

Homework 3: due by 11:59pm on Wednesday, March 29th, 2022

(Total: 100 points)
Mack Joyner

Commit all work to the github classroom hw3 repo at <https://classroom.github.com/a/4ir64HSD> that we created for you. In case of problems committing your files, please contact the teaching staff at comp322-staff@mailman.rice.edu before the deadline to get help resolving your issues.

Your solution to the written assignment should be submitted as a PDF file named `hw3_written.pdf` in your github classroom hw3 top level directory. This is important — you will be penalized 5 points if you place the file in some other folder or with some other name. The PDF file can be created however you choose. If you scan handwritten text, make sure that the writing is clearly legible in the scanned copy. Your solution to the programming assignment should be submitted in the appropriate location in the hw3 directory.

The slip day policy for COMP 322 is similar to that of COMP 321. All students will be given 3 slip days to use throughout the semester. When you use a slip day, you will receive up to 24 additional hours to complete the assignment. You may use these slip days in any way you see fit (3 days on one assignment, 1 day each on 3 assignments, etc.). If you plan to use a slip day, you need to say so in an github committed README.md file before the deadline. You should specifically mention how many slip days you plan to use. The README.md file should be placed in the top level directory (e.g. hw3). Other than slip days, no extensions will be given unless there are exceptional circumstances (such as severe sickness, not because you have too much other work). Such extensions must be requested and approved by the instructor (via e-mail, phone, or in person) before the due date for the assignment. Last minute requests are likely to be denied.

If you see ambiguity or inconsistency in a question, please seek clarification on Piazza (remember not to share homework solutions in public posts) or from the teaching staff. If it is not resolved through those channels, you should state the ambiguity/inconsistency that you see, as well as any assumptions that you make to resolve it.

Honor Code Policy: All submitted homework is expected to be the result of your individual effort. You are free to discuss course material and approaches to problems with your other classmates, the teaching assistants and the instructors, but you should never misrepresent someone else's work as your own. If you use any material from external sources, you must provide proper attribution.

1 Written Assignment (30 points total)

As mentioned earlier, your solution to the written assignment should be submitted as a PDF file named `hw3_written.pdf` in the hw3 directory.

1.1 Amdahl's Law (15 points)

In Lecture 13 (Topic 1.5), you will learn the following statement of Amdahl's Law:

If $q \leq 1$ is the fraction of WORK in a parallel program that must be executed sequentially, then the best speedup that can be obtained for that program, even with an unbounded number of processors, is $\text{Speedup} \leq 1/q$.

Now, consider the following generalization of Amdahl's Law. Let q_1 be the fraction of WORK in a parallel program that must be executed sequentially, q_2 be the fraction of WORK that can use at most 2 processors, and $(1 - q_1 - q_2)$ the fraction of WORK that can use an unbounded number of processors. The fractions of

WORK represented by q_1 , q_2 , and $(1 - q_1 - q_2)$ are disjoint, and cannot be overlapped with each other in a parallel execution. Your assignment is as follows:

1. (10 points) Provide the best possible (smallest) upper bound on the Speedup as a function of q_1 and q_2 .
2. (5 points) Explain your answer, and justify why it is a correct upper bound. Use the cases when $q_1 = 0$, $q_2 = 0$, $q_1 = 1$ or $q_2 = 1$ to explain why your bound is correct.

Hints:

- As with Amdahl's Law, your answer should not include the number of processors, P . It should be an upper bound that applies to all values of P .
- To check your answer, consider the cases when q_1 or q_2 are equal to 0 or 1.

