

Midterm Exam

CS 4364/5364

Spring 2021

Instructions: Please read all of the instructions below before you begin:

- Read all 6 of the questions in the exam before you begin.
- Questions marked with a dagger (†) are required for students in CS 5364, and bonus (optional) for those in CS 4364
- This document will be released to students by Midnight (12:00 am) on 11 March 2021, and will be due by 11:59pm the same day.
- Your submission should be sent to Dr. DeBlasio (dan@deblasiolab.org) by the deadline, please also send Dr. DeBlasio a private message on MS Teams (teamschat.deblasiolab.org) to inform him of your submission in case something happens with the email delivery.
- All submissions should be made as a single PDF file with all of your responses.
- All students are permitted to submit their assignments as either a typed document or as hand-written responses (scanned and clearly readable).
- Figured (pictograms) can be included if they help describe a solution, and are encouraged if they are clear.
- Remember that unanswered questions will receive 0 credit, any reasonable try at a response will receive at least half-credit. If you feel you're unable to provide a reasonable answer to a question, you can answer with "I cannot provide a reasonable attempt for this question", which will be provided quarter-credit.
- Dr. DeBlasio will have his normal office hours 1pm-2pm on test day (on teams as normal).
- The class period, 3-4:20pm, on the test day, will be replaced with an open question session using Zoom (the same class link).
- Additional questions can be posed using teams private messages, but note that questions outside the times above may receive delayed responses.
- Monitor the 'general' channel on the course team (specialtopics.deblasiolab.org/s21/teams) for errata corrections.
- **Warning:** Questions asked after 5pm on test day may not receive responses before the exam is due.

1. Use the *Needleman-Wunch* dynamic programming table for $S = \text{CTACTGTGT}$ and $T = \text{CACCCCTGT}$ below to the next questions.

		C	T	A	C	T	G	T	G	T
	0	$\leftarrow -0.5$	$\leftarrow -1$	$\leftarrow -1.5$	$\leftarrow -2$	$\leftarrow -2.5$	$\leftarrow -3$	$\leftarrow -3.5$	$\leftarrow -4$	$\leftarrow -4.5$
C	$\uparrow -0.5$	$\nearrow 5$	$\leftarrow 4.5$	$\nwarrow 4$	$\nwarrow 3.5$	$\leftarrow 3$	$\leftarrow 2.5$	$\leftarrow 2$	$\leftarrow 1.5$	$\leftarrow 1$
A	$\uparrow -1$	$\uparrow 4.5$	$\nwarrow 4$	$\nwarrow 9.5$	$\leftarrow 9$	$\leftarrow 8.5$	$\leftarrow 8$	$\leftarrow 7.5$	$\leftarrow 7$	$\leftarrow 6.5$
C	$\uparrow -1.5$	$\uparrow 4$	$\nwarrow 3.5$	$\uparrow 9$	$\nwarrow 14.5$	$\leftarrow 14$	$\leftarrow 13.5$	$\leftarrow 13$	$\leftarrow 12.5$	$\leftarrow 12$
C	$\uparrow -2$	$\nwarrow 3.5$	$\nwarrow 3$	$\uparrow 8.5$	$\nwarrow 14$	$\nwarrow 13.5$	$\nwarrow 13$	$\nwarrow 12.5$	$\nwarrow 12$	$\nwarrow 11.5$
C	$\uparrow -2.5$	$\nwarrow 3$	$\nwarrow 2.5$	$\uparrow 8$	$\nwarrow 13.5$	$\nwarrow 13$	$\nwarrow 12.5$	$\nwarrow 12$	$\nwarrow 11.5$	$\nwarrow 11$
C	$\uparrow -3$	$\nwarrow 2.5$	$\nwarrow 2$	$\uparrow 7.5$	$\nwarrow 13$	$\nwarrow 12.5$	$\nwarrow 12$	$\nwarrow 11.5$	$\nwarrow 11$	$\nwarrow 10.5$
T	$\uparrow -3.5$	$\uparrow 2$	$\nwarrow 7.5$	$\uparrow 7$	$\uparrow 12.5$	$\nwarrow 18$	$\leftarrow 17.5$	$\nwarrow 17$	$\leftarrow 16.5$	$\nwarrow 16$
G	$\uparrow -4$	$\uparrow 1.5$	$\uparrow 7$	$\nwarrow 6.5$	$\uparrow 12$	$\uparrow 17.5$	$\nwarrow 23$	$\leftarrow 22.5$	$\nwarrow 22$	$\leftarrow 21.5$
T	$\uparrow -4.5$	$\uparrow 1$	$\nwarrow 6.5$	$\nwarrow 6$	$\uparrow 11.5$	$\nwarrow 17$	$\uparrow 22.5$	$\nwarrow 28$	$\leftarrow 27.5$	$\nwarrow 27$
G	$\uparrow -5$	$\uparrow 0.5$	$\uparrow 6$	$\nwarrow 5.5$	$\uparrow 11$	$\uparrow 16.5$	$\nwarrow 22$	$\uparrow 27.5$	$\nwarrow 33$	$\leftarrow 32.5$

