



Module: The 2011 European E. coli Outbreak (Week 1 out of 3)  
Course: Genome Assembly Programming Challenge (Course 6 out of 6)  
Specialization: Data Structures and Algorithms

# Programming Assignment 1: Assembling phi X174 Using Overlap Graphs

Revision: October 25, 2016

## Introduction

Welcome to the first programming assignment of the [Genome Assembly Programming Challenge](#)! In this assignment, you will be practicing assembling the phi X174 genome using *overlap graphs*.

## Passing Criteria: 1 out of 2

Passing this programming assignment requires passing at least 1 out of 2 code problems from this assignment. In turn, passing a code problem requires implementing a solution that passes all the tests for this problem in the grader and does so under the time and memory limits specified in the problem statement.

## Contents

<b>1 Dataset Problem: Assembling the phi X174 Genome from Error-Free Reads Using Overlap Graphs</b>	<b>3</b>
<b>2 Dataset Problem: Assembling the phi X174 Genome from Error-Prone Reads Using Overlap Graphs</b>	<b>5</b>
<b>3 General Instructions and Recommendations on Solving Algorithmic Problems</b>	<b>7</b>
3.1 Reading the Problem Statement . . . . .	7
3.2 Designing an Algorithm . . . . .	7
3.3 Implementing Your Algorithm . . . . .	7
3.4 Compiling Your Program . . . . .	7
3.5 Testing Your Program . . . . .	9
3.6 Submitting Your Program to the Grading System . . . . .	9
3.7 Debugging and Stress Testing Your Program . . . . .	9
<b>4 Frequently Asked Questions</b>	<b>10</b>
4.1 I submit the program, but nothing happens. Why? . . . . .	10
4.2 I submit the solution only for one problem, but all the problems in the assignment are graded. Why? . . . . .	10
4.3 What are the possible grading outcomes, and how to read them? . . . . .	10
4.4 How to understand why my program fails and to fix it? . . . . .	11
4.5 Why do you hide the test on which my program fails? . . . . .	11
4.6 My solution does not pass the tests? May I post it in the forum and ask for a help? . . . . .	12

- 4.7 My implementation always fails in the grader, though I already tested and stress tested it a lot. Would not it be better if you give me a solution to this problem or at least the test cases that you use? I will then be able to fix my code and will learn how to avoid making mistakes. Otherwise, I do not feel that I learn anything from solving this problem. I am just stuck. . . . 12

# 1 Dataset Problem: Assembling the phi X174 Genome from Error-Free Reads Using Overlap Graphs

In dataset problems, your solution is going to be tested against a *single* dataset as opposed to problems in the previous classes in this specialization. For this reason, there is no Constraints section in the problem description below. The sample section shows just a similar dataset. Your program is not going to be tested on this dataset.

## Problem Introduction

In this challenge, you will be given the task of performing Genome Assembly on a simulated error-free sequencing dataset.

## Problem Description

**Task.** Given a list of error-free reads, perform the task of Genome Assembly and return the circular genome from which they came.

**Dataset.** Each of 1618 lines of the input contains a single *read*, that is, a string over {A, C, G, T}. The reads are given to you in alphabetical order because their true order is hidden from you. Each read is 100 nucleotides long and contains no sequencing errors. **Note that you are not given the 100-mer composition of the genome, i.e., some 100-mers may be missing.**

**Output.** Output the assembled genome on a single line.

### Time Limits.

language	C	C++	Java	Python	C#	Go	Haskell	JavaScript	Ruby	Scala
time (sec)	3	3	4.5	15	4.5	4.5	6	15	15	9

**Memory Limit.** 512MB.

### Sample 1.

Input:

```
AAC
ACG
GAA
GTT
TCG
```

Output:

```
ACGTTCGA
```

Explanation:

In this sample, the circular genome is ACGTTCGA (the sample output), and the reads were all generated from the genome: **ACGTTCGA**, **ACGTTCGA**, **ACGTTCGA**, **ACGTTCGA**, and **ACGTTCGA**.

## What to Do

Construct an overlap graph: two reads are joined by a directed edge of weight equal to the length of the maximum overlap of these two reads. Then construct a Hamiltonian path in this graph in a greedy fashion: for each read select an out-going edge of maximum weight. Then read a string spelled by this path. To avoid computing overlaps between all pairs of reads, you may want to first compute a list of all pairs of reads that share a  $k$ -mer (for, say,  $k = 12$ ).

## Need Help?

Ask a question or see the questions asked by other learners at [this forum](#) thread.

