

RNA-seq Bioinformatics: Format and QC of short reads

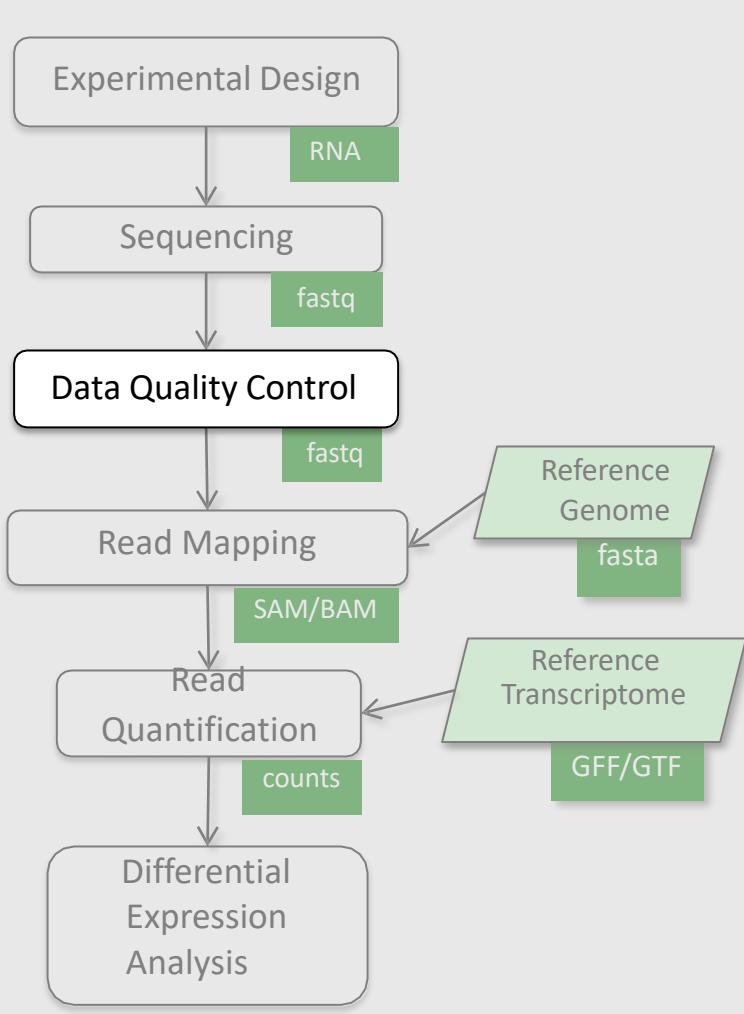
Falko Noé



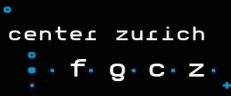
**University of
Zurich** UZH



Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich



- Short read file format
- Quality control of short reads

10
01
01
10101 10
01 01
101 10
10 10
01 01
.. f g c z ..
01 10
01 01

Properties of sequencing data

Technology	Read Length	Accuracy	Major error type
Sanger	400 to 900 bp	99.9%	Mismatch
Illumina	50 to 300 bp	98%	Mismatch
ONT	Limited only by the DNA molecules presented	65%-88% (traditionally) 99%-99.9% (current gen)	Indel
PacBio	10 kb to > 40 Kb	99.9999% circular consensus; 87% subread	Indel



10

01

101

010

01

10

101

10

10

10

10

10

10

10

01

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

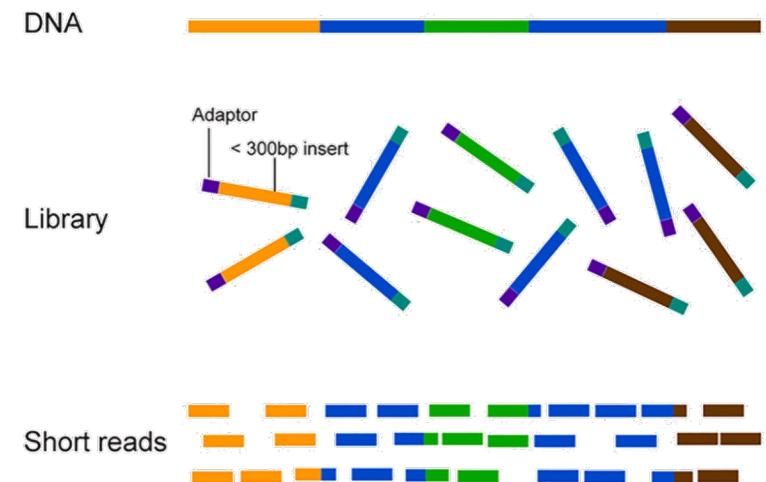
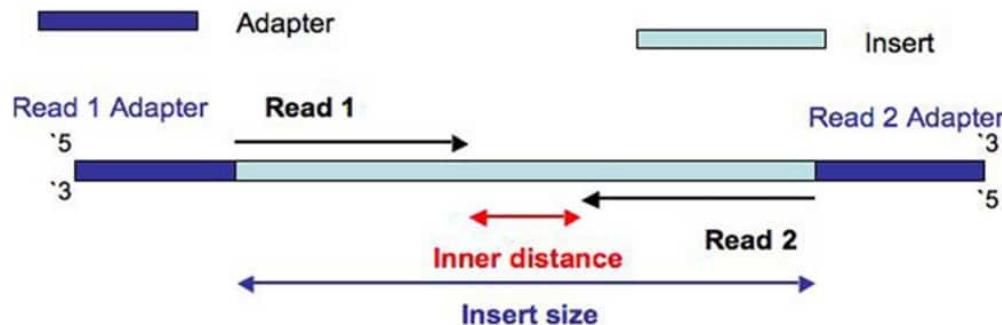
Sequence data – fastq files

- Data delivery (**fastq files**) from sequencing core
 - http://fgcz-gstore.uzh.ch/projects/pXXXX/NextSeq500_20160413_NS38_o2430/
- File names
 - Single end
 - 20160413.A-APC_mut_1_R1.fastq.gz
 - 20160413.A-APC_mut_2_R1.fastq.gz
 - Paired end
 - 20160413.A-APC_mut_1_R1.fastq.gz
 - 20160413.A-APC_mut_1_R2.fastq.gz
 - 20160413.A-APC_mut_2_R1.fastq.gz
 - 20160413.A-APC_mut_2_R2.fastq.gz

10
01
101

Sequence data – Illumina

- Single-end vs. paired-end reads



10
01
101

Fastq file format

- 4 lines per read

Machine ID

Read ID → @HWI-M00262:4:000000000-A0ABC:1:1:18376:2027 1:N:0:AGATC
Sequence → TTCAGAGAGAATGAATTGTACGTGCTTTTTTGT
+ → +
Quality score → =1:?7A7+?77+<<@AC<3<,33@A;<A?A=:4=
Phred+33

QC Filter flag
Y=bad
N=good
barcode
Read pair #

1. Header line for Read (starts with “@” and the sequence ID)
2. Sequence
3. Header line for Qualities (starts with “+”)
4. Quality score

