 |
| 146 | T | 0.984 |
| 147 | E | 0.283 |
| 148 | K | 0.454 |
| 149 | V | 0.930 |
| 150 | V | 0.692 |
| 151 | F | 0.440 |
| 152 | E | 0.562 |
| 153 | Q | 0.993 |
| 154 | T | 0.623 |
| 155 | K | 0.424 |
| 156 | V | 0.712 |
| 157 | I | 0.491 |
| 158 | A | 0.507 |
| 159 | D | 0.569 |
| 160 | N | 0.511 |
| 161 | V | 0.458 |
| 162 | : | 0.494 * |

**(D) 163-216 pos**

| Pos | AA | Score |
|---|---|---|
| 163 | : | 0.494 * |
| 164 | K | 0.568 |
| 165 | D | 0.718 |
| 166 | W | 1.013 |
| 167 | S | 0.617 |
| 168 | K | 0.558 |
| 169 | V | 0.815 |
| 170 | V | 0.930 |
| 171 | L | 0.518 |
| 172 | A | 0.905 |
| 173 | Y | 1.011 |
| 174 | E | 0.960 |
| 175 | P | 0.997 |
| 176 | V | 0.930 |
| 177 | W | 1.013 |
| 178 | A | 0.905 |
| 179 | I | 0.970 |
| 180 | G | 0.951 |
| 181 | T | 0.984 |
| 182 | G | 0.951 |
| 183 | K | 0.718 |
| 184 | T | 0.535 |
| 185 | A | 0.905 |
| 186 | T | 0.870 |
| 187 | P | 0.880 |
| 188 | Q | 0.514 |
| 189 | Q | 0.853 |
| 190 | A | 0.905 |
| 191 | Q | 0.852 |
| 192 | E | 0.689 |
| 193 | V | 0.788 |
| 194 | H | 1.017 |
| 195 | E | 0.641 |
| 196 | K | 0.351 |
| 197 | L | 0.592 |
| 198 | R | 1.000 |
| 199 | G | 0.319 |
| 200 | W | 0.792 |
| 201 | L | 0.808 |
| 202 | K | 0.476 |
| 203 | S | 0.525 |
| 204 | N | 0.690 |
| 205 | V | 0.471 |
| 206 | S | 0.656 |
| 207 | D | 0.532 |
| 208 | A | 0.319 |
| 209 | V | 0.633 |
| 210 | A | 0.796 |
| 211 | Q | 0.405 |
| 212 | S | 0.575 |
| 213 | T | 0.680 |
| 214 | R | 1.000 |
| 215 | I | 0.970 |
| 216 | I | 0.586 |

**(E) 216-264 pos**

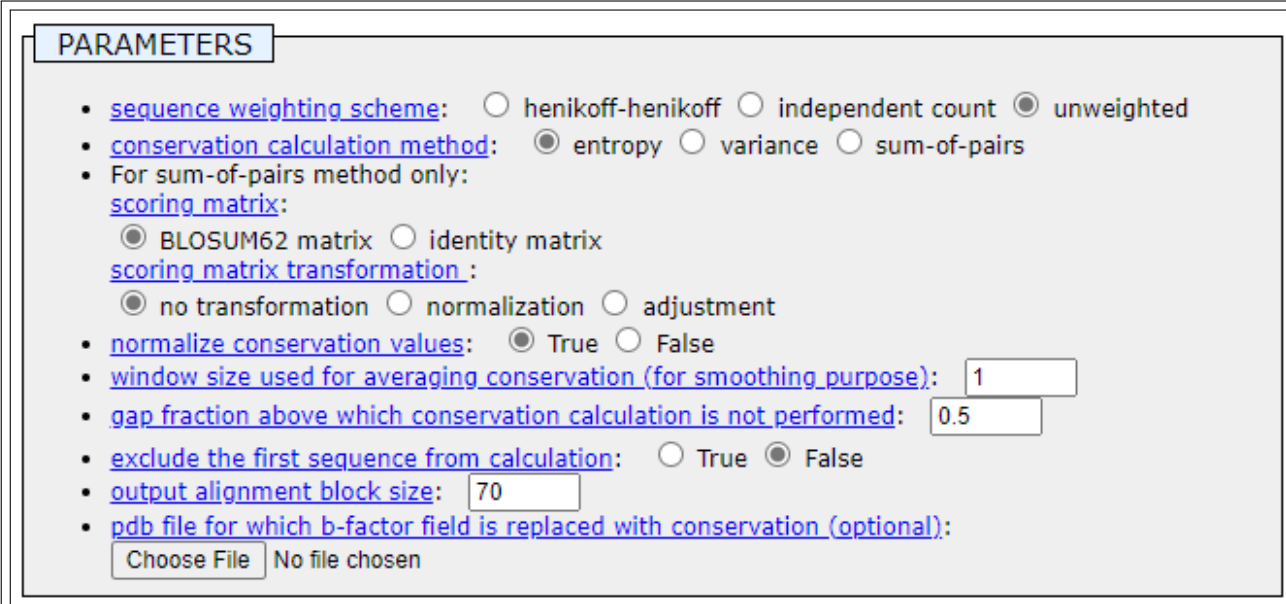| Pos | AA | Score |
|---|---|---|
| 217 | Y | 1.011 |
| 218 | G | 0.951 |
| 219 | G | 0.951 |
| 220 | S | 0.977 |
| 221 | V | 0.776 |
| 222 | T | 0.711 |
| 223 | G | 0.521 |
| 224 | A | 0.518 |
| 225 | T | 0.814 |
| 226 | C | 0.650 |
| 227 | K | 0.429 |
| 228 | E | 0.596 |
| 229 | L | 0.823 |
| 230 | A | 0.520 |
| 231 | S | 0.427 |
| 232 | Q | 0.601 |
| 233 | P | 0.707 |
| 234 | D | 0.852 |
| 235 | V | 0.638 |
| 236 | D | 0.847 |
| 237 | G | 0.951 |
| 238 | F | 1.006 |
| 239 | L | 0.959 |
| 240 | V | 0.930 |
| 241 | G | 0.951 |
| 242 | G | 0.951 |
| 243 | A | 0.905 |
| 244 | S | 0.977 |
| 245 | L | 0.829 |
| 246 | K | 0.966 |
| 247 | P | 0.997 |
| 248 | E | 0.713 |
| 249 | F | 1.006 |
| 250 | V | 0.678 |
| 251 | D | 0.628 |
| 252 | I | 0.970 |
| 253 | I | 0.824 |
| 254 | N | 0.714 |
| 255 | A | 0.663 |
| 256 | K | 0.457 |
| 257 | Q | 0.451 |
| 258 | : | 0.494 * |
| 259 | : | 0.494 * |
| 260 | : | 0.494 * |
| 261 | : | 0.494 * |
| 262 | : | 0.494 * |
| 263 | : | 0.494 * |
| 264 | : | 0.494 * |

FIGURE 22. Positional conservation scores for set 2

(v) Normalize the scores obtained with (i)

The parameters for this scenario are enclosed in the image below. The important parameters are as follows:

- sequence weighting scheme
- conservation calculation method
- scoring matrix (for sum of pairs method only)
- scoring matrix transformation (for sum of pairs method only)
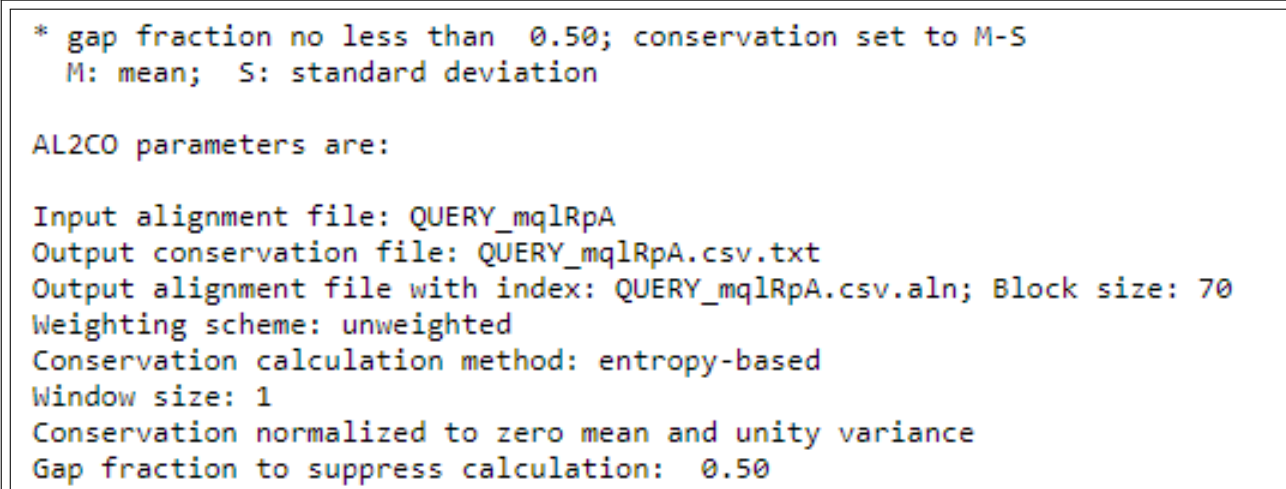- normalize conservation values



FIGURE 23. Parameters for the given scenario (v)

The AL2CO gives the list of positional conservation values and the alignment with integer conservation indices. The question asks to only calculate the positional conservation values.
The window of positional conservation values generates the following set of parameters which are the ones displayed in the above image.
It also displays some parameters taken into consideration to compute the desired positional conservation scores.



FIGURE 24. Parameter output display in window

## Figure 25. Positional conservation scores for set 1

### (A) 1-35 pos

| Pos | Res | Score |
|---|---|---|
| 1 | M | 0.943 |
| 2 | V | 0.253 |
| 3 | L | 0.943 |
| 4 | S | 0.943 |
| 5 | P | -1.597 |
| 6 | A | -1.685 |
| 7 | D | 0.943 |
| 8 | K | 0.943 |
| 9 | T | -1.705 |
| 10 | N | 0.253 |
| 11 | V | -0.385 |
| 12 | K | 0.253 |
| 13 | A | -1.399 |
| 14 | A | -1.685 |
| 15 | W | 0.253 |
| 16 | G | -1.134 |
| 17 | K | 0.943 |
| 18 | V | -0.618 |
| 19 | G | 0.253 |
| 20 | A | -1.399 |
| 21 | H | -1.064 |
| 22 | A | -0.131 |
| 23 | G | -1.399 |
| 24 | E | -1.114 |
| 25 | Y | -1.064 |
| 26 | G | 0.943 |
| 27 | A | -1.064 |
| 28 | E | 0.253 |
| 29 | A | 0.253 |
| 30 | L | 0.943 |
| 31 | E | -1.064 |
| 32 | R | 0.943 |
| 33 | M | 0.253 |
| 34 | F | 0.253 |
| 35 | L | -1.991 |

### (B) 36-70 pos

| Pos | Res | Score |
|---|---|---|
| 36 | S | -2.154 |
| 37 | F | -0.131 |
| 38 | P | 0.943 |
| 39 | T | -0.131 |
| 40 | T | 0.943 |
| 41 | K | 0.943 |
| 42 | T | 0.943 |
| 43 | Y | 0.943 |
| 44 | F | 0.943 |
| 45 | P | -0.131 |
| 46 | H | 0.943 |
| 47 | F | -0.417 |
| 48 | : | -1.000 * |
| 49 | D | 0.253 |
| 50 | L | -0.385 |
| 51 | S | -0.417 |
| 52 | H | -0.385 |
| 53 | G | 0.943 |
| 54 | S | 0.943 |
| 55 | A | -0.417 |
| 56 | Q | 0.253 |
| 57 | V | 0.253 |
| 58 | K | 0.943 |
| 59 | G | -1.006 |
| 60 | H | 0.943 |
| 61 | G | 0.943 |
| 62 | K | -0.417 |
| 63 | K | 0.943 |
| 64 | V | 0.943 |
| 65 | A | -0.778 |
| 66 | D | -0.778 |
| 67 | A | 0.253 |
| 68 | L | 0.253 |
| 69 | T | -1.006 |
| 70 | N | -2.848 |

