

# After sequencing QC

---

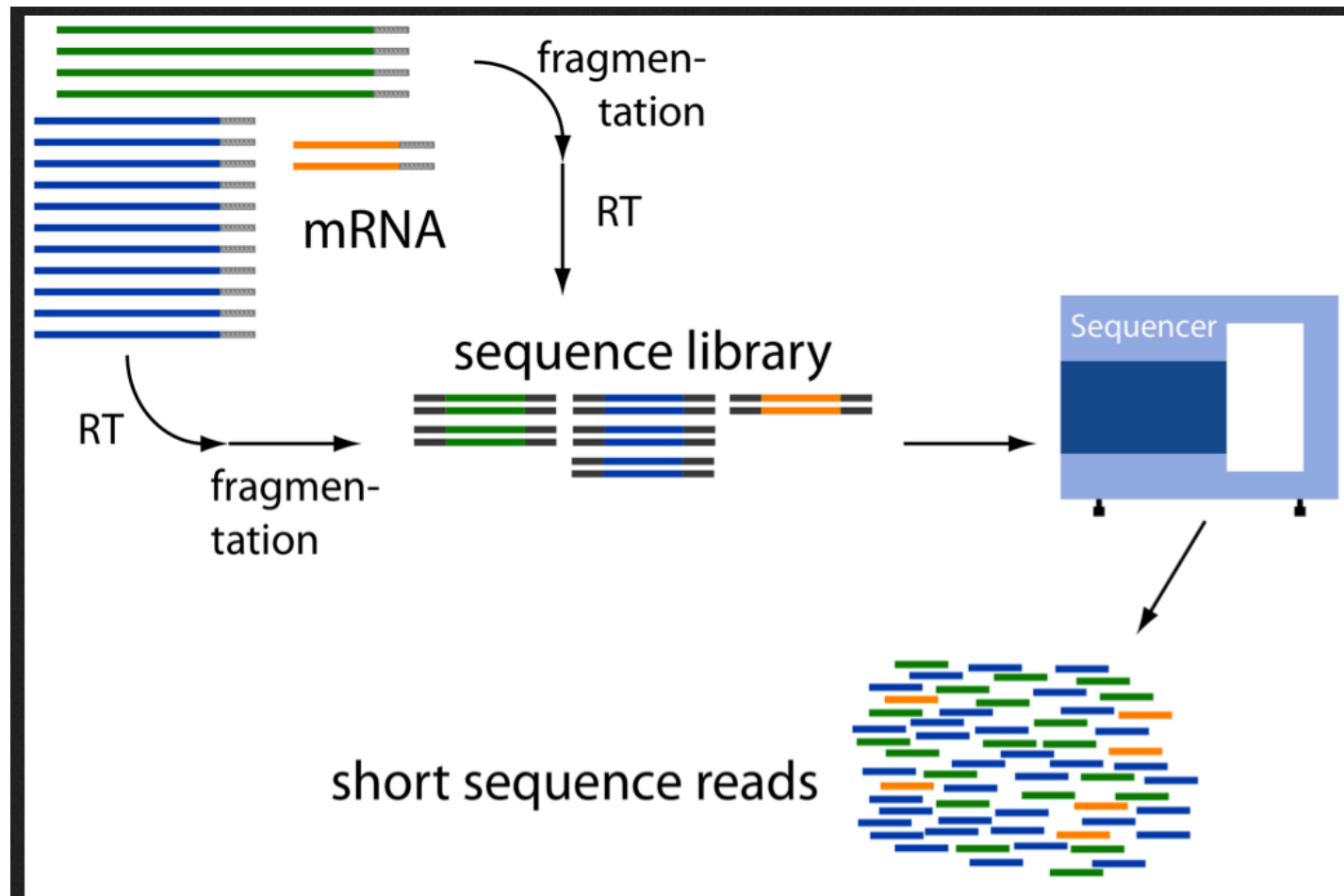
RNA-seq data analysis

**Johan Reimegård** | 13-May-2019

# Overview

- What can affect your data?
- FastQC – read based QC
- RSeQC – mapping based QC
- PCA
- Preventive measurements: spike-in controls, experimental design

# RNA-seq libraries



What could go wrong?

