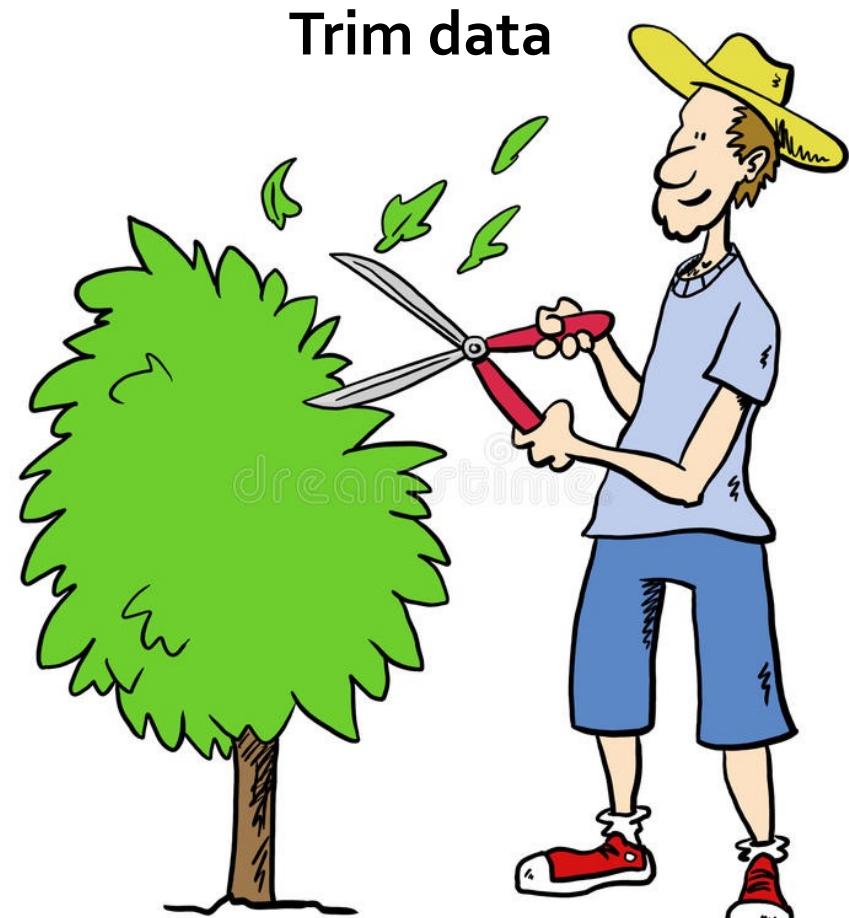


Trimming and filtering



Filter data

© Can Stock Photo - csp34642023



Trim data

<https://www.dreamstime.com/stock-illustration-man-trimming-bush-garden-image64030511>

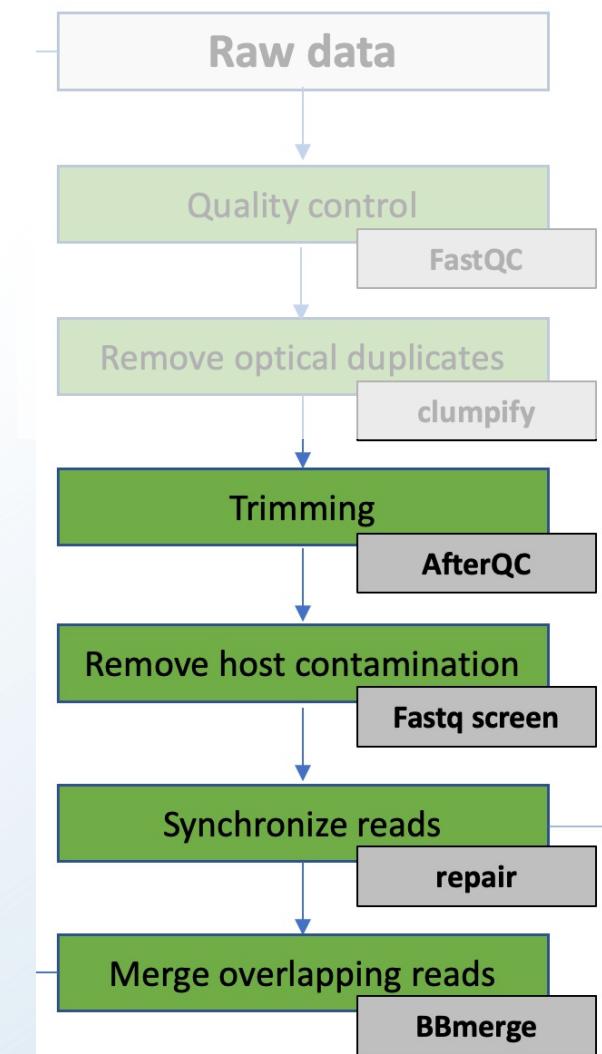
Overview of this talk

Removing poor quality data

Data decontamination

Synchronization of paired-end reads

Merging overlapping reads pairs



Why is it important to perform QC and filtering/trimming?

