

# Programming Project 1: Sequence Alignment

CS 181, Fall 2022

**Out:** Sep. 30

**Due:** Oct. 13, 11:59 PM

## 1 Task

You must implement the following three alignment algorithms you've learned in class:

1. Global Alignment
2. Local Alignment
3. Global Alignment with Affine Gap Penalties

Each algorithm should take in two strings  $u$  and  $v$ , a scoring matrix, and a gap penalty (or penalties in the case of affine gap). The scoring matrix will be a table that gives scores for the match or mismatch of two letters. More details on the precise form of the inputs will be given in the following section.

After implementing these algorithms, you will use them to perform alignments using real sequence data.

## 2 Program Specifications

### 2.1 Setup

To grab the support code, download from the website. We've provided example shell scripts and a few test cases and matrices.

### 2.2 Programming Language

- You may use any of the following programming languages to write your code for this class: Python 3, Java, R, and Julia. Our Gradescope autograder is configured to accept solutions written in these.
- Your solutions generally should not require the installation of any packages that do not come in the standard installations of your chosen programming language. However, if you are using Python, you will also have access to numpy and pandas.
- To facilitate anonymized and automated grading, each of your solutions must be accompanied by a shell script. Make sure each problem is able to output the correct result using the shell script provided. If you are using Python 3, your shell script must run your program using the command "python3" rather than "python", as "python" will run your code with Python 2.
- Your shell scripts should print exactly what is shown in the examples given for each problem. If you print any extra text, you will fail our autograder and lose points. This means that if you are coding in R, you may need to print text using "cat()" instead of "print()".

- Be sure to print any terminal output to stdout, which is the channel that the standard print functions write to in most programming languages. If you print to stderr, your solution will be interpreted as an error message and fail the autograder.

## 2.3 Alignment Characters

Use ‘-’ to represent gap characters. Your project should also be case-insensitive, so ‘a’ and ‘A’ should be equivalent.

## 2.4 Shell Scripts And Examples

Your handin should provide three shell scripts: `global.sh`, `local.sh`, and `affine.sh`. Each script should take three arguments: (1) a path (relative or absolute) to a file with the two sequences to align, (2) a path to a file with the scoring matrix, and (3) a gap penalty (or penalties, in the case of `affine.sh`). Your script should then print the first string aligned and then the second string aligned, followed by the alignment score, as an integer. Examples of the input files are shown below.

The sequences file just has two sequences on two separate lines. For example, here are the contents of `test_cases/01.txt`. Note that input files may or may not have newlines at the end of line 2.

```
AAAGAATTCA
AAATCA
```

The scoring matrix file has rows on separate lines and columns delimited by a single space. (The top left cell will be a simple ‘X’, to make parsing the matrix a little simpler.) An example scoring matrix can be found at `matrices/standard.m`. All test case matrices will follow the same format.

```
X C T A G
C 1 -1 -1 -1
T -1 1 -1 -1
A -1 -1 1 -1
G -1 -1 -1 1
```

Here are a few examples illustrating how we expect to be able to call your shell scripts and what we expect as output. The examples assume that you have run `cs1810_setup alignment`.

```
> sh global.sh test_cases/01.txt matrices/standard.m -1
AAAGAATTCA
AAA---T-CA
2
```

```
> sh local.sh test_cases/01.txt matrices/standard.m -1
AATCA
AAT-CA
4
```

```
> sh affine.sh test_cases/01.txt matrices/standard.m -1 0
AAAGAATCA
AAA----TCA
5
```

Note that `affine.sh` requires a gap opening penalty and a gap extension penalty (**in that order!**).

## 2.5 README

All projects must include a README file. Include the following information:

- A description of any known bugs.
- Anything you want the TAs to know about your project.

## 3 Patient Sample Collection and Consent

In April 2020, the UNESCO International Bioethics Committee (IBC) and the UNESCO World Commission on the Ethics of Scientific Knowledge and Technology (COMEST) released a [statement](#) addressing ethical concerns during the COVID-19 pandemic, one of them being the importance of responsible research practices.

Patient-collected samples play a crucial role in COVID-19 research- these samples enable the construction of COVID-19 reference genomes, one of which you will be working with in the next portion of this assignment.

When dealing with individual patients' samples, researchers must receive consent before moving forward. Read [this](#) paper about navigating consent in the context of the COVID-19 pandemic.

1. What are possible accessibility issues that may impact consent?
2. Would you use samples from each of the following patients for research purposes? What concerns might you have and how could they be addressed? Describe your thought process.
  - (a) A patient does not speak English so their next of kin relays their verbal consent, although it is not clear what the patient is saying.
  - (b) A patient of specific interest is unaware of their surroundings and physically unable to consent. They have no legal representative or next of kin.
  - (c) A patient is unable to consent and a distant relative is the only person accompanying them. Their personal relationship is not clear.

