

March 2016

An Introduction to Deep-Sequencing Data Analysis

Exercise #1

Dena Leshkowitz and Ester Feldmesser

Introduction

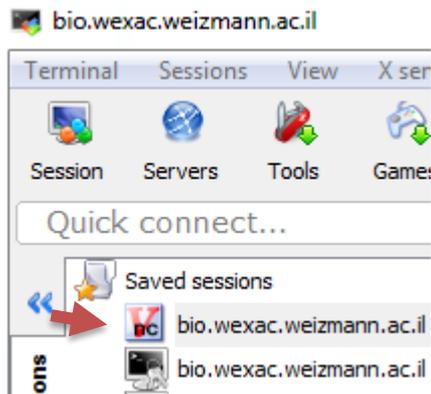
In this workshop we will learn how to evaluate the sequence quality and how to map the reads to a reference genome. The data set in this workshop is a collection of RNA-Seq data from mRNAs extracted from acute lymphoblastic leukemia (ALL) precursor B cell line.

We will use Chipster (Kallio et al. BMC Genomics 2011, 12:507), a user-friendly analysis software for high-throughput data. Its intuitive graphical user interface enables biologists to access a powerful collection of data analysis and integration tools, and to visualize data interactively. Users can collaborate by sharing analysis sessions and workflows. Chipster is an open source, and the server installation package is available for free.

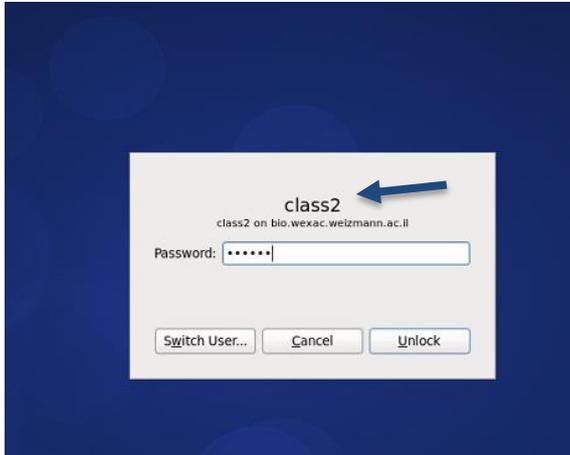
Instructions

1. Accessing the data

- a. Find the icon MobaXterm on your desktop and click to open.
- b. You will see on the left under “Saved sessions” the following:



Click on the VNC session (red arrow), it can take a while until you are prompted to insert the required password (ask the instructor), this depends on your user ID. The userID is class (see blue arrow below) and a certain number, insert the password and press enter in your keyboard. The user ID and password that you are now using is needed also to open Chipster.



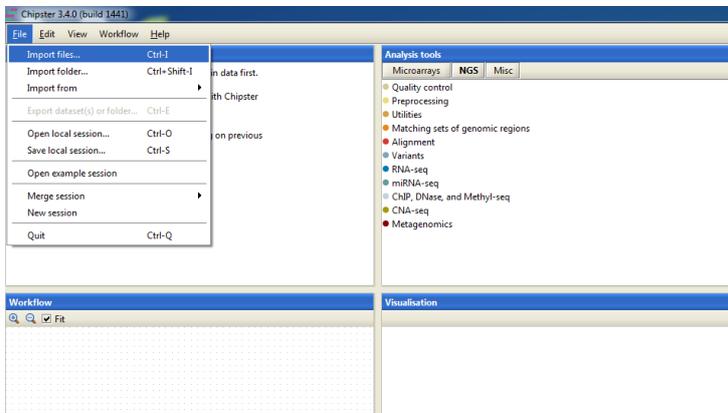
- c. Open the Chipster application found on the VNC session.



Enter your userID and password (the same as before).