Data analysis also costs money and time

Filtering data

Length related

Quality score related

GC content related

Ambiguity code related

Sequence complexity related

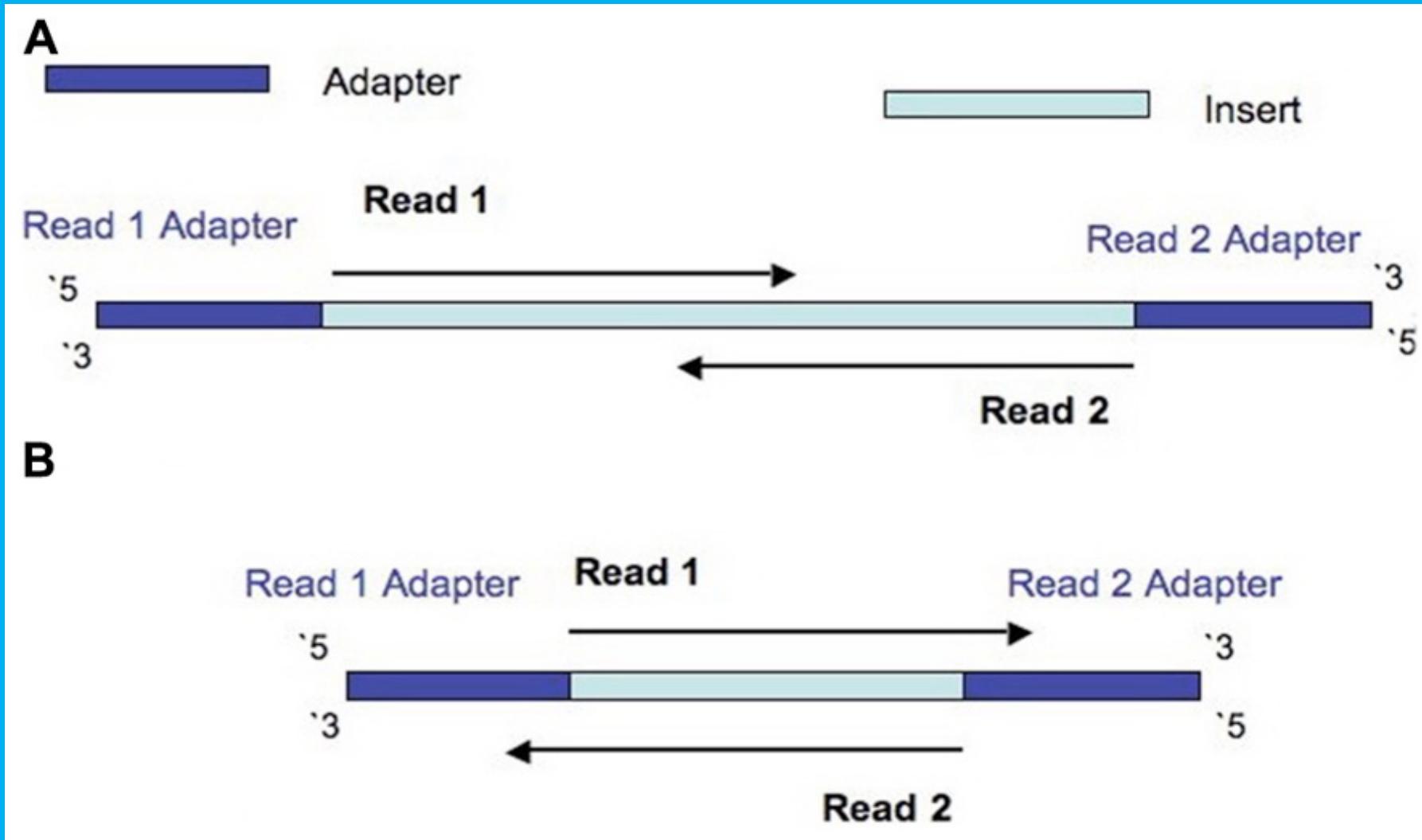
Trimming data

Trim by length/position – fixed, e.g. 20 bp

Trim tails

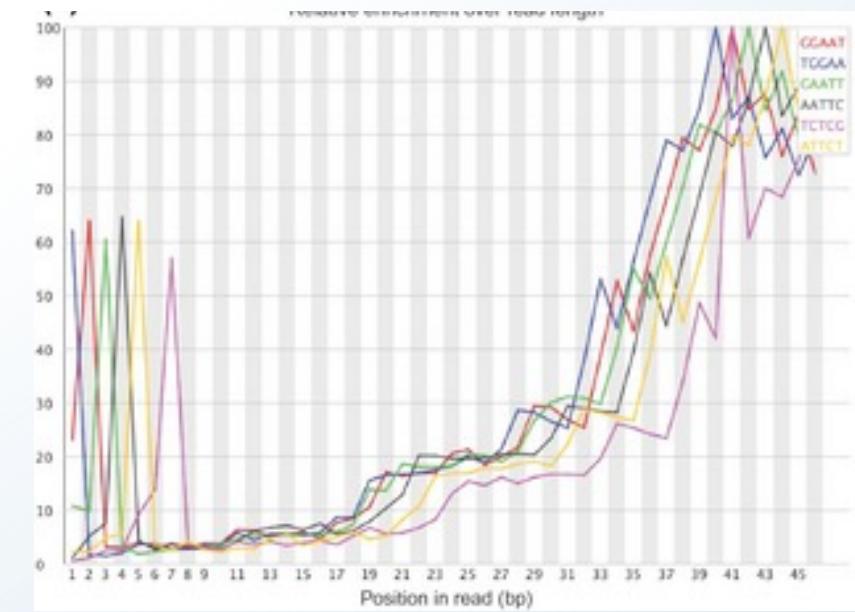
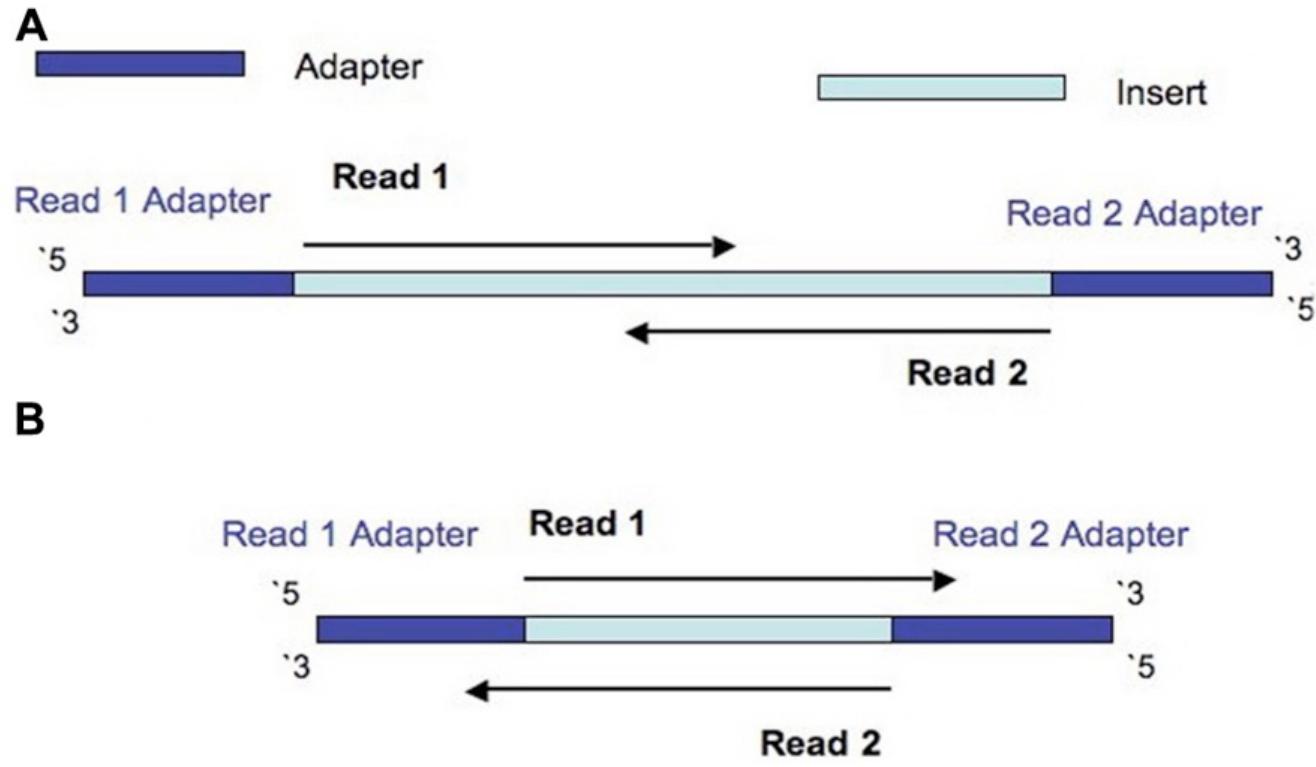
Trim ends by quality scores

Trimming - Discuss two and two: what is this figure showing?



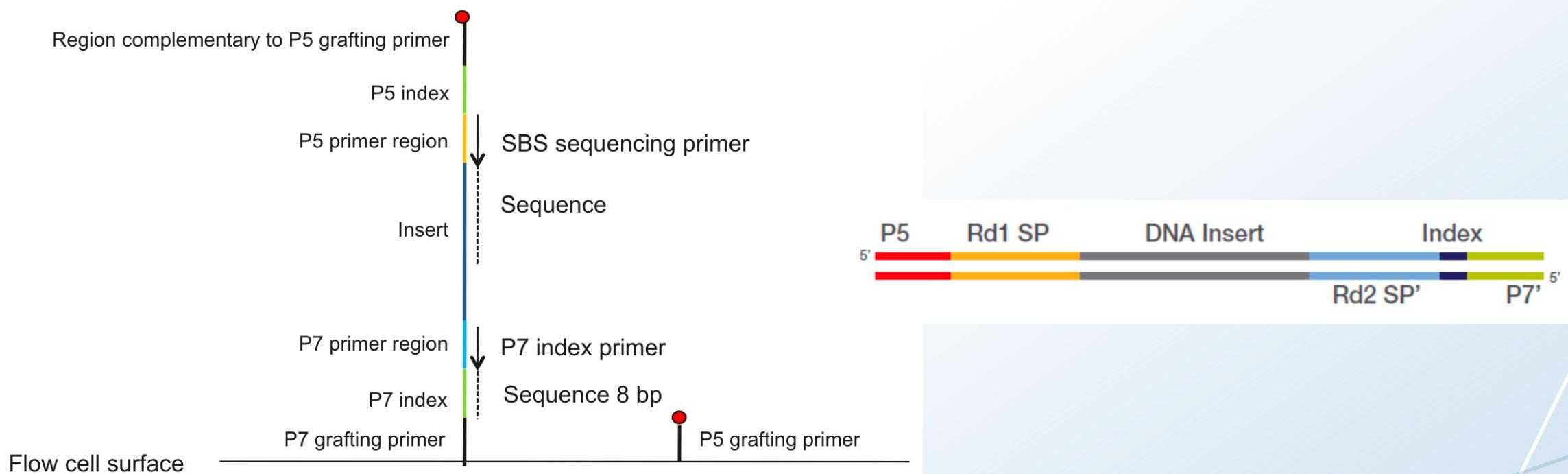
We trim the start and end of reads to remove poor quality data or adapter sequences

