

Industrial Training Report

Muhammad Farhan bin Mohd Tahir

3 January 2017

Contents


1. Company
2. Project's Background
3. Project's Aim and Objectives
4. Term's Definition
5. Programming Language and Tools Used
6. Program's FLOWchart
7. Method Implemented
8. Achievements
9. Demonstration
10. Conclusion

Company

Company Background

- ▶ Novocraft Technologies SDN BHD



- ▶ Culmination of over 8 years of experience in software development and performance tuning. 

Company

Company Background

► Mission

- To be a globally recognized provider of innovative and accurate bioinformatics tools and services for discovery, primarily in the next-generation genomics (NGS) space.

► Vision

- Intelligent software that automates the complex data processing pipelines required to conduct genomics analysis.
- Ensuring better predictive power from NGS data in personalized genomics & therapies.
- Intelligent software that automates the complex data processing pipelines required to conduct genomics analysis.

Company

Company Organization's Chart

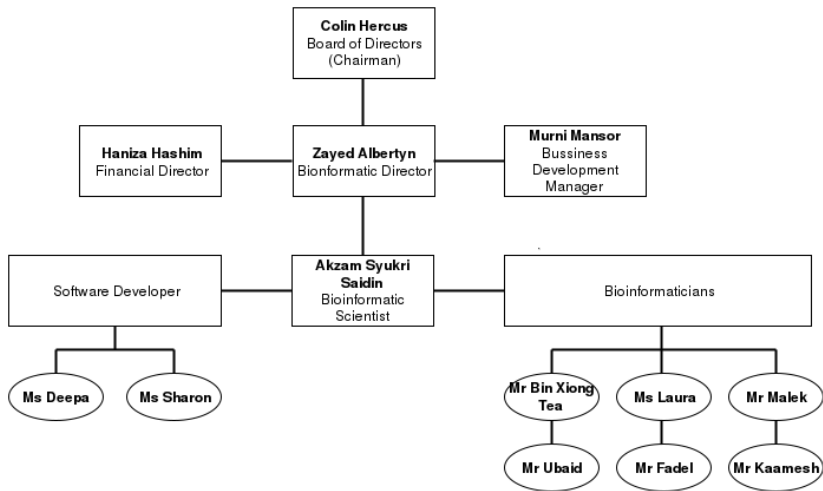


Figure 1:

Company

Company's Services and Product

- ▶ Services:
 - ▶ Bioinformatics Consultancy
 - ▶ Pipeline and Software Development Services
 - ▶ Bioinformatics Contract Research Services
- ▶ Products:
 - ▶ NovoSort
 - ▶ NovoAlignCS
 - ▶ NovoAlign
 - ▶ NovoWorx (Latest released tools)

Project's Background

- ▶ Problem Statement:

- ▶ Illumina Sequencing technology is a tool that used by the biologist/biological researcher to extract the genome sequence from any biological sample. However, the output from Illumina sequencing may contain sequences that are not originally found in the sample, but chemically synthesized, such as, the Adapter Sequence which is used to link the ends of two other DNA molecules. In order to get a clean genomic sequence, the adapter sequence needs to be identified and removed. A program has previously been developed in Novocraft to overcome this issue, Nevertheless, there are room for improvement in terms of the program's result accuracy and executional speed.

- ▶ Project's Assigned:

- ▶ To develop a program to find an adapter sequence from two paired-end FASTQ file (input file from user).

- ▶ Program's Name:

- ▶ PEAdapterFinder program (**Pair Ends Adapter Finder**)

- ▶ Project's Initiator:

Project's Aim and Objectives

- ▶ Project's Aim:
 - ▶ to develop program to find an adapter sequence from two paired-end FASTQ file (input file from user).
- ▶ Project's Objective:
 - ▶ To implement the code that can make the program's result more accurate than before.
 - ▶ To implement the code that can increase the speed of program's execution time.

Term's Definition

► FASTQ File

- Line 1: Begins with a '@' character and is followed by a sequence identifier and an optional description (like a FASTA title line).
- Line 2: Raw DNA sequence.
- Line 3: Begins with a '+' character and is optionally followed by the same sequence identifier (and any description) again.
- Line 4: Encodes the quality values for the sequence in Line 2, and must contain the same number of symbols as letters in the sequence.