1.2 Finish Accumulators (15 points)

Consider the pseudocode shown below in Listing 1 for a Parallel Search algorithm that is intended to compute Q , the number of occurrences of the pattern array in the text array. What possible values can variables $count0$, $count1$, and $count2$ contain at line 16? Write your answers in terms of M , N , and Q , and explain your answers.

```
1 // Assume that count0, count1, count2 are declared
2 // as object/static fields of type int
3 . . .
4 count0 = 0;
5 accumulator a = new accumulator(SUM, int.class);
6 finish (a) {
7     for (int i = 0; i <= N - M; i++)
8         async {
9             int j;
10            for (j = 0; j < M; j++) if (text[i+j] != pattern[j]) break;
11            if (j == M) { count0++; a.put(1); } // found at offset i
12            count1 = a.get();
13        } // for-async
14    } // finish
15    count2 = a.get();
16 // Print count0, count1, count2
```

Listing 1: Parallel Search using Finish Accumulators

Hints based on common errors from past years: Be sure to include exactly all possible values for the variables, not a subset or superset of the values. Remember to use Q in your answers, even though Q can have different values for different values of the `text[]` and `pattern[]` arrays (even for the same values of M and N .) Finally, don't forget to explain your answers.

2 Programming Assignment (70 points)

2.1 Habanero-Java Library (HJ-lib) Setup

See the Lab 1 handout for instructions on HJ-lib installation for use in this homework.

2.2 Pairwise Sequence Alignment

This homework focuses on computing the optimal score for pairwise sequence alignment, not on the alignment itself. Though a biologist is ultimately interested in seeing the alignment, there are many applications where the score alone is of interest. For example, in multiple sequence alignment, the most commonly used approach is called progressive alignment, where an evolutionary tree is first built based on the scores of pairwise alignments, and then the tree is used as a guide for doing the multiple sequence alignment. In this case, the pairwise alignments are performed solely for the sake of obtaining scores, and the alignments themselves are not needed. However, it is important to compute the scores as quickly as possible when exploring alignments of large DNA sequences.

Let X and Y be two sequences over alphabet Σ (for DNA sequences, $\Sigma = \{A, C, T, G\}$). An *alignment* of X and Y is two sequences X' and Y' over the alphabet $\Sigma \cup \{-\}$, where X' is formed from X by adding only dashes to it, and Y' is formed from Y by adding only dashes to it, such that

- 1 $|X'| = |Y'|$ i.e., X' and Y' have the same size,
- 2 there does not exist an i such that $X'[i] = Y'[i] = -$

This alignment is also referred to as *global pairwise alignment* (as opposed to *local pairwise alignment*, which is used to align selected regions of sequences X and Y).

Sequence alignment helps biologists make inferences about the evolutionary relationship between two DNA sequences. Aligning two sequences amounts to “reverse engineering” the evolutionary process that acted upon the two sequences and modified them so that their characters and their lengths differ. As an example, one possible alignment of the two sequences and $X = TACGGT$ and $Y = ACCT$ is as follows:

$$\begin{array}{rcccccc} X' & = & T & A & C & G & G & T \\ Y' & = & - & A & C & - & C & T \end{array}$$

As you may imagine, there may be multiple alignments for the same pair of sequences. For example, a trivial alternate alignment for X and Y is as follows:

$$\begin{array}{rcccccccc} X'' & = & - & - & - & - & T & A & C & G & G & T \\ Y'' & = & A & C & C & T & - & - & - & - & - & - \end{array}$$

2.3 Scoring in Pairwise Sequence Alignment: Optimality Criterion

As discussed above, a number of alignments exist for a given pair of sequences; therefore, we define a *scoring scheme* that gives higher scores to “better” alignments. Once the scoring scheme is defined, we seek an alignment with the highest score (among all feasible alignments). For DNA, a scoring scheme is given by a 5×5 matrix M , where for $p, q \in \{A, C, T, G\}$, $M_{p,q}$ specifies the score for aligning p in sequence X' with q in sequence Y' , $M_{p,-}$ denotes the penalty for aligning p in sequence X' with a dash in sequence Y' , and $M_{-,q}$ denotes the penalty for aligning q in sequence Y' with a dash in sequence X' . Assuming $|X'| = |Y'| = k$, the score of the alignment is

$$\sum_{i=1}^k M_{X'[i], Y'[i]}. \tag{1}$$

For this assignment, we will assume the following scoring scheme: $M_{p,p} = 5$, $M_{p,q} = 2$ (for $p \neq q$), $M_{p,-} = -4$ and $M_{-,q} = -2$.