If the DNA fragment is shorter than the read length, the sequence reaction will go through the read and into the adapter



Illumina adapter sequences

When preparing a (TrueSeq) library, adaptors are ligated to the DNA of interest



We filter reads to remove poor quality data

Common to set some quality thresholds and remove all reads that does no comply
For example remove all reads with a lower average Q score than 20, or reads with more than 5 Ns

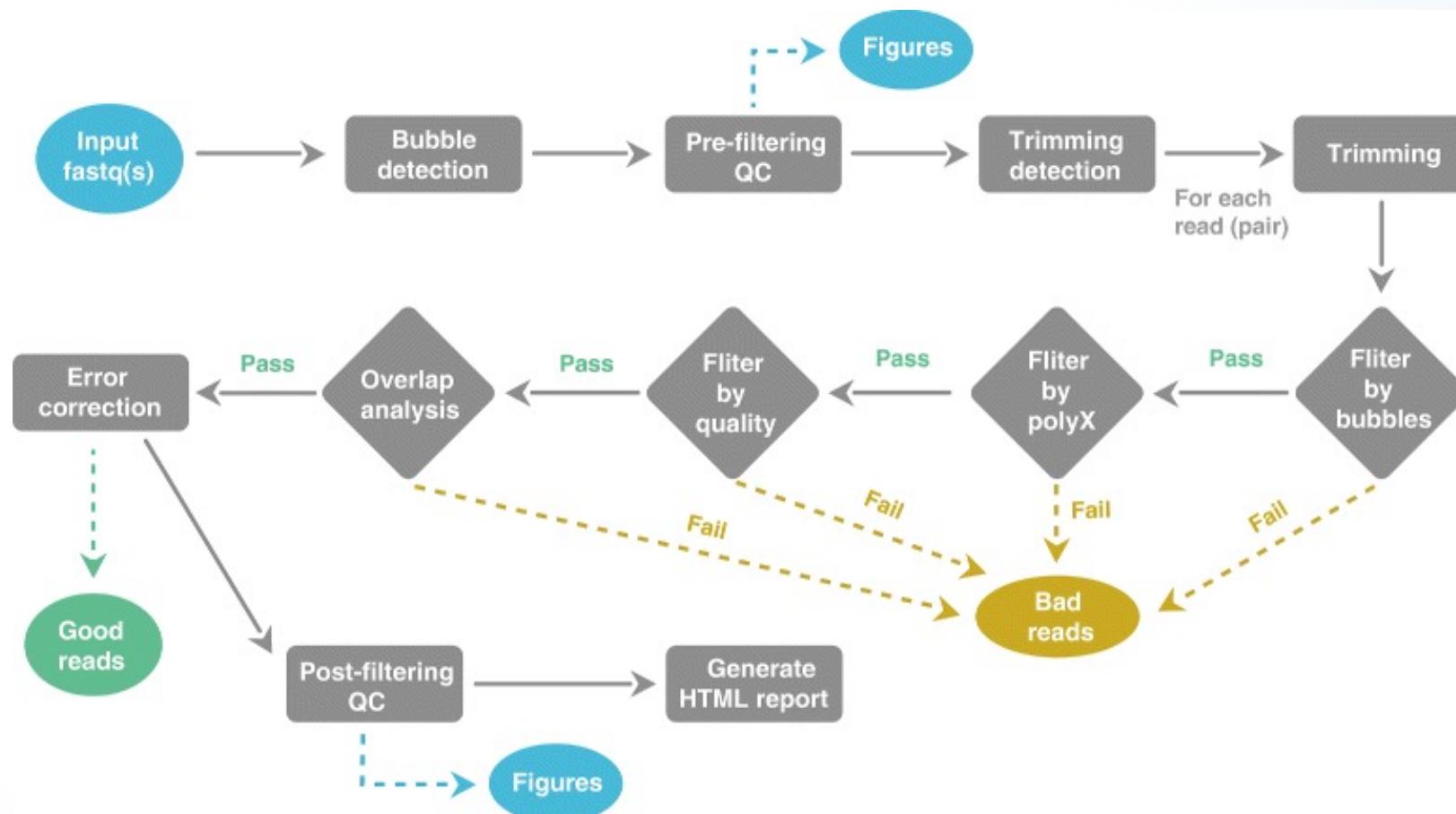
For example:
Remove all reads with a lower
average Q score than 20



For example:
Remove all reads with
more than 5 Ns

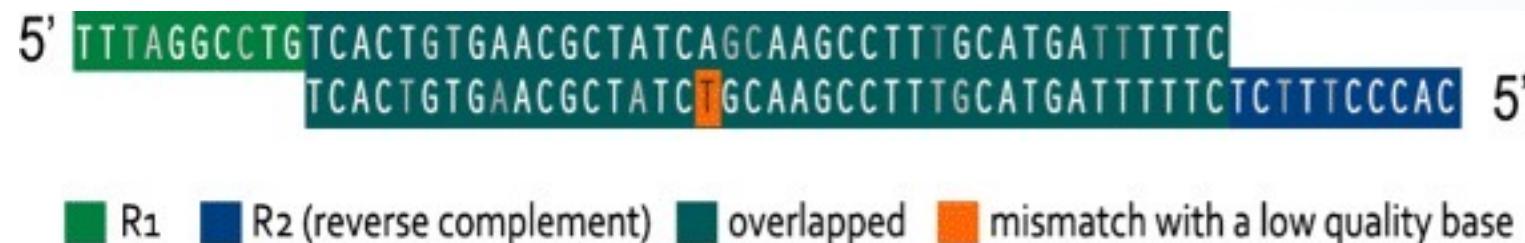
AfterQC - Automatic Filtering, Trimming, Error Removing and Quality Control for FASTQ data