# What could go wrong?

- RNA quality:

ies)

- Library prep:

- Contaminations

- Se

- Contaminations

- Sequence complexity

# What could go wrong?

- RNA quality:
  - Degradation
  - Contaminations (pathogens or other sources)
  - GC-bias
  - Nuclear vs organelle reads
- Library prep:
  - Contaminations
- Sequencing:
  - Contaminations
  - Sequence complexity

# What could go wrong?

- RNA quality:
  - Degradation
  - Contaminations (pathogens or other sources)
  - GC-bias
  - Nuclear vs organelle reads
- Library prep:
  - Failed reactions
  - RNA / Adapter ratios – primer dimers
  - Clonal duplicates
  - Chimeric reads
  - Contaminations
- Se
  - Contaminations
  - Sequence complexity

# What could go wrong?

- RNA quality:
  - Degradation
  - Contaminations (pathogens or other sources)
  - GC-bias
  - Nuclear vs organelle reads
- Library prep:
  - Failed reactions
  - RNA / Adapter ratios – primer dimers
  - Clonal duplicates
  - Chimeric reads
  - Contaminations
- Sequencing:
  - Base calling errors
  - Uncalled bases
  - Low quality bases (3' end)
  - Contaminations
  - Sequence complexity

# From samples to reads

- may not be what you think they are

- Mixing samples
  - 30 samples with 5 steps from samples to reads has 24 300 000 potential mix ups of samples
  - Error rate  $1/100$  with 5 steps suggest that one of every 20 sample is mislabeled
- Experiments go wrong
  - 30 samples with 5 steps from samples to reads has 150 potential steps for errors
  - Error rate  $1/100$  with 5 steps suggest that one of every 20 samples the reads does not represent the sample
- Combine the two error sources and approximately one in every 10 samples is wrong



# From samples to reads

- may not be what you think they are

- Mixing samples
- Experiments go wrong
- How do we understand what went wrong?

