

# Programming Assignment 1: Assembling phi X174 Using Overlap Graphs

Revision: November 15, 2018

## Introduction

Welcome to the first programming assignment of the [Algorithms and Data Structures Capstone](#)! 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</b>	<b>Dataset Problem: Assembling the phi X174 Genome from Error-Free Reads Using Overlap Graphs</b>	<b>2</b>
<b>2</b>	<b>Dataset Problem: Assembling the phi X174 Genome from Error-Prone Reads Using Overlap Graphs</b>	<b>3</b>
<b>3</b>	<b>Solving a Programming Challenge in Five Easy Steps</b>	<b>4</b>
3.1	Reading Problem Statement . . . . .	4
3.2	Designing an Algorithm . . . . .	4
3.3	Implementing an Algorithm . . . . .	4
3.4	Testing and Debugging . . . . .	4
3.5	Submitting to the Grading System . . . . .	5
<b>4</b>	<b>Appendix: Compiler Flags</b>	<b>5</b>

# 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	JavaScript	Scala
time (sec)	3	3	4.5	15	15	9

**Memory Limit.** 512MB.

### Sample 1.

Input:

```
AAC
ACG
GAA
GTT
TCG
```

Output:

```
ACGTTCGA
```

In this sample, the circular genome is ACGTTCGA (the sample output), and the reads were all generated from the genome: **ACGTTCGA**, **ACGTT**CGA, **ACGTT**CG**A**, **ACGTT**CGA, and **ACGT**TCGA.

## 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,  $k = 12$ ).

## 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	JavaScript	Scala
time (sec)	3	3	4.5	15	15	9

**Memory Limit.** 512MB.

**Sample 1.**

Input:

```
AAC
ACG
GAA
GTT
TCG
```

Output:

```
ACGTTCGA
```

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

## 3 Solving a Programming Challenge in Five Easy Steps

### 3.1 Reading Problem Statement

Start by reading the problem statement that contains the description of a computational task, time and memory limits, and a few sample tests. Make sure you understand how an output matches an input in each sample case.

If time and memory limits are not specified explicitly in the problem statement, the following default values are used.

#### Time Limits.

language	C	C++	Java	Python	JavaScript	Scala
time (sec)	1	1	1.5	5	5	3

**Memory limit:** 512 Mb.

### 3.2 Designing an Algorithm

After designing an algorithm, prove that it is correct and try to estimate its expected running time on the most complex inputs specified in the constraints section. If your laptop performs roughly  $10^8$ – $10^9$  operations per second, and the maximum size of a dataset in the problem description is  $n = 10^5$ , then an algorithm with quadratic running time is unlikely to fit into the time limit (since  $n^2 = 10^{10}$ ), while a solution with running time  $O(n \log n)$  will. However, an  $O(n^2)$  solution will fit if  $n = 1\,000$ , and if  $n = 100$ , even an  $O(n^3)$  solution will fit. Although polynomial algorithms remain unknown for some hard problems in this book, a solution with  $O(2^n n^2)$  running time will probably fit into the time limit as long as  $n$  is smaller than 20.

### 3.3 Implementing an Algorithm

Start implementing your algorithm in one of the following programming languages supported by our automated grading system: C, C++, Haskell, Java, JavaScript, Python2, Python3, or Scala. For all problems, we provide starter solutions for C++, Java, and Python3. For other programming languages, you need to implement a solution from scratch. The grading system detects the programming language of your submission automatically, based on the extension of the submission file.

We have reference solutions in C++, Java, and Python3 (that we don't share with you) which solve the problem correctly under the given constraints, and 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.

In the Appendix, we list compiler versions and flags used by the grading system. We recommend using the same compiler flags when you test your solution locally. 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).

### 3.4 Testing and Debugging

Submitting your implementation to the grading system without testing it first is a bad idea! Start with small datasets and make sure that your program produces correct results on all sample datasets. Then proceed to checking how long it takes to process a large dataset. To estimate the running time, 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  $10^5$ , pass it to your `solve()` function, and ensure that the program outputs the result quickly.

Check the boundary values to ensure that your program processes correctly both short sequences (e.g., with 2 elements) and long sequences (e.g., with  $10^5$  elements). If a sequence of integers from 0 to, let's 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  $10^6, 10^6, \dots, 10^6$ . Afterwards, check it also on randomly generated data. Check degenerate cases like an empty set, three points on a single line, a tree which consists of a single path of nodes, etc.

After it appears that your program works on all these tests, proceed to stress testing. Implement a slow, but simple and correct algorithm and check that two programs produce the same result (note however that this is not applicable to problems where the output is not unique). Generate random test cases as well as biased tests cases such as those with only small numbers or a small range of large numbers, strings containing a single letter “a” or only two different letters (as opposed to strings composed of all possible Latin letters), and so on. Think about other possible tests which could be peculiar in some sense. For example, if you are generating graphs, try generating trees, disconnected graphs, complete graphs, bipartite graphs, etc. If you generate trees, try generating paths, binary trees, stars, etc. If you are generating integers, try generating both prime and composite numbers.

### 3.5 Submitting to the Grading System

When you are done with testing, submit your program to the grading system! Go to the submission page, create a new submission, and upload a file with your program (make sure to upload a source file rather than an executable). The grading system then compiles your program and runs it on a set of carefully constructed tests to check that it outputs a correct result for all tests and that it fits into the time and memory limits. The grading usually takes less 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. You want to see the “Good job!” message indicating that your program passed all the tests. The messages “Wrong answer”, “Time limit exceeded”, “Memory limit exceeded” notify you that your program failed due to one of these reasons. If you program fails on one of the first two test cases, the grader will report this to you and will show you the test case and the output of your program. This is done to help you to get the input/output format right. In all other cases, the grader will *not* show you the test case where your program fails.

## 4 Appendix: Compiler 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`.

**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
```

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

```
scalac
```