### (C) 71-105 pos

| Pos | Res | Score |
|---|---|---|
| 71 | A | 0.943 |
| 72 | V | -0.385 |
| 73 | A | -2.670 |
| 74 | H | 0.253 |
| 75 | V | -1.312 |
| 76 | D | 0.253 |
| 77 | D | 0.943 |
| 78 | M | -0.778 |
| 79 | P | -0.417 |
| 80 | N | -0.778 |
| 81 | A | 0.253 |
| 82 | L | 0.943 |
| 83 | S | 0.253 |
| 84 | A | -2.262 |
| 85 | L | 0.943 |
| 86 | S | 0.943 |
| 87 | D | 0.253 |
| 88 | L | 0.253 |
| 89 | H | 0.943 |
| 90 | A | 0.943 |
| 91 | H | -0.778 |
| 92 | K | 0.943 |
| 93 | L | 0.943 |
| 94 | R | 0.943 |
| 95 | V | 0.943 |
| 96 | D | 0.943 |
| 97 | P | 0.943 |
| 98 | V | 0.253 |
| 99 | N | 0.943 |
| 100 | F | 0.943 |
| 101 | K | 0.943 |
| 102 | L | -0.417 |
| 103 | L | 0.943 |
| 104 | S | -0.417 |
| 105 | H | 0.253 |

### (D) 106-143 pos

| Pos | Res | Score |
|---|---|---|
| 106 | C | 0.253 |
| 107 | L | -0.417 |
| 108 | L | 0.253 |
| 109 | V | -0.131 |
| 110 | T | -0.131 |
| 111 | L | -0.417 |
| 112 | A | 0.253 |
| 113 | A | -2.955 |
| 114 | H | 0.253 |
| 115 | L | -1.175 |
| 116 | P | 0.943 |
| 117 | A | -3.241 |
| 118 | E | -1.134 |
| 119 | F | 0.253 |
| 120 | T | 0.253 |
| 121 | P | 0.943 |
| 122 | A | -0.778 |
| 123 | V | -1.064 |
| 124 | H | 0.943 |
| 125 | A | 0.253 |
| 126 | S | 0.943 |
| 127 | L | 0.943 |
| 128 | D | 0.943 |
| 129 | K | 0.943 |
| 130 | F | 0.943 |
| 131 | L | 0.253 |
| 132 | A | -1.006 |
| 133 | S | -2.262 |
| 134 | V | 0.943 |
| 135 | S | -0.417 |
| 136 | T | -0.417 |
| 137 | V | 0.253 |
| 138 | L | 0.943 |
| 139 | T | 0.253 |
| 140 | S | 0.253 |
| 141 | K | 0.253 |
| 142 | Y | 0.943 |
| 143 | R | 0.943 |

FIGURE 25. Positional conservation scores for set 1

Above is the image of the positional conservation scores for the set 1 proteins. Below is the image of the positional conservation scores for the set 2 proteins.

## Figure 26. Positional conservation scores for set 2

### (A) 1-54 pos

| Pos | Res | Score |
|---|---|---|
| 1 | M | -1.000 * |
| 2 | A | -1.000 * |
| 3 | P | -0.560 |
| 4 | S | -1.507 |
| 5 | R | 0.529 |
| 6 | K | -0.125 |
| 7 | F | 0.529 |
| 8 | F | -0.125 |
| 9 | V | 0.529 |
| 10 | G | 0.529 |
| 11 | G | 0.529 |
| 12 | N | 1.215 |
| 13 | W | 0.580 |
| 14 | K | 1.215 |
| 15 | M | -0.492 |
| 16 | N | 1.215 |
| 17 | G | 0.580 |
| 18 | R | -1.173 |
| 19 | K | -1.159 |
| 20 | Q | -1.560 |
| 21 | S | -1.666 |
| 22 | L | -0.718 |
| 23 | G | -1.841 |
| 24 | E | -0.612 |
| 25 | L | -0.036 |
| 26 | I | -0.543 |
| 27 | G | -1.946 |
| 28 | T | -0.612 |
| 29 | L | 0.580 |
| 30 | N | 0.580 |
| 31 | A | -0.612 |
| 32 | A | -0.031 |
| 33 | K | -2.227 |
| 34 | V | -1.666 |
| 35 | P | -0.879 |
| 36 | A | -1.841 |
| 37 | D | -1.560 |
| 38 | : | -1.000 * |
| 39 | : | -1.000 * |
| 40 | T | -0.036 |
| 41 | E | -0.331 |
| 42 | V | 0.580 |
| 43 | V | 1.215 |
| 44 | C | -0.718 |
| 45 | A | -0.612 |
| 46 | P | 0.580 |
| 47 | P | 0.580 |
| 48 | T | -0.998 |
| 49 | A | -1.841 |
| 50 | Y | -0.031 |
| 51 | I | -0.036 |
| 52 | D | -0.879 |
| 53 | F | -1.454 |
| 54 | A | -0.879 |

### (B) 55-108 pos

| Pos | Res | Score |
|---|---|---|
| 55 | R | -0.492 |
| 56 | Q | -0.998 |
| 57 | K | -1.385 |
| 58 | L | 0.580 |
| 59 | : | -1.000 * |
| 60 | D | -1.560 |
| 61 | P | -0.879 |
| 62 | K | -1.385 |
| 63 | I | -1.173 |
| 64 | A | -2.121 |
| 65 | V | -0.612 |
| 66 | A | -0.031 |
| 67 | A | 0.580 |
| 68 | Q | 1.215 |
| 69 | N | 1.215 |
| 70 | C | -0.036 |
| 71 | Y | -0.331 |
| 72 | K | -0.741 |
| 73 | V | -0.879 |
| 74 | T | -1.560 |
| 75 | N | -1.279 |
| 76 | G | 1.215 |
| 77 | A | 1.215 |
| 78 | F | 1.215 |
| 79 | T | 1.215 |
| 80 | G | 1.215 |
| 81 | E | 1.215 |
| 82 | I | -0.879 |
| 83 | S | 0.580 |
| 84 | P | -0.879 |
| 85 | G | -1.454 |
| 86 | M | -0.031 |
| 87 | I | -0.492 |
| 88 | K | 0.580 |
| 89 | D | 0.580 |
| 90 | C | -1.560 |
| 91 | G | 1.215 |
| 92 | A | -0.612 |
| 93 | T | -1.841 |
| 94 | W | 1.215 |
| 95 | V | 1.215 |
| 96 | V | -0.036 |
| 97 | L | 1.215 |
| 98 | G | 1.215 |
| 99 | H | 1.215 |
| 100 | S | 1.215 |
| 101 | E | 1.215 |
| 102 | R | 1.215 |
| 103 | R | 1.215 |
| 104 | H | -0.879 |
| 105 | V | -1.104 |
| 106 | : | -1.000 * |
| 107 | : | -1.000 * |
| 108 | F | -0.331 |

### (C) 109-162 pos

| Pos | Res | Score |
|---|---|---|
| 109 | G | -0.031 |
| 110 | E | 1.215 |
| 111 | S | -0.031 |
| 112 | D | 0.250 |
| 113 | E | -0.879 |
| 114 | L | -0.879 |
| 115 | I | 0.250 |
| 116 | G | -0.036 |
| 117 | Q | -1.173 |
| 118 | K | 1.215 |
| 119 | V | -0.598 |
| 120 | A | -1.159 |
| 121 | H | -0.879 |
| 122 | A | 1.215 |
| 123 | L | 0.580 |
| 124 | A | -1.385 |
| 125 | : | -1.000 * |
| 126 | E | -0.718 |
| 127 | G | 0.580 |
| 128 | L | -0.612 |
| 129 | G | -0.492 |
| 130 | V | 1.215 |
| 131 | I | 0.580 |
| 132 | A | -1.159 |
| 133 | C | 1.215 |
| 134 | I | 0.250 |
| 135 | G | 1.215 |
| 136 | E | 1.215 |
| 137 | K | -0.543 |
| 138 | L | 0.580 |
| 139 | D | 0.056 |
| 140 | E | 0.580 |
| 141 | R | 0.580 |
| 142 | E | 0.580 |
| 143 | A | 0.580 |
| 144 | G | 1.215 |
| 145 | I | -1.841 |
| 146 | T | 1.215 |
| 147 | E | -2.121 |
| 148 | K | -1.173 |
| 149 | V | 1.215 |
| 150 | V | -0.031 |
| 151 | L | -1.104 |
| 152 | E | -0.998 |
| 153 | Q | 1.215 |
| 154 | T | -0.543 |
| 155 | K | -1.385 |
| 156 | V | 0.056 |
| 157 | I | -0.879 |
| 158 | A | -0.612 |
| 159 | D | -0.598 |
| 160 | N | -0.998 |
| 161 | V | -0.998 |
| 162 | : | -1.000 * |

### (D) 163-216 pos

| Pos | Res | Score |
|---|---|---|
| 163 | : | -1.000 * |
| 164 | K | -0.612 |
| 165 | D | -0.031 |
| 166 | W | 1.215 |
| 167 | S | -0.612 |
| 168 | K | -0.998 |
| 169 | V | 0.580 |
| 170 | V | 1.215 |
| 171 | L | -0.787 |
| 172 | A | 1.215 |
| 173 | Y | 1.215 |
| 174 | E | 1.215 |
| 175 | P | 1.215 |
| 176 | V | 1.215 |
| 177 | W | 1.215 |
| 178 | A | 1.215 |
| 179 | I | 1.215 |
| 180 | G | 1.215 |
| 181 | T | 1.215 |
| 182 | G | 1.215 |
| 183 | K | -0.031 |
| 184 | T | -0.492 |
| 185 | A | 1.215 |
| 186 | T | 0.580 |
| 187 | P | 0.580 |
| 188 | Q | -0.998 |
| 189 | Q | 0.580 |
| 190 | A | 1.215 |
| 191 | Q | 0.580 |
| 192 | E | -0.031 |
| 193 | V | 0.580 |
| 194 | H | 1.215 |
| 195 | E | 0.056 |
| 196 | K | -1.841 |
| 197 | L | -0.331 |
| 198 | R | 1.215 |
| 199 | G | -1.841 |
| 200 | W | 0.250 |
| 201 | L | 0.580 |
| 202 | K | -0.879 |
| 203 | S | -1.159 |
| 204 | N | 0.056 |
| 205 | V | -0.879 |
| 206 | S | -0.331 |
| 207 | D | -0.998 |
| 208 | A | -1.841 |
| 209 | V | -0.031 |
| 210 | A | 0.580 |
| 211 | Q | -1.454 |
| 212 | S | -0.612 |
| 213 | T | 0.056 |
| 214 | R | 1.215 |
| 215 | I | 1.215 |
| 216 | I | -0.492 |

### (E) 216-264 pos

| Pos | Res | Score |
|---|---|---|
| 217 | Y | 1.215 |
| 218 | G | 1.215 |
| 219 | G | 1.215 |
| 220 | S | 1.215 |
| 221 | V | 0.580 |
| 222 | T | -0.031 |
| 223 | G | -0.543 |
| 224 | A | -0.612 |
| 225 | T | 0.056 |
| 226 | C | -0.036 |
| 227 | K | -1.159 |
| 228 | E | -0.331 |
| 229 | L | 0.580 |
| 230 | A | -0.612 |
| 231 | S | -1.454 |
| 232 | Q | -0.492 |
| 233 | P | -0.031 |
| 234 | D | 0.580 |
| 235 | V | -0.036 |
| 236 | D | 0.580 |
| 237 | G | 1.215 |
| 238 | F | 1.215 |
| 239 | L | 1.215 |
| 240 | V | 1.215 |
| 241 | G | 1.215 |
| 242 | G | 1.215 |
| 243 | A | 1.215 |
| 244 | S | 1.215 |
| 245 | L | 0.580 |
| 246 | K | 1.215 |
| 247 | P | 1.215 |
| 248 | E | -0.031 |
| 249 | F | 1.215 |
| 250 | V | -0.031 |
| 251 | D | -0.331 |
| 252 | I | 1.215 |
| 253 | I | 0.580 |
| 254 | N | -0.031 |
| 255 | A | 0.250 |
| 256 | K | -1.454 |
| 257 | Q | -1.385 |
| 258 | : | -1.000 * |
| 259 | : | -1.000 * |
| 260 | : | -1.000 * |
| 261 | : | -1.000 * |
| 262 | : | -1.000 * |
| 263 | : | -1.000 * |
| 264 | : | -1.000 * |

FIGURE 26. Positional conservation scores for set 2

**Question 2.** Tabulate the topmost 10 residues with highest and lowest conservation scores (in both Set1 and Set 2) obtained with method (i).

