

# CS612 – Algorithms in Bioinformatics

## Homework Assignment 1

Due Fri. Feb. 21 before midnight, on Gradescope

### Objectives

- Learn how to perform sequence alignment and analyze the results
- Understand and apply the theory of sequence alignment
- Practice programming BLAST with Biopython.

### The FASTA format

FASTA (pronounced "fast-ay") is a text-based format for representing either nucleotide sequences or peptide sequences, in a single letter code. It is very simple. The top line starts with a > followed by free text. The rest are the nucleotides or amino acid letter codes.

### Part 1 – Practice

1. **Sequence alignment hands-on exercise.** You are given the following protein sequence:

>protein

```
TCPFADPAALYSRQDTTSGQSPLAAYEVDDSTGYLTSVGGPIQDQTSLKAGIRGPTLLEDFMFRQKIQHFDHERVPERAV
```

The sequence is also available at the course webpage: <http://www.cs.umb.edu/nurith/cs612/protein.fasta>

- (a) Go to the Blast website at <http://blast.ncbi.nlm.nih.gov>. Select "Protein Blast" to get to the BlastP website. Paste the above sequence to the top window. Use the default parameters – nr database (non-redundant protein sequence database). Hit "BLAST" – the search may take up to a minute.

Answer the following questions about the best hit:

- i. What is the name of the sequence?
  - ii. what is the identifier (Accession) on the far right?
  - iii. what is the alignment score ("max score")?
  - iv. what is the percent identity and query coverage?
  - v. what is the E-value?
  - vi. are there any gaps in the alignment?
- (b) Repeat the search above with the UniprotKB/SwissProt database and answer the same questions as above.

- (c) Repeat search with nr as the database and PAM30 as a substitution matrix. This can be done in the blastP homepage by opening “algorithm parameters” at the bottom of the page. Observe the changes between the results of a and c due to the change in the substitution matrix: Look at the first entry that differs between a and c. What is its rank in a and c? What is the name of the protein sequence in this entry?

2. **DNA sequence alignment:** The following sequence was constructed by NCBI scientist Mark Boguski for Michael Chrichton’s “The Lost World” of the Jurassic Park series:

>DinoDNA from THE LOST WORLD p. 135

```
GAATTCGGAAGCGAGCAAGAGATAAGTCCTGGCATCAGATACAGTTGGAGATAAGGACG
GACGTGTGGCAGCTCCCGCAGAGGATTCACCTGGAAGTGCATTACCTATCCCATGGGAGCC
ATGGAGTTTCGTGGCGCTGGGGGGGCGGATGCGGGCTCCCCACTCCGTTCCCTGATGAA
GCCGGAGCCTTCTGGGGCTGGGGGGGGCGAGAGGACGGAGGCGGGGGGGCTGCTGGCC
TCCTACCCCCCTCAGGCCGCGTGTCCCTGGTGCCGTGGGCAGACACGGGTACTTTGGGG
ACCCCCCAGTGGGTGCCGCCGCCACCCAAATGGAGCCCCCCCCACTACCTGGAGCTGCTG
CAACCCCCCGGGCAGCCCCCCCCATCCCTCCTCCGGGCCCTACTGCCACTCAGCAGC
GGGCCCCACCCCTGCGAGGCCCGTGAGTGCATGGCCAGGAAGAACTGCGGAGCGACG
GCAACGCCGCTGTGGCGCCGGGACGGCACCGGGCATTACCTGTGCAACTGGGCCTCAGCC
TGCGGGCTCTACCACCGCCTCAACGGCCAGAACCGCCCGCTCATCCGCCCAAAAAGCGC
CTGCTGGTGAGTAAGCGCGCAGGCACAGTGTGCAGCCACGAGCGTGAAAAGTCCAGACA
TCCACCACCACTCTGTGGCGTCGCAGCCCCATGGGGGACCCCGTCTGCAACAACATTCAC
GCCTGCGGCCTCTACTACAAACTGCACCAAGTGAACCGCCCCCTCACGATGCGCAAAGAC
GGAATCCAAACCCGAAACCGCAAAGTTTCTCCAAGGGTAAAAAGCGGCGCCCCCGGGG
GGGGGAAACCCCTCCGCCACCGCGGGAGGGGGCGCTCCTATGGGGGAGGGGGGACCCC
TCTATGCCCCCCCCGCCGCCGCCCGCCGCCGCCGCCCTCAAAGCGACGCTCTGTAC
GCTCTCGCCCCGTGGTCTTTTCGGGCCATTTTCTGCCCTTTGGAAACTCCGGAGGGTTT
TTTGGGGGGGGGCGGGGGGTTACACGGCCCCCGGGGCTGAGCCCGCAGATTTAAATA
ATAACTCTGACGTGGGCAAGTGGGCCTTGCTGAGAAGACAGTGTAACATAATAATTTGCA
CCTCGGCAATTGCAGAGGGTCGATCTCCAATTTGGACACAACAGGGCTACTCGGTAGGAC
CAGATAAGCACTTTGCTCCCTGGACTGAAAAAGAAAGGATTTATCTGTTTGCTTCTTGCT
GACAAATCCCTGTGAAAGGTAAAAGTTCGGACACAGCAATCGATTATTTCTCGCCTGTGTG
AAATTACTGTGAATATTGTAATATATATATATATATATATATATATATCTGTATAGAACAGCC
TCGGAGGCGGCATGGACCCAGCGTAGATCATGCTGGATTTGTAAGTCCCGGAATTC
```

