

## MBI HOMEWORK 3 FOR CS STUDENTS

Autor: Marián Kravec

### 1. RNA structure and dynamic programming.

a)

To compute structure with maximal number of base pairs without any not nested cases we will use slightly modified Nussinov algorithm.

```
def nussin(seq):
    n = len(seq)
    A = np.zeros((n,n))
    B = np.zeros((n,n))
    for length in range(1, n):
        for i in range(n):
            j = i + length
            if j < n:
                pos_val = 0
                A[i, j], B[i, j] = max([A[i+1, j], 1),
                                       (A[i, j-1], 2),
                                       (A[i+1, j-1]+pair(seq[i], seq[j]), 3)
                                       ], key= lambda x: x[0])
    return(A[0, n-1], B)
```

This algorithm is filling out table  $A$  (starting from diagonal and getting closer to upper-top corner) where  $A[i, j]$  is number maximal number of nested base pair of sub-sequence starting from position  $i$  and ending on position  $j$  (final solution is in cell  $A[0, n]$  where  $n$  is length of sequence), it's a bit simple than standard Nussinov algorithm because it does not take into account situation where branching of secondary structure create more base pairs, so situation where we are looking for position  $k$  between positions  $i$  and  $j$  such that  $A[i, k] + A[k + 1, j]$  is bigger than any other solution. We omit this rules because breaching creates not nested pairs.

So to compute  $A[i, j]$  our algorithm choose maximal value out of three possibilities:  $A[i + 1, j]$ ,  $A[i, j - 1]$  or  $A[i + 1, j - 1] + pair(x_i, x_j)$ .

- $A[i + 1, j]$  - first base of sub-sequence in unpaired and rest is optimally paired
- $A[i, j - 1]$  - last base of sub-sequence in unpaired and rest is optimally paired
- $A[i + 1, j - 1] + pair(x_i, x_j)$  - interpretation of this depends on value of  $pair(x_i, x_j)$  if  $x_i$  and  $x_j$  creates pair then this is situation where we take optimal pairing of bases between them plus their own, if they do not create pair then this is situation where we take optimal pairing of bases between them and first and last base is unpaired

This algorithm compute half of values of matrix with size  $n \times n$  and for each position it finds maximum of constant amount of possibilities (only 3), so asymptotic running time of this algorithm is  $O(n^2)$ .

b)

Similarly to algorithm from part a), our stochastic context-free grammar needs only one nonterminal which we will call  $S$  and to have only three types of rules (three types does not mean three rules), and that is rules to create unpaired bases before nonterminal, rules to create unpaired bases after nonterminal and rules to create pairs (and lastly there will be rule to terminate sequence generation):

$$\begin{aligned}
S &\rightarrow aS|uS|cS|gS| \\
&Sa|Su|Sc|Sg| \\
&aSu|uSa|cSg|gSc| \\
&\epsilon
\end{aligned}$$

c)

## 2. Bioinformatics tools and databases.

a)

After running Blast with in this setup:

The screenshot shows the NCBI BLAST search interface. The 'Enter Query Sequence' section contains a text area with the following FASTA sequence:

```

CAGCCTAAAGCTTCTGCTGCTCCTGAGGCTCAAGCAAGTCAACAAGTGG
TTGATGTTCAAATCCCTGATA
TTGGTGTAGAAAAAGCCACTGTCGGTGAAATTCTGTTTCTGTCGGTGA
TGAAATCGAGGTAGATCAAAG
  
```

Below the text area are options to 'Or, upload file' (Choose File), 'Job Title' (a text input field), and a checkbox for 'Align two or more sequences'. The 'Choose Search Set' section shows the 'Database' dropdown set to 'Standard databases (nr etc.)', with a 'Try experimental taxonomic nt databases' button. The 'Organism' dropdown is set to 'RefSeq Genome Database (refseq\_genomes)'. The 'Exclude' section has checkboxes for 'Models (XM/XP)' and 'Uncultured/environmental sample sequences'. The 'Limit to' section has a checkbox for 'Sequences from type material'. The 'Entrez Query' section has a text input field and a 'Create custom database' button. The 'Program Selection' section shows the 'Optimize for' dropdown set to 'Highly similar sequences (megablast)'. The 'Choose a BLAST algorithm' link is also visible.