**Solution.** First the output obtained in the Q1 (i) is taken. The output has three columns as seen from its image. The position on sequence alignment, the residue, and its position conservation score.

Now, I have copied the above text from the AL2CO server's positional conservation score output to a text file. I have now written a code to read the file and then convert it into a dataframe, which has then been sorted to extract the top 10 and bottom 10 values as desired. The code is given below:

```python
import pandas as pd

new_line_list = []

# file.txt has the the positive conservation values the way it is shown in
    AL2CO server

with open("file.txt", 'r') as f:

    # Read the text file and store in list

    line = f.read()
    line_list = line.split("\n")

    # Remove the additional spaces and positions with gap alignments

    for i in range(len(line_list)):
        line_i = line_list[i].split(" ")
        line_i = [str(value) for value in line_i if value != ""]
        if line_i[1] != ":":
            new_line_list.append(line_i)

# The above obtained list is then converted into a dataframe for easy
    readability and easier to perform sorting techniques and observe

df = pd.DataFrame(new_line_list)

# The first and third columns are converted to numeric data type because they
    represent the alignment position and position conservation scores
    respectively

df[2] = pd.to_numeric(df[2])
df[0] = pd.to_numeric(df[0])

# Sort the values in the descending order as per the last column, which is
    the column for positional conservation scores
final_df = df.sort_values(by=[2], ascending=False)

# Display the top 10 and bottom 10 residues
# Top 10 will be the highest scores
# Bottom 10 will be the lowest scores
final_df.head(10)
final_df.tail(10)
```

The last two lines will print the top 10 residues with highest and lowest positional conservation scores.

(A) Highest scores          (B) Lowest scores

FIGURE 27. Top 10 residues with highest and lowest scores in set 1

### The topmost 10 residues with highest and lowest conservation scores in set 1

In the following table, I have enclosed the readings from the above image with the top 10 residues with the highest and the lowest residues.

The results in this table are with respect to set 1 only.

| Top 10 residues with the highest conservation scores | | | Top 10 residues with the lowest conservation scores | | |
|---|---|---|---|---|---|
| Position | Residue | Score | Position | Residue | Score |
| 1 | Methionine | 0.0 | 117 | Alanine | -1.846 |
| 42 | Threonine | 0.0 | 113 | Alanine | -1.720 |
| 92 | Lysine | 0.0 | 70 | Asparagine | -1.673 |
| 90 | Alanine | 0.0 | 73 | Alanine | -1.594 |
| 89 | Histidine | 0.0 | 133 | Serine | -1.414 |
| 86 | Serine | 0.0 | 84 | Alanine | -1.414 |
| 85 | Leucine | 0.0 | 36 | Serine | -1.367 |
| 82 | Leucine | 0.0 | 35 | Leucine | -1.295 |
| 77 | Aspartate | 0.0 | 9 | Threonine | -1.169 |
| 71 | Alanine | 0.0 | 14 | Alanine | -1.160 |

(A) Highest scores      (B) Lowest scores

FIGURE 28. Top 10 residues with highest and lowest scores in set 2

## The topmost 10 residues with highest and lowest conservation scores in set 2

In the following table, I have enclosed the readings from the above image with the top 10 residues with the highest and the lowest residues.

The results in this table are with respect to set 2 only.

| Top 10 residues with the highest conservation scores | | | Top 10 residues with the lowest conservation scores | | |
|---|---|---|---|---|---|
| Position | Residue | Score | Position | Residue | Score |
| 240 | Valine | 0.0 | 33 | Lysine | -1.889 |
| 172 | Alanine | 0.0 | 147 | Glutamate | -1.831 |
| 78 | Phenylalanine | 0.0 | 64 | Alanine | -1.831 |
| 77 | Alanine | 0.0 | 27 | Glycine | -1.735 |
| 76 | Glycine | 0.0 | 49 | Alanine | -1.677 |
| 237 | Glycine | 0.0 | 93 | Threonine | -1.677 |
| 220 | Serine | 0.0 | 145 | Isoleucine | -1.677 |
| 166 | Tryptophan | 0.0 | 199 | Glycine | -1.677 |
| 69 | Asparagine | 0.0 | 196 | Lysine | -1.677 |
| 68 | Glutamine | 0.0 | 208 | Alanine | -1.677 |

**Question 3.** Write a program to compute the conservation score from MSA using unweighted frequency, and entropy, variance and sum of pairs-based measures.

**Solution.** The code for computing the conservation score from the MSA using the above mentioned techniques is given below:

```
1 # Creating a blosum62 matrix for sum of pairs measure
2 # The blosum matrix is read from a text file stored locally
3 # blosum_dict initializes the index for each amino acid in matrix
4 blosum_dict = {"A":0,"R":1,"N":2,"D":3,"C":4,"Q":5,"E":6,"G":7,"H":8,"I":9,
5 "L":10,"K":11,"M":12,"F":13,"P":14,"S":15,"T":16,"W":17,"Y":18,"V":19,"-":20}
6
7 blosum_matrix = []
8 with open("blosum62.txt", 'r') as f:
9     line = f.read()
10    line_list = line.split("\n")
11
12    for i in range(len(line_list)):
13        line_i = line_list[i].split(" ")
14        line_i = [int(value) for value in line_i if value != ""]
15        blosum_matrix.append(line_i)
```

LISTING 1. Blosum62 Matrix

```
1 import math
2
3 amino_acids_track =
4 { "0": "G",   "1": "A",   "2": "V",   "3": "L",   "4": "I",
5   "5": "T",   "6": "S",   "7": "M",   "8": "C",   "9": "P",
6  "10": "F",  "11": "Y",  "12": "W",  "13": "H",  "14": "K",
7  "15": "R",  "16": "D",  "17": "E",  "18": "N",  "19": "Q",
8  "20": "-"}
9
10 amino_acids =
11 {"G":  0, "A":  1, "V":  2, "L":  3, "I":  4, "T":  5,
12  "S":  6, "M":  7, "C":  8, "P":  9, "F": 10, "Y": 11,
13  "W": 12, "H": 13, "K": 14, "R": 15, "D": 16, "E": 17,
14  "N": 18, "Q": 19, "-": 20}
15
16 # Creating lists to store the conservation scores using the weighing scheme
       of unweighted amino acid frequencies with the three conservation
       calculation methods given in the question, entropy, variance and sum-of-
       pairs
17
18 entropy_score_list = []
19 variance_score_list = []
20 sum_pairs_score_list = []
21
22 # Iterating through entire aligned sequence
23 for i in range(len(seqs[0])):
24     entropy_score = 0
25     variance_score = 0
26     sum_pairs_score = 0
27
28     each_freq = [0 for k in range(21)]
29     overall_freq = [0 for k in range(21)]
30     total_non_aligns = 0
```

```
31
32    # Iterating through each column of sequences
33    for j in range(len(seqs)):
34        each_freq[amino_acids[seqs[j][i]]] += 1
35
36    for j in range(len(each_freq)):
37        each_freq[j] /= (len(seqs)-each_freq[20])
38
39    for j in range(len(seqs)):
40        for k in range(len(seqs[0])):
41            overall_freq[amino_acids[seqs[j][k]]] += 1
42            if seqs[j][k] == "-":
43                total_non_aligns += 1
44
45    # For the sake of variance based conservation score
46    for j in range(len(overall_freq)):
47        overall_freq[j] /= (len(seqs) * len(seqs[0]) - total_non_aligns)
48
49    # Calculation of conservation score via the entropy method
50    for j in range(len(each_freq)-1):
51        if each_freq[j] != 0:
52            entropy_score += each_freq[j] * math.log(each_freq[j])
53
54    # Calculation of conservation score via the variance method
55    for j in range(len(each_freq)-1):
56        variance_score += abs(each_freq[j] - overall_freq[j])**2
57    variance_score = variance_score**0.5
58
59    # Calculation of conservation score via the sum-of-pairs method
60    for j in range(len(each_freq)-1):
61        for k in range(len(each_freq)-1):
62            sum_pairs_score += (each_freq[j]*each_freq[k]*
63            blosum_matrix[blosum_dict[amino_acids_track[str(j)]]]
64                        [blosum_dict[amino_acids_track[str(k)]]])
65
66    # Rounding off the values to 4 decimal places (just like AL2CO)
67    entropy_score = round(entropy_score, 4)
68    variance_score = round(variance_score, 4)
69    sum_pairs_score = round(sum_pairs_score, 4)
70
71    # Appending the scores for each column to the list
72    entropy_score_list.append(entropy_score)
73    variance_score_list.append(variance_score)
74    sum_pairs_score_list.append(sum_pairs_score)
75
76 # Print the output lists
77 print(entropy_score_list)
78 print(variance_score_list)
79 print(sum_pairs_score_list)
```

LISTING 2. Computing the different metric conservation scores

The above code generates the lists of entropy based, variance based, and sum-of-pairs based scores. The input sequences used were same as the one that is used in the Set 1 of the first question. The value of output generated by the code is compared against the ones generated by the AL2CO server. They match.

The input passed to the above code is the list of 11 sequences that were passed as input in the first question. They are: P69905, P01946, P01942, P01966, P01958, P01959, P01965, P06635, P60529, P80043 and P01980. These correspond to HBA_HUMAN, HBA_RAT, HBA_MOUSE, HBA_BOVIN, HBA_HORSE, HBA_EQUAS, HBA_PIG, HBA_PONPY, HBA_CANLF, HBA_TREBE, HBA_APTFO respectively. The inputs and their outputs are given below:

**Input given to the above code**

```
 1 seqs[0] =
 2 "MVLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHF-
 3 DLSHGSAQVKGHGKKVADALTNAVAHVDDMPNALSALSDLHAHKLRVDP
 4 VNFKLLSHCLLVTLAAHLPAEFTPAVHASLDKFLASVSTVLTSKYR"
 5 seqs[1] =
 6 "MVLSADDKTNIKNCWGKIGGHGGEYGEEALQRMFAAFPTTKTYFSHI-
 7 DVSPGSAQVKAHGKKVADALAKAADHVEDLPGALSTLSDLHAHKLRVDP
 8 VNFKFLSHCLLVTLACHHPGDFTPAMHASLDKFLASVSTVLTSKYR"
 9 seqs[2] =
10 "MVLSGEDKSNIKAAWGKIGGHGAEYGAEALERMFASFPTTKTYFPHF-
11 DVSHGSAQVKGHGKKVADALASAAGHLDDLPGALSALSDLHAHKLRVDP
12 VNFKLLSHCLLVTLASHHPADFTPAVHASLDKFLASVSTVLTSKYR"
13 seqs[3] =
14 "MVLSAADKGNVKAAWGKVGGHAAEYGAEALERMFLSFPTTKTYFPHF-
15 DLSHGSAQVKGHGAKVAAALTKAVEHLDDLPGALSELSDLHAHKLRVDP
16 VNFKLLSHSLLVTLASHLPSDFTPAVHASLDKFLANVSTVLTSKYR"
17 seqs[4] =
18 "MVLSAADKTNVKAAWSKVGGHAGEYGAEALERMFLGFPTTKTYFPHF-
19 DLSHGSAQVKAHGKKVGDALTLAVGHLDDLPGALSNLSDLHAHKLRVDP
20 VNFKLLSHCLLSTLAVHLPNDFTPAVHASLDKFLSSVSTVLTSKYR"
21 seqs[5] =
22 "MVLSAADKTNVKAAWSKVGGNAGEFGAEALERMFLGFPTTKTYFPHF-
23 DLSHGSAQVKAHGKKVGDALTLAVGHLDDLPGALSNLSDLHAHKLRVDP
24 VNFKLLSHCLLSTLAVHLPNDFTPAVHASLDKFLSTVSTVLTSKYR"
25 seqs[6] =
26 "-VLSAADKANVKAAWGKVGGQAGAHGAEALERMFLGFPTTKTYFPHF-
27 NLSHGSDQVKAHGQKVADALTKAVGHLDDLPGALSALSDLHAHKLRVDP
28 VNFKLLSHCLLVTLAAHHPDDFNPSVHASLDKFLANVSTVLTSKYR"
29 seqs[7] =
30 "MVLSPADKTNVKTAWGKVGAHAGDYGAEALERMFLSFPTTKTYFPHF-
31 DLSHGSAQVKDHGKKVADALTNAVAHVDDMPNALSALSDLHAHKLRVDP
32 VNFKLLSHCLLVTLAAHLPAEFTPAVHASLDKFLASVSTVLTSKYR"
33 seqs[8] =
34 "-VLSPADKTNIKSTWDKIGGHAGDYGGEALDRTFQSFPTTKTYFPHF-
35 DLSPGSAQVKAHGKKVADALTTAVAHLDDLPGALSALSDLHAYKLRVDP
36 VNFKLLSHCLLVTLACHHPTEFTPAVHASLDKFFAAVSTVLTSKYR"
37 seqs[9] =
38 "-SLSDKDKAAVRALWSKIGKSADAIGNDALSRMIVVYPQTKTYFSHWP
39 DVTPGSPHIKAHGKKVMGGIALAVSKIDDLKTGLMELSEQHAYKLRVDP
40 ANFKILNHCILVVISTMFPKEFTPEAHVSLDKFLSGVALALAERYR"
41 seqs[10] =
42 "MVLSANDKSNVKSIFSKISSHAEEYGAETLERMFTTYPQTKTYFPHF-
43 DLHHGSAQVKAHGKKVAAALIEAANHIDDIAGALSKLSDLHAEKLRVDP
44 VNFKLLGQCFMVVVAIHHPSALTPEIHASLDKFLCAVGNVLTSKYR"
```