The sequence is also available at the course webpage: <http://www.cs.umb.edu/nurith/cs612/dino.fasta>

Perform a Blast search using blastn (nucleotide search) and the default non-redundant (nt) nucleotide database.

- (a) What are the two main species used to construct the dinosaur DNA sequence?
- (b) Repeat the search with blastx (DNA vs. protein sequence) using the default non-redundant protein sequence database. Look at the top sequence alignment and retrieve the hidden message there (hint: look at the gaps...).

## Part 2 – Theory

1. Lesk book question 5.3: The edit distance between the strings `agtcc` and `cgctca` is 3, consistent with the following alignment:

```
ag-tcc
cgctca
```

Find the sequence of three edit operations that convert `agtc` to `cgctca`.

## 2. Dynamic programming:

- Use the Needleman Wunsch global alignment Dynamic programming formula in slide set no. 2 to find the sequence alignment score of the two DNA sequences `ATCGAACTGCC` and `TACGCACTCCA`. Show the filled dynamic programming matrix using +1 for a match, -1 for a mismatch and -1 for a gap penalty in a way similar to the slide sets. Additionally, show the alignment itself.
- Repeat with the Smith-Waterman local alignment algorithm and the same scoring scheme.
- The semi-global alignment is a variant of global alignment that does not penalize for gaps at the beginning and/or the end of one of the sequences but penalizes in the middle. In other words, the conditions are the same as in global alignment:

$$S_{i,j} = \max \begin{cases} S_{i-1,j-1} + s(x_i, y_j) \\ S_{i-1,j} + g \\ S_{i,j-1} + g \end{cases}$$

but the boundary conditions are different:  $S_{i,0} = S_{0,j} = 0$ , and the alignment backtrack starts at the maximum column of the bottom row (instead of always at the bottom right like in global alignment). Calculate the semi-global alignment of the two sequences above.

- We defined the longest common subsequence (LCS) in class. Let us now define the Shortest common supersequence (SCS) of two sequences X and Y as the shortest sequence which has X and Y as subsequences. For example – The shortest common supersequence of `ABLE` and `BLUE` is `ABLUE`.
  - Describe how to obtain the SCS of two sequences, X and Y, given their LCS.
  - What is the SCS of `TACGGGTAT` and `GGACGTACG`?

## 4. Substitution matrices:

- Given the BLOSUM-62 matrix (see sequence class notes or [http://www.ncbi.nlm.nih.gov/IEB/ToolBox/C\\_DOC/lxr/source/data/BLOSUM62](http://www.ncbi.nlm.nih.gov/IEB/ToolBox/C_DOC/lxr/source/data/BLOSUM62)), find the score of the following alignment (assume this is the optimal alignment, so no need to align further):  
`THISSEQ`  
`THATSEQ`
- Repeat with the PAM-250 matrix. It can be found in Lesk, page 257, or here: [http://www.ncbi.nlm.nih.gov/IEB/ToolBox/C\\_DOC/lxr/source/data/PAM250](http://www.ncbi.nlm.nih.gov/IEB/ToolBox/C_DOC/lxr/source/data/PAM250)

- Multiple Sequence Alignment:** Extend the dynamic programming formula to 3 dimensions. What is the run time in this case? How many cases do we have to compare this time?

Hint: this time the matrix is cubic since instead of a 2-dimensional matrix we need to run on a cube of  $m \times n \times k$  where m,n, and k are the lengths of the three sequences. Every path goes from one vertex of the cube and traveling inside the cube to another vertex. Try to count how many such paths there can be.

## Part 3 – Programming

For this part you need the biopython module installed (<https://biopython.org/docs/1.75/api/index.html>) and in particular - BioBlast. In particular – look at the NCBIWWW module.

Here is an extensive tutorial: <https://biopython.org/docs/dev/Tutorial/index.html> Write a python program that implements questions 1 and 2 using BioBlast. Using NCBIWWW it should not be difficult. Use the web run

and do not try to install and run standalone BLAST! It will take a lot of space to install the databases locally. Try to have your program obtain the results in a human-readable format. By default the results are in xml, but you can change the format or use the parse tools from the NCBI XML library. **Notice:** The BLAST run may take a while (30–60 seconds or more). Attach the code as part of your answer.