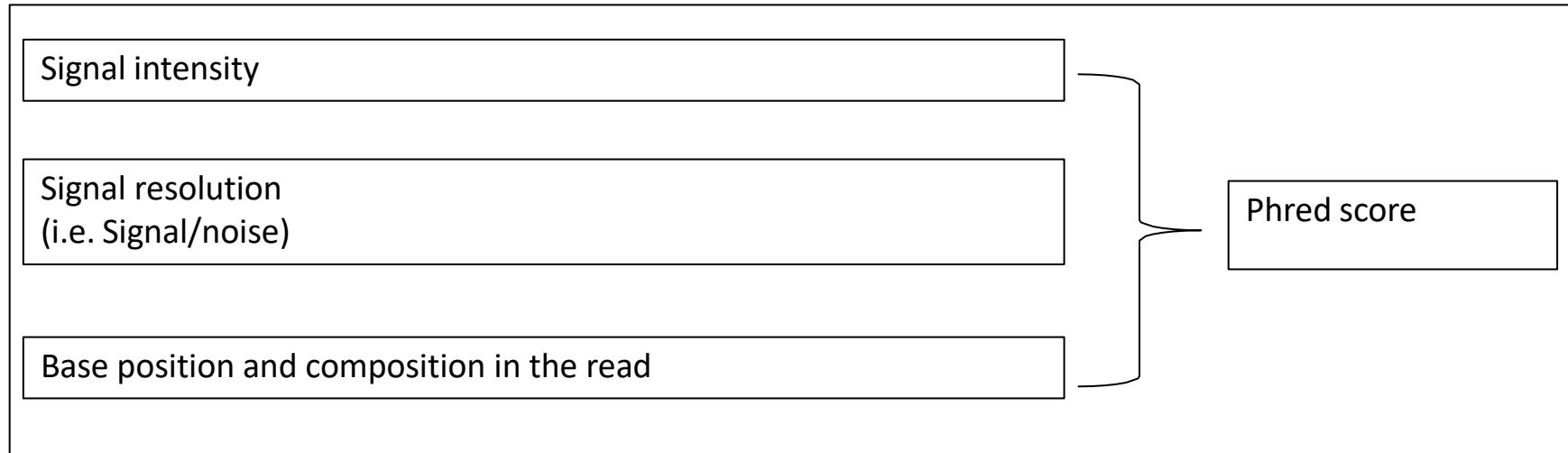
Phred scores

- Measure base calling accuracy
Accuracy of assigning **bases**
(nucleobases) to signal peaks
- P
error probability of a given base call
- $Q = -10 \log_{10} P$
- Assign to each base
- Range from 0-41 for Illumina sequencing

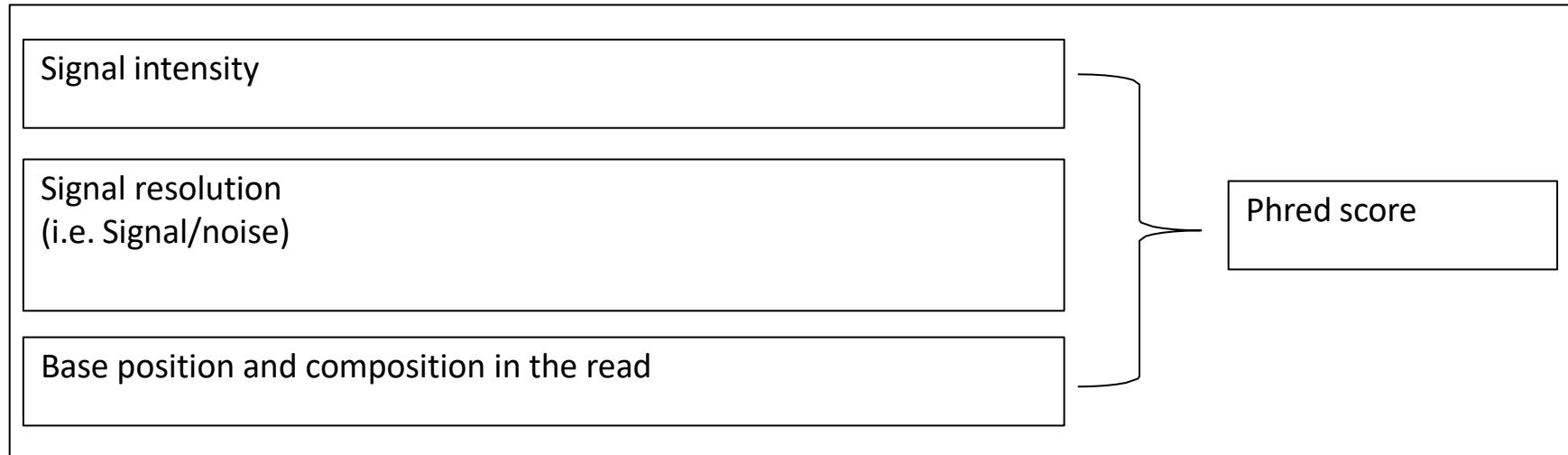


Phred Quality Score	Probability of Incorrect Base Call	Base Call Accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1,000	99.9%
40	1 in 10,000	99.99%
50	1 in 100,000	99.999%

How are Phred scores generated?

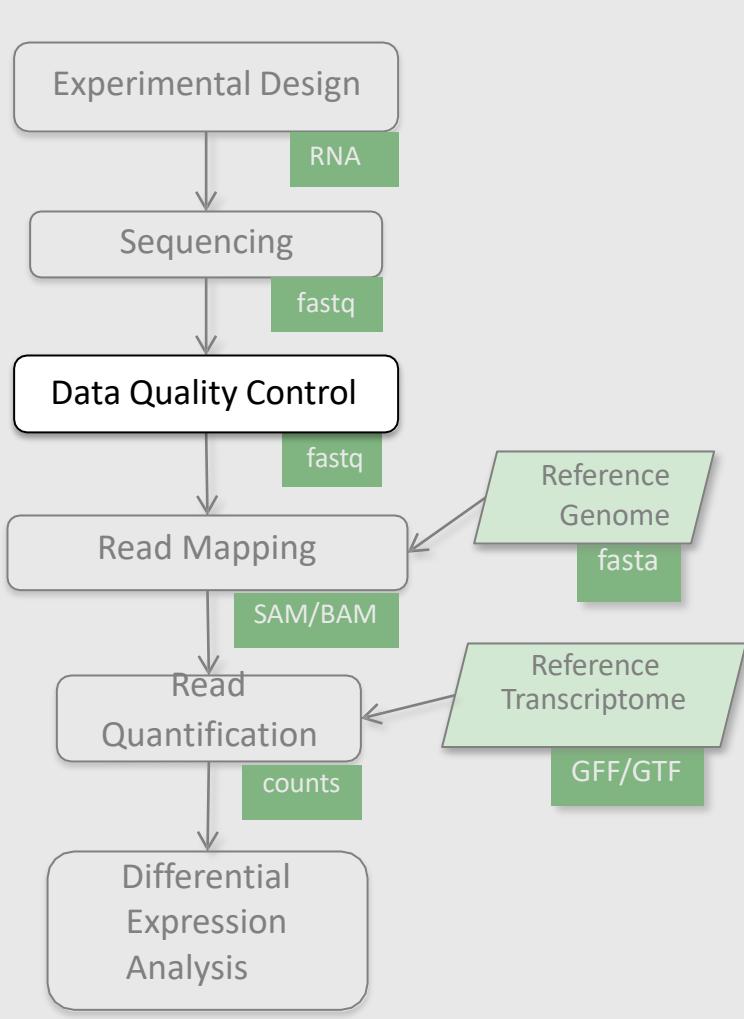


How are Phred scores generated?



- Parameters measured during sequencing **real samples**
 - Exact sequences are not known
 - **Scores assigned** by searching the look- up table

Million of reads



- Short read file format
- **Quality control of short reads**

Million of reads

FastQC

- Can be embedded in workflows as a data analysis module
 - Can also be run as an independent application with GUI using a laptop
 - <http://www.youtube.com/watch?v=bz93ReOv87Y>



Babraham Bioinformatics

About | People | Services | Projects | Training

FastQC

Function	A quality control tool for high throughput sequence data.
Language	Java
Requirements	A suitable Java Runtime Environment The Picard BAM/SAM Libraries (included in download)
Code Maturity	Stable. Mature code, but feedback is appreciated.
Code Released	Yes, under GPL v3 or later .
Initial Contact	Simon Andrews

