

# INTRODUCTION TO UNIX-LINUX – QUALITY CONTROL

*Klebsiella Workshop*

*Sep 2024*



UNIVERSITY of the  
WESTERN CAPE



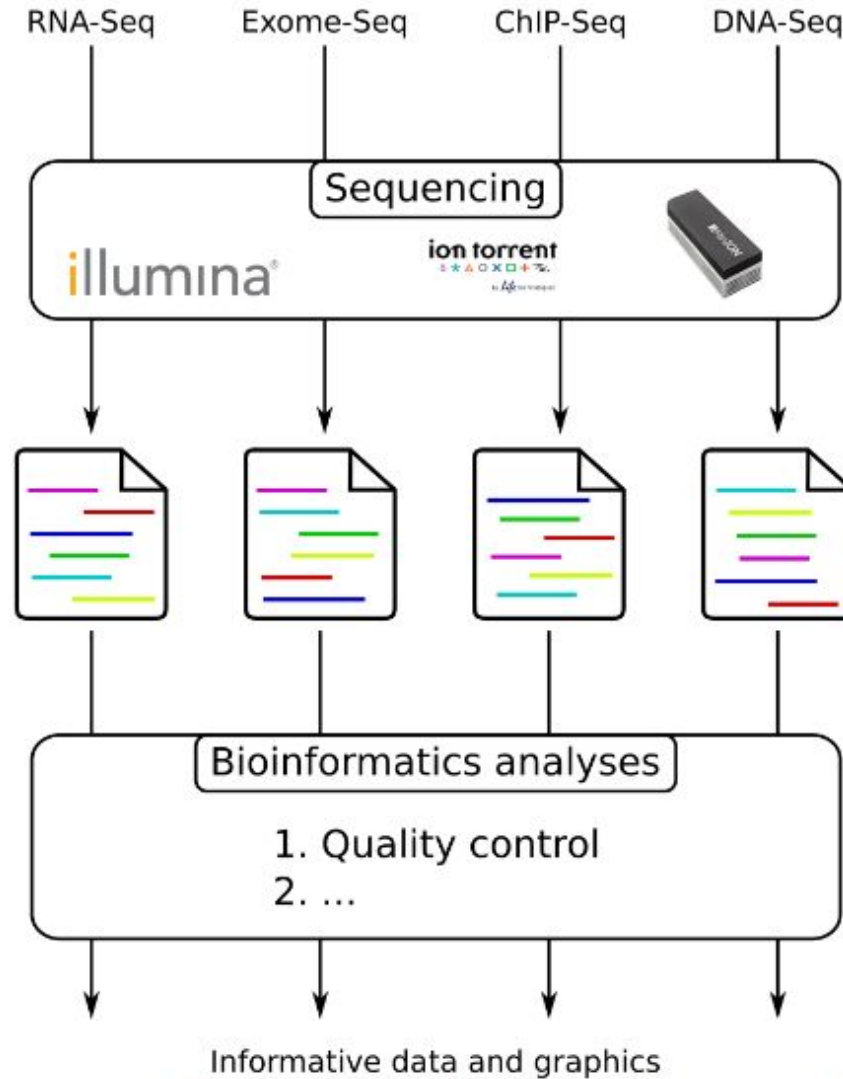
SANBI

South African National  
Bioinformatics Institute



# From experiments to data

we are here



Quality control = First step of the bioinformatics analyses

Quality  
control  
(raw reads)

Trimming  
(adapters/low  
quality reads)

Alignment  
/mapping  
(QC)

Variant  
calling  
(QC)

Tree  
Building  
(QC/caveats)

# Different Quality Control tools

we are here

## Illumina

- FastQC (most widely used)
- MultiQC (compresses dataset)
- Fastp (all-in-one)

## DeNovo

- Metaplan
- BLAST

Quality  
control

Adaptor /  
Trimming

Alignment  
/mapping

Variant  
calling

Tree  
Building

# Quality Control: FastQC

we are here

- A quality control tool for high throughput sequence data.
- FastQC aims to provide a simple way to do some quality control checks on raw sequence data coming from high throughput sequencing pipelines. It provides a modular set of analyses which you can use to give a quick impression of whether your data has any problems of which you should be aware before doing any further analysis.

The main functions of FastQC are:

- Import of data from BAM, SAM or FastQ files (any variant)
- Providing a quick overview to tell you in which areas there may be problems
- Summary graphs and tables to quickly assess your data
- Export of results to an HTML based permanent report
- Offline operation to allow automated generation of reports without running the interactive application

Quality  
control

Adaptor /  
Trimming

Alignment  
/mapping

Variant  
calling

Tree  
Building

# Installation & Usage: FastQC

we are here

- Installation:

- Wget

<https://github.com/s-andrews/FastQC/archive/refs/tags/v0.12.1.zip>

- Usage:

- Run either interactive graphical application in which you can dynamically load FastQ files and view their results OR
  - In a non-interactive mode where you specify the files you want to process on the command line and FastQC will generate an HTML report for each file without launching a user interface.