LISTING 3. Input given to the above code

```
 1  # Entropy based conservation score
 2  [0.0, -0.305, 0.0, 0.0, -1.121, -1.16, 0.0, 0.0, -1.169, -0.305,
 3  -0.586, -0.305, -1.034, -1.16, -0.305, -0.916, 0.0, -0.689, -0.305,
 4  -1.034, -0.886, -0.474, -1.034, -0.908, -0.886, 0.0, -0.886,
 5  -0.305, -0.305, 0.0, -0.886, 0.0, -0.305, -0.305, -1.295, -1.367,
 6  -0.474, 0.0, -0.474, 0.0, 0.0, 0.0, 0.0, 0.0, -0.474, 0.0, -0.6,
 7  0.0, -0.305, -0.586, -0.6, -0.586, 0.0, 0.0, -0.6, -0.305, -0.305,
 8  0.0, -0.86, 0.0, 0.0, -0.6, 0.0, 0.0, -0.76, -0.76, -0.305, -0.305,
 9  -0.86, -1.673, 0.0, -0.586, -1.594, -0.305, -0.995, -0.305, 0.0,
10  -0.76, -0.6, -0.76, -0.305, 0.0, -0.305, -1.414, 0.0, 0.0, -0.305,
11  -0.305, 0.0, 0.0, -0.76, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, -0.305, 0.0,
12  0.0, 0.0, -0.6, 0.0, -0.6, -0.305, -0.305, -0.6, -0.305, -0.474,
13  -0.474, -0.6, -0.305, -1.72, -0.305, -0.935, 0.0, -1.846, -0.916,
14  -0.305, -0.305, 0.0, -0.76, -0.886, 0.0, -0.305, 0.0, 0.0, 0.0,
15  0.0, 0.0, -0.305, -0.86, -1.414, 0.0, -0.6, -0.6, -0.305, 0.0,
16  -0.305, -0.305, -0.305, 0.0, 0.0]
17
18  # Variance based conservation score
19  [1.021, 0.865, 0.909, 0.953, 0.527, 0.561, 0.973, 0.956, 0.541,
20  0.917, 0.744, 0.871, 0.567, 0.56, 0.943, 0.612, 0.956, 0.691,
21  0.886, 0.613, 0.713, 0.74, 0.624, 0.665, 0.757, 0.977, 0.65, 0.922,
22  0.817, 0.909, 0.749, 1.015, 0.933, 0.904, 0.47, 0.473, 0.83, 0.993,
23  0.821, 0.977, 0.956, 0.977, 1.016, 0.988, 0.821, 0.969, 0.823,
24  0.993, 0.888, 0.674, 0.776, 0.745, 0.977, 0.953, 0.73, 0.936,
25  0.872, 0.956, 0.596, 0.969, 0.977, 0.777, 0.956, 0.956, 0.66,
26  0.704, 0.817, 0.825, 0.643, 0.36, 0.906, 0.701, 0.406, 0.879,
27  0.543, 0.888, 0.973, 0.675, 0.806, 0.733, 0.817, 0.909, 0.869,
28  0.451, 0.909, 0.953, 0.888, 0.826, 0.969, 0.906, 0.731, 0.956,
29  0.909, 1.015, 0.956, 0.973, 0.993, 0.86, 1.012, 0.988, 0.956,
30  0.739, 0.909, 0.781, 0.886, 0.937, 0.739, 0.825, 0.784, 0.806,
31  0.735, 0.814, 0.345, 0.885, 0.56, 0.993, 0.302, 0.633, 0.893,
32  0.891, 0.993, 0.66, 0.695, 0.969, 0.815, 0.953, 0.909, 0.973,
33  0.956, 0.988, 0.822, 0.595, 0.459, 0.956, 0.769, 0.797, 0.86,
34  0.909, 0.881, 0.868, 0.871, 1.016, 1.015]
35
36  # Sum-of-pairs based conservation score
37  [5.0, 3.008, 4.0, 4.0, 1.149, 1.157, 6.0, 5.0, 1.868, 4.661, 3.603,
38  4.504, 1.909, 1.521, 9.306, 2.264, 5.0, 3.504, 4.992, 2.405, 4.372,
39  2.876, 2.24, 2.455, 4.314, 6.0, 1.826, 4.512, 3.347, 4.0, 3.298,
40  5.0, 4.008, 4.992, 0.959, 1.132, 5.14, 7.0, 3.215, 5.0, 5.0, 5.0,
41  7.0, 6.0, 4.521, 8.0, 4.24, 7.0, 5.174, 2.81, 2.752, 3.959, 6.0,
42  4.0, 2.322, 4.198, 3.835, 5.0, 1.835, 8.0, 6.0, 3.405, 5.0, 4.0,
43  2.124, 2.694, 3.355, 3.669, 2.19, 0.124, 4.0, 2.413, 0.934, 6.488,
44  2.612, 5.331, 6.0, 3.14, 4.446, 3.149, 3.355, 4.0, 3.182, 0.678,
45  4.0, 4.0, 5.331, 3.017, 8.0, 4.0, 4.967, 5.0, 4.0, 5.0, 4.0, 6.0,
46  7.0, 3.339, 6.0, 6.0, 5.0, 3.058, 4.0, 2.926, 6.653, 7.306, 3.058,
47  3.678, 2.215, 3.479, 3.24, 3.504, 0.702, 6.322, 1.207, 7.0, 0.562,
48  3.008, 4.992, 4.182, 7.0, 2.182, 2.736, 8.0, 3.339, 4.0, 4.0, 6.0,
49  5.0, 6.0, 3.355, 2.289, 1.496, 4.0, 2.909, 3.231, 3.339, 4.0,
50  4.165, 3.347, 4.504, 7.0, 5.0]
```

LISTING 4. The conversation scores output generated by the above code

**Question 4.** Using the program written in Q3 (unweighted frequency and entropy-based measure), compare the MSA from Clustal Omega, MAFFT, and MUSCLE. Identify the residues with (i) similar and (ii) different conservation scores among the three alignment methods.

**Solution.** Below are the MSA obtained from the Clustal Omega, MAFFT, and MUSCLE for the protein sequences in set 1.

```
CLUSTAL O(1.2.4) multiple sequence alignment


HBA_TREBE      -SLSDKDKAAVRALWSKIGKSADAIGNDALSRMIVVYPQTKTYFSHWPDVTPGSPHIKAH      59
HBA_APTFO      MVLSANDKSNVKSIFSKISSHAEEYGAETLERMFTTYPQTKTYFPHF-DLHHGSAQVKAH      59
HBA_CANLF      -VLSPADKTNIKSTWDKIGGHAGDYGGEALDRTFQSFPTTKTYFPHF-DLSPGSAQVKAH      58
HBA_RAT        MVLSADDKTNIKNCWGKIGGHGGEYGEEALQRMFAAFPTTKTYFSHI-DVSPGSAQVKAH      59
HBA_PIG        -VLSAADKANVKAAWGKVGGQAGAHGAEALERMFLGFPTTKTYFPHF-NLSHGSDQVKAH      58
HBA_MOUSE      MVLSGEDKSNIKAAWGKIGGHGAEYGAEALERMFASFPTTKTYFPHF-DVSHGSAQVKGH      59
HBA_HORSE      MVLSAADKTNVKAAWSKVGGHAGEYGAEALERMFLGFPTTKTYFPHF-DLSHGSAQVKAH      59
HBA_EQUAS      MVLSAADKTNVKAAWSKVGGNAGEFGAEALERMFLGFPTTKTYFPHF-DLSHGSAQVKAH      59
HBA_BOVIN      MVLSAADKGNVKAAWGKVGGHAAEYGAEALERMFLSFPTTKTYFPHF-DLSHGSAQVKGH      59
HBA_HUMAN      MVLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHF-DLSHGSAQVKGH      59
HBA_PONPY      MVLSPADKTNVKTAWGKVGAHAGDYGAEALERMFLSFPTTKTYFPHF-DLSHGSAQVKDH      59
                        **   **   ::   :.*:.   .    * ::*.* :   :* ***** *   ::   ** ::* *


HBA_TREBE      GKKVMGGIALAVSKIDDLKTGLMELSEQHAYKLRVDPANFKILNHCILVVISTMFPKEFT     119
HBA_APTFO      GKKVAAALIEAANHIDDIAGALSKLSDLHAEKLRVDPVNFKLLGQCFMVVVAIHHPSALT     119
HBA_CANLF      GKKVADALTTAVAHLDDLPGALSALSDLHAYKLRVDPVNFKLLSHCLLVTLACHHPTEFT     118
HBA_RAT        GKKVADALAKAADHVEDLPGALSTLSDLHAHKLRVDPVNFKFLSHCLLVTLACHHPGDFT     119
HBA_PIG        GQKVADALTKAVGHLDDLPGALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHHPDDFN     118
HBA_MOUSE      GKKVADALASAAGHLDDLPGALSALSDLHAHKLRVDPVNFKLLSHCLLVTLASHHPADFT     119
HBA_HORSE      GKKVGDALTLAVGHLDDLPGALSNLSDLHAHKLRVDPVNFKLLSHCLLSTLAVHLPNDFT     119
HBA_EQUAS      GKKVGDALTLAVGHLDDLPGALSNLSDLHAHKLRVDPVNFKLLSHCLLSTLAVHLPNDFT     119
HBA_BOVIN      GAKVAAALTKAVEHLDDLPGALSELSDLHAHKLRVDPVNFKLLSHSLLVTLASHLPSDFT     119
HBA_HUMAN      GKKVADALTNAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPAEFT     119
HBA_PONPY      GKKVADALTNAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPAEFT     119
                        * **   .:   *.  :::*:   .*   **: ** ******.***:*.:.:: .::    *  :.


HBA_TREBE      PEAHVSLDKFLSGVALALAERYR   142
HBA_APTFO      PEIHASLDKFLCAVGNVLTSKYR   142
HBA_CANLF      PAVHASLDKFFAAVSTVLTSKYR   141
HBA_RAT        PAMHASLDKFLASVSTVLTSKYR   142
HBA_PIG        PSVHASLDKFLANVSTVLTSKYR   141
HBA_MOUSE      PAVHASLDKFLASVSTVLTSKYR   142
HBA_HORSE      PAVHASLDKFLSSVSTVLTSKYR   142
HBA_EQUAS      PAVHASLDKFLSTVSTVLTSKYR   142
HBA_BOVIN      PAVHASLDKFLANVSTVLTSKYR   142
HBA_HUMAN      PAVHASLDKFLASVSTVLTSKYR   142
HBA_PONPY      PAVHASLDKFLASVSTVLTSKYR   142
                        *   *.*****:. *.  .*:.:**
```