- (d) A patient's sample is collected as part of routine treatment. Their nurse mentions as an aside that the sample might also be used for research. They are unable to elaborate further, so the patient reluctantly says yes.
- (e) A patient is shown a 10-minute long video giving an overview of how their sample might be used, although the exact research pipeline is not yet clear. They can choose to stop the video at any point to sign a consent.

## 4 Application

Before the advent of rapid COVID testing, scientists had to rely on genome alignment to figure out if individuals were exposed to *SARS-CoV-2*. Imagine you are a computational biologist working in an early clinical trial for a vaccine. While recruiting patients, you encounter three individuals who all of a sudden feel 'under the weather'. While you suspect COVID-19, you must confirm their disease status using a saliva sample. How will you do it? In this part of the assignment, you will use your very own code to do so!

To be completely sure that you have a match, you should probably align the whole genome of each patient's samples to the whole reference genome of *SARS-CoV-2*. But running such large alignments on the possibly hundreds of different microbes in a person's mouth can be computationally prohibitive. So instead, we will focus on a single gene: *Spike or surface glycoprotein protein (S)*. The spike glycoprotein is quite characteristic of a coronavirus, allowing it to adhere to and enter healthy cells, and a prime target for initiating an immune response.

The setup script provides you with a `data` directory containing the three sample protein sequences from our sick patients. We also provide you with three BLOSUM matrices for amino acids, although you are welcome to edit these scoring matrices or use another scoring scheme that you deem fit (e.g. a PAM matrix). Your task is to determine which of these individuals might have been infected by *SARS-CoV-2* by running a global alignment on the samples. You should use the [protein sequence of the reference S gene](#) to complete this section. Please report and submit the following in a PDF file called `application.pdf`:

- A reproducible description of your workflow. You should provide enough detail for us to trace your steps exactly from start to finish. Please tell us if you wrote any new scripts or scoring matrices. If you save your outputs to a file, mention it. If you run any commands in the shell, tell us which commands exactly. If you create new files or edit files by hand, we should know that too. The point is that it should be clear exactly how you generate outputs from data.
- The percent identity for each of the three sequences with the reference gene. Percent identity is defined as  $\frac{\# \text{ matches}}{\text{length of longer string}}$ . (This quantity will of course depend on your scoring scheme.)
- A justification for your chosen scoring scheme. Discuss the pros and cons of your scoring scheme.
- Finally, whether you believe any of the three individuals may have been exposed to *SARS-CoV-2* or not. Tell us why in each case.

## 5 Handin

- To hand in this assignment, upload your shell scripts along with all files needed to run your code on Gradescope. If you have any subdirectories or folders that you would like to preserve in your handin, you will need to compress your submission into a zip file and upload the zip file to Gradescope. However, to ensure that the autograder runs your handin properly, make sure that all shell scripts are present at the root of your handin; in other words, do not place your shell scripts inside any folders.
- When you hand in your solution, our autograder will immediately run on your submission. We have made the test cases provided in the project handout immediately visible to you so that you can ensure your solution runs properly with the autograder. If you fail the autograder and cannot determine why, notify a TA and we will help you to diagnose the problem. If your final submission is not compatible with our autograder, it will be very difficult for us to give you credit for this assignment.

## 6 Optional: Computational Marine Biology Research!

After their visit to the New England Aquarium, the CS 181 TAs wanted to apply their computational biology knowledge to make novel discoveries about marine animal biology. After countless hours of beautiful algorithm development and manuscript drafting, the great Sorin Istrail reminds them of the dangers plagiarizing in the world of scientific research. For extra credit, add a few sentences to your application.pdf file explaining to the CS 181 TAs how they could use alignment to help prove their paper is not similar to previously published papers.

## 7 Grading

We will grade your handin using pre-generated test cases, with a certain number of points allocated per test case. Your code should not throw exceptions for any valid inputs. Test edge cases extensively! We will be checking for common gotchas, including but not limited to the following:

1. Global alignment of two very short strings, e.g. ‘C’ and ‘G’, using various scoring schemes.
2. Local alignment of two very short strings, e.g. ‘C’ and ‘G’, using various scoring schemes.

You will be graded solely on the correctness of your implementation. You will not be graded on code performance, as long as it is within reasonable bounds. For this project, we will give a lenient limit of one minute for aligning two 500bp sequences. (If implemented wisely, these alignments should take less than half a second to perform.)

## 8 Tips

- We do not provide you with pseudocode for the algorithms in this handout. Instead, you should refer to the lecture notes/slides.

- If you are having trouble with the shell scripts or with parsing the input files, try temporarily hard-coding some test cases and scoring schemes and returning to parsing later or come to office hours. We want you to focus on the implementation of the algorithm, not the technical details of I/O.
- In the past, students have had more problems with backtracking than with filling out the dynamic programming table. Think carefully about what you need to do to backtrack. Again, feel free to come to office hours if you need help checking your logic.