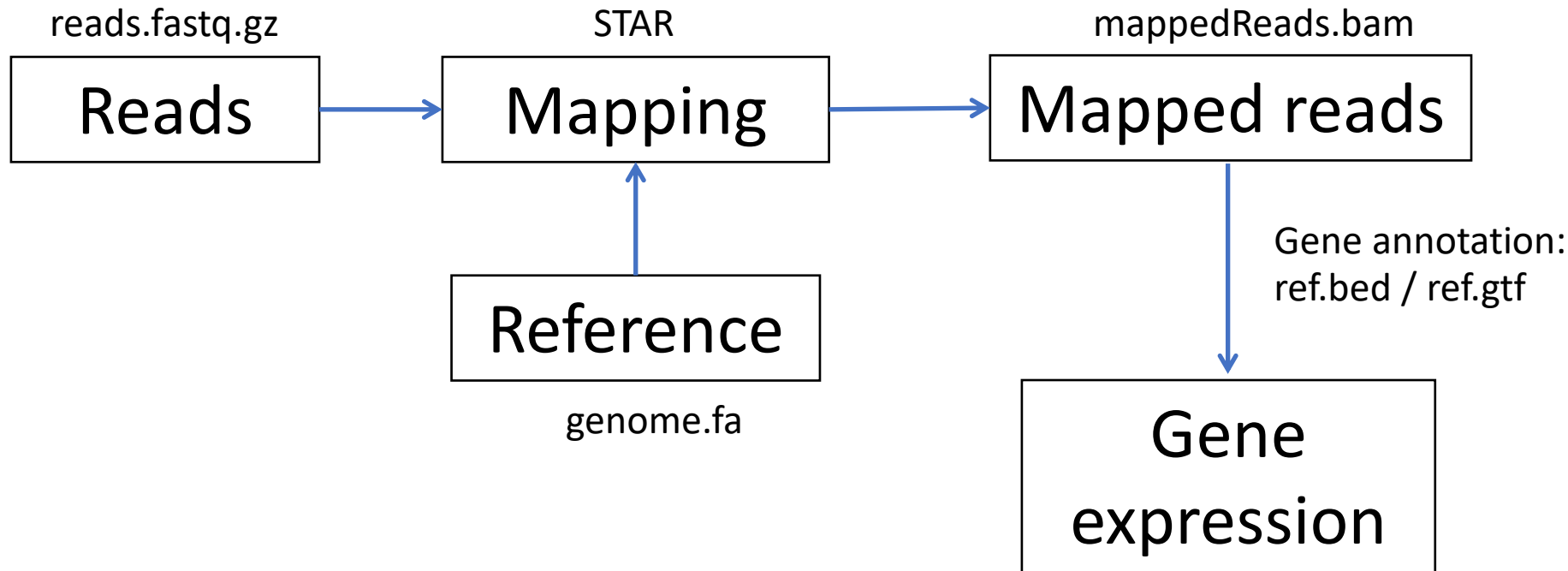
# From samples to reads

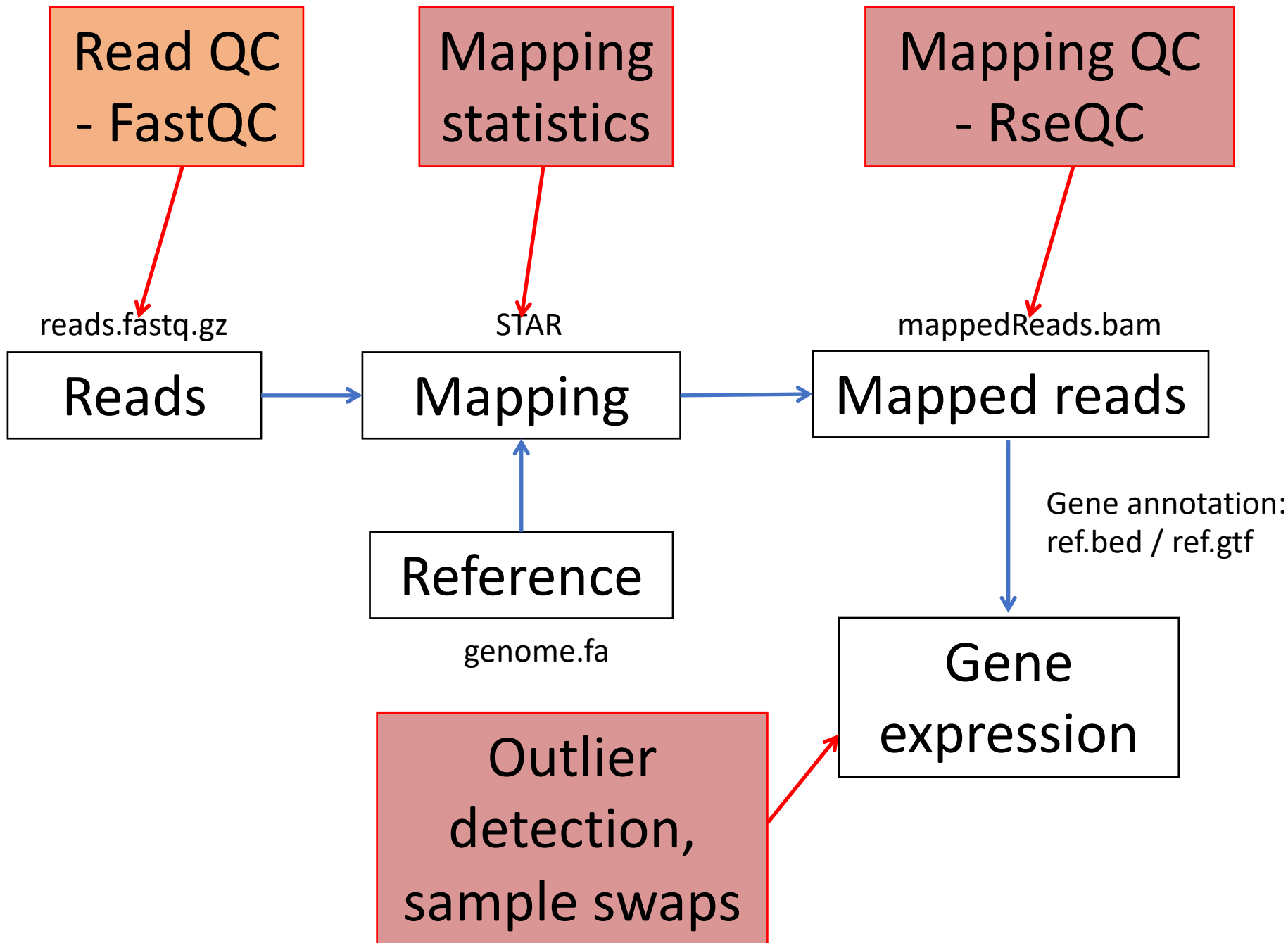
- may not be what you think they are
- Mixing samples
- Experiments go wrong
- How do we understand what went wrong?

BIOINFORMATICS!



# RNA-seq analysis workflow





# Fastq – read file format

```
@SEQ_ID
GATTGTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT
+
!''*((( (***) )%%%++) (%%%) .1***-+*'' ) **55CCF>>>>>CCCCCCC65
```