(As specified in assignment)

We got these alignments:

	Description	Scientific Name	Max Score	Total Score	Query Cover	E value	Per. Ident	Acc. Len	Accession
✓	<a href="#">Acinetobacter Iwoffii strain FDAARGOS_1393 chromosome, complete genome</a>	<a href="#">Acinetobacter I...</a>	4894	4894	100%	0.0	100.00%	3166595	<a href="#">NZ_CP077336.1</a>
✓	<a href="#">Acinetobacter Iwoffii NCTC 5886 = CIP 64.10 = NIPH 512 adgTz-supercont2.2, whole genome s...</a>	<a href="#">Acinetobacter I...</a>	4894	4894	100%	0.0	100.00%	1072087	<a href="#">NZ_KI530565.1</a>
✓	<a href="#">Acinetobacter sp. CIPA162 aclZB-supercont1.4, whole genome shotgun sequence</a>	<a href="#">Acinetobacter...</a>	4894	4894	100%	0.0	100.00%	1807690	<a href="#">NZ_KB849159.1</a>
✓	<a href="#">Acinetobacter Iwoffii strain WU_MDCI_AI101 NODE_6_length_95976_cov_75.071242, whole ge...</a>	<a href="#">Acinetobacter I...</a>	4894	4894	100%	0.0	100.00%	95976	<a href="#">NZ_IAHPSO010000006.1</a>
✓	<a href="#">Acinetobacter Iwoffii strain NCTC5886, whole genome shotgun sequence</a>	<a href="#">Acinetobacter I...</a>	4894	4894	100%	0.0	100.00%	3264872	<a href="#">NZ_CAADHN010000002.1</a>
✓	<a href="#">Acinetobacter Iwoffii strain NCTC5867, whole genome shotgun sequence</a>	<a href="#">Acinetobacter I...</a>	4894	4894	100%	0.0	100.00%	3108475	<a href="#">NZ_UFSE01000003.1</a>
✓	<a href="#">Acinetobacter Iwoffii strain DSM 2403 chromosome, complete genome</a>	<a href="#">Acinetobacter I...</a>	4894	4894	100%	0.0	100.00%	3166595	<a href="#">NZ_CP118963.1</a>
✓	<a href="#">Acinetobacter Iwoffii strain GTC 00071 sequence05, whole genome shotgun sequence</a>	<a href="#">Acinetobacter I...</a>	4894	4894	100%	0.0	100.00%	116176	<a href="#">NZ_BJLI010000005.1</a>
✓	<a href="#">Acinetobacter Iwoffii ATCC 9957 = CIP 70.31 aclsr-supercont1.9, whole genome shotgun seque...</a>	<a href="#">Acinetobacter I...</a>	4894	4894	100%	0.0	100.00%	213534	<a href="#">NZ_KB849831.1</a>
✓	<a href="#">Acinetobacter Iwoffii strain NBRC 109760, whole genome shotgun sequence</a>	<a href="#">Acinetobacter I...</a>	4894	4894	100%	0.0	100.00%	116175	<a href="#">NZ_BBSQ01000006.1</a>
✓	<a href="#">Acinetobacter Iwoffii strain AMA23 AMA23_NODE_21, whole genome shotgun sequence</a>	<a href="#">Acinetobacter I...</a>	4894	4894	100%	0.0	100.00%	46679	<a href="#">NZ_VYTK01000021.1</a>
✓	<a href="#">Acinetobacter sp. RRD8.6, whole genome shotgun sequence</a>	<a href="#">Acinetobacter...</a>	4700	4700	100%	0.0	98.68%	186557	<a href="#">NZ_JAUZUJ010000006.1</a>