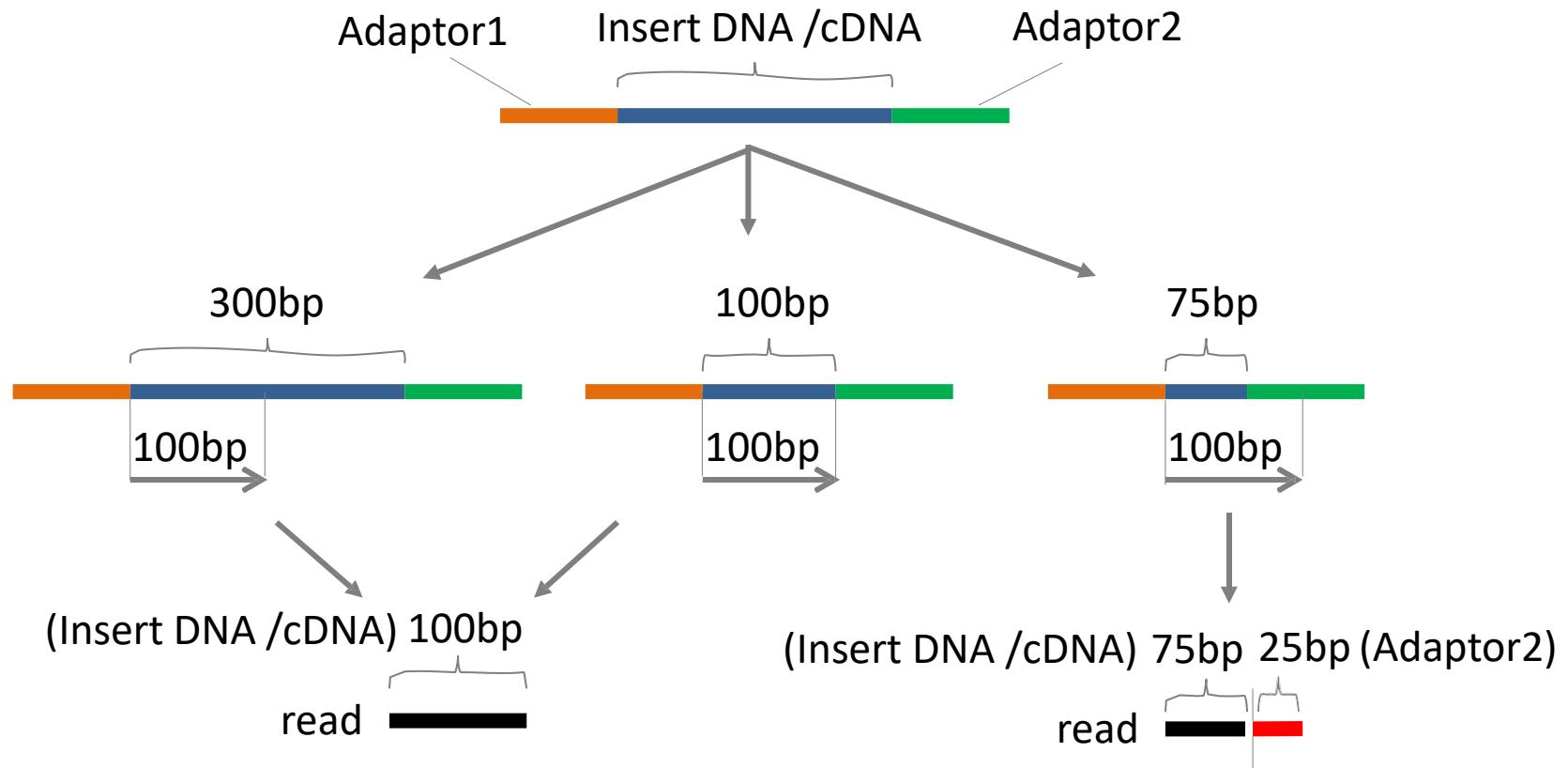
[Download Now](#)



Boxplots, histograms, heatmaps

- Fastqc uses boxplots and histograms for the analysis of phred scores and GC content
- Heatmaps can also be used. They are more visual and colors can be used to represent different information

Different scenarios



Bias and errors

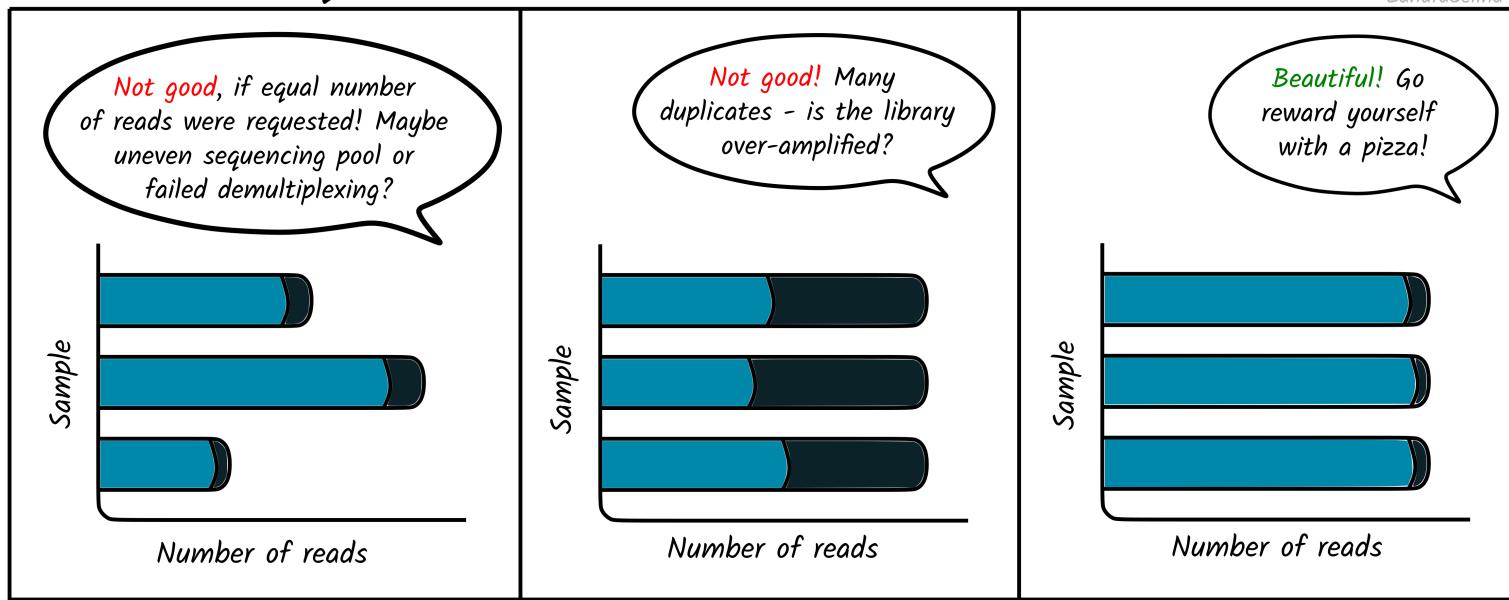
- Library construction could introduce bias
 - Fragmentation, ligation, amplification
 - GC bias
 - Over-amplification
 - Contamination
 - Sequencing errors
 - Chemical, optical, computational

FastQC

- Sequence counts and basic statistics
- Do I have enough read counts/coverage?

