# Galaxy for virologist training Exercise 2: Quality control and trimming

Despite the improvement of sequencing methods, there is no error-free technique. A correct measuring of the sequencing quality is essential for identifying problems in the sequencing, thus, this must be the first step in every sequencing analysis. Once the quality control is finished, it's important to remove those low quality reads, or short reads, for which a trimming step is mandatory. After the trimming step it is recommended to perform a new quality control step to be sure that trimming worked.

## 1. Illumina Quality control and trimming

| Title | Pre-processing |
|---|---|
| **Training dataset:** | PRJEB43037 - In August 2020, an outbreak of West Nile Virus affected 71 people with meningoencephalitis in Andalusia and 6 more cases in Extremadura (south-west of Spain), causing a total of eight deaths. The virus belonged to the lineage 1 and was relatively similar to previous outbreaks occurred in the Mediterranean region. Here, we present a detailed analysis of the outbreak, including an extensive phylogenetic study. This is one of the outbreak samples. |
| **Questions:** | • How do I check whether my Illumina data was correctly sequenced?<br>• How can I improve the quality of my data? |
| **Objectives**: | • Perform a quality control in raw Illumina reads<br>• Perform a quality trimming in raw Illumina reads<br>• Perform a quality control in trimmed Illumina reads |
| **Estimated time**: | 25 min |

### 1.1. Quality control

#### 1.1.1. Upload data

To run the quality control over the samples, follow these steps:

1. Create a new history, as we explained yesterday named **Illumina preprocessing**
2. Upload data as seen yesterday, copy and paste the following URLs:

```
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/ERR531/002/ERR5310322/ERR5310322_1.fastq.gz
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/ERR531/002/ERR5310322/ERR5310322_2.fastq.gz
```

3. Add some tags to the files. *It is mandatory that the tag starts with # to be propagated to the processes*.
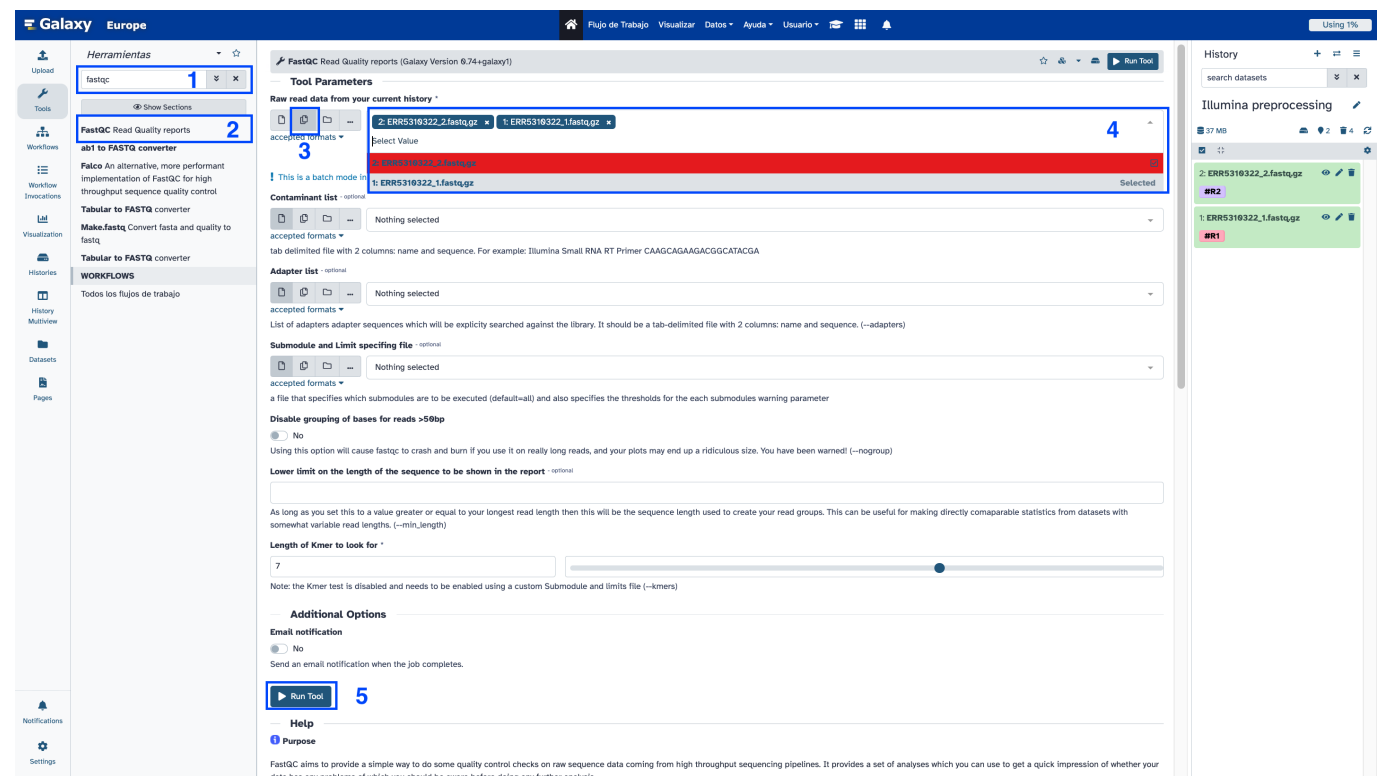