## 2 Dataset Problem: Assembling the phi X174 Genome from Error-Prone Reads Using Overlap Graphs

In dataset problems, your solution is going to be tested against a *single* dataset as opposed to problems in the previous classes in this specialization. For this reason, there is no Constraints section in the problem description below. The sample section shows just a similar dataset. Your program is not going to be tested on this dataset.

### Problem Introduction

In this challenge, you will be given the task of performing Genome Assembly on a simulated error-prone sequencing dataset.

### Problem Description

**Task.** Given a list of error-prone reads, perform the task of Genome Assembly and return the circular genome from which they came.

**Dataset.** Each of 1618 lines of the input contains a single *read*. The reads are given to you in alphabetical order because their true order is hidden from you. Each *read* is 100 nucleotides long and contains a single sequencing error (i.e., one mismatch per read) in order to simulate the 1% error rate of Illumina sequencing machines. Note that you are not given the 100-mer composition of the genome (i.e., some 100-mers may be missing).

**Output.** Output the assembled genome on a single line.

#### Time Limits.

language	C	C++	Java	Python	C#	Go	Haskell	JavaScript	Ruby	Scala
time (sec)	3	3	4.5	15	4.5	4.5	6	15	15	9

**Memory Limit.** 512MB.

#### Sample 1.

Input:

```
AAC
ACG
GAA
GTT
TCG
```

Output:

```
ACGTTCGA
```

Explanation:

In this sample, the circular genome is ACGTTCGA (the sample output), and the reads were all generated from the genome: **ACGTTCGA**, **ACGTTCGA**, **ACGTTCGA**, **ACGTT**CGA, and **ACGTT**CGA. Note that we did not put mismatches in these reads because they are extremely short, so introducing the sequencing errors on such short reads would only make the solution difficult to see. On the real dataset, however, each read (of length 100) will have exactly one error.

### What to Do

You may want to declare two error-prone reads as overlapping if they have at most 5% error rate in their overlap (since error rate in each read is 1%, 5% should account for possible clustering of random errors.)

## Need Help?

Ask a question or see the questions asked by other learners at [this forum thread](#).

### 3 General Instructions and Recommendations on Solving Algorithmic Problems

Your main goal in an algorithmic problem is to implement a program that solves a given computational problem in just few seconds even on massive datasets. Your program should read a dataset from the standard input and write an answer to the standard output.

Below we provide general instructions and recommendations on solving such problems. Before reading them, go through readings and screencasts in the first module that show a step by step process of solving two algorithmic problems: [link](#).

#### 3.1 Reading the Problem Statement

You start by reading the problem statement that contains the description of a particular computational task as well as time and memory limits your solution should fit in, and one or two sample tests. In some problems your goal is just to implement carefully an algorithm covered in the lectures, while in some other problems you first need to come up with an algorithm yourself.

#### 3.2 Designing an Algorithm

If your goal is to design an algorithm yourself, one of the things it is important to realize is the expected running time of your algorithm. Usually, you can guess it from the problem statement (specifically, from the subsection called constraints) as follows. Modern computers perform roughly  $10^8$ – $10^9$  operations per second. So, if the maximum size of a dataset in the problem description is  $n = 10^5$ , then most probably an algorithm with quadratic running time is not going to fit into time limit (since for  $n = 10^5$ ,  $n^2 = 10^{10}$ ) while a solution with running time  $O(n \log n)$  will fit. However, an  $O(n^2)$  solution will fit if  $n$  is up to  $10^3 = 1000$ , and if  $n$  is at most 100, even  $O(n^3)$  solutions will fit. In some cases, the problem is so hard that we do not know a polynomial solution. But for  $n$  up to 18, a solution with  $O(2^n n^2)$  running time will probably fit into the time limit.

To design an algorithm with the expected running time, you will of course need to use the ideas covered in the lectures. Also, make sure to carefully go through sample tests in the problem description.

#### 3.3 Implementing Your Algorithm

When you have an algorithm in mind, you start implementing it. Currently, you can use the following programming languages to implement a solution to a problem: C, C++, C#, Haskell, Java, JavaScript, Python2, Python3, Ruby, Scala. For all problems, we will be providing starter solutions for C++, Java, and Python3. If you are going to use one of these programming languages, use these starter files. For other programming languages, you need to implement a solution from scratch.

#### 3.4 Compiling Your Program

For solving programming assignments, you can use any of the following programming languages: C, C++, C#, Haskell, Java, JavaScript, Python2, Python3, Ruby, and Scala. However, we will only be providing starter solution files for C++, Java, and Python3. The programming language of your submission is detected automatically, based on the extension of your submission.