FASTQC - Sequence count

ZandraSelina

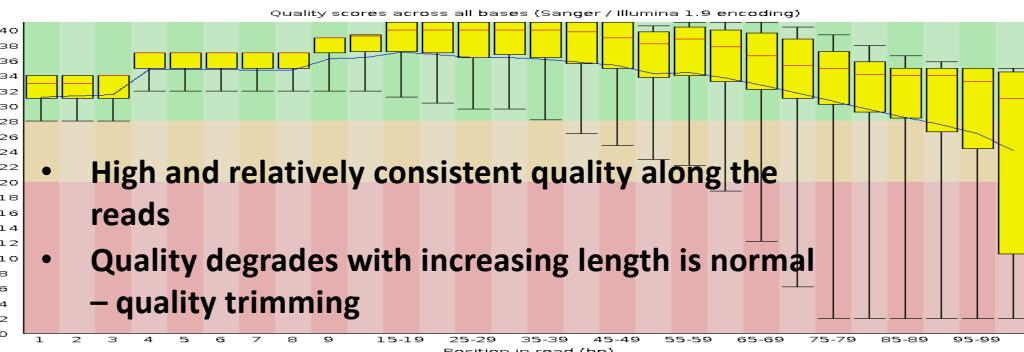


CC-BY 4.0

10
01
101010
01
101
10
010
01f
g
c
z01
10
01
10
01
10
01
10

Per base sequence quality - FastQC

- Range of quality values across all bases at each position

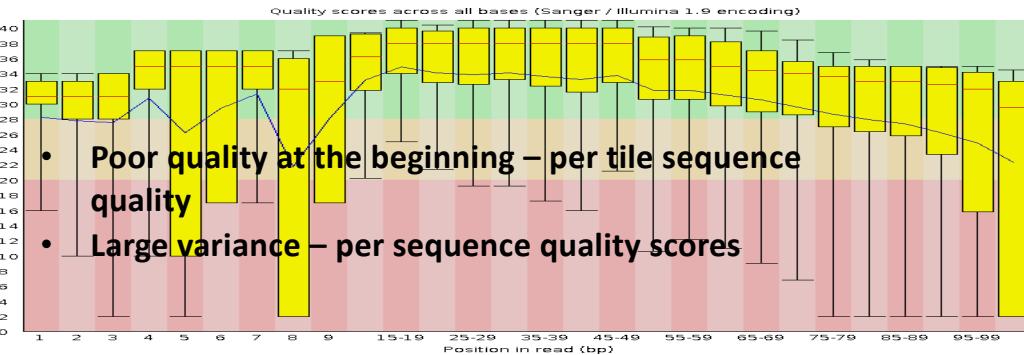


Green: >Q28, good

Orange: >Q20, reasonable

Red: <Q20, poor

Median > Q25

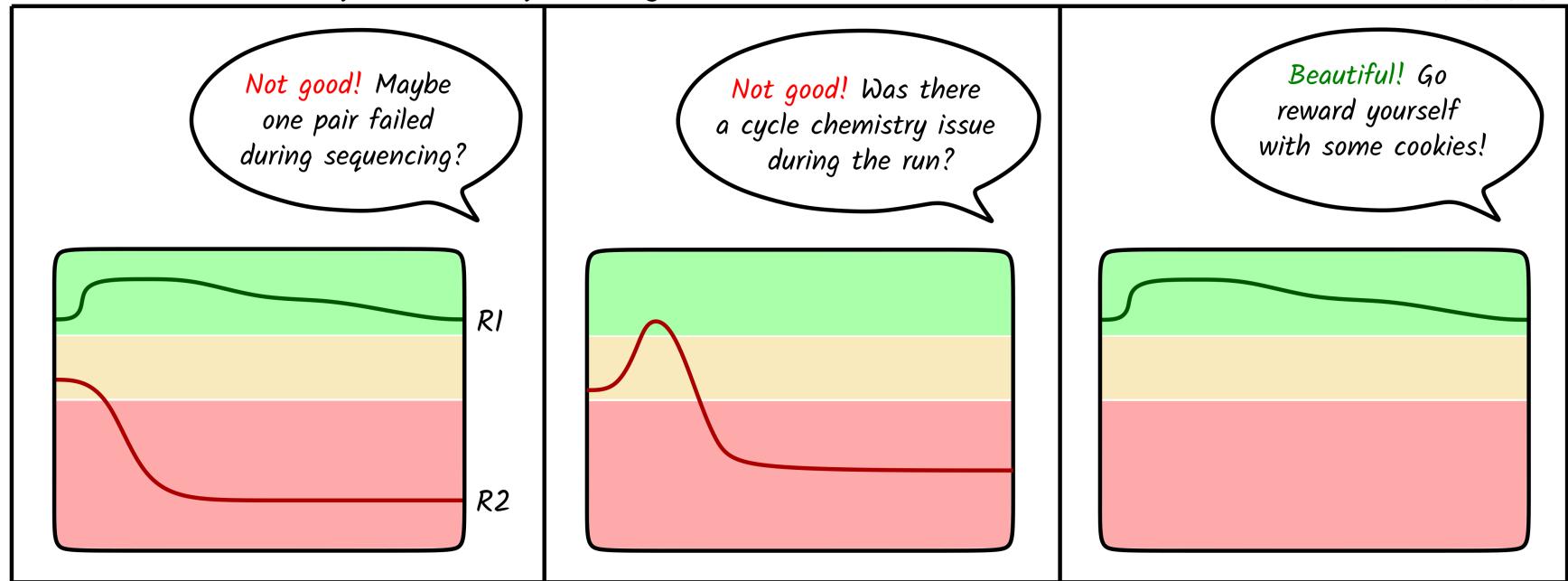


Median < Q20

Per base sequence quality - FastQC

FASTQC - Sequence quality

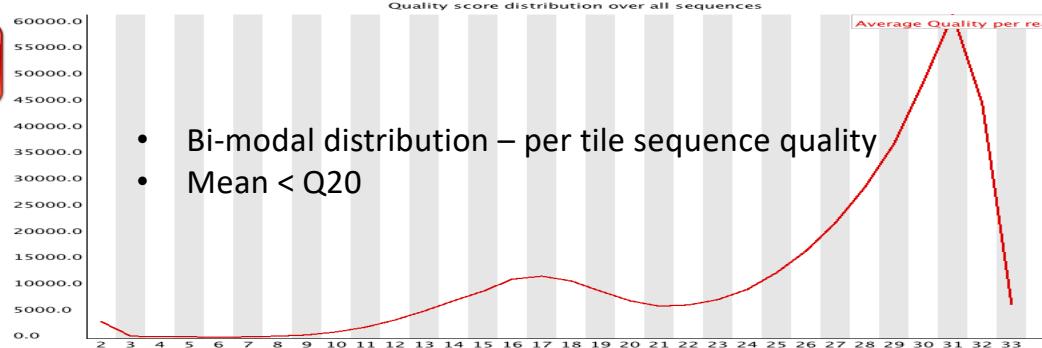
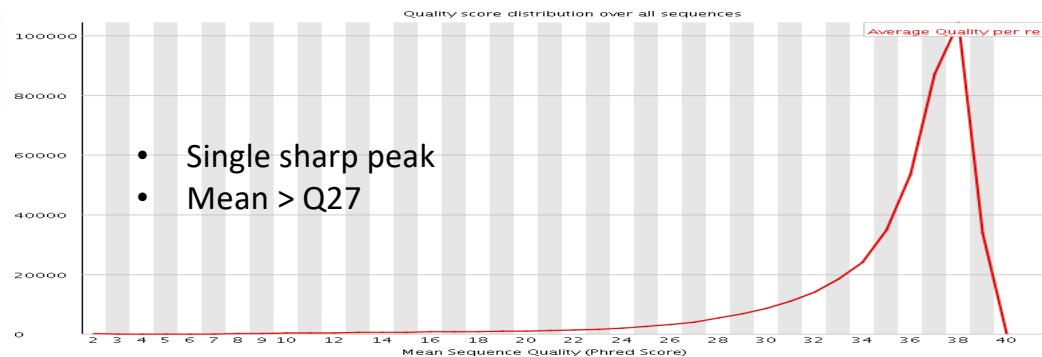
ZandraSelina