Performs quality control and filtering/trimming of the sequence reads

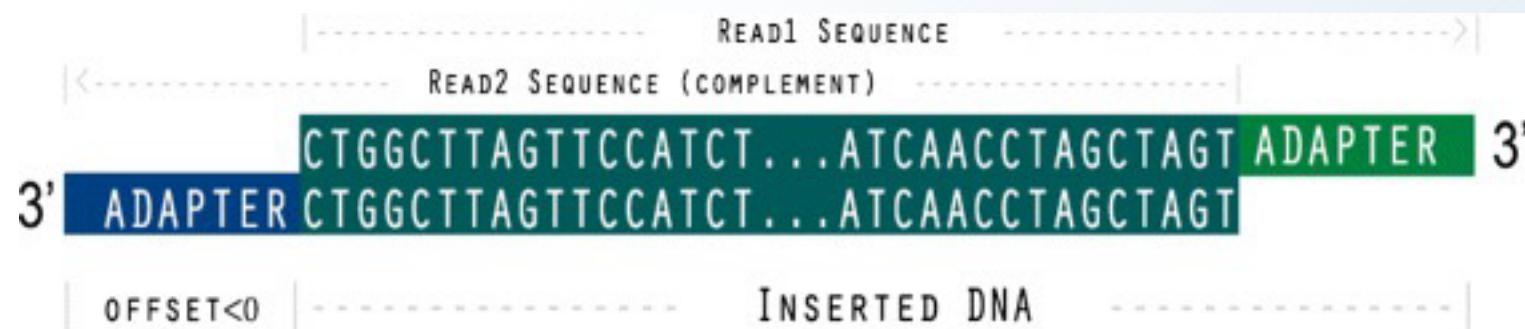


AfterQC analyses the overlap of paired sequences for pair-end sequencing data

AfterQC will correct the low quality base according to its high quality mate



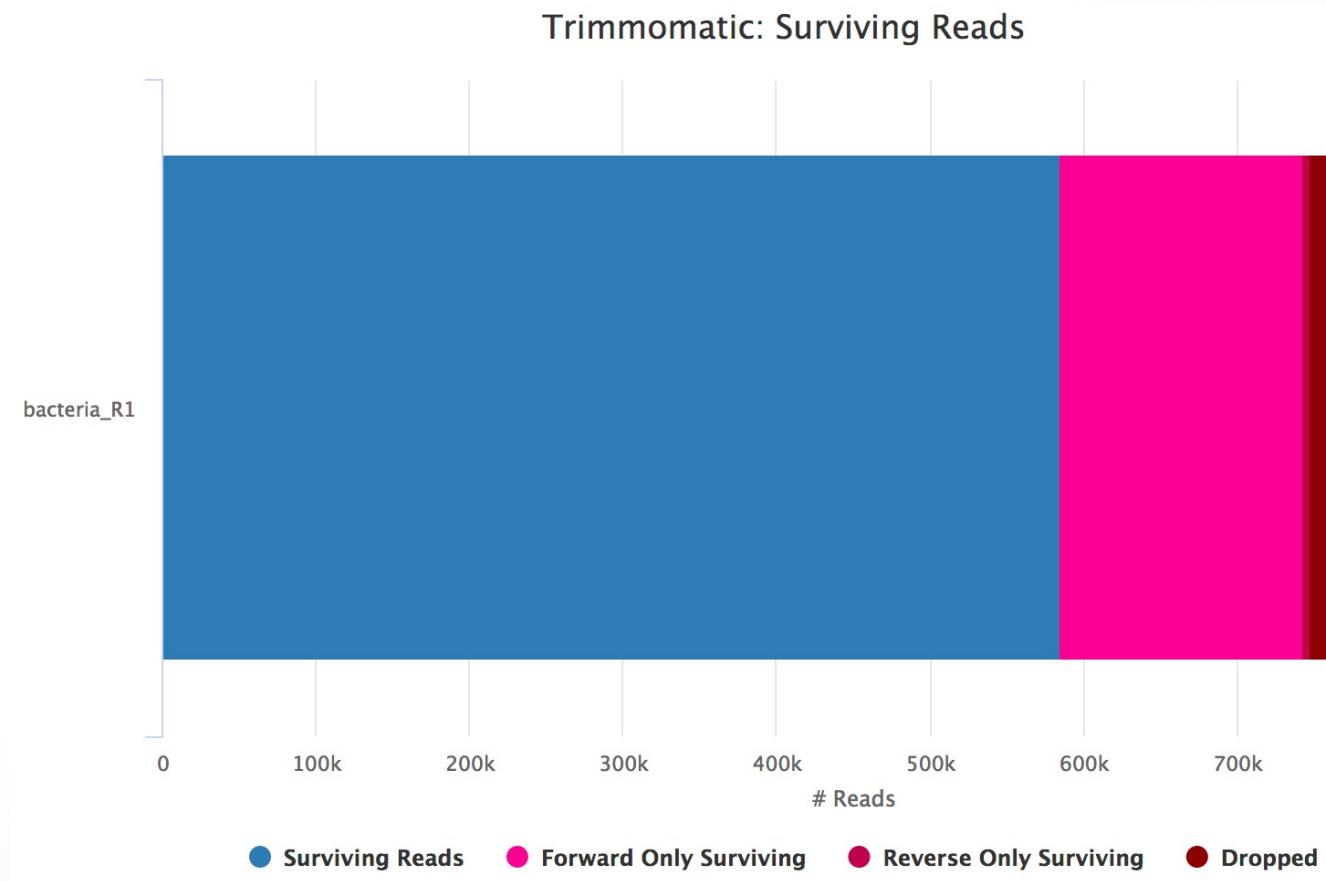
Removal of 3' adapter in the tail



Popular trimming tools

Many tools available – here are some :

- AfterQC
- Fastp
- Trimmomatic
- CutAdapt
- AlienTrimmer
- Sickle
- Trim Galore
- Sythe
- Prinseq



It is important to remove sequence contaminations as early as possible

There can be many sources of contamination in the final sequence library

For example:
PhiX sequences from
the sequencing kit



For example:
Metagenomic samples may
contain sequences from the host

PhiX control

PhiX is used as a quality and calibration control for Illumina sequencing runs
10% of the genomes that are published in literature are contaminated with PhiX



Commentary | [Open Access](#) | Published: 30 March 2015

Large-scale contamination of microbial isolate genomes by Illumina PhiX control

[Supratim Mukherjee](#) [Marcel Huntemann](#), [Natalia Ivanova](#), [Nikos C Kyprides](#) & [Amrita Pati](#)

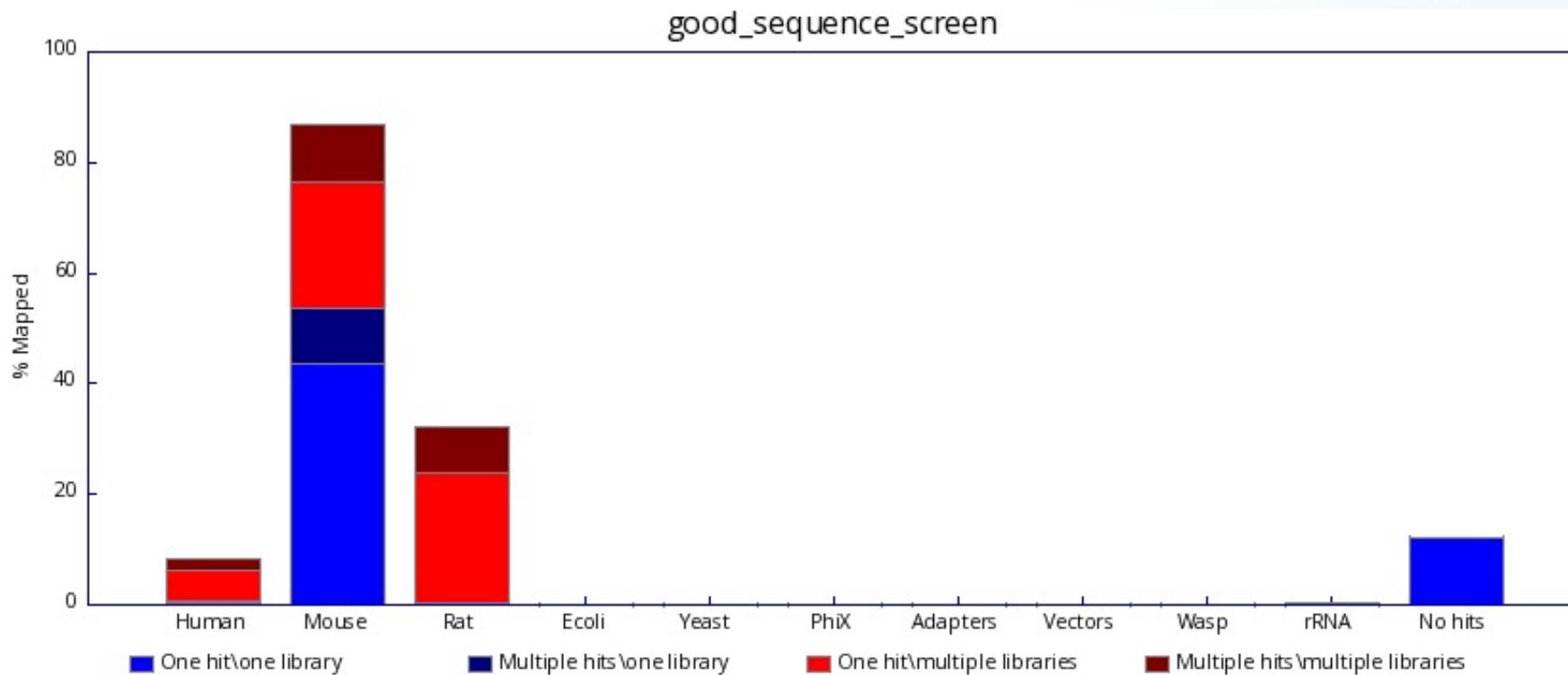
[Standards in Genomic Sciences](#) 10, Article number: 18 (2015) | [Cite this article](#)