- (a) How many co-optimal alignments of the two strings are there?
- (b) What is the optimal alignment of $S[1\dots 3]$ and $T[1\dots 5]$? (note these are prefixes CTA and CACCC)
- (c) What is the mismatch penalty used to construct the table? match score? indel penalty?
- (d) [†] Using *only the scores in the table above* is it possible to determine the score of the optimal alignment of $S[4\dots 9]$ and $T[6\dots 10]$? Why or why not?

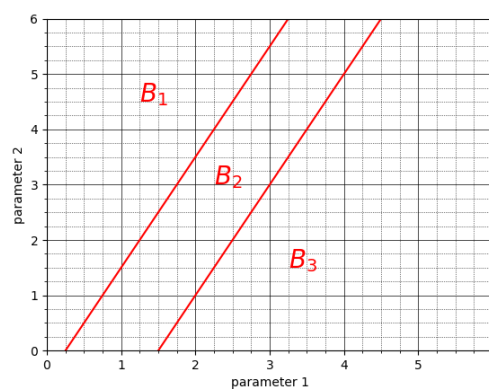
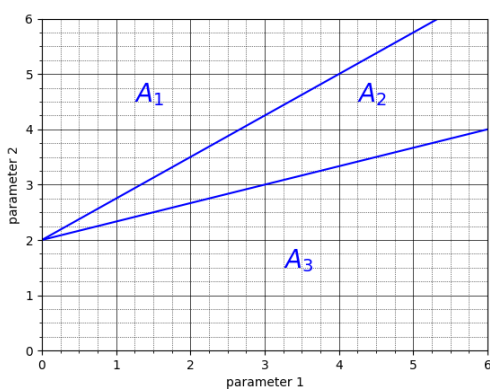
2. Use the alignments and plots below to answer the following questions:

- Calculate the accuracies for the following groups of alignments, the reference alignment is provided at the top of each column.
- Given the accuracies and the parameter decompositions shown in the figures, what is the region of the parameter space (identify the corners of the polygon) that provides the best alignments on average across these two pairs of sequences.

Remember accuracy, with respect to a reference, is the fraction of columns from the reference that are recovered in a computed alignment; and that each region of the plots corresponds to a set of parameters that produce the labeled alignment.

Reference A	<div>ATG-CTGGAT</div> <div>-TGA-TCGAT</div>	Reference B	<div>TTGTGTCC--</div> <div>TT-T-TCCAA</div>
A_1	<div>ATG-CT-GGAT</div> <div>-TGA-TC-GAT</div>	B_1	<div>TTGTGTCC</div> <div>TTTTCCAA</div>
A_2	<div>ATGCTGGAT</div> <div>-TGATCGAT</div>	B_2	<div>TTGTGTCC--</div> <div>TTTT--CCAA</div>
A_3	<div>ATG--CTGGAT</div> <div>-TGATC--GAT</div>	B_3	<div>TTGTGTCC--</div> <div>--TTTTCCAA</div>

	A_1	A_2	A_3	B_1	B_2	B_3
Accuracies:						