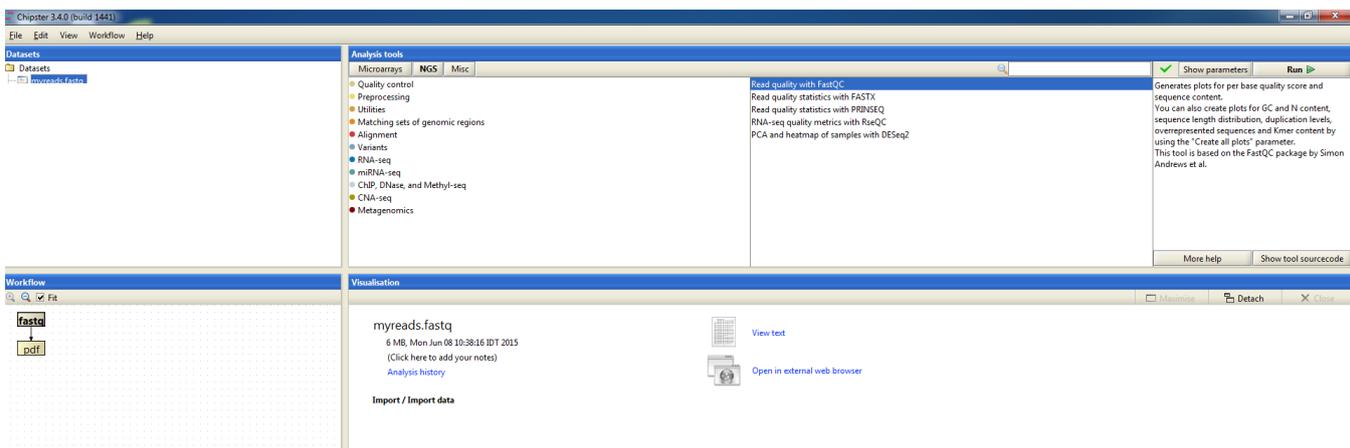
2. Running FASTQC on Chipster

- a. Import the fastq file (File->importFile->select the folder files_for_course and in it myreads.fastq file)



- b. Select the imported file from the “Workflow” left panel by clicking on it.

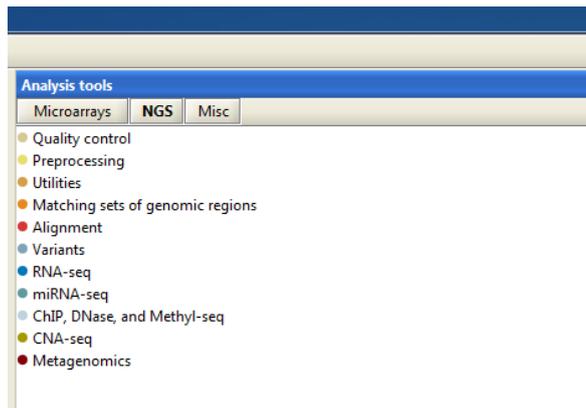
- c. In the “Visualization” pane (bottom right), click on View text:



- d. Look at the quality value of the first base from the first sequence. Convert the character into a numeric value based on the Supplementary converter below. Our sequences comply with the Sanger convention.

Question 1: What is the probability of an error in the first base?

- e. To create a QC report for myreads.fastq with the fastqc program, select the imported file from the left panel by clicking on it.
- f. In the “Analysis” tools pane, click on **NGS** tab, and then on “Quality control” and then on “Read quality with FastQC”.



- g. Click on the **“Run”** button at the right part. Wait for the results, an html result will appear in the left window called Workflow. Select the html result and from the visualization pane select **“Open in an external browser”**”

Question 2:

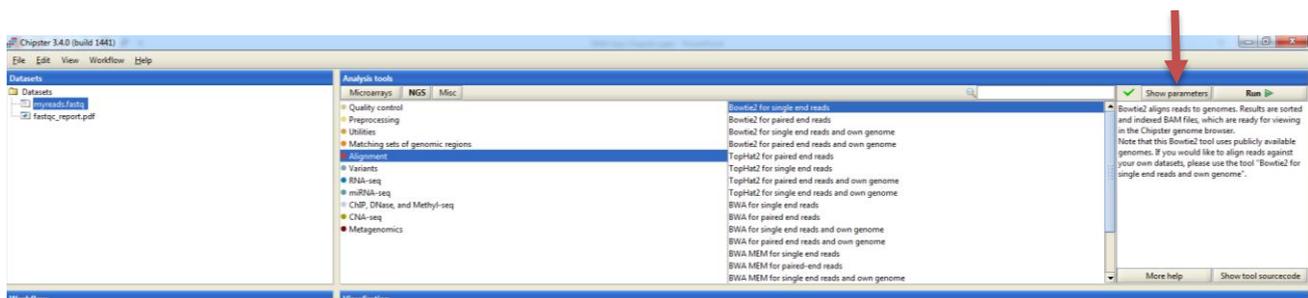
- How many sequences does the fastq file contain?
- Is the base quality the same for all the cycles?
- Do all the cycles have an equal base content?