Term's Definition

FASTQ File

► 4-Line FASTQ file example

```
@HWI-EAS179_0001:5:1:7:1350#0/1
CGTGATTTTGTATGTAAGTTTTTCTTTGAAGTNTTGTTCAGTGAAACAANAATTGNNNNNNNNNNNNNNNN
+HWI-EAS179_0001:5:1:7:1350#0/1
HHDH@HHHHHHHHHHFF>FFEEEEFFFFFEE@E>##A?AFA:GC15C#####
@HWI-EAS179_0001:5:1:7:72#0/1
GGTGATCCTCAAAGTCAAACCTGCAGTAGTTCGGAANNNGGTCTTATTCTGTCAATNTCATANNNNNNNNNNNNNN
+HWI-EAS179_0001:5:1:7:72#0/1
HHHEEDGHHHCAGE4HHHHHHHAGE8:HEDFBF3<##A7.58=,1BECDF0.B#####
@HWI-EAS179_0001:5:1:7:746#0/1
CGTGATCTTATATCATGGAGTTGGGCGAGTCATACNNTTCTCGTGATTTATTNATTTNNNNNNNNNNNNNNNN
+HWI-EAS179_0001:5:1:7:746#0/1
HHGHHHHGHHCAD:5DGGADGGGGFAFFFC>CC9##=<2CCCC<D>7D<<5@#####
@HWI-EAS179_0001:5:1:7:1657#0/1
CCTGATACGCAACAATTTAATGCTAGAAATTTTCNNTAATGGTCCTTGATCGTNNNNNNNNNNNNNNNNNN
+HWI-EAS179_0001:5:1:7:1657#0/1
HHDEHHFHHHHHHH5GGG<DFFFF=D>:ACACCC##@6:>@C:FFCF9@?#>;A#####
@HWI-EAS179_0001:5:1:7:1002#0/1
CGTGATGATCAAACAATTACAGATGACTCTTTGGANNTCCATACAGTCCCTTNCGAATNNNNNNNNNNNNNNNN
+HWI-EAS179_0001:5:1:7:1002#0/1
CGTTATHBHF&G;ADDFDGGGEE.G8GD=9:C5##<=<9><CFF?FC7?#####
@HWI-EAS179_0001:5:1:7:14#0/1
AATTATTGGACTTTGTGGTGAATTTATAACAAGTCNNCCATTTTCTGTAGGNTATTNNNNNNNNNNNNNNNN
+HWI-EAS179_0001:5:1:7:14#0/1
GFGAGA>DDDDGGDGD?FFF.7?<??7478;155+5##=0@=94+CEC6A<C@#A3?#####
@HWI-EAS179_0001:5:1:7:1350#0/1
CGTGATTTTGTATGTAAGTTTTTCTTTGAAGTNTTGTTCAGTGAAACAANAATTGNNNNNNNNNNNNNNNN
+HWI-EAS179_0001:5:1:7:1350#0/1
HHDH@HHHHHHHHHHFF>FFEEEEFFFFFEE@E>##A?AFA:GC15C#####
@HWI-EAS179_0001:5:1:7:72#0/1
GGTGATCCTCAAAGTCAAACCTGCAGTAGTTCGGAANNNGGTCTTATTCTGTCAATNTCATANNNNNNNNNNNNNN
+HWI-EAS179_0001:5:1:7:72#0/1
HHHEEDGHHHCAGE4HHHHHHHAGE8:HEDFBF3<##A7.58=,1BECDF0.B#####
@HWI-EAS179_0001:5:1:7:746#0/1
CGTGATCTTATATCATGGAGTTGGGCGAGTCATACNNTTCTCGTGATTTATTNATTTNNNNNNNNNNNNNNNN
+HWI-EAS179_0001:5:1:7:746#0/1
HHGHHHHGHHCAD:5DGGADGGGGFAFFFC>CC9##=<2CCCC<D>7D<<5@#####
```

Line 1:
identifier

Line 2: Raw
DNA sequence

Line 3:
separator '+'

Line 4:
Qualities Value

Term's Definition

- ▶ **Adapter Sequence**

- ▶ Adapter Sequence is a short, chemically synthesized, double stranded DNA molecule. which is used to link the ends of two other DNA molecules.

- ▶ **Needleman-Wunsch Algorithm (NW)**

- ▶ NW is an algorithm used in bioinformatics to align two DNA sequences.

- ▶ **Consensus Sequence (CS)**

- ▶ CS is used to find the highest frequency of DNA base(A/T/G/C) that need to be assign at certain place.

Programming Language and Tools Used

- ▶ Programming Language

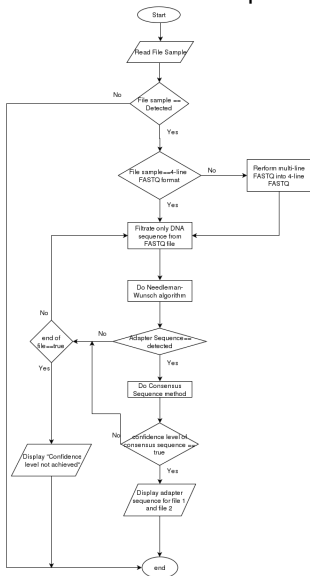
- ▶ C++ language

- ▶ Tools

Tools	Description
NetBeans IDE -version 8.1	Provide framework to develop C++ program.
GPREFTOOLS (PPOF) -version 2.0	Provide tools that can help which line in the code that consume the most time during execution of the program.
GNU gdb -version 7.4.1-debian	Tools that can help in trace if there any error during execution of program such as <i>Segmentation Fault</i>

Program's Flowchart

- Flowchart of PEAdapterFinder in finding adapter sequence.