CC BY 4.0

Per sequence quality scores - FastQC

- Subset of sequences with universally low quality values

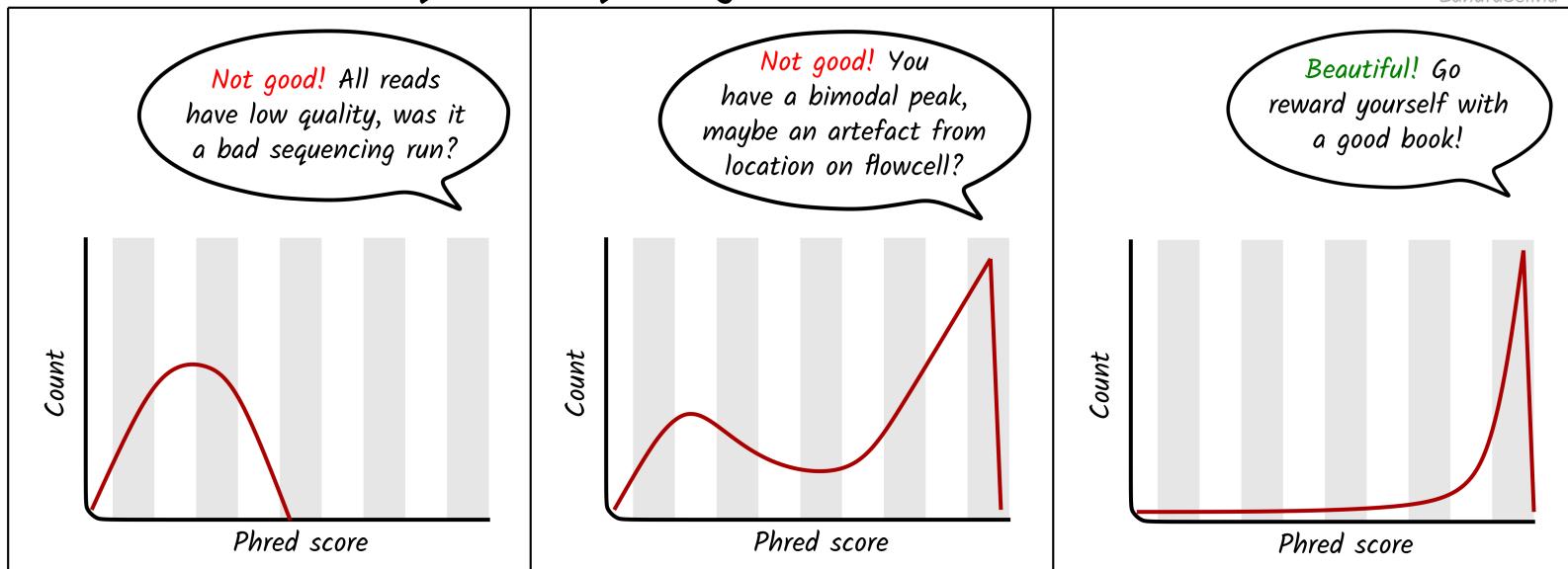


Per sequence quality scores - FastQC

- Subset of sequences with universally low quality values

FASTQC - Per sequence quality score

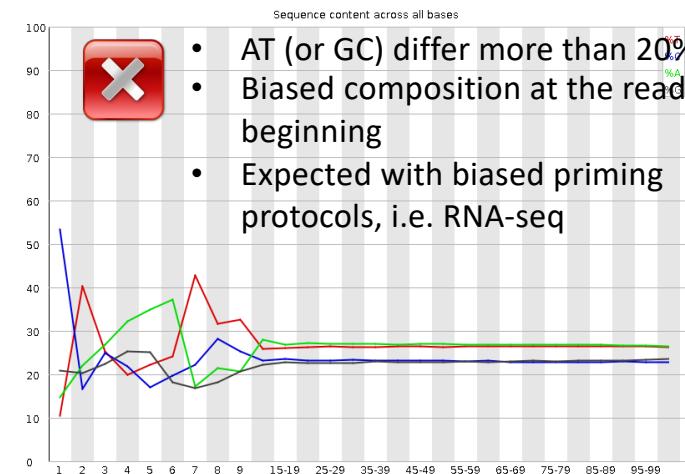
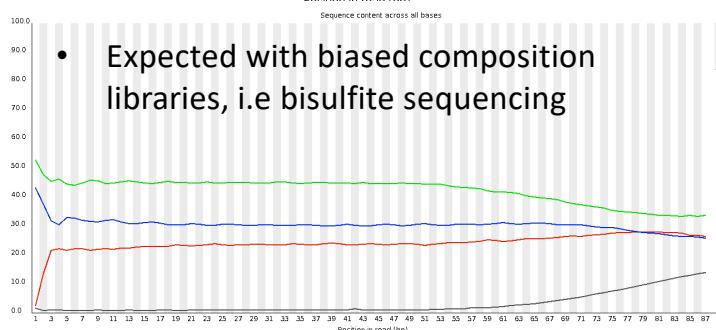
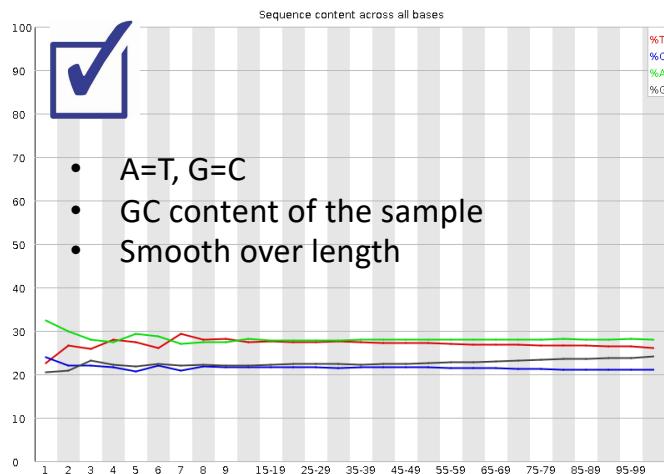
ZandraSelina



CC BY 4.0

Per base sequence content - FastQC

- The portion of A, T, G, and C at each position



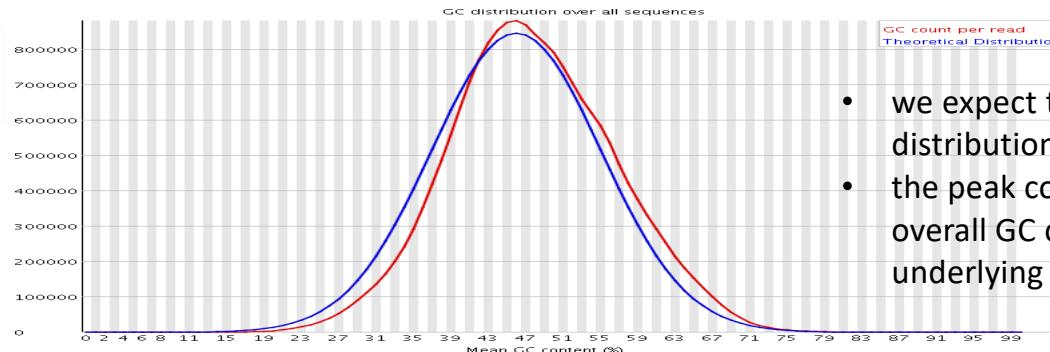
Biases in Illumina transcriptome sequencing caused by random hexamer priming

Kasper D. Hansen^{1,*}, Steven E. Brenner² and Sandrine Dudoit^{1,3}

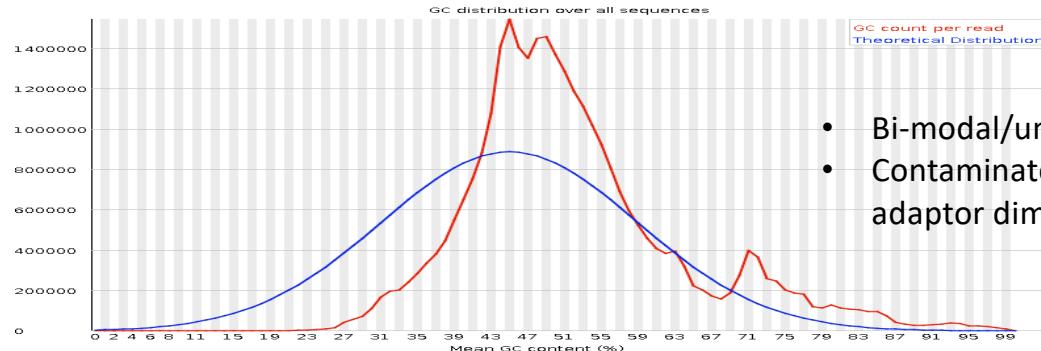
Treatment of DNA with bisulfite converts cytosine to uracil, but leaves methylated cytosine unaffected. Therefore, DNA that has been treated with bisulfite retains only methylated cytosines.