Quality  
control

Adaptor /  
Trimming

Alignment  
/ mapping

Variant  
calling

Tree  
Building

# Installation of fastQC using conda

we are here

```
conda install bioconda::fastqc
conda install
bioconda/label/broken::fastqc
conda install
bioconda/label/cf201901::fastqc
```

Quality  
control

Adaptor /  
Trimming

Alignment  
/mapping

Variant  
calling

Tree  
Building



# Quality Scoring

we are here

Quality  
control

Measure of the quality of the identification of the nucleobases  
generated by automated DNA sequencing

Phred Quality Score	Probability of incorrect base call	Base call accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1000	99.9%
40	1 in 10,000	99.99%
50	1 in 100,000	99.999%
60	1 in 1,000,000	99.9999%

Adaptor /  
Trimming

Alignment  
/mapping

Variant  
calling

Tree  
Building

# Sequence: FASTA

we are here

# Quality control

# Adaptor / Trimming Alignmen t / mapping Variant calling

# Tree Building

```
>Identifier1 (comment)
```

XX  
XX  
XX

```
>Identifier2 (comment)
```

[illegible]



# Sequence: FASTQ

we are here

Quality  
control

Adaptor /  
Trimming

Alignmen  
t /  
mapping

Variant  
calling

Tree  
Building

@Identifier1 (comment)

XX

+ ←

QQ

@Identifier2 (comment)

XX

+

QQ



Quality control

Adaptor /  
Trimming

Alignmen  
t /

mapping

Variant  
calling

Tree  
Building

# FastQC Report

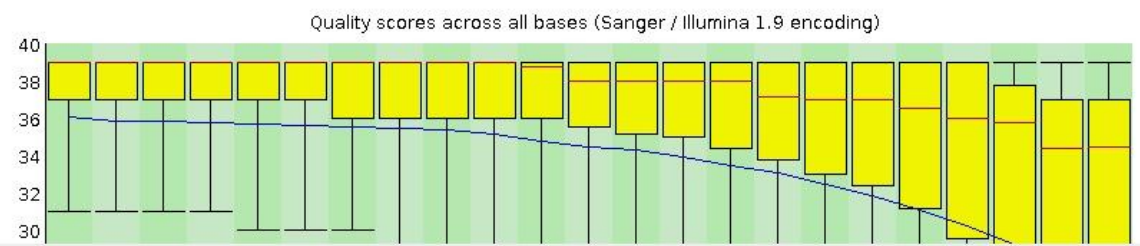
## Summary

- ✓ Basic Statistics
- ✓ [Per base sequence quality](#)
- ✓ [Per sequence quality scores](#)
- ✓ [Per base sequence content](#)
- ✓ [Per base GC content](#)
- ! [Per sequence GC content](#)
- ✓ [Per base N content](#)
- ✓ [Sequence Length Distribution](#)
- ✓ [Sequence Duplication Levels](#)
- ✓ [Overrepresented sequences](#)
- ! [Kmer Content](#)

## ✓ Basic Statistics

Measure	Value
Filename	WES_human_Illumina.pe_1.fastq
File type	Conventional base calls
Encoding	Sanger / Illumina 1.9
Total Sequences	4942814
Filtered Sequences	0
Sequence length	76
%GC	47