The sequences are from a RNA-Seq experiment. During the course we will discuss the reason for the unequal base content in the beginning of the sequences.

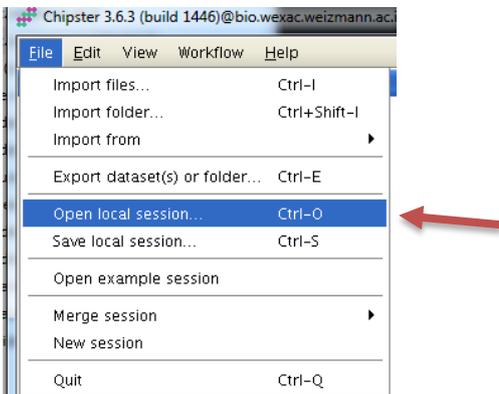
3. Running Bowtie and understanding the output

We are going to run bowtie:

- Select the imported fastq file from the left panel by clicking on it.
- In the Analysis tools, click on **NGS tab**, then on **“Alignment”** and then on **“Bowtie2 for single end reads”**”.
- Click on **“Show parameters”** on the right (Red arrow below). Be sure to use the most recent version of the human genome: the GRCh38 version.

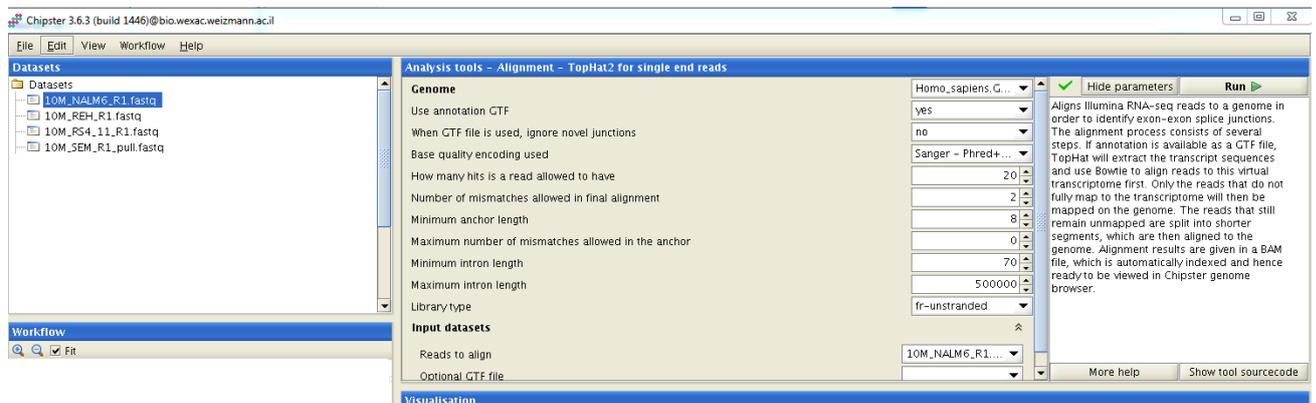


Open the local session (see below; File->Open local session) named Exercise-fastq-session.zip found in folder files_for_course.



Select the first fastq file from the “Database” pane and align it with “Tophat for single end reads” against the Homo_sapiens.GRCh38.81 genome, the other parameters should be as shown below. We will use the option “use annotation GTF” yes. Tomorrow we will discuss what this option means.

Repeat this for all 4 fastq files. You should see in the bottom of the Chipster window that all 4 jobs are running. This will take more than an hour, so you can close the VNC using the disconnect button in top of the mobaXterm screen and close the computer, the jobs will still run.



Good job. You completed this exercise.

M alignment match

I insertion to the reference

D deletion from the reference

N skipped region from the reference

S soft clipping (clipped sequences present in SEQ)

H hard clipping (clipped sequences NOT present in SEQ)

P padding (silent deletion from padded reference)

= sequence match

X sequence mismatch

THE END.