## 1.1.2. Run FastQC

1. Search for the **fastqc** tool
2. Select **FastQC Read Quality reports** and set the following parameters:
3. Select multiple file data set in Raw read data from your current history
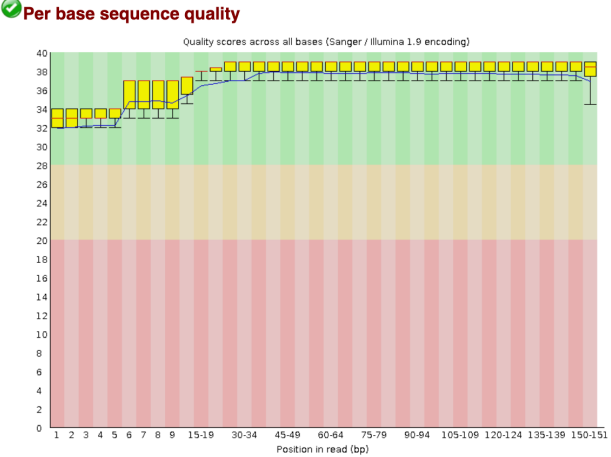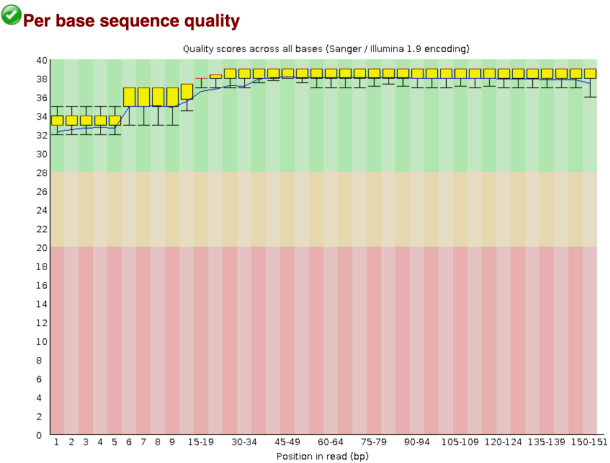4. Select the two datasets

5. Then go down and select **Run tool**

Galaxy Europe — FastQC Read Quality reports (Galaxy Version 0.74+galaxy1)

To see the results we are going to open the jobs with **Web page** in their name for both data 1 and data 2.