✓	<a href="#">Acinetobacter sp. RG5.5, whole genome shotgun sequence</a>	<a href="#">Acinetobacter...</a>	4700	4700	100%	0.0	98.68%	186557	<a href="#">NZ_JAUZUJ010000005.1</a>
✓	<a href="#">Acinetobacter Iwoffii strain VE196-10 contig00004, whole genome shotgun sequence</a>	<a href="#">Acinetobacter I...</a>	4519	4519	99%	0.0	97.47%	153900	<a href="#">NZ_IAMXXP010000004.1</a>
✓	<a href="#">Acinetobacter Iwoffii SH145 supercont1.10, whole genome shotgun sequence</a>	<a href="#">Acinetobacter I...</a>	4348	4348	100%	0.0	96.34%	167894	<a href="#">NZ_GG705084.1</a>
✓	<a href="#">Acinetobacter sp. isolate CTOTU50698 NODE_24_length_96746_cov_3.622829, whole genome...</a>	<a href="#">Acinetobacter...</a>	4348	4348	100%	0.0	96.34%	96746	<a href="#">NZ_DAMDPK010000001.1</a>
✓	<a href="#">Acinetobacter Iwoffii strain S252-2 contig00014, whole genome shotgun sequence</a>	<a href="#">Acinetobacter I...</a>	4340	4340	100%	0.0	96.23%	75191	<a href="#">NZ_IAMXXL010000014.1</a>
✓	<a href="#">Acinetobacter sp. CIP 101966 aclst-supercont1.12, whole genome shotgun sequence</a>	<a href="#">Acinetobacter...</a>	4324	4324	100%	0.0	96.12%	570054	<a href="#">NZ_KB850158.1</a>
✓	<a href="#">Acinetobacter Iwoffii strain S246-3 contig00002, whole genome shotgun sequence</a>	<a href="#">Acinetobacter I...</a>	4313	4313	100%	0.0	96.01%	105902	<a href="#">NZ_IAMXXK010000002.1</a>
✓	<a href="#">Acinetobacter Iwoffii strain FDAARGOS_620 unitig_0_quiver_quiver_pilon, whole genome shotg...</a>	<a href="#">Acinetobacter I...</a>	4307	4307	100%	0.0	95.98%	3341993	<a href="#">NZ_JAAXYZ010000003.1</a>
✓	<a href="#">Acinetobacter sp. NEB149 chromosome, complete genome</a>	<a href="#">Acinetobacter...</a>	4289	4289	99%	0.0	95.89%	3218664	<a href="#">NZ_CP051208.1</a>
✓	<a href="#">Acinetobacter Iwoffii strain WU_MDCI_AI262 NODE_9_length_63957_cov_63.179195, whole ge...</a>	<a href="#">Acinetobacter I...</a>	4263	4263	100%	0.0	95.62%	63957	<a href="#">NZ_JAHPWU010000009.1</a>
✓	<a href="#">Acinetobacter Iwoffii strain WU_MDCI_AI83 NODE_11_length_54607_cov_20.378746, whole ge...</a>	<a href="#">Acinetobacter I...</a>	4261	4261	99%	0.0	95.67%	54607	<a href="#">NZ_JAHPRY010000011.1</a>
✓	<a href="#">Acinetobacter Iwoffii NIPH 715 acl rE-supercont1.19, whole genome shotgun sequence</a>	<a href="#">Acinetobacter I...</a>	4255	4255	99%	0.0	95.64%	646333	<a href="#">NZ_KB849264.1</a>

We can see that in 11 sequences we got same highest score 4894, in all 11 cases algorithm identified 100% of input sequence and E-value is approximately 0.