FIGURE 29. Output MSA by Clustal Omega for set 1

```
CLUSTAL format alignment by MAFFT FFT-NS-i (v7.487)


HBA_HUMAN      MVLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHF-DLSHGSAQVKGH
HBA_PONPY      MVLSPADKTNVKTAWGKVGAHAGDYGAEALERMFLSFPTTKTYFPHF-DLSHGSAQVKDH
HBA_BOVIN      MVLSAADKGNVKAAWGKVGGHAAEYGAEALERMFLSFPTTKTYFPHF-DLSHGSAQVKGH
HBA_HORSE      MVLSAADKTNVKAAWSKVGGHAGEYGAEALERMFLGFPTTKTYFPHF-DLSHGSAQVKAH
HBA_EQUAS      MVLSAADKTNVKAAWSKVGGNAGEFGAEALERMFLGFPTTKTYFPHF-DLSHGSAQVKAH
HBA_MOUSE      MVLSGEDKSNIKAAWGKIGGHGAEYGAEALERMFASFPTTKTYFPHF-DVSHGSAQVKGH
HBA_RAT        MVLSADDKTNIKNCWGKIGGHGGEYGEEALQRMFAAFPTTKTYFSHI-DVSPGSAQVKAH
HBA_PIG        -VLSAADKANVKAAWGKVGGQAGAHGAEALERMFLGFPTTKTYFPHF-NLSHGSDQVKAH
HBA_CANLF      -VLSPADKTNIKSTWDKIGGHAGDYGGEALDRTFQSFPTTKTYFPHF-DLSPGSAQVKAH
HBA_APTFO      MVLSANDKSNVKSIFSKISSHAEEYGAETLERMFTTYPQTKTYFPHF-DLHHGSAQVKAH
HBA_TREBE      -SLSDKDKAAVRALWSKIGKSADAIGNDALSRMIVVYPQTKTYFSHWPDVTPGSPHIKAH
                        **    **   ::    :.*:.   .    * ::*.* :   :* *****.*   ::   ** ::* *


HBA_HUMAN      GKKVADALTNAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPAEFT
HBA_PONPY      GKKVADALTNAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPAEFT
HBA_BOVIN      GAKVAAALTKAVEHLDDLPGALSELSDLHAHKLRVDPVNFKLLSHSLLVTLASHLPSDFT
HBA_HORSE      GKKVGDALTLAVGHLDDLPGALSNLSDLHAHKLRVDPVNFKLLSHCLLSTLAVHLPNDFT
HBA_EQUAS      GKKVGDALTLAVGHLDDLPGALSNLSDLHAHKLRVDPVNFKLLSHCLLSTLAVHLPNDFT
HBA_MOUSE      GKKVADALASAAGHLDDLPGALSALSDLHAHKLRVDPVNFKLLSHCLLVTLASHHPADFT
HBA_RAT        GKKVADALAKAADHVEDLPGALSTLSDLHAHKLRVDPVNFKFLSHCLLVTLACHHPGDFT
HBA_PIG        GQKVADALTKAVGHLDDLPGALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHHPDDFN
HBA_CANLF      GKKVADALTTAVAHLDDLPGALSALSDLHAYKLRVDPVNFKLLSHCLLVTLACHHPTEFT
HBA_APTFO      GKKVAAALIEAANHIDDIAGALSKLSDLHAEKLRVDPVNFKLLGQCFMVVVAIHHPSALT
HBA_TREBE      GKKVMGGIALAVSKIDDLKTGLMELSEQHAYKLRVDPANFKILNHCILVVISTMFPKEFT
                        * **   .:   *. :::*:   .*   **: ** ******.***:*.:.:: .::   *  :.


HBA_HUMAN      PAVHASLDKFLASVSTVLTSKYR
HBA_PONPY      PAVHASLDKFLASVSTVLTSKYR
HBA_BOVIN      PAVHASLDKFLANVSTVLTSKYR
HBA_HORSE      PAVHASLDKFLSSVSTVLTSKYR
HBA_EQUAS      PAVHASLDKFLSTVSTVLTSKYR
HBA_MOUSE      PAVHASLDKFLASVSTVLTSKYR
HBA_RAT        PAMHASLDKFLASVSTVLTSKYR
HBA_PIG        PSVHASLDKFLANVSTVLTSKYR
HBA_CANLF      PAVHASLDKFFAAVSTVLTSKYR
HBA_APTFO      PEIHASLDKFLCAVGNVLTSKYR
HBA_TREBE      PEAHVSLDKFLSGVALALAERYR
                        *   *.*****:. *.  .*:.:**
```

FIGURE 30.  Output MSA by MAFFT for set 1

```
CLUSTAL multiple sequence alignment by MUSCLE (3.8)


HBA_TREBE       -SLSDKDKAAVRALWSKIGKSADAIGNDALSRMIVVYPQTKTYFSHWPDVTPGSPHIKAH
HBA_APTFO       MVLSANDKSNVKSIFSKISSHAEEYGAETLERMFTTYPQTKTYFPHF-DLHHGSAQVKAH
HBA_CANLF       -VLSPADKTNIKSTWDKIGGHAGDYGGEALDRTFQSFPTTKTYFPHF-DLSPGSAQVKAH
HBA_RAT         MVLSADDKTNIKNCWGKIGGHGGEYGEEALQRMFAAFPTTKTYFSHI-DVSPGSAQVKAH
HBA_PIG         -VLSAADKANVKAAWGKVGGQAGAHGAEALERMFLGFPTTKTYFPHF-NLSHGSDQVKAH
HBA_MOUSE       MVLSGEDKSNIKAAWGKIGGHGAEYGAEALERMFASFPTTKTYFPHF-DVSHGSAQVKGH
HBA_HORSE       MVLSAADKTNVKAAWSKVGGHAGEYGAEALERMFLGFPTTKTYFPHF-DLSHGSAQVKAH
HBA_EQUAS       MVLSAADKTNVKAAWSKVGGNAGEFGAEALERMFLGFPTTKTYFPHF-DLSHGSAQVKAH
HBA_BOVIN       MVLSAADKGNVKAAWGKVGGHAAEYGAEALERMFLSFPTTKTYFPHF-DLSHGSAQVKGH
HBA_HUMAN       MVLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHF-DLSHGSAQVKGH
HBA_PONPY       MVLSPADKTNVKTAWGKVGAHAGDYGAEALERMFLSFPTTKTYFPHF-DLSHGSAQVKDH
                   **   **   :.   :.*:.   .    * ::*.* :   :* *****.*   ::   ** ::* *


HBA_TREBE       GKKVMGGIALAVSKIDDLKTGLMELSEQHAYKLRVDPANFKILNHCILVVISTMFPKEFT
HBA_APTFO       GKKVAAALIEAANHIDDIAGALSKLSDLHAEKLRVDPVNFKLLGQCFMVVVAIHHPSALT
HBA_CANLF       GKKVADALTTAVAHLDDLPGALSALSDLHAYKLRVDPVNFKLLSHCLLVTLACHHPTEFT
HBA_RAT         GKKVADALAKAADHVEDLPGALSTLSDLHAHKLRVDPVNFKFLSHCLLVTLACHHPGDFT
HBA_PIG         GQKVADALTKAVGHLDDLPGALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHHPDDFN
HBA_MOUSE       GKKVADALASAAGHLDDLPGALSALSDLHAHKLRVDPVNFKLLSHCLLVTLASHHPADFT
HBA_HORSE       GKKVGDALTLAVGHLDDLPGALSNLSDLHAHKLRVDPVNFKLLSHCLLSTLAVHLPNDFT
HBA_EQUAS       GKKVGDALTLAVGHLDDLPGALSNLSDLHAHKLRVDPVNFKLLSHCLLSTLAVHLPNDFT
HBA_BOVIN       GAKVAAALTKAVEHLDDLPGALSELSDLHAHKLRVDPVNFKLLSHSLLVTLASHLPSDFT
HBA_HUMAN       GKKVADALTNAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPAEFT
HBA_PONPY       GKKVADALTNAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPAEFT
                * **   .:   *. :::*:   .*   **: ** ******.***:*.:.:: .::    *   :.


HBA_TREBE       PEAHVSLDKFLSGVALALAERYR
HBA_APTFO       PEIHASLDKFLCAVGNVLTSKYR
HBA_CANLF       PAVHASLDKFFAAVSTVLTSKYR
HBA_RAT         PAMHASLDKFLASVSTVLTSKYR
HBA_PIG         PSVHASLDKFLANVSTVLTSKYR
HBA_MOUSE       PAVHASLDKFLASVSTVLTSKYR
HBA_HORSE       PAVHASLDKFLSSVSTVLTSKYR
HBA_EQUAS       PAVHASLDKFLSTVSTVLTSKYR
HBA_BOVIN       PAVHASLDKFLANVSTVLTSKYR
HBA_HUMAN       PAVHASLDKFLASVSTVLTSKYR
HBA_PONPY       PAVHASLDKFLASVSTVLTSKYR
                *   *.*****:. *.  .*:..**
```

FIGURE 31. Output MSA by MUSCLE for set 1

Now, I will be putting these MSA aligned sequences (in the CLUSTAl format) to the code from question 3 to compute the given three score for each of these MSA alignments.

I used the following code to check for covert the MSA obtained by CLUSTAL OMEGA, MAFFT, and MUSCLE into their respective dataframes.

```python
import pandas as pd

# Reading the CLUSTAL Omega MSA file
clustal = []
with open("clustal.txt", 'r') as f:
    line = f.read()
    line_list = line.split("\n")

    for i in range(len(line_list)):
        line_i = line_list[i].split(" ")
        line_i = [str(value) for value in line_i if value != ""]
        clustal.append(line_i)

df1 = pd.DataFrame(clustal)
df1[2] = pd.to_numeric(df1[2])
df1[0] = pd.to_numeric(df1[0])

# Reading the MAFFT MSA file
mafft = []
with open("mafft.txt", 'r') as f:
    line = f.read()
    line_list = line.split("\n")

    for i in range(len(line_list)):
        line_i = line_list[i].split(" ")
        line_i = [str(value) for value in line_i if value != ""]
        mafft.append(line_i)

df2 = pd.DataFrame(mafft)
df2[2] = pd.to_numeric(df2[2])
df2[0] = pd.to_numeric(df2[0])

# Reading the MUSCLE MSA file
muscle = []
with open("muscle.txt", 'r') as f:
    line = f.read()
    line_list = line.split("\n")

    for i in range(len(line_list)):
        line_i = line_list[i].split(" ")
        line_i = [str(value) for value in line_i if value != ""]
        muscle.append(line_i)

df3 = pd.DataFrame(muscle)
df3[2] = pd.to_numeric(df3[2])
df3[0] = pd.to_numeric(df3[0])
```

LISTING 5. Code to create the dataframes for different MSAs

Now that I have the three dataframes for the MSA obtained by each of the alignment algorithms CLUSTAL OMEGA, MAFFT, and MUSCLE, I use the following code to list the similarities and differences.

```
1  similarities_3 = [] # All three same residue same position
2  differences_2 = [] # Any two different residue same position
3  differences_3 = [] # All three different residue same position
4
5  for i in range(len(df1)):
6    # Condition for checking any two sequences with different residues at same
         position
7    if df1[2][i]!=df2[2][i] or df1[2][i]!=df3[2][i] or df2[2][i]!=df3[2][i]:
8      print("Differences_2")
9      print(df1[2][i],df1[3][i],df2[2][i],df2[3][i],df3[2][i],df3[3][i],i)
10     differences_2.append([df1[2][i], df2[2][i], df3[2][i], i+1])
11
12   # Condition for checking all three sequences with different residues at
         same position
13   if df1[2][i]!=df2[2][i] and df1[2][i]!=df3[2][i] and df2[2][i]!=df3[2][i]:
14     print("Differences_3")
15     print(df1[2][i],df1[3][i],df2[2][i],df2[3][i],df3[2][i],df3[3][i],i)
16     differences_3.append([df1[2][i], df2[2][i], df3[2][i], i+1])
17
18   # Condtion for checking all three sequences with similar residues at same
         position
19   if df1[2][i] == df2[2][i] and df1[2][i] == df3[2][i]:
20     print("Similarities_3")
21     print(df1[2][i],df1[3][i],df2[2][i],df2[3][i],df3[2][i],df3[3][i],i)
22     similarities_3.append([df1[2][i], df2[2][i], df3[2][i], i+1])
23
24 # Converting all of these lists into dataframes for easy visualization
25 df1 = pd.DataFrame(similarities_3, columns=["Clustal", "Mafft", "Muscle", "
       Position"])
26 df2 = pd.DataFrame(differences_3, columns=["Clustal", "Mafft", "Muscle",  "
       Position"])
27 df3 = pd.DataFrame(differences_2, columns=["Clustal", "Mafft", "Muscle",  "
       Position"])
```

LISTING 6. Code to find the similarities and differences in MSA

When I passed the above set 1 MSA, obtained by all three methods, to the code, unfortunately, all the differences dataframes were empty. This is indicative of the fact that all three MSA algorithms of CLUSTAL OMEGA, MAFFT, and MUSCLE aligned the given 11 sequences in set 1 in the exact same manner.