For this scoring scheme, the score of the (X', Y') alignment in Section 2.2 is

$$M_{T,-} + M_{A,A} + M_{C,C} + M_{G,-} + M_{G,C} + M_{T,T} = (-4) + 5 + 5 + (-4) + 2 + 5 = 9$$

and the score of the (X'', Y'') alignment is $6 \times M_{p,-} + 4 \times M_{-,q} = -32$.

2.4 Sequential Algorithm to compute the Optimal Scoring for Pairwise Sequence Alignment

In this problem, we introduce a sequential dynamic programming algorithm (called the Smith-Waterman algorithm) to compute the Optimal Scoring for Pairwise Sequence Alignment. For two sequences X and Y of lengths m and n , respectively, denote by $S[i, j]$, $0 \leq i \leq m$ and $0 \leq j \leq n$, the score of the best alignment of the first i characters of X with the first j characters of Y . The boundary values are, $S[i, 0] = i * M_{p,-}$ and $S[0, j] = j * M_{-,p}$. It has been shown that this optimal scoring can be defined as follows $\forall i, j \geq 1$:

$$S[i, j] = \max \begin{cases} S[i-1, j-1] + M_{X[i],Y[j]} \\ S[i-1, j] + M_{X[i],-} \\ S[i, j-1] + M_{-,Y[j]} \end{cases} . \quad (2)$$

The above definition directly leads to a sequential dynamic programming algorithm that can be implemented as shown in Listing 2. Assume that the input sequences are represented as Java strings, and the scoring matrix, S , is represented as a 2-dimensional array of size $(X.length()+1) \times (Y.length()+1)$. After the algorithm terminates, the final score is available in $S[X.length()][Y.length()]$.

The dependence structure of the iterations in Listing 2 is shown in Figure 1. The cells in the figure correspond to $S[i, j]$ values, and the arrows show the dependences among the $S[i, j]$ computations.

```

1  for ( i = 1; i <= xLength; i++)
2  for ( j = 1; j <= yLength; j++) {
3  char xChar = X.charAt(i-1);
4  char YChar = Y.charAt(j-1);
5  int diagScore = S[i-1][j-1] + M[charMap(xChar)][charMap(YChar)];
6  int topColScore = S[i-1][j] + M[charMap(xChar)][0];
7  int leftRowScore = S[i][j-1] + M[0][charMap(YChar)];
8  S[i][j] = Math.max(diagScore, Math.max(leftRowScore, topColScore));
9  }
10 int finalScore = S[xLength][yLength];

```

Listing 2: Sequential implementation of Smith-Waterman Algorithm for Optimal Scoring for Pairwise Sequence Alignment

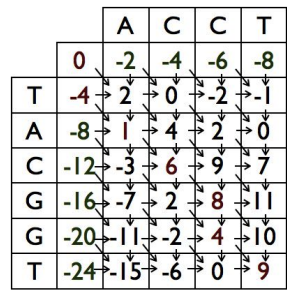


Figure 1: Dependences in Pairwise Sequence Alignment

This homework focuses on computing the optimal score for pairwise sequence alignment, not on the alignment itself. Though a biologist is ultimately interested in seeing the alignment, there are many applications where the score alone is of interest. For example, in multiple sequence alignment, the most commonly used approach is called progressive alignment, where an evolutionary tree is first built based on the scores of pairwise alignments, and then the tree is used as a guide for doing the multiple sequence alignment. In this case, the pairwise alignments are performed solely for the sake of obtaining scores, and the alignments themselves are not needed. However, it is important to compute the scores as quickly as possible when exploring alignments of large DNA sequences.