<https://www.illumina.com>

FastQ Screen allows you to screen sequences against a set of sequence databases

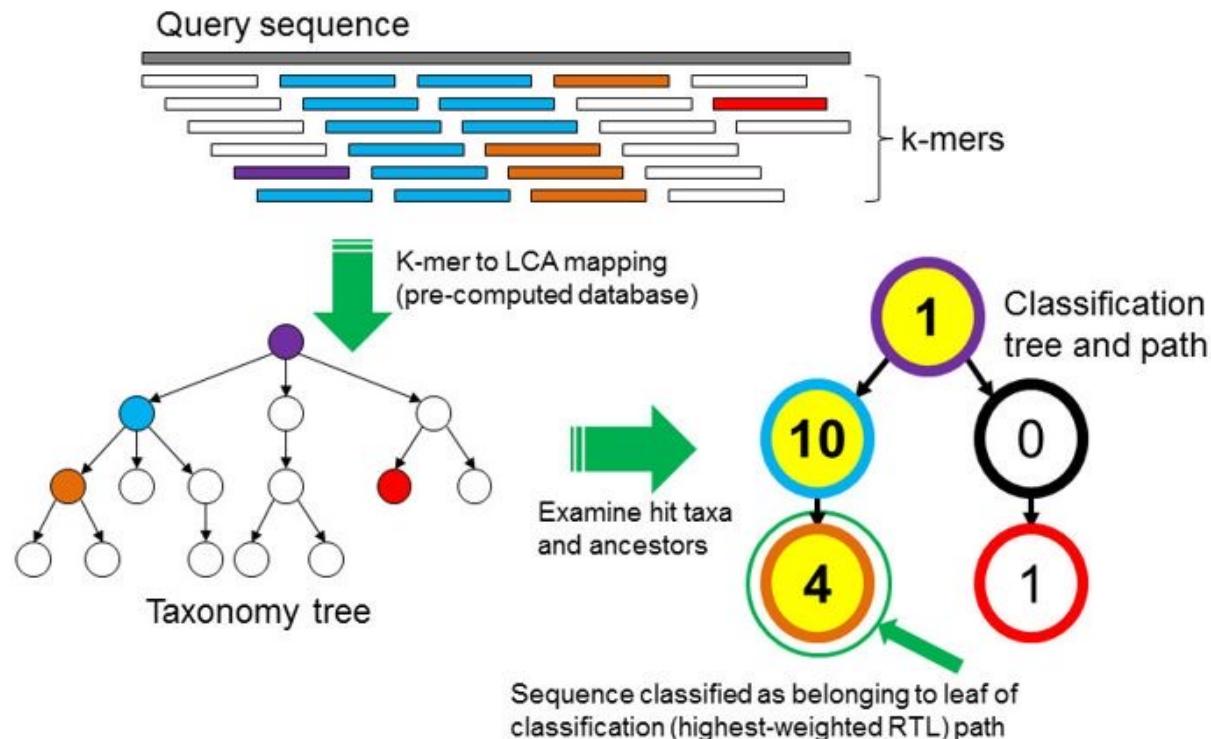
Normally used to screen for host contaminations

Label sequences that match sequences in the database you provide



Kraken for contamination analysis

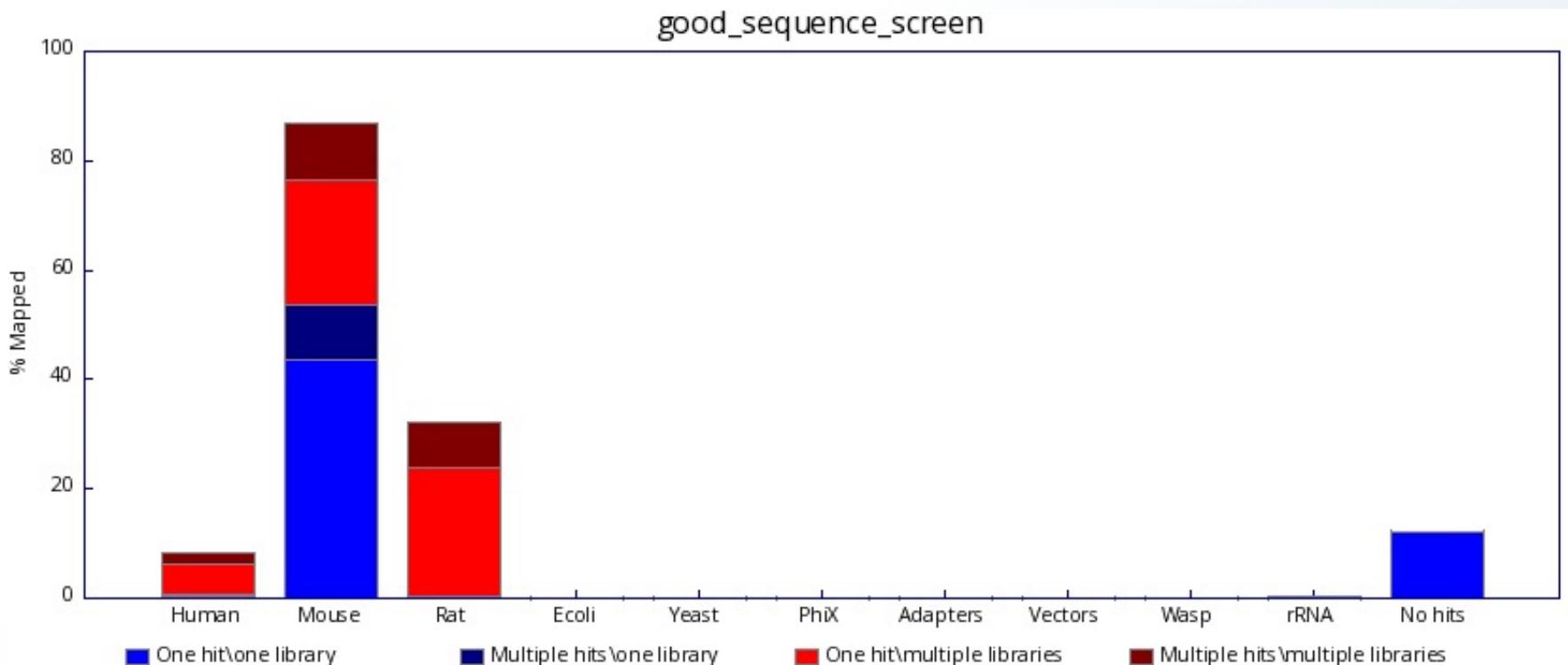
Use exact alignment of k -mers against a reference database to assign taxonomic labels to DNA sequences