Per sequence GC content - FastQC

- Distribution of average GC in all reads



- we expect to see a roughly normal distribution of GC content
- the peak corresponds to the overall GC content of the underlying genome



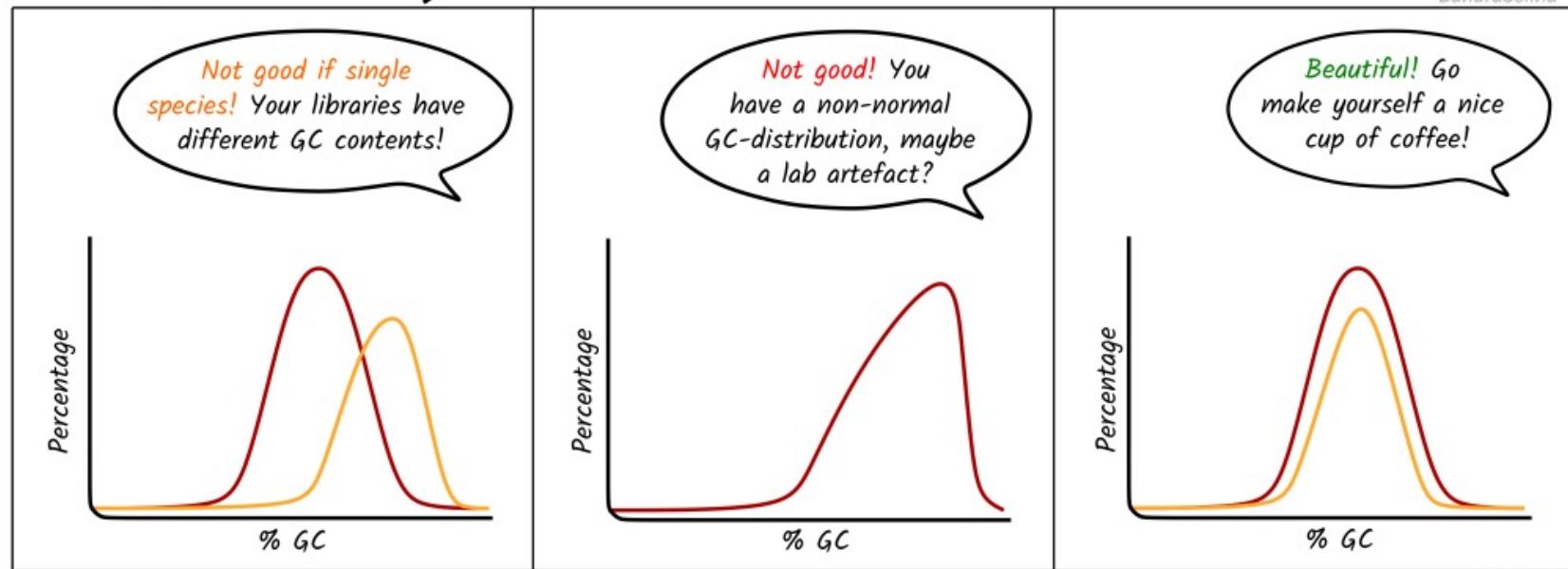
- Bi-modal/unusual distribution
- Contaminated/biased subset, i.e. adaptor dimmers, rRNA etc

Per sequence GC content - FastQC

- Distribution of average GC in all reads

FASTQC - Per sequence GC-content

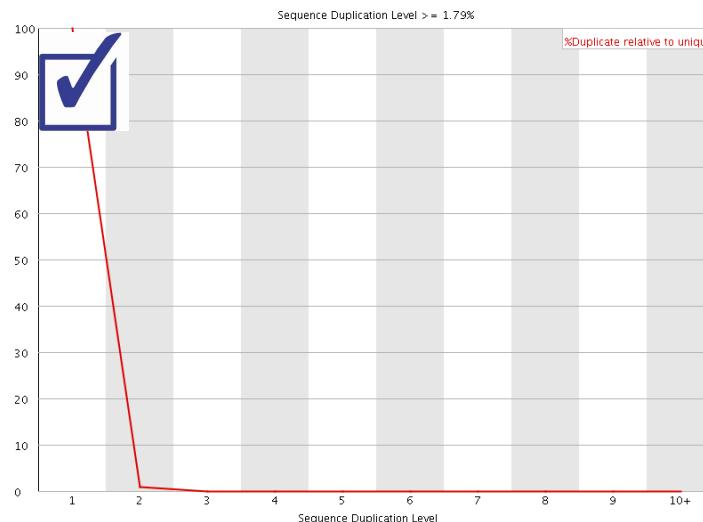
ZandraSelina



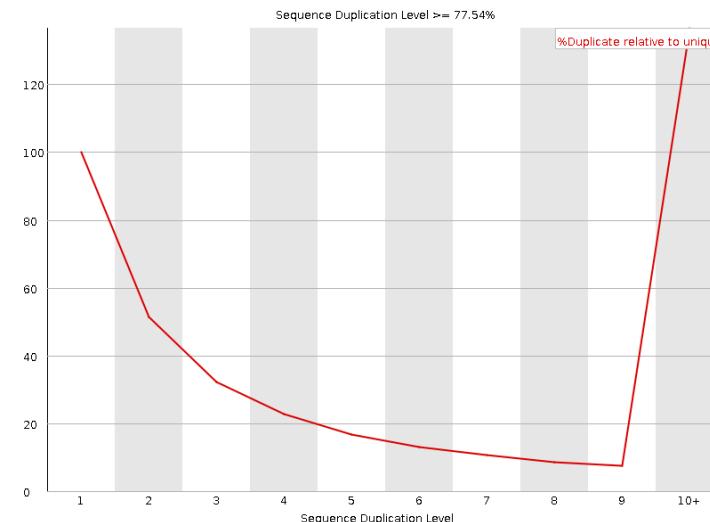
CC BY 4.0

Sequence duplication - FastQC

- Relative number of sequences with different degrees of duplication



- Low level duplication is expected for a diverse library



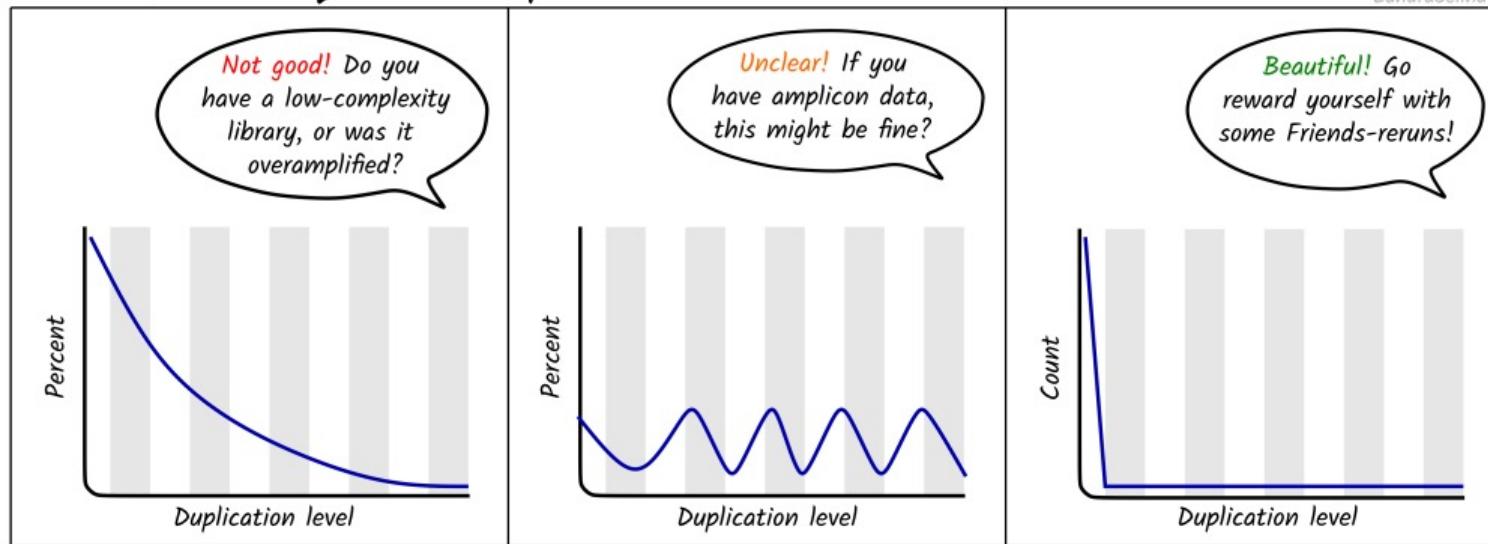
- High level duplication: enrichment bias, saturated sequencing depth
- Normal for RNA-seq (high sequencing depth) and ChIP-seq (enriched libraries)

Sequence duplication - FastQC

- Relative number of sequences with different degrees of duplication

FASTQC - Sequence duplication levels

ZandraSelina



CC BY 4.0

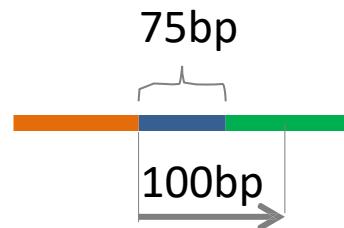
Overrepresented sequences - FastQC

- Sequences make up >0.1 % of the total
 - Compare those with a contamination database for finding contamination (i.e. adaptor dimmers)

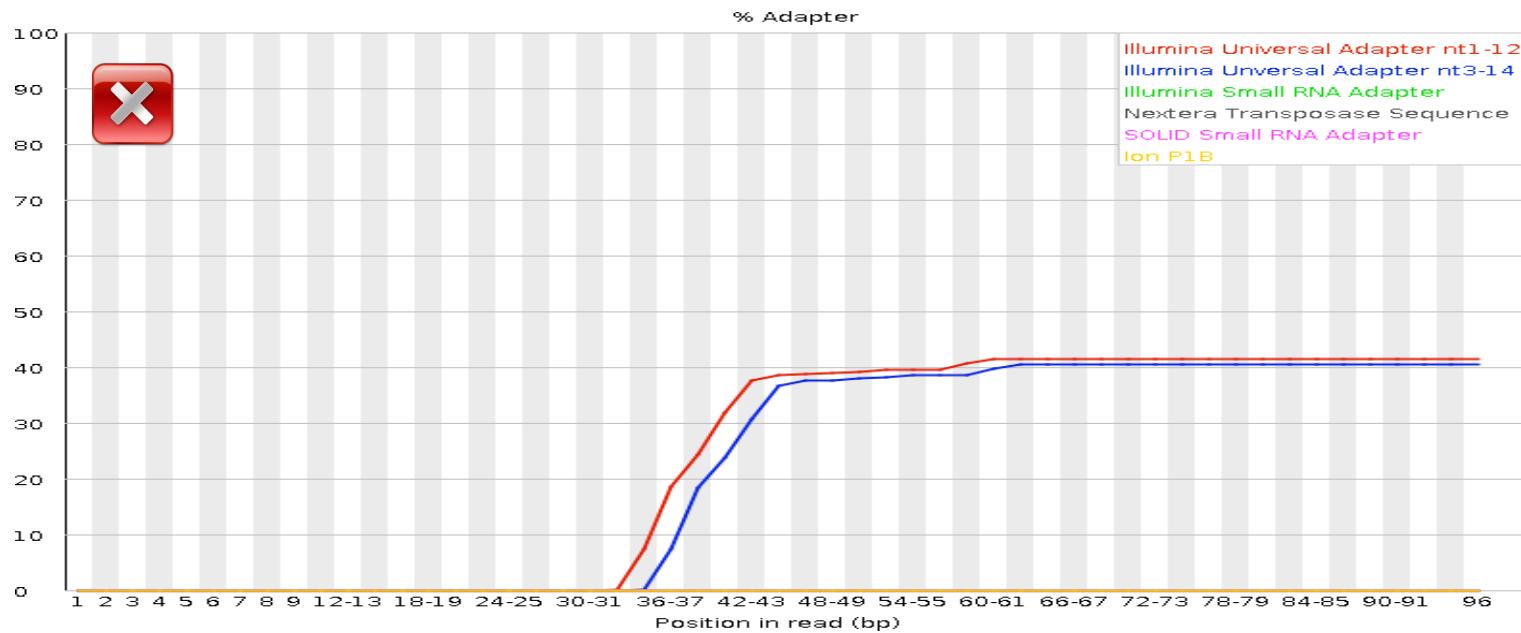
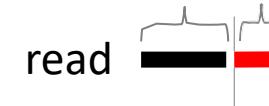
Overrepresented sequences				
Sequence	Count	Percentage	Possible Source	
GGAAGAGCACACGCTCTGAACTCCAGTCACCGATCATCTCGTATGCCGTC	75874	1.5613887498682963	TruSeq Adapter, Index 7	(100% over 50bp)
GGAAGAGCACACGCTCTGAACTCCAGTCACCGATGTATCTGTATGCCGTC	7636	0.15713900010536297	TruSeq Adapter, Index 2	(100% over 50bp)
GGAAGAGCACACGCTCTGAACTCCAGTCACACAGTGATCTGTATGCCGTC	7539	0.1551428656095248	TruSeq Adapter, Index 5	(100% over 50bp)
GGAAGAGCACACGCTCTGAACTCCAGTCACGCCAATATCTGTATGCCGTC	5117	0.10530123933199874	TruSeq Adapter, Index 6	(100% over 50bp)

- Can be normal and biologically meaningful
 - highly expressed transcripts
 - high copy number repeats
 - Less diverse library (amplicons)

Adapter Content - FastQC



(Insert DNA /cDNA) 75bp 25bp (Adaptor2)



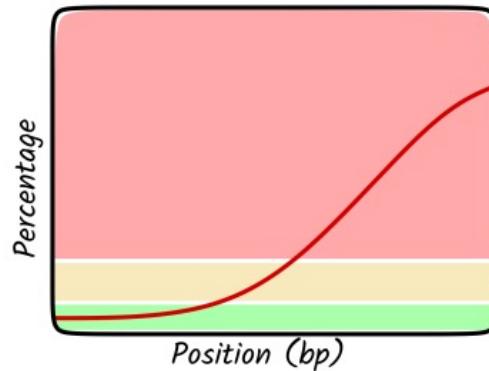
10
01
101101 1
010 0
0101 10010 01
101 10
10 00
01 01
f g c z10
01
101

Adapter Content - FastQC

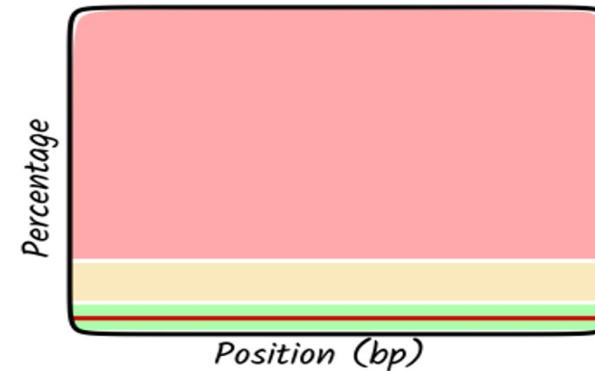
FASTQC - Adapter content

ZandraSelina

not good!
You have short inserts!