FIGURE 32. Output of differences_3 and differences_2

This is indicative of the fact that all the positions are same in all the three alignment techniques for the proteins in set 1. The similarities are:

| Position | Clustal Residue | Mafft Residue | Muscle Residue | Score |
|---|---|---|---|---|
| 3 | Leucine | Leucine | Leucine | 0.0 |
| 5 | Proline | Proline | Proline | -1.121 |
| 143 | Argenine | Argenine | Argenine | 0.0 |

When I passed the above set 2 MSA, obtained by all three methods, to the code, unfortunately, the differences_3 dataframe was empty, however, the dataframes differences_2 and similarities_3 were not empty, indicating that differences and similarities exist upto some extent. This is indicative of the fact that all three MSA algorithms of CLUSTAL OMEGA, MAFFT, and MUSCLE do not align the given 9 sequences in set 2 in the exact same manner.



(A) similarities_3



(B) differences_2

FIGURE 33. The two non empty dataframes in MSA comparison for set 2

Another observation is that the differences observed in differences_2 is only the differences between either **mafft** or **muscle** and **clustal**. The mafft and muscle scores are exactly same throughout indicating their same alignment.

The similarities_3 are:

| Position | Clustal Residue | Mafft Residue | Muscle Residue | Score |
|---|---|---|---|---|
| 6 | Lysine | Lysine | Lysine | -0.736 |
| 7 | Phenylalanine | Phenylalanine | Phenylalanine | -0.377 |
| 255 | Alanine | Alanine | Alanine | -0.530 |

The differences_2 are:

| Position | Clustal | Score | Mafft | Score | Muscle | Score |
|---|---|---|---|---|---|---|
| 71 | Tyrosine | -0.849 | Tyrosine | -0.562 | Tyrosine | -0.562 |
| 72 | Lysine | -1.074 | Lysine | -1.149 | Lysine | -1.149 |
| 162 | Gap | * | Lysine | -0.849 | Lysine | -0.849 |

The differences_2 are a total of 20 in number. The remaining are basically due to gap scores of clustal technique not aligning perfectly with the gap score of mafft and muscle techniques. Also, there are no occurences of differences_3. The dataframe is empty.



FIGURE 34. Output of differences_3 for seq 2

**Question 5.** Check the scores manually at positions 9, 11, 20, 22 and 30 (use MSA from Clustal Omega)

**Solution.** Below are the images of manual calculation of scores. The scores calculated are unweighted frequencies, with scoring computational methods, entropy based, variance based, and sum-of-pairs based measures.

Also, in the sum-of-pairs method, after the entire summation, it needs to be squared and then square root should to be taken. This measure is to ensure that the value is positive.

Also, I saw a research paper on AL2CO that also had omitted the squaring and square rooting component in the sum-of-pairs method. Along with that, none of the values I computed showed a negative sign, hence I have dropped the squaring and square rooting component and computed the values. Below are the images of the manual calculations.



FIGURE 35. Calculations for 9th position of set 1

$11^{th}$ position: $[V, I, I, V, V, V, V, V, I, V, V]$

$n(Total) = 11 \quad n(V) = 8 \quad n(I) = 3$

Unweighted frequency: $f_a^u(i) = n_a(i)/n(i)$

$f_V(11) = [8/11] \quad f_I(11) = [3/11]$

Entropy based: $C_i^e(i) = \Sigma f_a(i) \cdot \ln(f_a(i)), \quad a = 1,20$

For all $a \notin \{v, I\}$, $f_a(i) = 0 \Rightarrow$ No impact on summation

For $a \in \{v, I\}$

$\quad C^e(11) = (8/11) \cdot \ln(8/11) + (3/11) \cdot \ln(3/11)$

$\qquad \approx -0.586$

$\qquad \Rightarrow \boxed{C^e(11) \approx -0.586}$

Variance based: $C_V(i) = \{ \Sigma [f_a(i) - f(i)]^2 \}^{0.5}, \quad a = 1,20$

The $f_a(i)$ is given above and $f(i)$ previous question

$\Rightarrow C^V(11) = [(0.7272 - 0.0807)^2 + \ldots + (0.2727 - 0.014)^2]$

$\qquad = \approx 0.744$

$\qquad \boxed{C^V(11) \approx 0.744}$

Sum-of-pairs based: $C^P(i) = \{ \Sigma_a \Sigma_b f_a(i) \cdot f_b(i) \cdot S_{ab} \} \quad a,b = 1,20$

For all $a, b \notin \{v, I\}$, $f_a(i)$ or $f_b(i) = 0 \Rightarrow$ No impact on sum.

|       | V     | I     |
|-------|-------|-------|
| V     | 0.727 | —     |
| I     | —     | 0.272 |

$f_a(i)$

|       | V   | I   |
|-------|-----|-----|
| V     | 4   | —   |
| I     | 3   | 4   |

BLossum 62
matrix

$\Rightarrow C^P(11) = [0.727 * 0.727 * 4 + \ldots + 0.272 * 0.272 * 4]$

$\qquad \approx 3.603$

$\qquad \boxed{C^P(11) \approx 3.603}$

FIGURE 36. Calculations for 11th position of set 1

**20th position** : $[A, G, G, G, G, G, G, A, G, K, S]$

$n(Total) = 11 \quad n(A) = 2 \quad n(G) = 7 \quad n(K) = 1 \quad n(S) = 1$

**Unweighted frequency** : $f_a^u(i) = n_a(i)/n(i)$

$$f_A(20) = [2/11] \quad f_G(20) = [7/11] \quad f_K(20) = [1/11] \quad f_S(20) = [1/11]$$

**Entropy based** : $c^e(i) = \sum f_a(i) . \ln(f_a(i)) , \quad a = 1,20.$

For all $a \notin \{A, G, K, s\}$, $f_a(i) = 0 \Rightarrow$ No impact on sum

For $a \in \{A, G, K, s\}$

$$c^e(20) = (2/11).\ln(2/11) + (7/11).\ln(7/11) + (1/11).\ln(1/11) + (1/11).\ln(1/11)$$

$$\approx -1.034$$

$$\boxed{c^e(20) \approx -1.034}$$

**Variance based** : $c^v(i) = \{ \sum [f_a(i) - f(i)]^2 \}^{0.5} , \quad a = 1,20$

The $f_a(i)$ is above and $f(i)$ in first part.

$$\Rightarrow c^v(20) = [(0.6363 - 0.0608)^2 + \ldots + (0.1818 - 0.0833)^2 + \ldots + (0.09 - 0.0807)^2]^{0.5}$$

$$\approx 0.613 \quad \Rightarrow \boxed{c^v(20) \approx 0.613}$$

**Sum-of-pairs based** : $c^p(i) = \{ \sum_a \sum_b f_a(i).f_b(i).S_{ab} \} \quad a,b = 1,20$

For all $ba \notin \{A, G, K, s\}$, $f_a(i)$ or $f_b(i) = 0 \Rightarrow$ No impact on sum

$f_a(i)$

|   | G | A | K | S |
|---|---|---|---|---|
| G | 0.636 | — | — | — |
| A | — | 0.1818 | — | — |
| K | — | — | 0.0909 | — |
| S | — | — | — | 0.0909 |

|   | G | A | K | S |
|---|---|---|---|---|
| G | 6 | — | — | — |
| A | 0 | 4 | — | — |
| K | -2 | -1 | 5 | — |
| S | 0 | 1 | 0 | 4 |

Blosum62 matrix

$$c^p(20) = [(0.636 * 0.636 * 6) + (0.636 * 0.0909 * 2) + \ldots\ldots$$
$$+ (0.0909 * 0.0909 * 4)]$$

$$\approx 2.405$$

$$\boxed{c^p(20) \approx 2.405}$$

FIGURE 37. Calculations for 20th position of set 1

$22^{nd}$ position : $[A, G, G, A, A, A, A, A, A, A, A]$

$n(Total) = 11$ $\quad n(A) = 9$ $\quad n(G) = 2$

Unweighted frequency : $f_a^u(i) = n_a(i) / n(i)$

$\quad f_A(22) = [9/11]$ $\quad f_G(22) = [2/11]$

Entropy based : $C^e(i) = \sum f_a(i) \cdot \ln f_a(i)$ , $a = 1, 20$.

$\quad$ For all $a \notin \{A, G\}$ , $f_a(i) = 0$ ⟹ No impact on sum.

For $a \in \{A, G\}$

$\quad C^e(22) = (2/11) \cdot \ln(2/11) + (9/11) \cdot \ln(9/11) \approx -0.474$

$\boxed{C^e(22) = -0.474}$

Variance based : $C^v(i) = \{\sum (f_a(i) - f(i))^2\}^{0.5}$ , $a = 1, 20$.

The $f_a(i)$ is above and $f(i)$ was computed for first Q.

⟹ $\quad C^v(22) = [(0.1818 - 0.0608)^2 + (0.818 - 0.1275)^2 + \dots + (0 - 0.0802)^2]^{0.5}$

$\quad \approx 0.74$

$\boxed{C^v(22) \approx 0.74}$

Sum-of-pairs based: $C^p(i) = \{\sum_a \sum_b f_a(i) \cdot f_b(i) \cdot S_{ab}\}$ $\quad a, b = 1, 20$

For all $b, a \notin \{A, G\}$ , $f_a(i) = 0$ or $f_a(i) = 0$ ⟹ No impact on sum

| | A | G |
|---|---|---|
| A | 0.8181 | — |
| G | — | 0.1818 |

$f_a(i)$

| | A | G |
|---|---|---|
| A | 4 | — |
| G | 0 | 6 |

Blosum62 matrix

⟹ $C^p(22) = [(0.8181 * 0.8181 * 4) + \dots + (0.1818 * 0.1818 * 6)]$

$\quad \approx 2.876$

$\boxed{C^p(22) \approx 2.876}$

FIGURE 38. Calculations for 22th position of set 1

$\underline{30^{th} \text{ position}}$ : $[L, L, L, L, L, L, L, L, L, L, L]$

$n(\text{Total}) = 11 \qquad n(L) = 11$

Unweighted frequency : $f_a^u(i) = n_a(i) / n(i)$

$\qquad f_L(30) = [11/11] = 1$

---

$\underline{\text{Entropy based}}$ : $C^e(i) = \sum f_a(i) . \ln(f_a(i))$ , $a = 1, 20$.

For all $a \notin \{L\}$, $f_a(i) = 0 \Rightarrow$ No impact on sum.