Popular decontamination tools

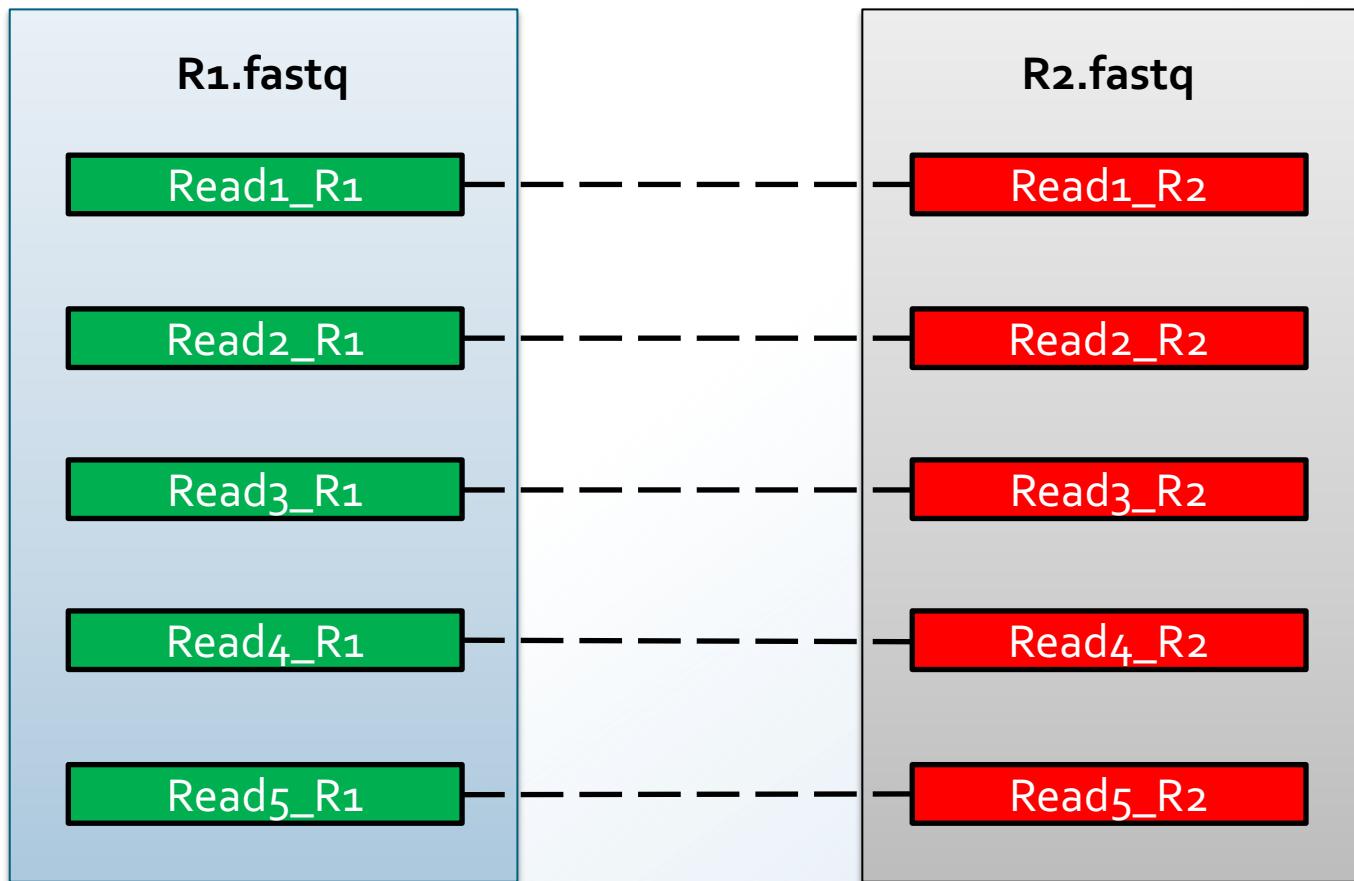
Removal of host contamination

- Fastq Screen
- DeconSeq
- CS-SCORE
- VecScreen
- Kraken



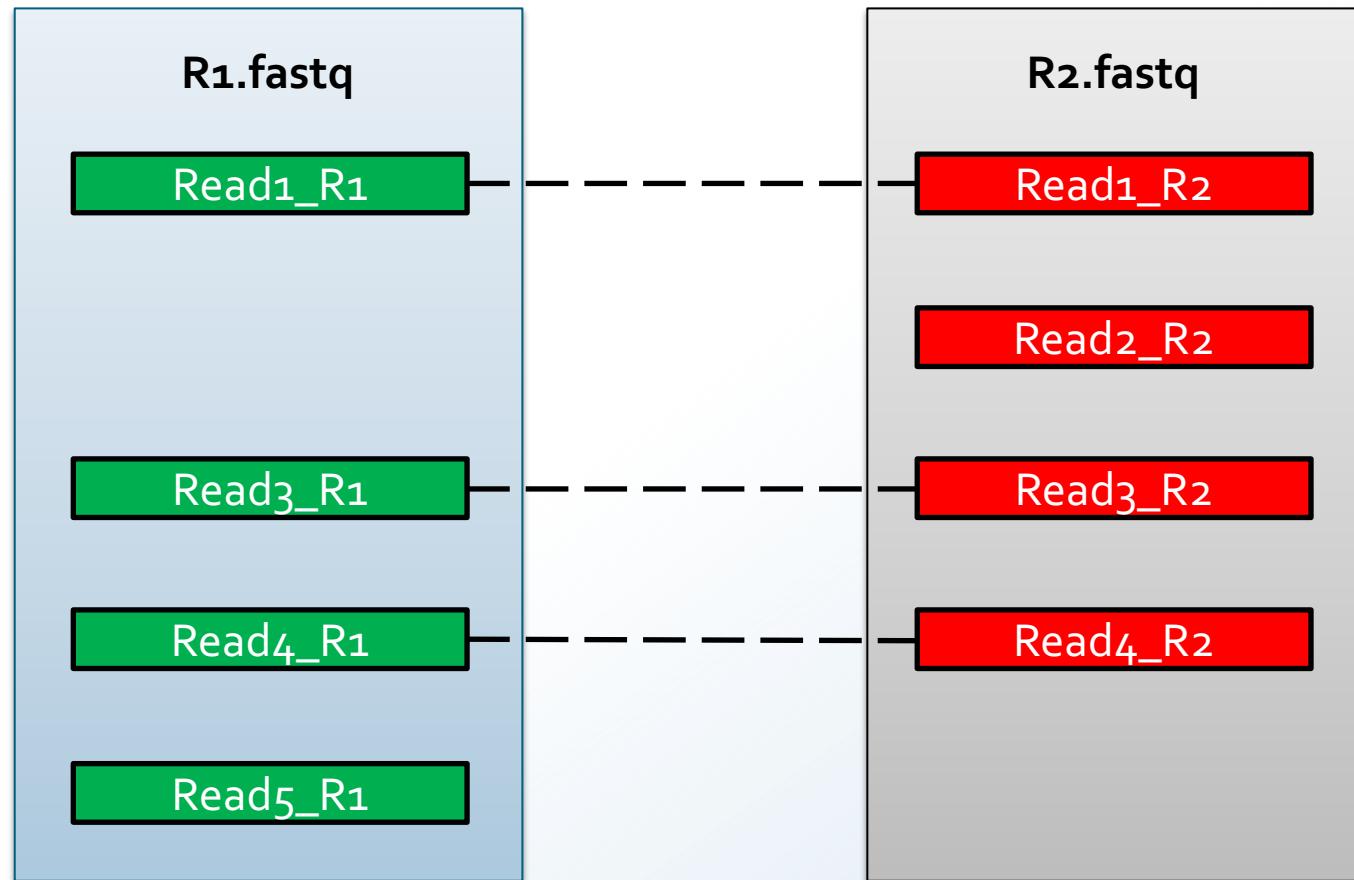
Synchronized FASTQ files are important for many tools

Some tools may remove one of the reads in a read-pair, hence the FASTQ files gets out of sync