We have reference solutions in C++, Java and Python3 which solve the problem correctly under the given restrictions, and in most cases spend at most 1/3 of the time limit and at most 1/2 of the memory limit. You can also use other languages, and we've estimated the time limit multipliers for them, however, we have no guarantee that a correct solution for a particular problem running under the given time and memory constraints exists in any of those other languages.

Your solution will be compiled as follows. We recommend that when testing your solution locally, you use the same compiler flags for compiling. This will increase the chances that your program behaves in the

same way on your machine and on the testing machine (note that a buggy program may behave differently when compiled by different compilers, or even by the same compiler with different flags).

- C (gcc 5.2.1). File extensions: `.c`. Flags:

```
gcc -pipe -O2 -std=c11 <filename> -lm
```

- C++ (g++ 5.2.1). File extensions: `.cc`, `.cpp`. Flags:

```
g++ -pipe -O2 -std=c++14 <filename> -lm
```

If your C/C++ compiler does not recognize `-std=c++14` flag, try replacing it with `-std=c++0x` flag or compiling without this flag at all (all starter solutions can be compiled without it). On Linux and MacOS, you most probably have the required compiler. On Windows, you may use your favorite compiler or install, e.g., `cygwin`.

- C# (mono 3.2.8). File extensions: `.cs`. Flags:

```
mcs
```

- Haskell (ghc 7.8.4). File extensions: `.hs`. Flags:

```
ghc -O
```

- Java (Open JDK 8). File extensions: `.java`. Flags:

```
javac -encoding UTF-8  
java -Xmx1024m
```

- JavaScript (Node v6.3.0). File extensions: `.js`. Flags:

```
nodejs
```

- Python 2 (CPython 2.7). File extensions: `.py2` or `.py` (a file ending in `.py` needs to have a first line which is a comment containing “python2”). No flags:

```
python2
```

- Python 3 (CPython 3.4). File extensions: `.py3` or `.py` (a file ending in `.py` needs to have a first line which is a comment containing “python3”). No flags:

```
python3
```

- Ruby (Ruby 2.1.5). File extensions: `.rb`.

```
ruby
```

- Scala (Scala 2.11.6). File extensions: `.scala`.

```
scalac
```



### 3.5 Testing Your Program

When your program is ready, you start testing it. It makes sense to start with small datasets — for example, sample tests provided in the problem description. Ensure that your program produces a correct result.

You then proceed to checking how long does it take your program to process a massive dataset. For this, it makes sense to implement your algorithm as a function like `solve(dataset)` and then implement an additional procedure `generate()` that produces a large dataset. For example, if an input to a problem is a sequence of integers of length  $1 \leq n \leq 10^5$ , then generate a sequence of length exactly  $10^5$ , pass it to your `solve()` function, and ensure that the program outputs the result quickly.

Also, check the boundary values. Ensure that your program processes correctly sequences of size  $n = 1, 2, 10^5$ . If a sequence of integers from 0 to, say,  $10^6$  is given as an input, check how your program behaves when it is given a sequence  $0, 0, \dots, 0$  or a sequence  $10^6, 10^6, \dots, 10^6$ . Check also on randomly generated data. For each such test check that you program produces a correct result (or at least a reasonably looking result).

In the end, we encourage you to stress test your program to make sure it passes in the system at the first attempt. See the readings and screencasts from the first week to learn about testing and stress testing: [link](#).

### 3.6 Submitting Your Program to the Grading System

When you are done with testing, you submit your program to the grading system. For this, you go the submission page, create a new submission, and upload a file with your program. The grading system then compiles your program (detecting the programming language based on your file extension, see Subsection 3.4) and runs it on a set of carefully constructed tests to check that your program always outputs a correct result and that it always fits into the given time and memory limits. The grading usually takes no more than a minute, but in rare cases when the servers are overloaded it might take longer. Please be patient. You can safely leave the page when your solution is uploaded.

As a result, you get a feedback message from the grading system. The feedback message that you will love to see is: **Good job!** This means that your program has passed all the tests. On the other hand, the three messages **Wrong answer**, **Time limit exceeded**, **Memory limit exceeded** notify you that your program failed due to one these three reasons. Note that the grader will not show you the actual test you program have failed on (though it does show you the test if your program have failed on one of the first few tests; this is done to help you to get the input/output format right).