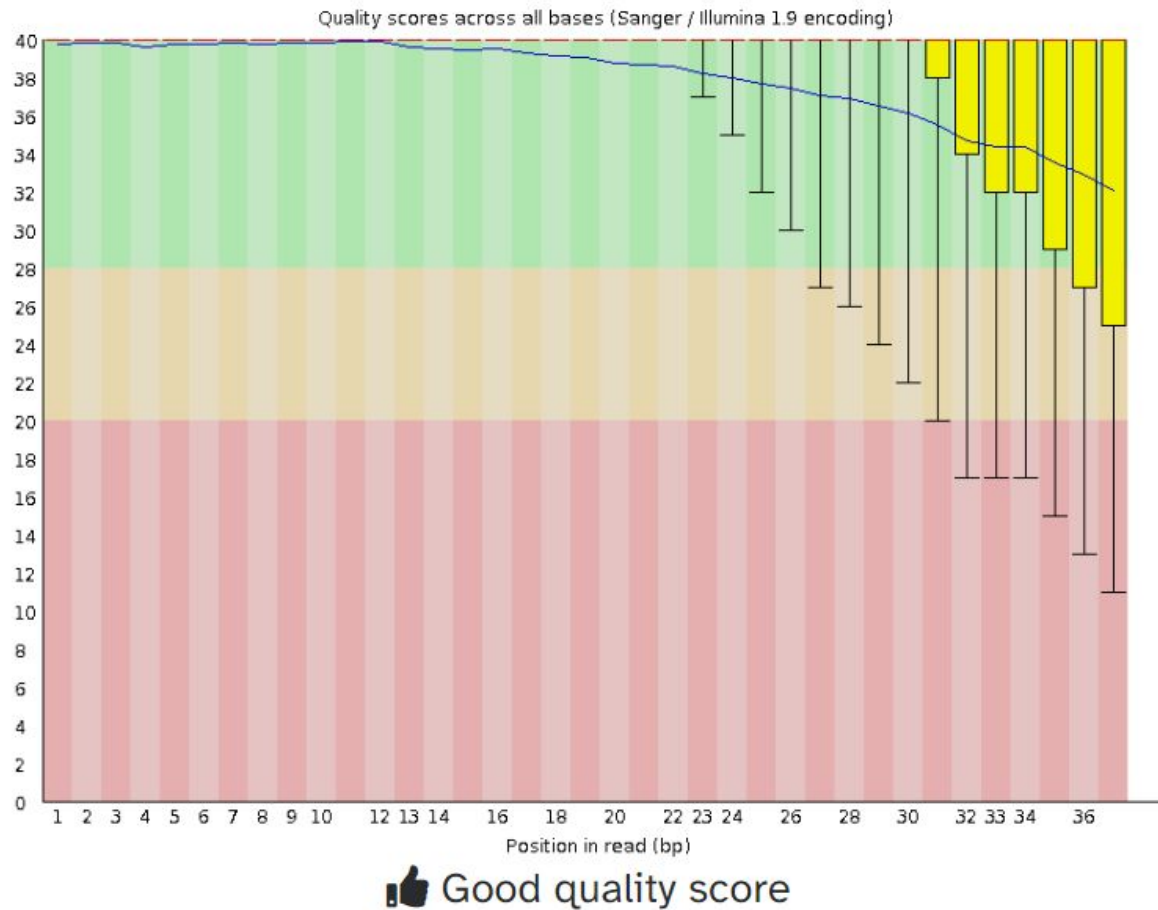
## ✓ Per base sequence quality



Produced by [FastQC](#) (version 0.10.1)