For $a \in \{L\}$

$\qquad C^e(30) = (11/11) \ln(11/11) = 0$

$\qquad \boxed{C^e(30) = 0}$

---

$\underline{\text{Variance based}}$ : $C_v^v(i) = \{ [\sum [f_a(i) - f(i)]^2 \}^{0.5}$

The $f_a(i)$ is above and $f(i)$ was computed for $1^{st}$ Q.

$\Rightarrow C^v(30) = [(1 - 0.12435)^2 + (0 - 0.0807)^2 + .. + (0 - 0.0115)^2]^{0.5}$

$\qquad \approx 0.909$

$\qquad \boxed{C^v(30) \approx 0.909}$

---

$\underline{\text{Sum-of-pairs method}}$ : $C^p(i) = \{ \sum_a \sum_b f_a(i) . f_b(i) . S_{ab} \} \quad a = 1, 20$

for all $a, b \notin \{L\}$, $f_a(i)$ or $f_b(i) = 0 \Rightarrow$ No impact on sum.

|   | L |
|---|---|
| L | 1 |

$f_a(i)$

|   | L |
|---|---|
| L | 4 |

Blosum62 matrix

$\Rightarrow C^p(30) = [(1 * 1 * 4)]$

$\qquad \Rightarrow \boxed{C^p(30) = 4}$

FIGURE 39. Calculations for 30th position of set 1

**Question 6.** Obtain the conservation score of 1BTM, A-chain using Consurf server https://consurf.tau.ac.il/

**Solution.** First thing I need to find is the UniProt ID. I searched for the 1BTM, A chain. The below is the structure I found along with the accession number of the protein. The accession number of the protein is **P00943**. I searched for the same on the UniProtKB and found the similarity in structure, and other features and structural elements. Below is the structure I found initially.
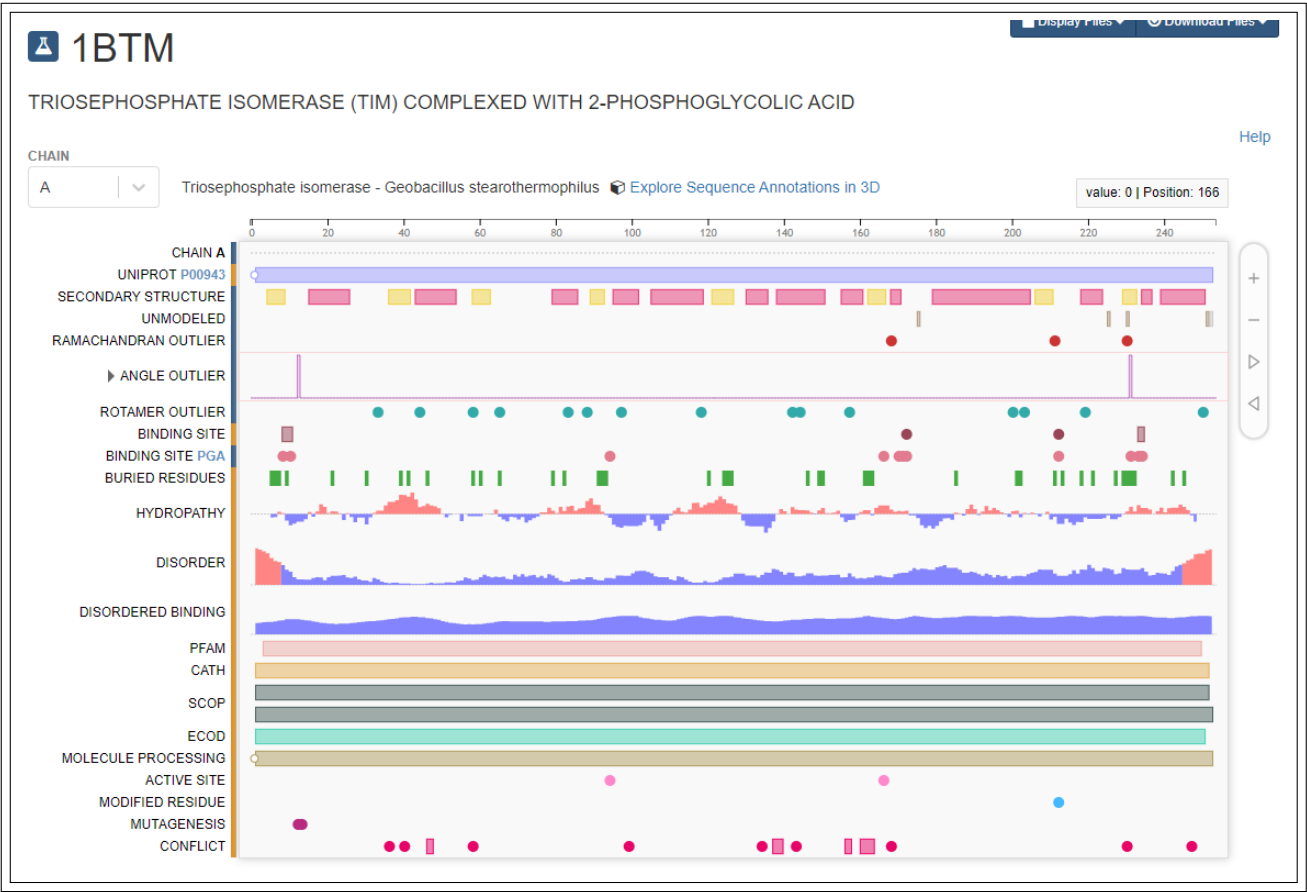


FIGURE 40. Finding the protein sequence UniProt ID

---

**Protein sequence of 1BTM A-chain obtained from UniProtKB**

```
1  >P00943
2  MRKPIIAGNWKMHKTLAEAVQFVEDVKGHVPPADEVISVVCAPFLFLDRLVQAAD
3  GTDLKIGAQTMHFADQGAYTGEVSPVMLKDLGVTYVILGHSERRQMFAETDETVN
4  KKVLAAFTRGLIPIICCGESLEEREAGQTNAVVASQVEKALAGLTPEQVKQAVIA
5  YEPIWAIGTGKSSTPEDANSVCGHIRSVVSRLFGPEAAEAIRIQYGGSVKPDNIR
6  DFLAQQQIDGPLVGGASLEPASFLQLVEAGRHE
```

LISTING 7. The protein sequence of the protein mentioned in the question

---

The above FASTA file is given as input to the **Consurf server**.

There also exists an alternative approach to execute the above given task. In the input page of the Consurf server, there is an option to enter the **PDB or UniProt ID**. Enter the term **1BTM** there. Then under the section of **Select the chain identifier**, choose **Chain A**. Now enter the job tite and email ID and **Run with default parameters**.

Below are two images, one for the input screen on the Consurf server, the other is the running parameters on Consurf server.
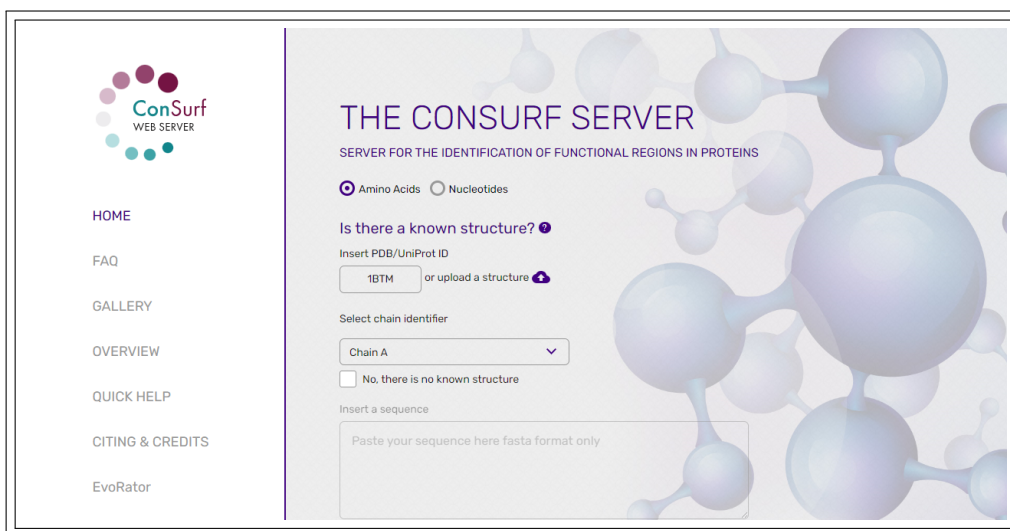


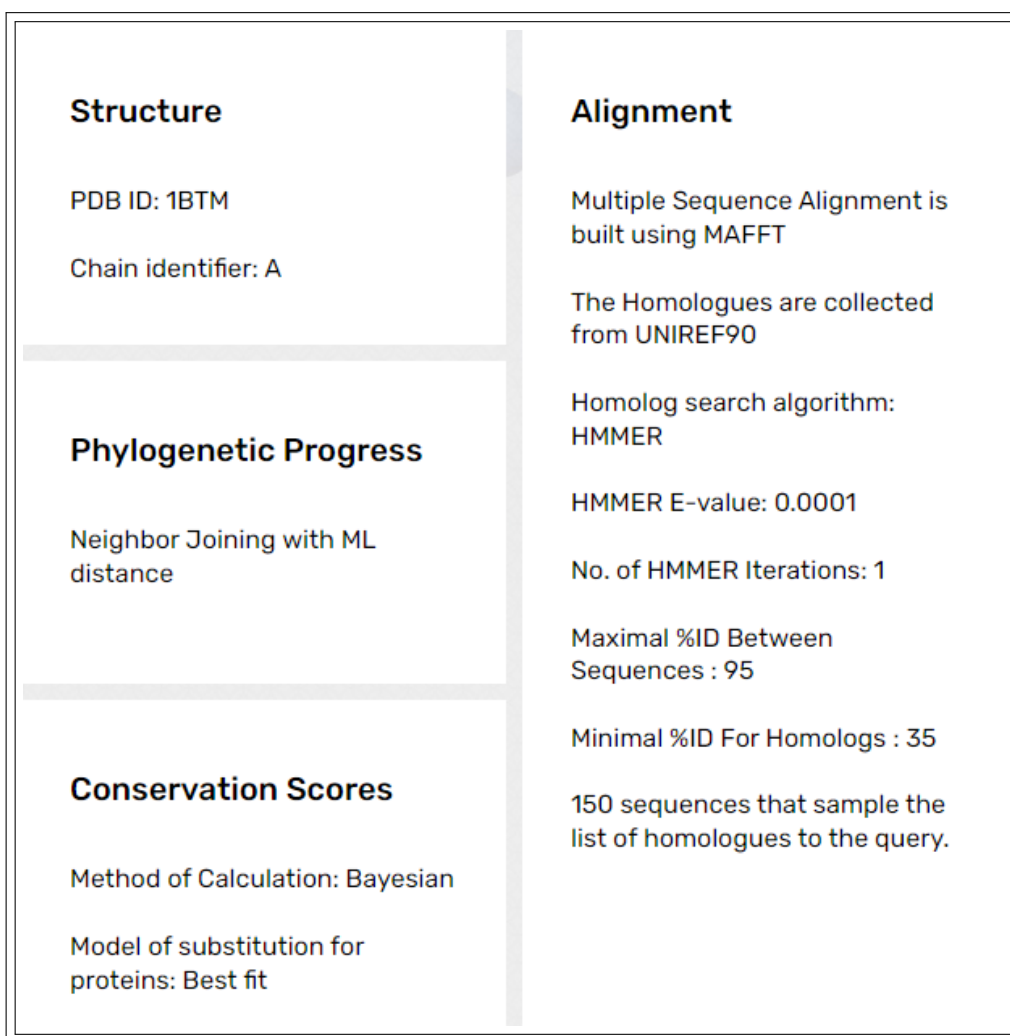FIGURE 41. Input screen on the Consurf server



**Structure**

PDB ID: 1BTM

Chain identifier: A

**Phylogenetic Progress**

Neighbor Joining with ML distance

**Conservation Scores**

Method of Calculation: Bayesian

Model of substitution for proteins: Best fit

**Alignment**

Multiple Sequence Alignment is built using MAFFT

The Homologues are collected from UNIREF90

Homolog search algorithm: HMMER

HMMER E-value: 0.0001

No. of HMMER Iterations: 1

Maximal %ID Between Sequences : 95

Minimal %ID For Homologs : 35

150 sequences that sample the list of homologues to the query.