2.5 Your Assignment: Parallel Optimal Scoring for Pairwise Sequence Alignment

Your assignment is to design and implement parallel algorithms for optimal scoring for pairwise sequence alignment. We have provided a sequential implementation of the algorithm in `SeqScoring.java` that you can use as a starting point. You should not use phasers for checkpoint 1. You should not use any parallel data constructs inside a constructor. Your homework deliverables are as follows:

1. **[Checkpoint 1 due on Monday, March 6th, 2023: Ideal parallelism with abstract execution metrics (15 points)]** Examine the dependence structure for $S[i, j]$ defined in Section 2.4 and create an ideal parallel version called `IdealParScoring.java` that computes the same output as `SeqScoring.java`, and delivers the maximum ideal parallelism ignoring all overheads. For analysis of ideal parallelism, assume that computing a single element of $S[i, j]$ takes one unit of time, *e.g.*, by inserting a call to `doWork(1)` between lines 7 and 8 of Listing 2. You will need to insert this `doWork()` call at the appropriate location in the parallel solution you create in `IdealParScoring`. Your solution will be evaluated using HJlib's abstract metrics. There's no cost to create an HJlib task. You should not use phasers for Checkpoint 1. You may use phasers after Checkpoint 1.

For this checkpoint, we have some unit tests in `Homework3Checkpoint1CorrectnessTest.java`.

Hints based on common errors/omissions from past years: Remember to check that your solution passes all unit tests, and that you don't have any checkstyle errors.

2. **[Checkpoint 2 due on Wednesday, March 22nd, 2023: Useful parallelism on NOTS compute nodes (20 points)]** Create a new parallel version of `SeqScoring.java` that is designed to achieve the smallest execution time using 16 cores on a dedicated NOTS compute node. Note that this is real execution time, not abstract metrics. Your code for this part will need to go into the `UsefulParScoring.java` file.

For this part of the assignment, we recommend first debugging your solution on small strings for correctness (which can be done on any platform) using `Homework3Checkpoint2CorrectnessTest`, and then evaluating the performance of your implementation with pairs of strings of length $O(10^4)$ on dedicated NOTS compute nodes using `Homework3PerformanceTest`. Lab 5 will provide instructions on submitting jobs to the NOTS cluster.

Even though each NOTS compute node has 32GB (or more) of memory, we will evaluate all homeworks using a maximum heap size of 8GB. The JVM heap size for tests running on NOTS is set in the provided `pom.xml`. Your submission will be evaluated with 8GB of heap, so changing this value in your `pom.xml` may result in incorrect test results. If you are running the JUnit tests locally through IntelliJ rather than using the provided `pom.xml`, you will want to add the following JVM command line argument to your Run Configurations to ensure that the JVM launched by IntelliJ is allowed to allocate up to 8GB of memory (in addition to the `-javaagent` argument that should already be there): `-Xmx8192m`

However, note that most laptops do not have 8GB of physical memory, so running some of the larger tests locally may be prohibitively slow as your machine swaps memory pages out to disk as you exceed physical memory capacity.

NOTE: a solution to the sparse memory solution for the final homework submission is also acceptable as a solution for Checkpoint 2. However, we feel that many students will benefit from first completing Checkpoint 2 for a dense memory version before starting on the sparse memory version below.

For this checkpoint, we have provided a set of unit tests in `Homework3Checkpoint2CorrectnessTest.java`. The `testUsefulParScoring` and `testUsefulParScoring2` tests in `Homework3PerformanceTest.java` also evaluate the performance of your `UsefulParScoring` implementation against the sequential version. We have provided a SLURM file under `src/main/resources` that can be used on NOTS to submit `Homework3PerformanceTest` for testing on a compute node. Note that you will need to edit this SLURM file to supply your e-mail for notification and to provide the correct path to your `hw3` folder on NOTS.