### 3.7 Debugging and Stress Testing Your Program

If your program failed, you will need to debug it. Most probably, you didn't follow some of our suggestions from the section 3.5. See the readings and screencasts from the first week to learn about debugging your program: [link](#).

You are almost guaranteed to find a bug in your program using stress testing, because the way these programming assignments and tests for them are prepared follows the same process: small manual tests, tests for edge cases, tests for large numbers and integer overflow, big tests for time limit and memory limit checking, random test generation. Also, implementation of wrong solutions which we expect to see and stress testing against them to add tests specifically against those wrong solutions.

**Go ahead, and we hope you pass the assignment soon!**

## 4 Frequently Asked Questions

### 4.1 I submit the program, but nothing happens. Why?

You need to create submission and upload the file with your solution in one of the programming languages C, C++, Java, or Python (see Subsections 3.3 and 3.4). Make sure that after uploading the file with your solution you press on the blue “Submit” button in the bottom. After that, the grading starts, and the submission being graded is enclosed in an orange rectangle. After the testing is finished, the rectangle disappears, and the results of the testing of all problems is shown to you.

### 4.2 I submit the solution only for one problem, but all the problems in the assignment are graded. Why?

Each time you submit any solution, the last uploaded solution for each problem is tested. Don’t worry: this doesn’t affect your score even if the submissions for the other problems are wrong. As soon as you pass the sufficient number of problems in the assignment (see in the pdf with instructions), you pass the assignment. After that, you can improve your result if you successfully pass more problems from the assignment. We recommend working on one problem at a time, checking whether your solution for any given problem passes in the system as soon as you are confident in it. However, it is better to test it first, please refer to the reading about stress testing: [link](#).

### 4.3 What are the possible grading outcomes, and how to read them?

Your solution may either pass or not. To pass, it must work without crashing and return the correct answers on all the test cases we prepared for you, and do so under the time limit and memory limit constraints specified in the problem statement. If your solution passes, you get the corresponding feedback "Good job!" and get a point for the problem. If your solution fails, it can be because it crashes, returns wrong answer, works for too long or uses too much memory for some test case. The feedback will contain the number of the test case on which your solution fails and the total number of test cases in the system. The tests for the problem are numbered from 1 to the total number of test cases for the problem, and the program is always tested on all the tests in the order from the test number 1 to the test with the biggest number.

Here are the possible outcomes:

**Good job! Hurrah!** Your solution passed, and you get a point!

**Wrong answer.** Your solution has output incorrect answer for some test case. If it is a sample test case from the problem statement, or if you are solving Programming Assignment 1, you will also see the input data, the output of your program and the correct answer. Otherwise, you won’t know the input, the output, and the correct answer. Check that you consider all the cases correctly, avoid integer overflow, output the required white space, output the floating point numbers with the required precision, don’t output anything in addition to what you are asked to output in the output specification of the problem statement. See this reading on testing: [link](#).

**Time limit exceeded.** Your solution worked longer than the allowed time limit for some test case. If it is a sample test case from the problem statement, or if you are solving Programming Assignment 1, you will also see the input data and the correct answer. Otherwise, you won’t know the input and the correct answer. Check again that your algorithm has good enough running time estimate. Test your program locally on the test of maximum size allowed by the problem statement and see how long it works. Check that your program doesn’t wait for some input from the user which makes it to wait forever. See this reading on testing: [link](#).

**Memory limit exceeded.** Your solution used more than the allowed memory limit for some test case. If it is a sample test case from the problem statement, or if you are solving Programming Assignment 1,

you will also see the input data and the correct answer. Otherwise, you won't know the input and the correct answer. Estimate the amount of memory that your program is going to use in the worst case and check that it is less than the memory limit. Check that you don't create too large arrays or data structures. Check that you don't create large arrays or lists or vectors consisting of empty arrays or empty strings, since those in some cases still eat up memory. Test your program locally on the test of maximum size allowed by the problem statement and look at its memory consumption in the system.

**Cannot check answer. Perhaps output format is wrong.** This happens when you output something completely different than expected. For example, you are required to output word "Yes" or "No", but you output number 1 or 0, or vice versa. Or your program has empty output. Or your program outputs not only the correct answer, but also some additional information (this is not allowed, so please follow exactly the output format specified in the problem statement). Maybe your program doesn't output anything, because it crashes.