Method Implemented

Improved Code

► Accuracy

Original Code	Function	Improved code
There is no code to convert multi-line FASTQ file into 4-line FASTQ file. So when the input file is multi-line FASTQ file, the output may incorrect.	Convert/Reform multi-line FASTQ file into 4-line FASTQ file	Code to convert from multi-line FASTQ file into 4-line FASTQ file is developed. So then, when the program want to filtrate only DNA sequence from FASTQ file, it just choose the string from line 2, 6, 10, ..., (NBEFORE+4)
Original code compare each character in each line whether the line only had bases (A/T/G/C/N) to be assign as DNA sequence. But if in case when the input file is multi-line FASTQ file, for example, there is two line which represent only ONE DNA sequence, the program will assume it is two different	Filtration only DNA sequence from 4 line FASTQ format	Improved code use line count to filtrate the only DNA sequence line from 4-line FASTQ file. This may be accurate than original code because the program will assign the string at line 2, 6, 10, ..., (NBEFORE+4) as DNA sequence which it is confirmed as DNA sequence's line.

Method Implemented

Improved Code

► Speed

Original Code	Function	Improved Code
<p>Original code comparing each character in each line for both file to assign whether the line is DNA sequence or not. This may increase in program's executional time since the program need to do comparison two times for each line (because have 2 file input):</p> <p>First comparison: Each character in each line in file 1.</p> <p>Second comparison: Each character in each line in file 2.</p>	<p>Filtration only DNA sequence from 4-line FASTQ file</p>	<p>Improved code just choose only line 2, 6, 10, ..., ($N_{\text{BEFORE}} + 4$) as DNA sequence at both file. <i>(If the input file is multi-line FASTQ file, the program will reformat it into 4-line FASTQ file. Then the program use above method in order to filtrate DNA sequence from FASTQ file)</i></p> <p>This may reduce program's execution time since the program do not need to compare each character for each string for both file.</p>

Method Implemented

Improved Code

► Speed

Original code implementing matrices and switch statement in order to determine the match/mismatch between 2 bases from file 1 and file 2 respectively. Since matrices need the usage of an array, this may increase the program's executional time.	Determination of match or mismatch between two bases	Improved code implementing if-else statement rather than switch statement. This code just store the result of match/mismatch in a single variable rather than use array and matrices. This may reduce program's executional time.
In original code, there is the function (user-define function) that code by the developer (Rayan) to find the maximum score between 3 score. User-define function may increase the program's executional time.	Finding maximum score to be filled in score-matrix table	In improved code, built-in function are implemented rather than using user-define function to find which score is higher between 3 score. Built-in function: <i>Max(max(SCORE1, SCORE2), SCORE3);</i> This may reduce program's

Achievements

Result's Accuracy

► File Sample

- Four-line FASTQ file: fourLine1.fastq & fourLine2.fastq
- Multi-line FASTQ file: multiLine1.fastq & multiLine2.fastq

Characteristics	Original Code	Improved Code
4-line FASTQ file Expected Output: Adapter 1: ATTACGAAATAATCGGATTC Adapter2: ATTACGAAATAATCGGATTC	Output: Adapter1: ATTACGAAATAATCGGATTC Adapter2: ATTACGAAATAATCGGATTC • Percentage of Accuracy: 100 %	Output: Adapter1: ATTACGAAATAATCGGATTC Adapter2: ATTACGAAATAATCGGATTC • Percentage of Accuracy: 100%
Multi-line FASTQ file Expected output: Adapter Sequence from file 1: ATTACGAAATAATCGGATTC	Output: Adapter1: <u>AGCCTAATCGGGATCTCATC</u> Adapter2: <u>CAGTAATCTTACTGTGATCC</u> *Base with line means the base is different from expected base	Output: Adapter1: ATTACGAAATAATCGGATTC Adapter2: ATTACGAAATAATCGGATTC • Percentage of

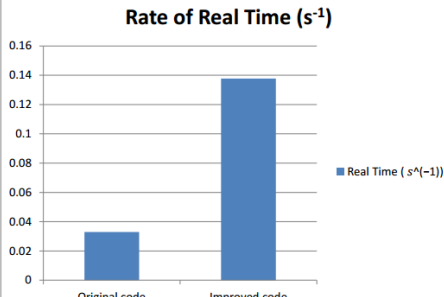
Achievements

Program's Speed

- File Sample: sample1.fastq & sample2.fastq

Characteristics	Original Code	Improved Code
Real time of program's executional (s)	30.315	7.261
Rate of Real Time of program's executional (s^{-1})	0.032987	0.1377221

Graph of code pattern vs rate of real time of program's executional



Demonstration

Conclusion

- ▶ Aim of this project had been achieved
- ▶ Objective of this project had been achieved

Conclusion



Figure 2:

Conclusion



Figure 3:

Conclusion



Figure 4: