



## Homework 2: Anchored Global Sequence Alignment

CSCI 5481, Computational Techniques for Genomics  
University of Minnesota  
Instructor: Dan Knights

---

### Instructions

- Please turn this assignment in on the course web page.
- There are multiple files to turn in. All text and code should be placed into a single folder with a name like *lastname\_exerciseXX*. The folder should then be compressed and submitted as a single archive (.zip or .tgz)
- You must do this work on your own, although you are encouraged to have general discussions with other students. The work you turn in must be your own. Your code will be checked for overlap and for surprising idiosyncrasies in common with other submissions.
- Please write the names of all students with whom you discussed the assignment at the top of your code.
- Please include copious comments in your code. Full credit will only be given for code that is fully commented, meaning that every line that is not completely obvious needs a comment. Partial credit may be given for broken/non-functioning code if the code is well-commented.
- You may use any programming language you wish.

### Background

This homework assignment is implementation of an anchored version of the standard Needleman-Wunsch algorithm and application of the algorithm to align PAX and HOX proteins from human and fruit fly. The anchored global sequence alignment assumes known matched regions between two sequences and applies Needleman-Wunsch algorithm to align the unaligned regions between the matched ones. Implement the anchored global sequence alignment algorithm and align the given sequences. (Hint: write Needleman-Wunsch first -- then the anchored algorithm is very simple extension of the Needleman-Wunsch algorithm and you only need to implement a wrapper function that calls your other function).

### Dataset

Download and extract the folder containing the sequences:

<https://www.dropbox.com/s/clbbgkq0tem66fj/Homework2-sequences.zip?dl=1>. For the matches files, i.e. "Match\_HOX.txt" and "Match\_PAX.txt", the first 2 columns represent the start and end positions of matched regions for the human sequence, whereas the last 2 columns represent the start and end positions of the matched regions for the fruit fly sequence.

### Input and Output Format:

You can hardcode your substitution matrix and gap penalty values (-3 for mismatches, 1 for a match, -2 for a gap; ignore the affine gaps). The command line for calling your program should be of the form:

`programname seq1.fasta seq2.fasta [matches.txt]`. Note that [matches.txt] means the third file is optional. If the matches.txt is not provided, your program should run standard Needleman-Wunsch. Output should be both the alignment score for this pair of sequences and the actual alignment itself printed

with gaps. You may also use command-line flags to label your parameters, e.g. `python programname.py -q eq1.fasta -r seq2.fasta [-m matches.txt]`.

Treat any special characters the same as the ones in alphabet, i.e. use the same match and mismatch costs.

## Tasks

1. (25 points) : Implement the Needleman-Wunsch algorithm (don't forget your comments).
2. (25 points) : Based on your Needleman-Wunsch algorithm to implement the anchored version (don't forget your comments).
3. (25 points) : Use your algorithm to align the provided two pairs of sequences. Report the alignment and the alignment score.
4. (25 points) : For both pairs of sequences, permute the amino acids in the sequences (use random library in your chosen language) and repeat the alignment 10,000 times. Report the distribution with a histogram of the random alignment score and mark the alignment score of the original sequences. Do this only for the non-anchored version.
5. (5 bonus points) : You are required to align Titin human sequence and Titin mouse sequence using anchored version you implemented in part 2. Titin is the largest known protein; its human variant consists of 34,350 amino acids, and its mouse homologue is even larger, comprising 35,213 amino acids. Hence, global alignment of these two sequences using Needleman-Wunsch will take very long time and use a lot of memory. However, using anchored-alignment approach will significantly improve the performance. Download and extract the bonus folder (<https://www.dropbox.com/s/50fx4ovd3zoeps9/Homework2-bonus-TITIN-sequences.zip?dl=1>) to get the Titin human sequence ("TITIN\_Human.fa"), Titin mouse sequence ("TITIN\_Mouse.fa"), and a match file "TITIN\_Match.txt", in which the first 2 columns represent the start and end positions of matched regions for the human sequence, and the last 2 columns represent the start and end positions of the matched regions for the mouse sequence. Report the alignment and alignment score.

## Deliverables

1. Source file (your code). Please note in your code the names of the people who worked on it.
2. Readme file (text). The readme file should contain instructions on how to compile and run the program.
3. Alignment results in a single file (text).
4. A pdf file giving the plots of the permuted and observed alignment scores in task 4.