So most likely our sequence in part of genome of bacteria called *Acinetobacter Iwoffii*.

b)

Now we will look at what kind of proteins might be encoded in our sequence, when we look for open reading frame with at least 180 codons we get 4 ORFs:

```

>orf1
QPKASAAPEAQASSQVVDVQIPDIGVEKATVGEILVSVGDEIEVDQSIWVVEDSKATVEV
PSTVSGTVESIEIEKEGDTIKEGVVILKVKTAVSAAQVQTEAPQAPVAQAATQEKAVEAPQ
TPAAPAGDVEVKVPDLGVDKAAVAEILVQVGDTVEKDQSIIVVEDSKATVEVPSTTAGVI
KAIHVELGQNVSQGIALMTIEAEAQAAAAPVAAKAEAPKAPAAKAPAPAASSTQTVAAS
DNADKLTKEQNVANSKVYAGPAVRKRLARELGVVLAADV KASGPHARVMKEDLKAYVKTRLT
TPQAAPVAAAAQVAGLPKLPDFSAGGGVEEKALTRLQQVSIPQLSLNNFIPQVTQFDAAD
ITELEAWRNEELKGNFKKEGLSLTIMAFIIKAVAHLLKEEREFAGHLADDGKSVLLRNEIH
MGIATVATPDGLTVPVLRHPDQKSIKQIATELGTLGQKARDKKLSPKDLQGANFTITSLGS