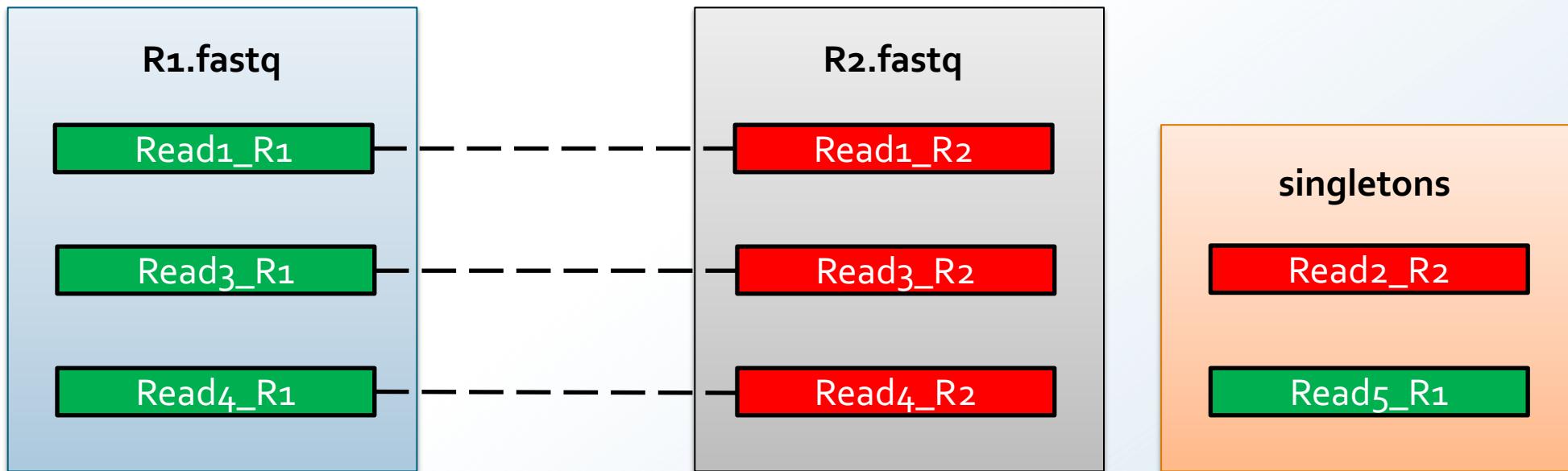
Synchronized FASTQ files are important for many tools

Some tools may remove one of the reads in a read-pair, hence the FASTQ files gets out of sync

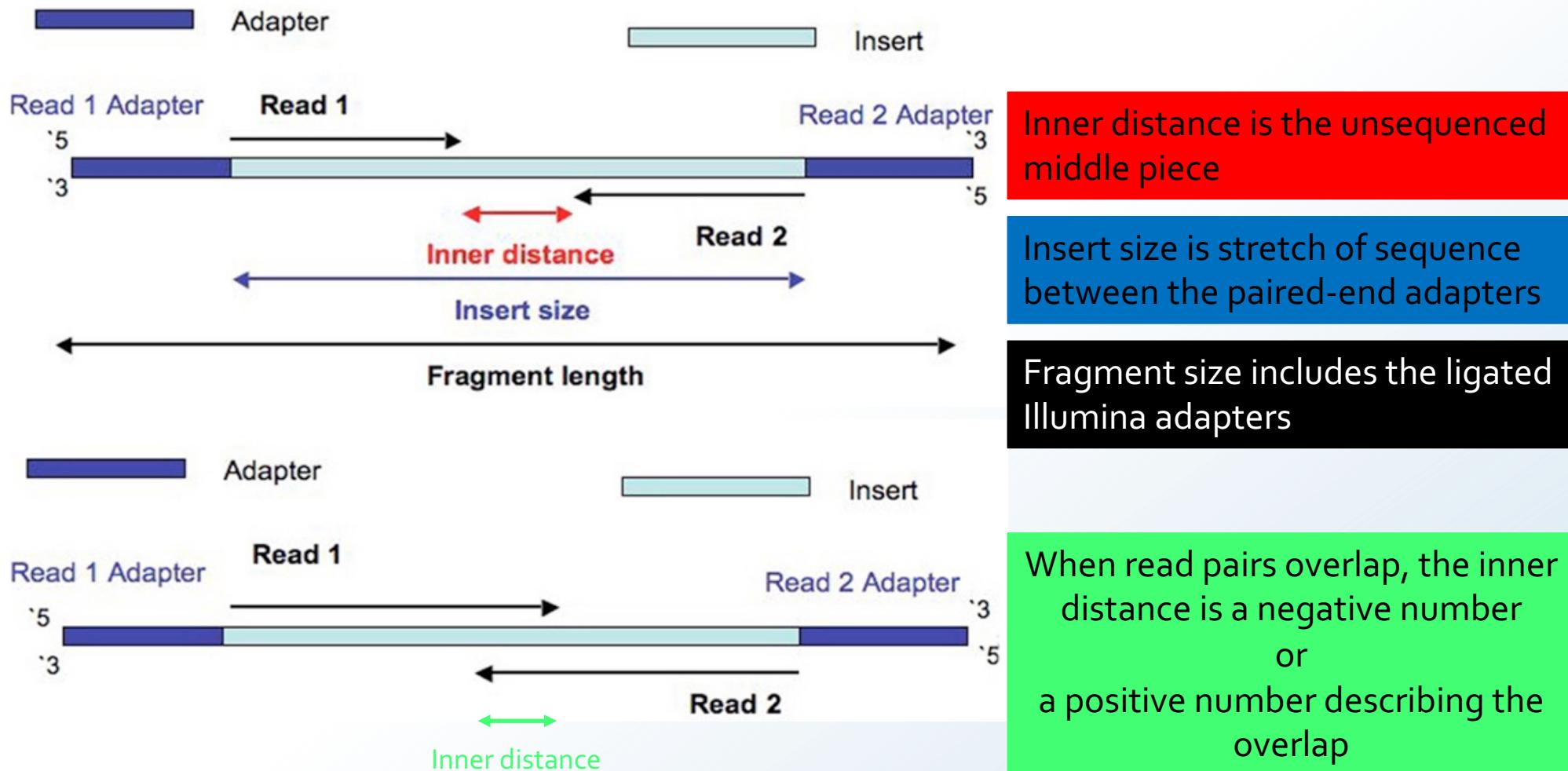


Synchronize FASTQ files with Repair

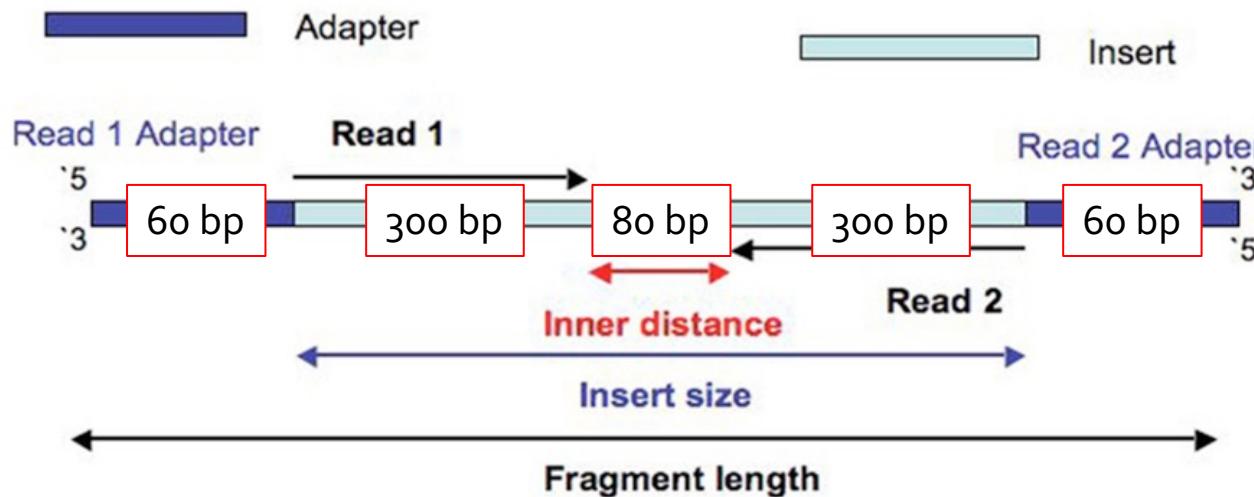
Repair pulls out the “singletons” and produce synchronized paired-end FASTQ files