**Unknown signal 6 (or 7, or 8, or 11, or some other).** This happens when your program crashes. It can be because of division by zero, accessing memory outside of the array bounds, using uninitialized variables, too deep recursion that triggers stack overflow, sorting with contradictory comparator, removing elements from an empty data structure, trying to allocate too much memory, and many other reasons. Look at your code and think about all those possibilities. Make sure that you use the same compilers and the same compiler options as we do. Try different testing techniques from this reading: [link](#).

**Internal error: exception...** Most probably, you submitted a compiled program instead of a source code.

**Grading failed.** Something very wrong happened with the system. Contact Coursera for help or write in the forums to let us know.

## 4.4 How to understand why my program fails and to fix it?

If your program works incorrectly, it gets a feedback from the grader. For the Programming Assignment 1, when your solution fails, you will see the input data, the correct answer and the output of your program in case it didn't crash, finished under the time limit and memory limit constraints. If the program crashed, worked too long or used too much memory, the system stops it, so you won't see the output of your program or will see just part of the whole output. We show you all this information so that you get used to the algorithmic problems in general and get some experience debugging your programs while knowing exactly on which tests they fail.

However, in the following Programming Assignments throughout the Specialization you will only get so much information for the test cases from the problem statement. For the next tests you will only get the result: passed, time limit exceeded, memory limit exceeded, wrong answer, wrong output format or some form of crash. We hide the test cases, because it is crucial for you to learn to test and fix your program even without knowing exactly the test on which it fails. In the real life, often there will be no or only partial information about the failure of your program or service. You will need to find the failing test case yourself. Stress testing is one powerful technique that allows you to do that. You should apply it after using the other testing techniques covered in this reading.

## 4.5 Why do you hide the test on which my program fails?

Often beginner programmers think by default that their programs work. Experienced programmers know, however, that their programs almost never work initially. Everyone who wants to become a better programmer needs to go through this realization.

When you are sure that your program works by default, you just throw a few random test cases against it, and if the answers look reasonable, you consider your work done. However, mostly this is not enough. To

make one's programs work, one must test them really well. Sometimes, the programs still don't work although you tried really hard to test them, and you need to be both skilled and creative to fix your bugs. Solutions to algorithmic problems are one of the hardest to implement correctly. That's why in this Specialization you will gain this important experience which will be invaluable in the future when you write programs which you really need to get right.

It is crucial for you to learn to test and fix your programs yourself. In the real life, often there will be no or only partial information about the failure of your program or service. Still, you will have to reproduce the failure to fix it (or just guess what it is, but that's rare, and you will still need to reproduce the failure to make sure you have really fixed it). When you solve algorithmic problems, it is very frequent to make subtle mistakes. That's why you should apply the testing techniques described in this reading to find the failing test case and fix your program.

#### **4.6 My solution does not pass the tests? May I post it in the forum and ask for a help?**

No, please do not post any solutions in the forum or anywhere on the web, even if a solution does not pass the tests (as in this case you are still revealing parts of a correct solution). Recall the third item of the Coursera Honor Code: "I will not make solutions to homework, quizzes, exams, projects, and other assignments available to anyone else (except to the extent an assignment explicitly permits sharing solutions). This includes both solutions written by me, as well as any solutions provided by the course staff or others" ([link](#)).

#### **4.7 My implementation always fails in the grader, though I already tested and stress tested it a lot. Would not it be better if you give me a solution to this problem or at least the test cases that you use? I will then be able to fix my code and will learn how to avoid making mistakes. Otherwise, I do not feel that I learn anything from solving this problem. I am just stuck.**

First of all, you always learn from your mistakes.

The process of trying to invent new test cases that might fail your program and proving them wrong is often enlightening. This thinking about the invariants which you expect your loops, ifs, etc. to keep and proving them wrong (or right) makes you understand what happens inside your program and in the general algorithm you're studying much more.

Also, it is important to be able to find a bug in your implementation without knowing a test case and without having a reference solution. Assume that you designed an application and an annoyed user reports that it crashed. Most probably, the user will not tell you the exact sequence of operations that led to a crash. Moreover, there will be no reference application. Hence, once again, it is important to be able to locate a bug in your implementation yourself, without a magic oracle giving you either a test case that your program fails or a reference solution. We encourage you to use programming assignments in this class as a way of practicing this important skill.

If you have already tested a lot (considered all corner cases that you can imagine, constructed a set of manual test cases, applied stress testing), but your program still fails and you are stuck, try to ask for help on the forum. We encourage you to do this by first explaining what kind of corner cases you have already considered (it may happen that when writing such a post you will realize that you missed some corner cases!) and only then asking other learners to give you more ideas for tests cases.