IGGTAFPTLVNWPQVAILGISPATMQPVWNGEGFDPRLMLPLSLSYDHRVINGADAARFT
NKLTKLLKDIRLLI*

>orf2
PEKDLMEYFLLSTAIVALAEMGDKTQLLALLLAARFRKVPVILVAILLATLINHGLSAVL
GQWVTTVIGPEVLLWIVSIGFIAMAVWMLIPDELGDENASINKWQKFGVFGATFILFFLA
EIGDKTQIATVALAARFDSVFVWMLGTTLGMMLANAPAVFLGDKLANKLPISLIHKIGAA
IFLVIGVATLVQYYFF*

>orf3
TRVATPMTRKMAAPILWIRLIGSLFASLSPKNTAGALASIMPSVVPNITQKTL SKRAASA
TVAIWVLSPI SARKNMKVAPNTPNFCHLLILAFSSPSSSGISIQTAIAIKPIETIHSST
SGPITVVT HCPNTAERPWLISVANKMATKIGTGLRKRAASSNASSWVLSPI SATIAVD
SKNSYMRSFSG*

>orf4
QVRNCFDDERHNGQAQAFFLEIAFQFITPCFQFRNICCIKLGHLRNEIIQRQLWNRYLLQ
TRQGFLFHTTKSAKVWQFWQTCDLSSSGYRSRLWRSKTRFDIRFQIFFHDARVWARCFDV
CQHNAQFTCQFTHSRTSIDFRVSDILLFGQLIGIIRRRDCLSRRC SRCRCFCFSRRFRF
CFG CNRCSSCLRFSFNCHQCNALRDILTQFHMNGFNHAGCSA WHFNSRFI*

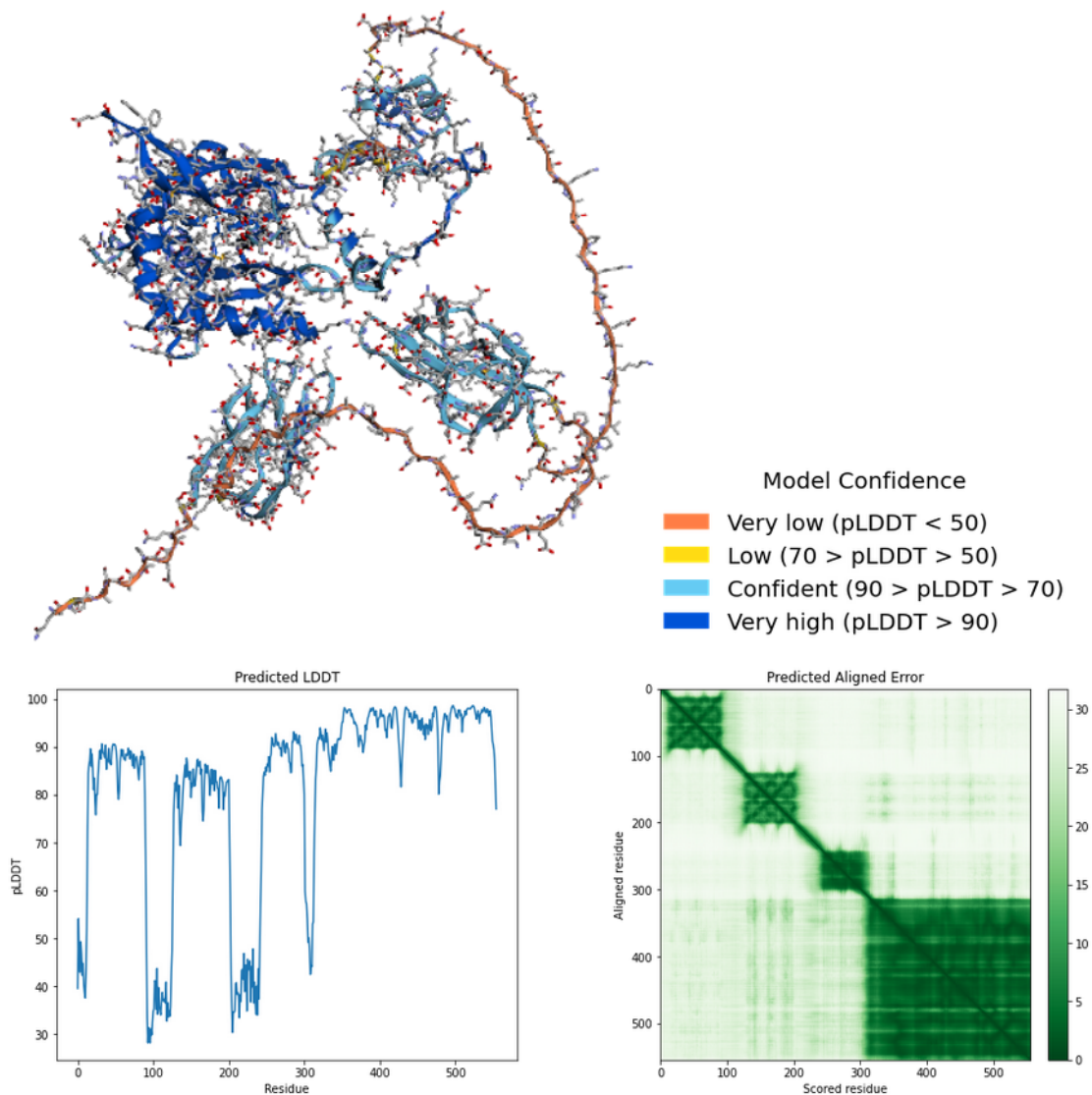
```