DNA insert sizes and overlapping read pairs



DNA insert sizes and overlapping read pairs

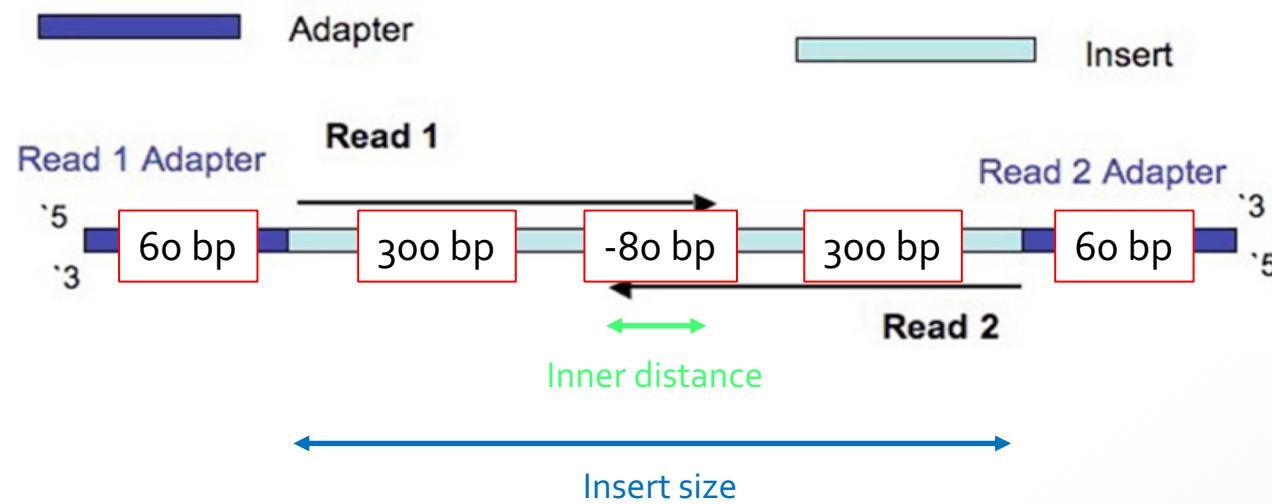


Inner distance = 80 bp

Insert size = $R_1 + R_2 + \text{inner distance} = 680 \text{ bp}$

Fragment size = $R_1 + R_2 + \text{inner distance} + \text{adapters (x2)} = 800 \text{ bp}$

Insert size in BBMerge is a bit confusing since the inner distance for overlapping reads is negative



Inner distance = BBMerge Insert size = -80

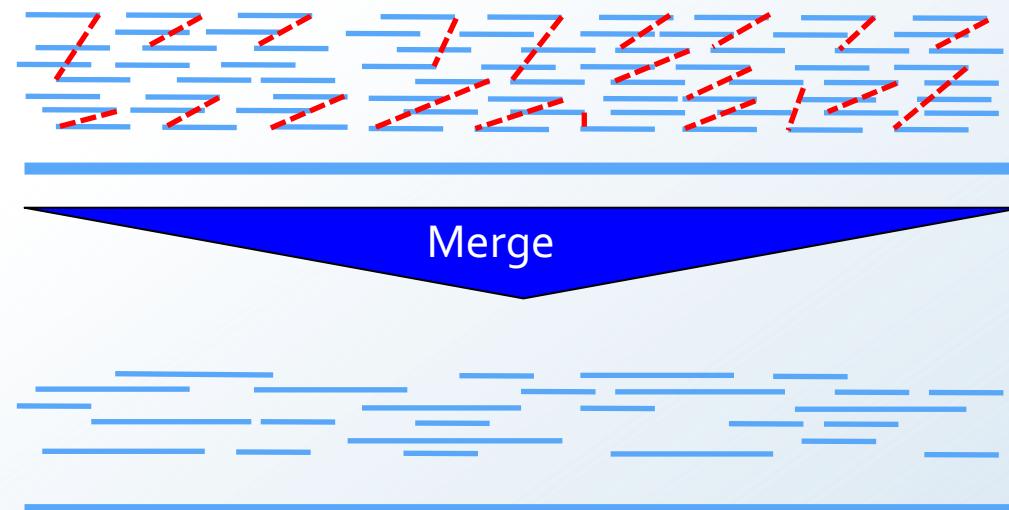
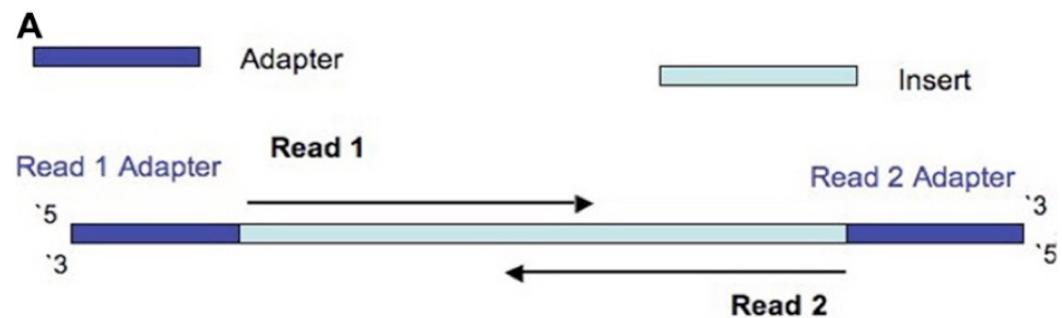
Insert size = R₁ + R₂ + BBMerge Insert size = 520 bp

Fragment size = R₁ + R₂ + inner distance + adapters (x2) = 640 bp

#InsertSize	Count (reads)
...	...
101	10248
102	10619
103	10397
104	10357
105	10218
106	10435
107	10288
108	9973
...	...

Generate longer reads by overlapping and merging read pairs before assembling a genome using BBmerge

Merging reads will reduce computational costs and improve the assembly



Generate longer reads by overlapping and merging read pairs before assembling a genome using BBmerge

Merging reads will improve the quality of the reads and generate longer reads
Longer reads allow the use of longer k-mers or fewer comparisons

Program	NA50 (bp)	Total Misassemblies	Indels/ 100 kbp	Genome Completeness (%)
Raw Data	60007	119	1.13	84.5
BBMerge	102577	127	0.84	84.88
BBMerge-REM	119328	117	0.81	85.18
BBMerge-RSEM	104441	115	0.84	84.88
COPE	89603	294	1.52	85.17
COPE-M3	98240	227	1.24	83.92
fastq-join	80672	183	1.17	84.74
FLASH	94846	282	1.41	85.20
leeHom	101992	290	1.1	84.91
PEAR	60937	660	1.46	84.28
Stitch	5623	20986	47.78	68.38
USEARCH	102156	131	0.88	84.77
XORRO	97403	158	1.08	84.85

Generate report and show to your boss😊

MultiQC is a reporting tool that parses summary statistics from results and log files generated by other bioinformatics tools

Parses relevant information from log files to a HTML report file

The screenshot displays the MultiQC reporting tool interface. On the left is a sidebar with a navigation menu:

- General Stats
- QUAST
- Assembly Statistics
- Number of Contigs
- FastQC
- Sequence Quality Histograms
- Per Sequence Quality Scores
- Per Base Sequence Content
- Per Sequence GC Content
- Per Base N Content
- Sequence Length Distribution
- Sequence Duplication Levels
- Overrepresented sequences
- Adapter Content

The main content area features the MultiQC logo at the top, followed by a brief description: "A modular tool to aggregate results from bioinformatics analyses across many samples in one report." Below this is a message: "Report generated on 2017-12-24, 14:22 based on data in: /Users/service/Box Sync/ELIXIR/Excellerate/MultiQC/reduced_data/multiqc". A welcome message "Welcome! Not sure where to start?" includes a "Watch a tutorial video" button. The central part of the interface shows a "General Statistics" table:

Sample Name	N50 (Kbp)	Length (Mbp)	% Dups
clean_megahit	3.4bp	29.6bp	
clean_metaspades	9.4bp	30.2bp	
sample_R1			0.0%
sample_R2			0.0%
sample_megahit	3.4bp	29.7bp	
sample_metaspades	4.0bp	29.7bp	
sample_trim_megahit	3.7bp	29.5bp	
sample_trim_metaspades	4.0bp	29.4bp	

To the right is a "MultiQC Toolbox" panel with the following sections:

- Rename Samples: Includes fields for "From" and "To", a "+" button, and a "Clear" button.
- Toolbox: Includes a "Click here for bulk input" link, a "Regex mode" switch, and a "Clear" button.
- A: A large letter "A" icon.