Unique identifier

Sequence

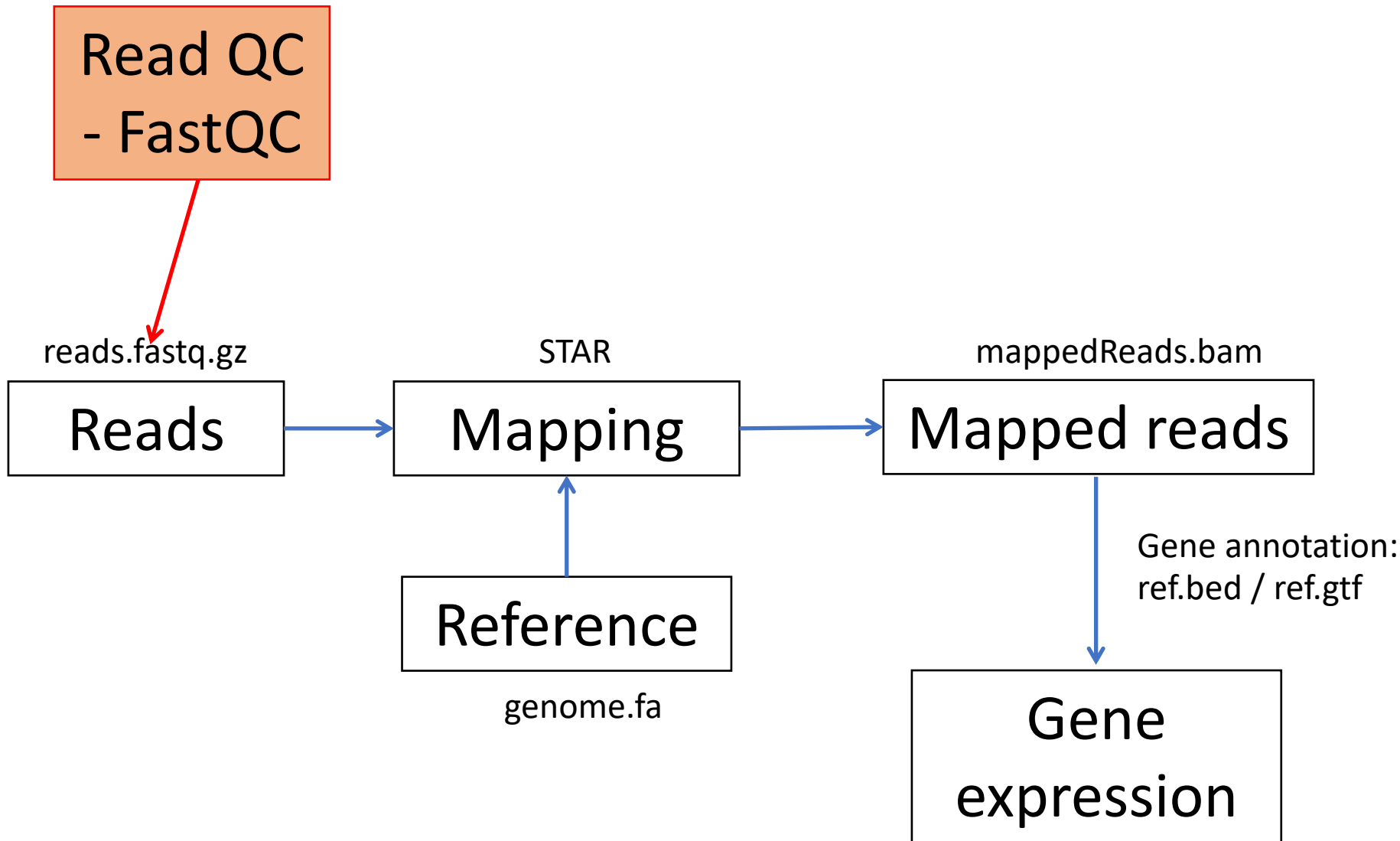
Sequence quality

Paired end data usually in format sampleX\_1.fastq and sampleX\_2.fastq with same SEQ\_ID for both mate pairs, followed by /1 and /2 (or \_f and \_r)

# Fastq – read file format

[illegible]

```
S - Sanger           Phred+33,  raw reads typically (0, 40)
X - Solexa           Solexa+64, raw reads typically (-5, 40)
I - Illumina 1.3+    Phred+64,  raw reads typically (0, 40)
J - Illumina 1.5+    Phred+64,  raw reads typically (3, 40)
                        with 0=unused, 1=unused, 2=Read Segment Quality Control Indicator (bold)
                        (Note: See discussion above).
L - Illumina 1.8+    Phred+33,  raw reads typically (0, 41)
```



# Basic read metrics with FastQC

A program that analyses some of the basic metrics on fastq raw read files.

- Quality
- Length
- Sequence bias
- GC content
- Repeated sequences
- Adapter contamination

## Code

```
$ module load bioinfo-tools
$ module load FastQC/0.11.2

$ fastqc -o outdir seqfile.fastq
# multiple files:
$ fastqc -o outdir seqfile_*.fastq
```

<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>





**Thank you. Questions?**

---

**Johan Reimegård | 13-May-2019**