## Quality score: Per-base

we are here



Quality  
control

Adaptor /  
Trimming

Alignment  
/ mapping

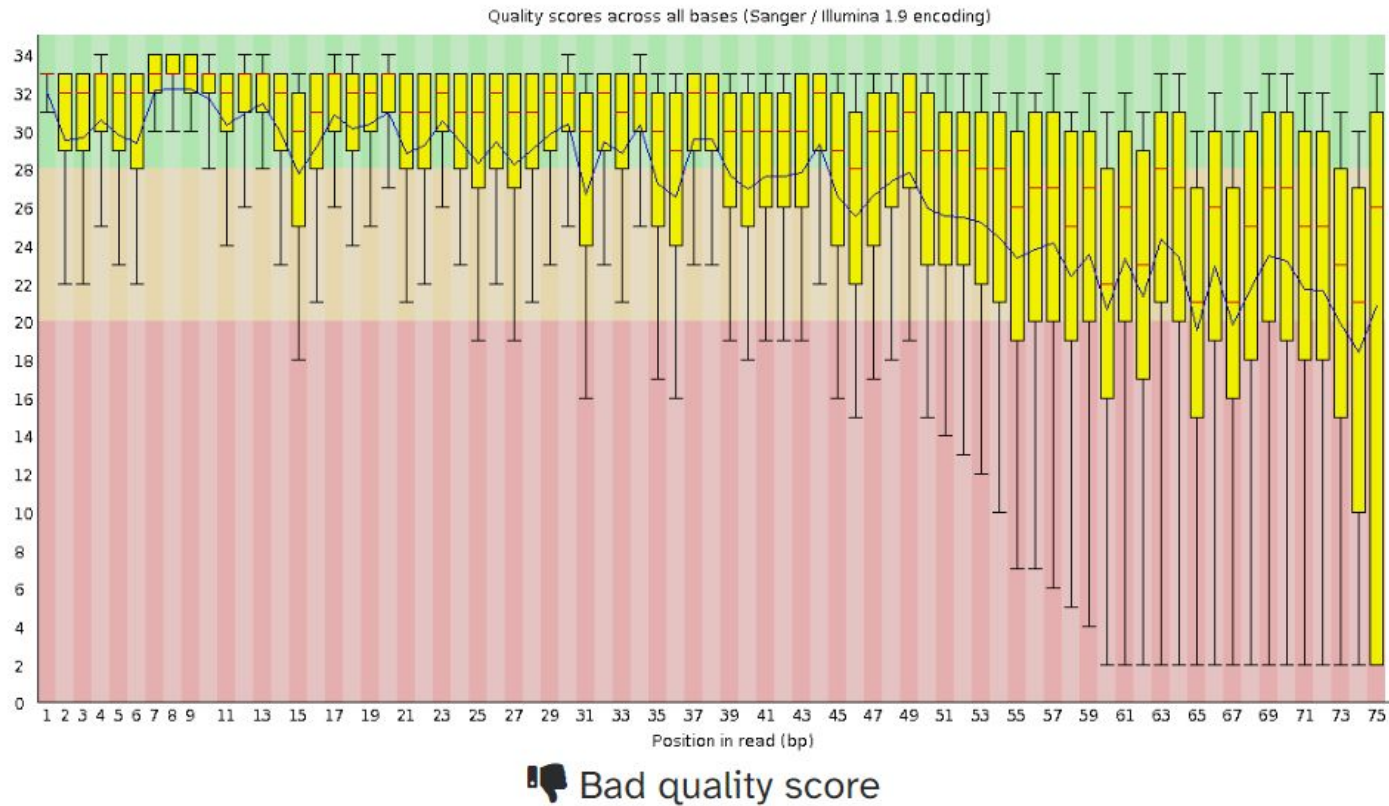
Variant  
calling

Tree  
Building

we are here

Quality  
control

## Per-base Quality



Adaptor /  
Trimming

Alignment  
/ mapping

Variant  
calling

Tree  
Building



# Improving Quality

- Filtering of sequences
  - with small mean quality score
  - too short
  - with too many N bases
  - based on their GC content
- Cutting/Trimming sequences
  - from low quality score parts
  - tails

we are here

Quality  
control

Adaptor /  
Trimming

Alignment  
/ mapping

Variant  
calling

Tree  
Building



# Different trimming tools

Illumina – ref data

- Trimmomatic
- Cutadapt
- Fastp
- Flexar
- Trimalore

DeNovo – no ref

- Kraken
- Blobology
- BLAST

we are here



Quality  
control

Adaptor /  
Trimming

Alignment  
/ mapping

Variant  
calling

Tree  
Building

# Trimmomatic

- Trimmomatic performs a variety of useful trimming tasks for Illumina paired-end and single ended data. The selection of trimming steps and their associated parameters are supplied on the command line.
- Requires Java installation check: `run java -version`
- **Inputs** - Single-end or Paired-end FASTQ or FASTQ.gz reads
- <https://github.com/usadellab/Trimmomatic?tab=readme-ov-file>

we are here

Quality  
control

Adaptor /  
Trimming

Alignment  
/ mapping

Variant  
calling

Tree  
Building

# Installation & Usage: Trimmomatic

- Installation:

- `wget https://github.com/usadellab/Trimmomatic.git`

- Usage:

- `java -jar trimmomatic-0.39.jar PE input_forward.fq.gz  
input_reverse.fq.gz output_forward_paired.fq.gz  
output_forward_unpaired.fq.gz output_reverse_paired.fq.gz  
output_reverse_unpaired.fq.gz  
ILLUMINACLIP:TruSeq3-PE.fa:2:30:10:2:True LEADING:3 TRAILING:3  
MINLEN:36`

- The above code is for reference only

we are here

Quality  
control

Adaptor /  
Trimming

Alignment  
/ mapping

Variant  
calling

Tree  
Building

# Fastp

- [fastp](#) is a tool designed to provide fast all-in-one preprocessing for FASTQ files. This tool is developed in C++ with multithreading supported to afford high performance.
- **Inputs** - Single-end or Paired-end FASTQ or FASTQ.GZ reads
- <https://github.com/OpenGene/fastp?tab=readme-ov-file#features>

we are here

Quality  
control

Adaptor /  
Trimming

Alignment  
/ mapping

Variant  
calling

Tree  
Building

# Installation & Usage: Fastp

- Installation:

- `wget http://opengene.org/fastp/fastp`
- `chmod a+x ./fastp`

- Usage:

- SE code: `fastp -i in.fq -o out.fq`
- PE code: `fastp -i in.R1.fq.gz -I in.R2.fq.gz -o out.R1.fq.gz -O out.R2.fq.gz`

- The above code is for reference only



we are here

Quality  
control

Adaptor /  
Trimming

Alignment  
/ mapping

Variant  
calling

Tree  
Building