3. Use the suffix tree below to answer the next questions.

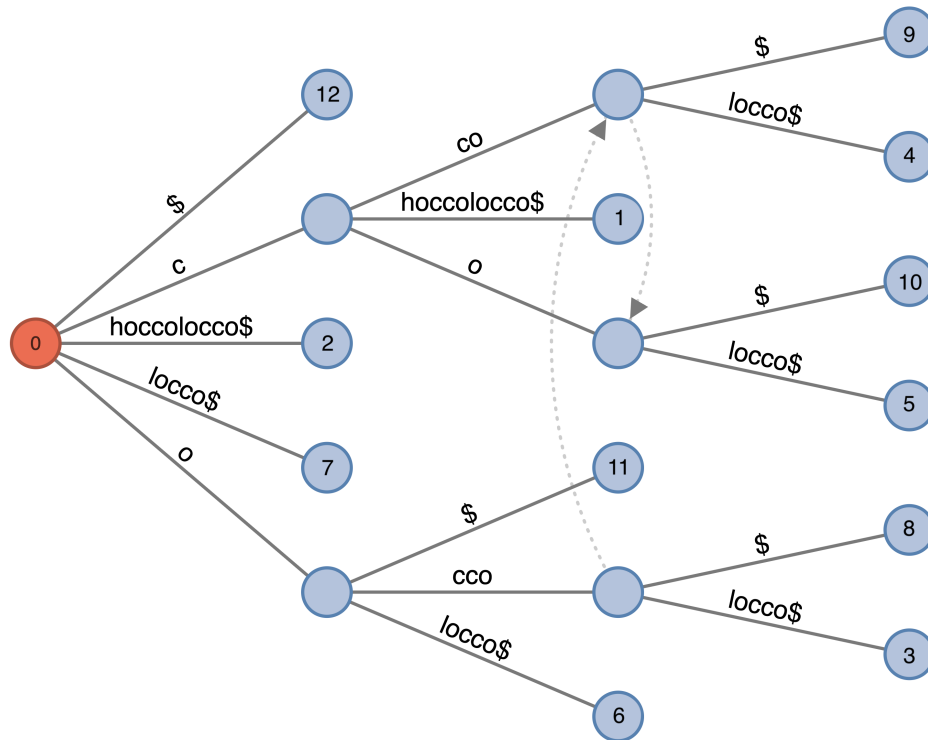


Figure 1: Suffix tree for question 3

- What is the full string for which this is the suffix tree?
- What is the longest sub-string that occurs 2 times?
- What is the lexicographically smallest suffix that is strictly longer than 1 character? (that is, its more than just \$)

4. Below is an ILP for pairwise global sequence alignment with several constraints missing. The scoring uses a match score of α , a mismatch penalty of β and an indel penalty of γ . Define these constraints.

$$\begin{array}{ll}
 \text{maximize} & \alpha \sum_{i,j} X_{ij} - \beta \sum_{i,j} Y_{ij} - \gamma \left(\sum_i Z_i^S + \sum_j Z_j^T \right) \\
 \text{subject to} & \sum_j X_{ij} + \sum_j Y_{ij} + Z_i^S = 1 \quad \forall i \\
 & \text{-----} \quad \forall j \\
 & \text{-----} \quad \forall i, j : S[i] = T[j] \\
 & \text{-----} \quad \forall i, j : S[i] \neq T[j] \\
 & \text{-----} \quad \forall i < i', j > j' \\
 & X_{ij} \in \{0, 1\}, \quad \forall i, j \\
 & Y_{ij} \in \{0, 1\}, \quad \forall i, j \\
 & Z_i^S \in \{0, 1\}, \quad \forall i \\
 & Z_j^T \in \{0, 1\}, \quad \forall j
 \end{array}$$

Hints:

- The first constraint enforces that each position i in S can only be a match, a mismatch or a deletion, the second constraint will do something similar but for T .
- The 3rd and 4th constraints decide if a position is a match or mismatch, remember when defining an ILP we usually say what a value *can't be* based on the information in the problem.
- The 5th constraint will be similar to the one we defined for LCS, but here we have two variables that can define a match between any two indexes.

5. [†] True or False: When computing the sum of pairs score of a multiple sequence alignment

$$id + 2mt + 2ms = L \binom{k}{2},$$

where id is the total number of indels, ms is the total number of mismatches, mt is the total number of matches, k is the number of strings aligned, and L is the length of the alignment itself. Justify your answer.

6. (bonus for all) What is the item that is the answer to question 3a? (You can use google for this one)