First two (ORF1 and ORF2) is on direct strand while last two (ORF3 and ORF4) are on reverse strand.

c)

Lastly we would like to see 3D structure of found proteins. We will look at ORF1 and ORF3. We used AlphaFold2 for this task.

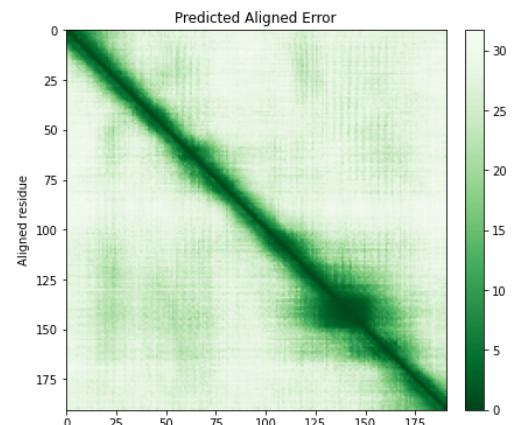
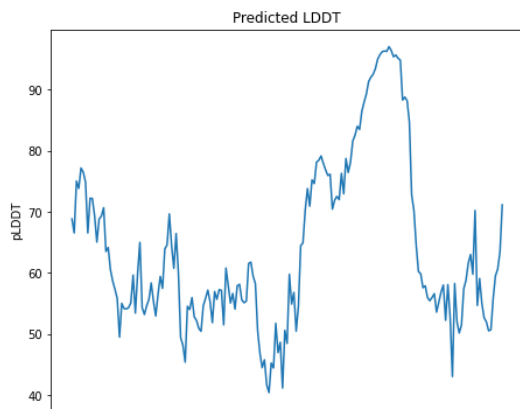
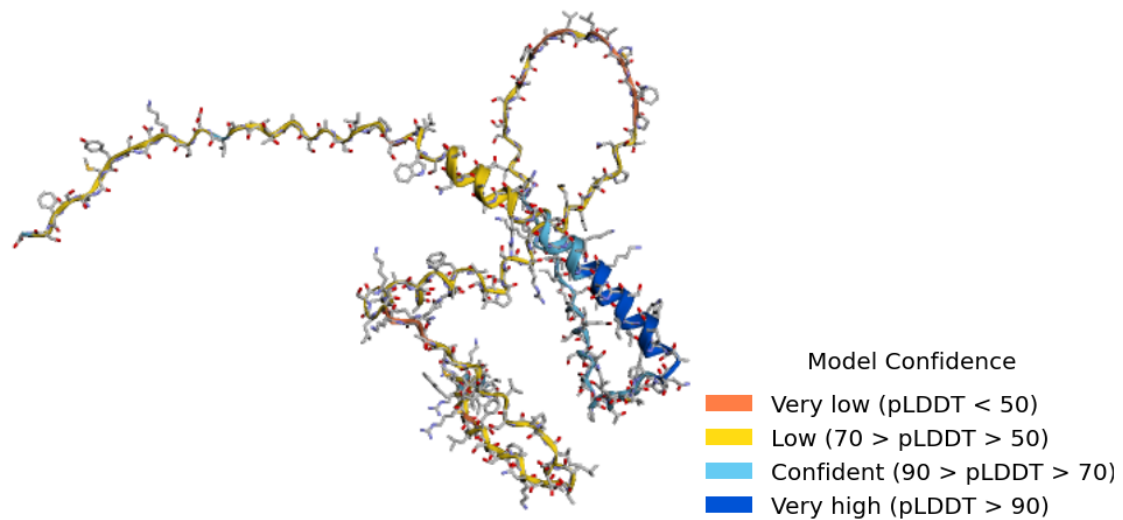
Now let's look firstly at ORF1:



We can see that most of the structure (mostly complex parts) is colored blue which signifies pretty high accuracy, only parts with low confidence are structures connecting highly accurate complex structures. This is confirmed by "Predicted LDDT" plot where we can see most values are higher than 80 so highly accurate.

On "Predicted aligned error" plot we can see 4 much darker squares which corresponds to 4 blue complex structure in 3D visualization. Dark cell mean that prediction error for relative positions of two those two amino acids (one from row and one from column) is low. So we can see that this protein seems to contain 4 domains.

Now let's look at ORF3:



This time we don't see any highly complicated parts. Also we can see that most of structure is yellow which indicates low confidence, so AlphaFold is not sure if this is correct structure. "Predicted LDDT" confirms it as we can see that most values are lower than 50 with one notable spike which corresponds to small blue part in 3D visualization.

"Predicted aligned error" plot also shows that other than amino acids right next to each other in sequence we have high error which means that relative position and/or orientation of almost any two amino acids is uncertain. One more information that we can get that algorithm was not able to find any notable domains.

If we compare results for those two ORF we can say that prediction for ORF1 is much more reliable because we can be almost certain shape of most of a structure.