**Basic Statistics**

| Measure | Value |
| --- | --- |
| Filename | ERR5310322_1_fastq_gz.gz |
| File type | Conventional base calls |
| Encoding | Sanger / Illumina 1.9 |
| Total Sequences | 265989 |
| Sequences flagged as poor quality | 0 |
| Sequence length | 35-151 |
| %GC | 51 |

**Basic Statistics**

| Measure | Value |
| --- | --- |
| Filename | ERR5310322_2_fastq_gz.gz |
| File type | Conventional base calls |
| Encoding | Sanger / Illumina 1.9 |
| Total Sequences | 265989 |
| Sequences flagged as poor quality | 0 |
| Sequence length | 35-151 |
| %GC | 51 |

**Per base sequence quality**

Quality scores across all bases (Sanger / Illumina 1.9 encoding)

**Per base sequence quality**

Quality scores across all bases (Sanger / Illumina 1.9 encoding)

Here, you can see the number of reads in each file, the maximum and minimum length of all reads in the sample, and the quality plots for both R1 and R2. They look quite good, but we are going to run trimming over the samples.

▶ How many reads do the samples have?

**First question**

▶ How do I check whether my Illumina data was correctly sequenced?

## 1.2. Trimming

Once we have performed the quality control, we have to perform the quality and read length trimming:

**1.2.1. Run Fastp**

1.Search for **fastp** in the tools

2.Then select **fastp - fast all-in-one preprocessing for FASTQ files**

-Select custom parameters:

```
  3.Single-end or paired reads > Paired

     4.Input 1 > Browse datasets (right folder icon) > Select
  ERR5310322_1.fastq.gz

     5.Input 2 > Browse datasets > Select ERR5310322_2.fastq.gz

  6.Display Filter Options

     -Quality Filtering options

        7.Qualified Quality Phred = 30

        8.Unqualified percent limit = 10

     -Length Filtering Options

        9.Length required = 50

  10.Read modification options

     11.PoliX tail trimming > Enable polyX tail trimming

     -Per read cutting by quality options

        12.Cut by quality in front (5') > Yes

        13.Cut by quality in tail (3') > Yes

        14.Cutting mean quality = 30
```

15.Finally, click on **Run tool**

To see the trimming stats, have a look at the **fastp on data 2 and data 1: HTML report** file. You should see something like that.

# fastp report for ERR5310322_1_fastq_gz.fastq.gz

## Summary

### General

| | |
|---|---|
| fastp version: | 0.20.1 (https://github.com/OpenGene/fastp) |
| sequencing: | paired end (151 cycles + 151 cycles) |
| mean length before filtering: | 105bp, 105bp |
| mean length after filtering: | 113bp, 113bp |
| duplication rate: | 19.977989% |
| Insert size peak: | 84 |

### Before filtering

| | |
|---|---|
| total reads: | 531.978000 K |
| total bases: | 56.257825 M |
| Q20 bases: | 54.842431 M (97.484094%) |
| Q30 bases: | 54.605191 M (97.062393%) |
| GC content: | 50.644494% |

### After filtering

| | |
|---|---|
| total reads: | 433.314000 K |
| total bases: | 49.003611 M |
| Q20 bases: | 48.876432 M (99.740470%) |
| Q30 bases: | 48.825481 M (99.636496%) |
| GC content: | 51.087943% |

### Filtering result

▶ How many reads have we lost?

**1.2.2. Other trimming tools: Trimmomatic**

1.Search for **trimmomatic** in the tools

2.Select **Trimmomatic flexible read trimming tool for Illumina NGS data**

-Select custom parameters:

```
3.Single-end or paired-end reads? = Paired-end (two separated files)

4.Input FASTQ file (R1/first of pair) = ERR5310322_1.fastq.gz

5.Input FASTQ file (R2/second of pair) = ERR5310322_2.fastq.gz

6.Average quality required = 30

7.Insert Trimmomatic Operation:

    8.Select Trimmomatic operation to perform: **MINLEN**
```

```
    9.Minimum length of reads to be kept = 50
```

10.Select **Run tool**





Trimmomatic does not perform statistics over trimmed reads, so we need to perform FastQC again over the Trimmomatic results.

▶ Try to do it on your own.

**Second question**

▶ How can I improve the quality of my data?