Hints based on common errors/omissions from past years: Remember to check that your solution passes all unit tests, and that you don't have any checkstyle errors. Also, you should aim to get a speedup of $\geq 8\times$ for this part; you will get a 1-point deduction if the speedup is in the $[7, 8)$ range, a 2-point deduction if it is in the $[6, 7)$ range, etc.

3. [Final submission: Sparse memory version and useful parallelism on NOTS compute nodes (25 points)]

The sequential algorithm outlined in Listing 2 and `SeqScoring.java` allocates and uses a dense two-dimensional matrix which requires $O(n^2)$ space when processing strings of size $O(n)$. The goal of this part of the assignment is to create a *sparse* memory version of the program that can process strings of length $O(10^5)$ or greater by using space that is less than $O(n^2)$. The key idea to think about is what data really needs to be retained as the computation advances. For example, in the sequential version, row 1 of the S matrix can be freed (set to null and garbage-collected) when the computation reaches row 3, since computation of row 3 only needs row 2 and not row 1. As a reminder, since the JVM uses automatic memory management through garbage collection, an object cannot be freed unless there are no remaining references to it. To produce a space-efficient version, you will need to ensure that no unnecessary data is retained, i.e. that it is not referenced by any variables in your implementation.

You will need to design and implement an analogous approach to reducing the space requirements of the parallel version. This will require reworking the data structure for matrix S , and may even require using a different algorithm (with different parallel constructs) from `UsefulParScoring.java`. Your code for this part will need to go into the `SparseParScoring.java` file.

As before, we recommend first debugging your solution on small strings for correctness, and then testing with pairs of strings of length $O(10^5)$ on dedicated NOTS compute nodes, with the heap size set to 8GB. For reliable timing measurements, *be sure to run these computations only on NOTS compute nodes, and not on the login node*. Also, there may be some impact of Java's garbage collection (GC) on the performance you observe. Please contact a teaching staff member if you believe that GC overheads are interfering with your performance measurements.

For the final submission, we have provided a small set of unit tests in `Homework3Checkpoint3CorrectnessTest.java`. The `testSparseParScoring` and `testSparseParScoring2` tests in `Homework3PerformanceTest.java` also evaluate the performance of your `SparseParScoring` implementation against the sequential version. We have provided a SLURM file under `src/main/resources` that can be used on NOTS to submit `Homework3PerformanceTest` for testing on a compute node.

Hints based on common errors/omissions from past years: Remember to check that your solution passes all unit tests, and that you don't have any checkstyle errors. Also, you should aim to get a speedup of $\geq 3\times$ for this part; you will get a 1-point deduction if the speedup is in the $[2.5, 3)$ range, a 2-point deduction if it is in the $[2, 2.5)$ range, etc. Finally, remember to include all the information listed below (summarize design of all three parallel versions, include performance results, etc.)

4. [Homework report (10 points)] With the final submission you should submit a report file, formatted as a PDF file named `hw3_report.pdf`, summarizing the design of your parallel algorithms in `IdealParScoring.java`, `UsefulParScoring.java`, and `SparseParScoring.java` explaining why you believe that each implementation is correct and data-race-free. Your report should also include the following measurements for `UsefulParScoring.java` and `SparseParScoring.java`:

- (a) Execution time of `SeqScoring.java` and `UsefulParScoring.java` on a NOTS compute node with inputs of length 10,000. You can get these numbers from a run of the `Homework3PerformanceTest.testUsefulParScoring` test on NOTS (manually).
- (b) Execution time of sequential and parallel versions of `SparseParScoring.java` with inputs of length 100,000. Note that you will need a single-threaded execution of `SparseParScoring` for this item, not a run of `SeqScoring` as it will run out of memory. You can get these measurements from the `Homework3PerformanceTest.testSparseParScoring` test on NOTS (manually).