FIGURE 42. Running parameters of Consurf server

The results page for the above run on the Consurf server shows up as below. It consists of the structure of the protein sequence, along with its sequence and highlighted are highly and lowly conserved regions:



FIGURE 43. Results page

Given below are some of the parameters that play a role in assigning conservation scores, like, layers for assigning grades, confidence interval colors, residue variety.
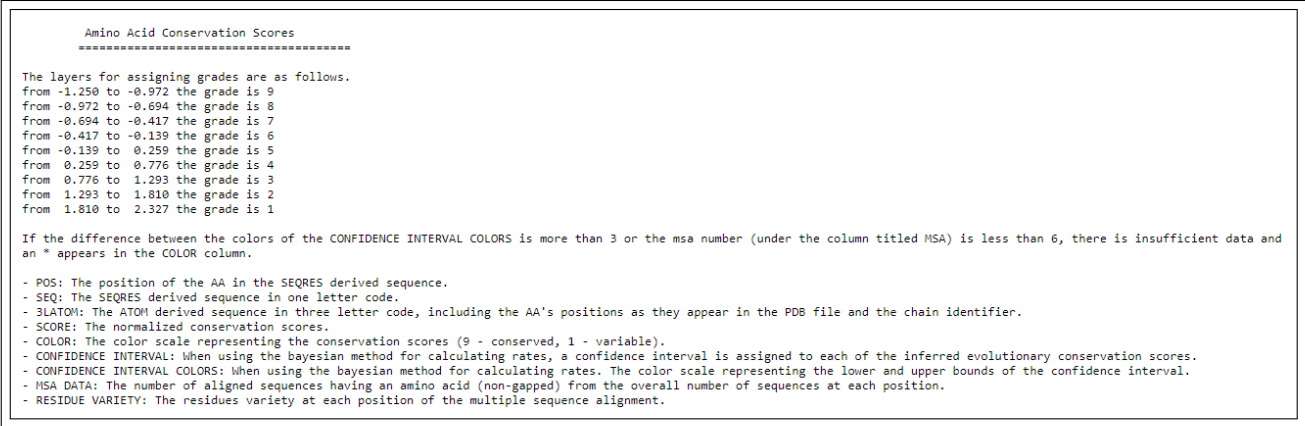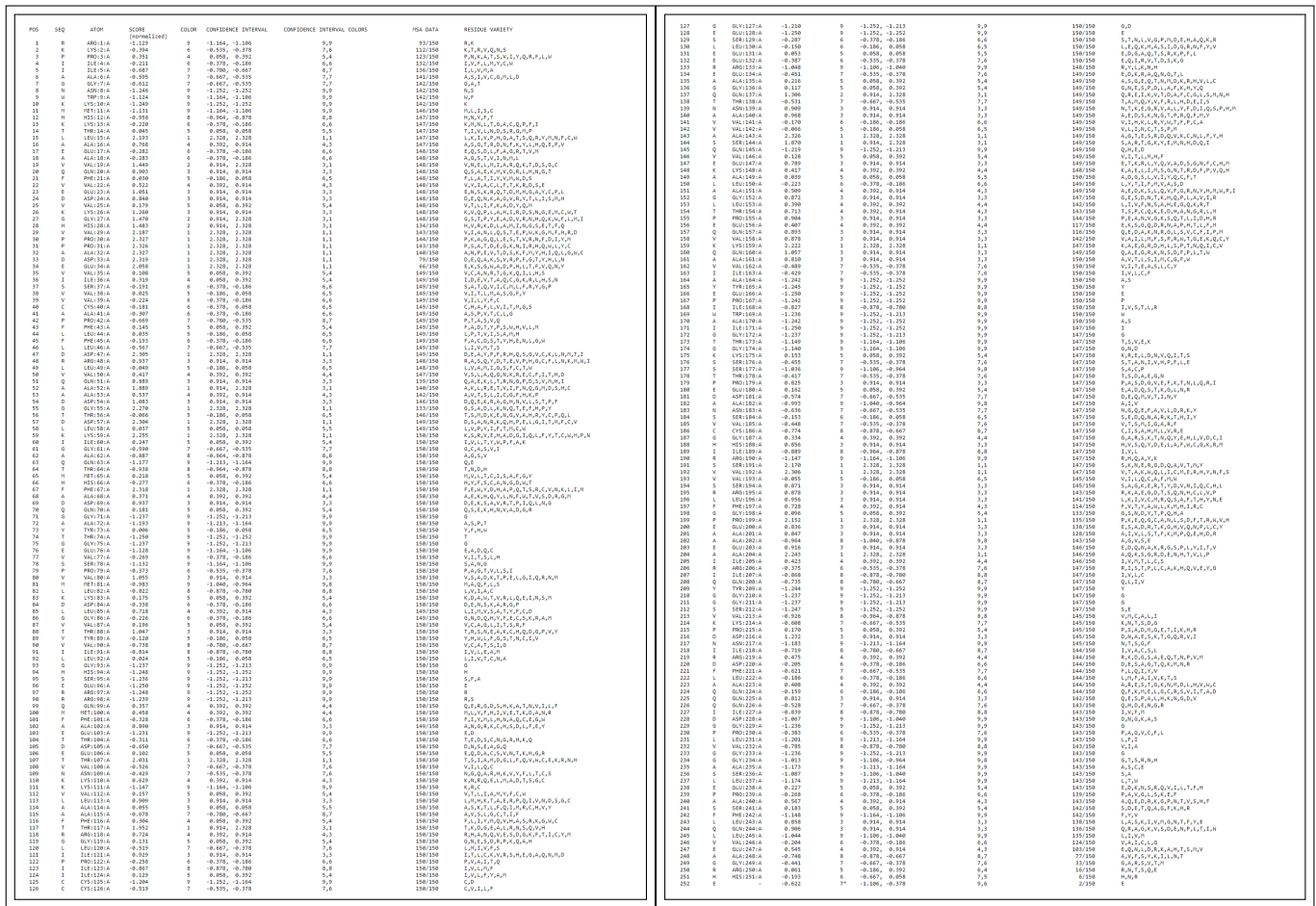


FIGURE 44. Parameters to compute the conservation scores

Now, given below is the set of all 252 residues in the given protein and also all properties related to each residue position, like, SEQ (amino acid), COLOR (indicating conservation), CONFIDENCE INTERVAL (indicate the range of conservation values for COLOR), RESIDUE VARIETY (presence of multiple residues at same position), etc.

(A) 1-126 positions

(B) 126-252 positions

FIGURE 45. Each residue-wise properties for the given protein

Names and Taxonomy related to protein from UniProtKB:

(1) Recommended name is **Triosephosphate isomerase**
(2) EC number is **EC:5.3.1.1**
(3) Its optimal temperature is 60 degrees Celsius. It is thermostable
(4) It plays a role in the following pathways: Carbohydrate biosynthesis; gluconeogenesis; Carbohydrate degradation; glycolysis; D-glyceraldehyde 3-phosphate from glycerone phosphate

The color scale in all of the above images are indicative of the conservation scores, with 9 indicating **conserved**, while 1 indicating **variable**. These numeric values have been assigned to different confidence interval ranges for the conservation scores. The split up is like this:

- from -1.250 to -0.972 the grade is 9
- from -0.972 to -0.694 the grade is 8
- from -0.694 to -0.417 the grade is 7
- from -0.417 to -0.139 the grade is 6
- from -0.139 to 0.259 the grade is 5
- from 0.259 to 0.776 the grade is 4
- from 0.776 to 1.293 the grade is 3
- from 1.293 to 1.810 the grade is 2
- from 1.810 to 2.327 the grade is 1

The Consurf can also perform multiple sequence alignment with some predefined sequences based on their closeness to the given sequence. Below image is an example of one such MSA.
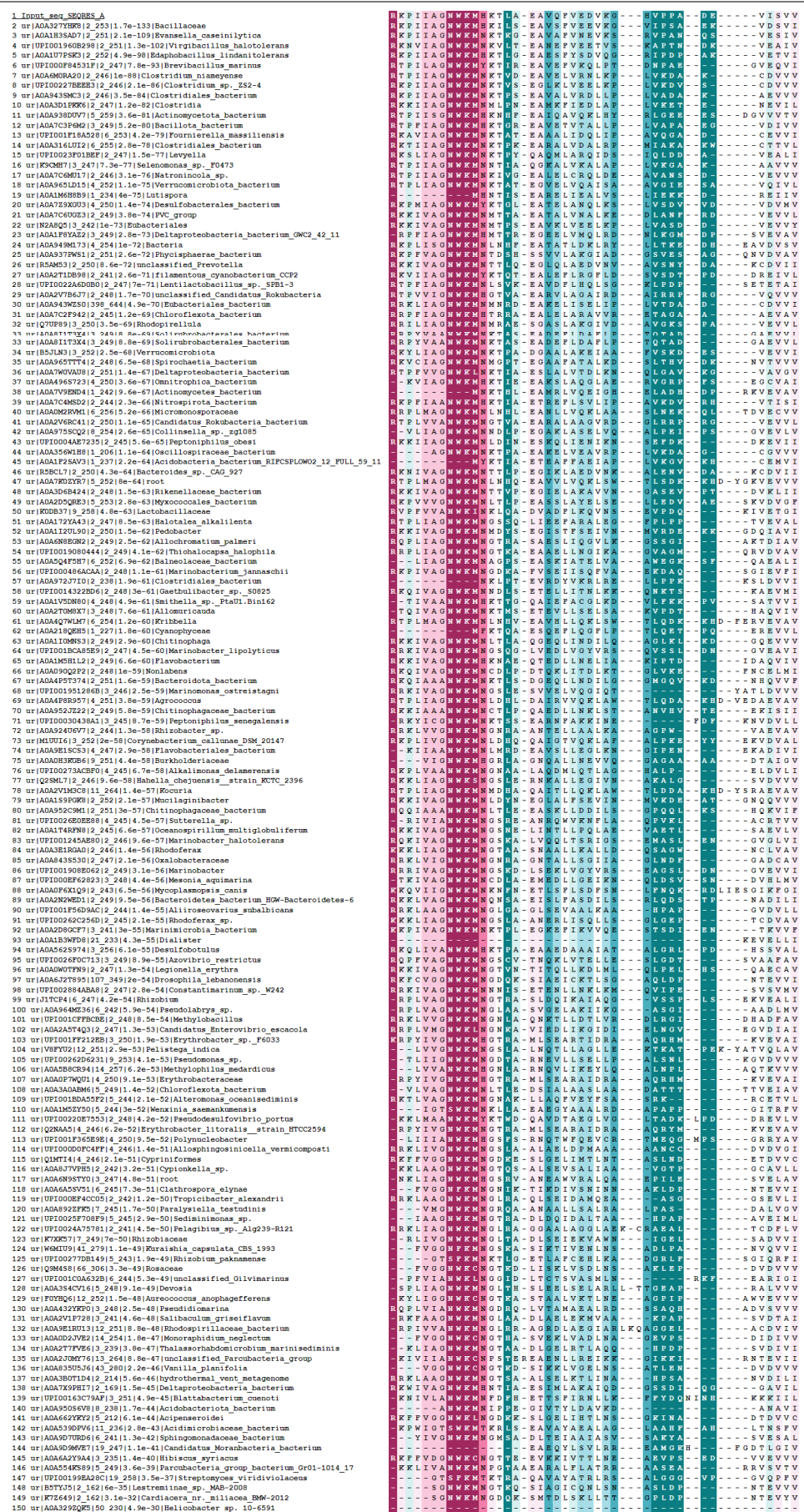
FIGURE 46. MSA on Consurf Server

Some additional results obtained from the Consurf server are given below in the image:

There are 28566 HMMER hits. 27844 of them are unique, including the query.
The calculation is performed on a sample of 150 sequences that represent the list of homologues to the query.

Here is the list of sequences that produced significant alignments, but were not chosen as hits.

The best evolutionary model was selected to be: WAG. (details).

FIGURE 47. Additional Homology related results