- This hands-on history URL: https://usegalaxy.eu/u/svarona/h/illumina-preprocessing

# 2. Nanopore Quality control and trimming

| Title | Galaxy |
|---|---|
| **Training dataset:** | The data we are going to manage corresponds to Nanopore amplicon sequencing data using ARTIC network primers por SARS-CoV-2 genome. From the Fast5 files generated by the ONT software, we are going to select the pass reads, so they are already filtered by quality. |
| **Questions:** | • How do I know if my Nanopore data was correctly sequenced? |
| **Objectives**: | • Perform a quality control in raw Illumina reads<br>• Perform a quality trimming in raw Nanopore reads<br>• Perform a quality control in trimmed Nanopore reads |
| **Estimated time**: | 15 min |

## 2.1. Quality control

To run the quality control over the samples, follow these steps:

1. Create a new history has explained yesterday named **Nanopore quality**
2. Upload data as seen yesterday, copy and paste the following URLs:

```
https://raw.githubusercontent.com/nf-core/test-
datasets/viralrecon/nanopore/minion/fastq_pass/barcode01/FA093606_pass_bar
code01_7650855b_0.fastq
https://raw.githubusercontent.com/nf-core/test-
datasets/viralrecon/nanopore/minion/fastq_pass/barcode01/FA093606_pass_bar
code01_7650855b_1.fastq
https://raw.githubusercontent.com/nf-core/test-
datasets/viralrecon/nanopore/minion/fastq_pass/barcode01/FA093606_pass_bar
code01_7650855b_2.fastq
```

### 2.1.1. PycoQC

To use PycoQC we need to use the `sequencing_summary.txt` provided by de Nanopore sequencing machine.

Upload data as seen yesterday, copy and paste the following URL:

```
https://raw.githubusercontent.com/nf-core/test-
datasets/viralrecon/nanopore/minion/sequencing_summary.txt
```

1. Search for the **Pycoqc** tool
2. Select **Pycoqc quality control for Nanopore sequencing data**
3. In *A sequencing_summary file*: Select the `sequencing_summary.txt` we just uploaded
4. Select **Run tool**



Then inspect the resulting PycoQC HTML Report:

# General run summary

| Status | Run Duration (h) | Active Channels | Number of Runids | Number of Barcodes |
|---|---|---|---|---|
| All Reads | 7.98 | 496 | 1 | 9 |
| Pass Reads | 7.98 | 494 | 1 | 9 |

# Basecall summary

| Status | Reads | Bases | N50 | Median Read Length | Median PHRED score |
|---|---|---|---|---|---|
| All Reads | 2.449000e+4 | 1.284332e+7 | 515 | 514 | 13.055 |
| Pass Reads | 2.435200e+4 | 1.276470e+7 | 515 | 514 | 13.067 |

**Question**

▶ How many reads do the samples have?

▶ Do you understand all the plots?

**Number of reads per barcode**:

This plot shows the number of reads per barcode, which means de number of reads per sample to be demultiplexed. In a goog experiment, all the barcodes should have the same number of reads. In this training we only used reads from barcode01 sample but we can see that barcode08 couldn't be correctly sequenced.
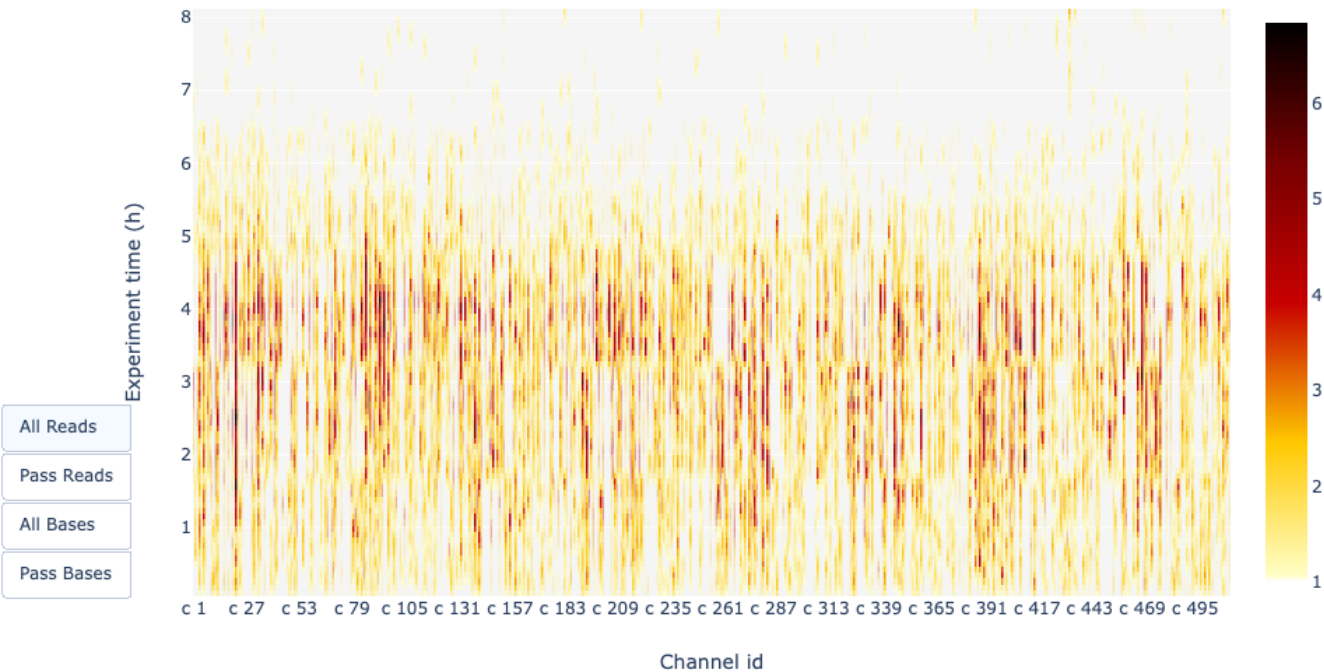
**Channel activity over time:**



It gives an overview of available pores, pore usage during the experiment, inactive pores and shows if the loading of the flow cell is good (almost all pores are used). In this case, the vast majority of channels/pores are inactive (white) after the 6h of experiment, so the run should have been dinished at that time. You would

hope for a plot that it is dark near the X-axis, and with higher Y-values (increasing time) doesn't get too light/white. Depending if you chose "Reads" or "Bases" on the left the colour indicates either number of bases or reads per time interval.

▶ How do I check whether my Nanopore data was correctly sequenced?

## 2.2. Trimming

When Nanopore reads are being sequenced, the MinKnown software splits Fast5 reads into quality pass and quality fail. As we will select only Fast5 pass reads, we won't need to perform a quality trimming, so even if we see that the reads have a bad Phred score, we know that the ONT software considered the reads as "good quality".

Then we will only be performing a read length trimming. As we are using amplicon sequencing data, we won't be expecting reads smaller than 400 nucleotides, nor higher than 600, which would obviously correspond to chimeric reads.

### 2.2.1. Artic

1. Search for **artic** tool
2. Select **ARTIC guppyplex Filter Nanopore reads by read length and (optionally) quality**
3. Structure of your input data: Multiple input datasets per sample
4. While pressing the *Ctrl* key, select the three samples
5. Remove reads longer than = 600
6. Remove reads shorter than = 300
7. Do not filter on quality score (speeds up processing) = Yes (we had already select pass reads)

### 2.2.2. Nanoplot

Now we are going to run NanoPlot on filtered data:

1. Search for the **Nanoplot** tool and select **NanoPlot Plotting suite for Oxford Nanopore sequencing data and alignments**
2. Run the tool as follows:
   - In the *files* part, select ARTIC output file.
   - Display **Options for customizing the plots created**:
     - **Specify the bivariate format of the plots** > *Select all*
     - **Show the N50 mark in the read length histogram** > *Yes*
   - Select **Execute**



### Questions

▶ Did our data length and quality improve?

▶ How many reads did we lost during trimming step?

- This hands-on history URL: https://usegalaxy.eu/u/svarona/h/nanopore-quality

> **NOTE:** We can't use nanofilt because it is not installed in Galaxy