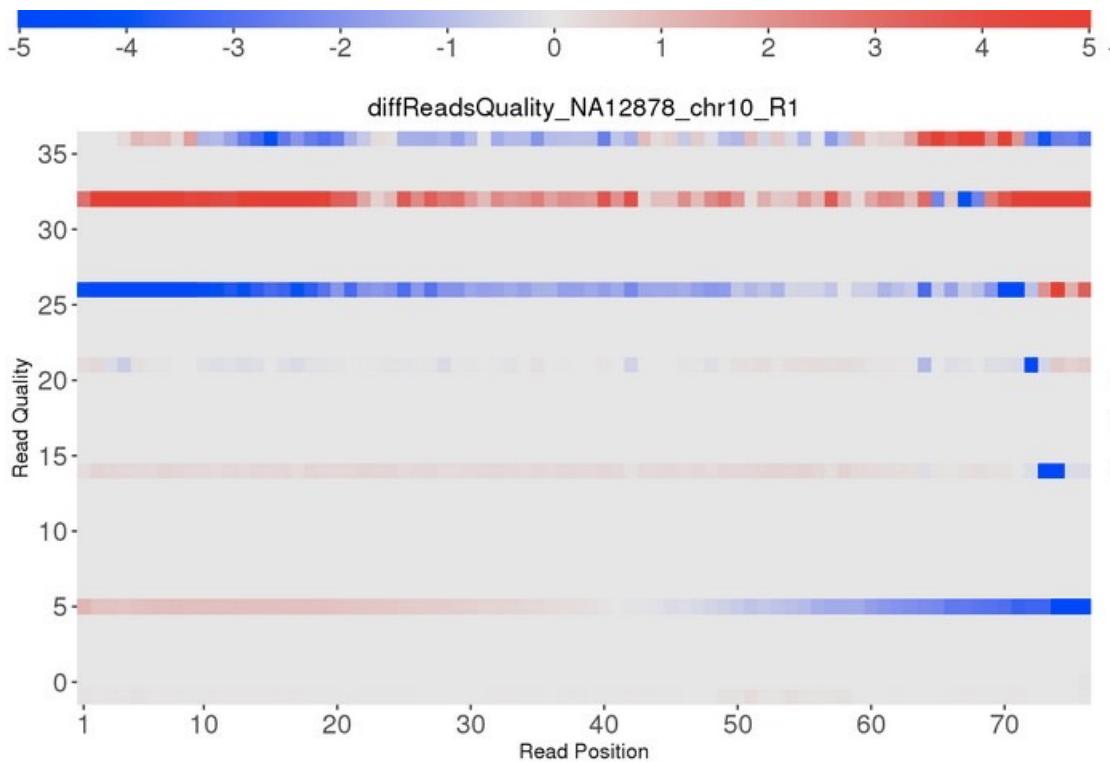
Beautiful
Go for a coffee
with your friends!



CC BY 4.0

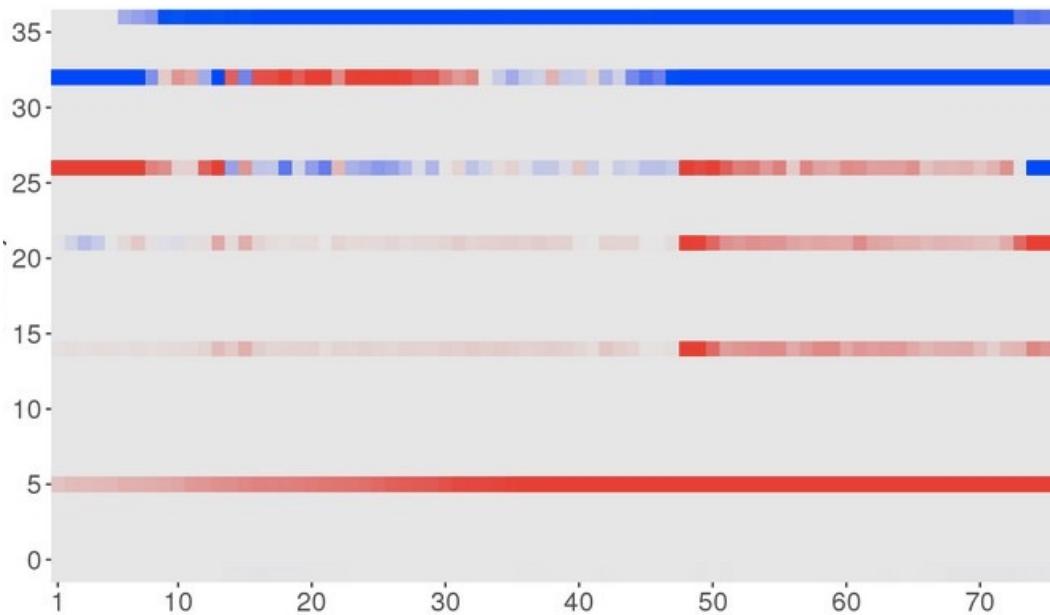
Comparative heatmap of per base phred score

- Reads in one sample /Average of reads in all samples
 - Sample with better quality than average



Comparative heatmap of per base phred score

- Sample with lower quality than average



10
01
101101 1
010 0
101 10
010 10

MultiQC

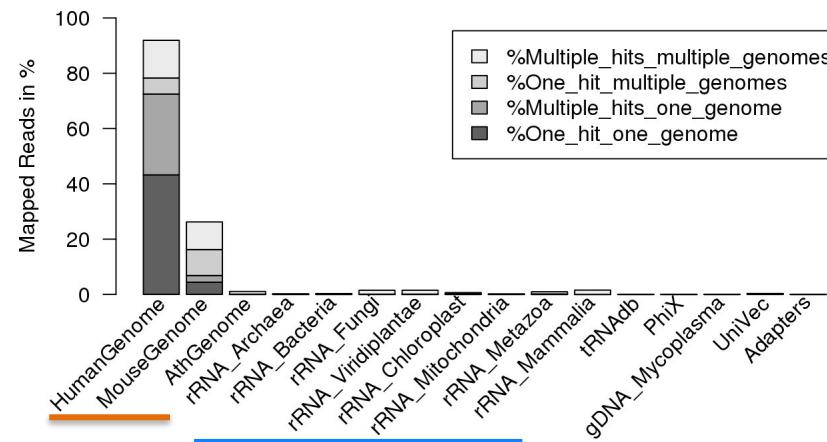
- <https://seqera.io/multiqc/>
- Aggregates bioinformatics analyses results from different samples and sources, including FastQC

The screenshot shows the MultiQC web interface. On the left, there's a sidebar with links to various analysis modules: General Stats, featureCounts, STAR, Summary Statistics, Alignment Scores, Cutadapt, Filtered Reads, Trimmed Sequence Lengths (3'), FastQC: trimmed, Sequence Counts, 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, and Overrepresented sequences by sample. A blue button labeled "Read documentation >" is also visible. The main content area has two tabs: "General Statistics" (selected) and "Feature Counts". The "General Statistics" tab displays a table of sample metrics. The table includes columns for Sample Name, Assigned, Aligned, Uniq aligned, Trimmed bases, Dups, GC, and Seqs. The table lists several samples, such as SRR3192396, SRR3192396 R1, SRR3192396 R2, SRR3192397, SRR3192398, SRR3192399, SRR3192400, SRR3192401, SRR3192657, and SRR3192658. Each row provides detailed statistics for that sample, such as 67.5% assigned reads and 97.2% aligned reads for SRR3192396. The "Feature Counts" tab is partially visible on the right.

Sample Name	Assigned	Aligned	Uniq aligned	Trimmed bases	Dups	GC	Seqs
SRR3192396	67.5%	97.2%	93.7%	3.2%	75.9%	50.5%	208.8 M
SRR3192396 R1				2.5%	72.8%	50.0%	104.4 M
SRR3192396 R2				4.0%	78.9%	51.0%	104.4 M
SRR3192397	66.8%	98.1%	94.7%	2.8%	74.7%	48.5%	183.9 M
SRR3192398	50.9%	95.9%	88.2%	5.0%	57.6%	47.0%	133.1 M
SRR3192399	52.3%	96.3%	88.2%	5.0%	59.3%	47.0%	148.7 M
SRR3192400	70.3%	83.1%	77.3%	6.2%	75.1%	45.0%	189.9 M
SRR3192401	71.2%	82.2%	76.4%	6.2%	75.9%	45.0%	190.5 M
SRR3192657	73.1%	94.1%	91.2%	2.5%	81.4%	50.5%	186.3 M
SRR3192658	71.2%	92.5%	89.7%	2.8%	81.4%	52.0%	194.1 M

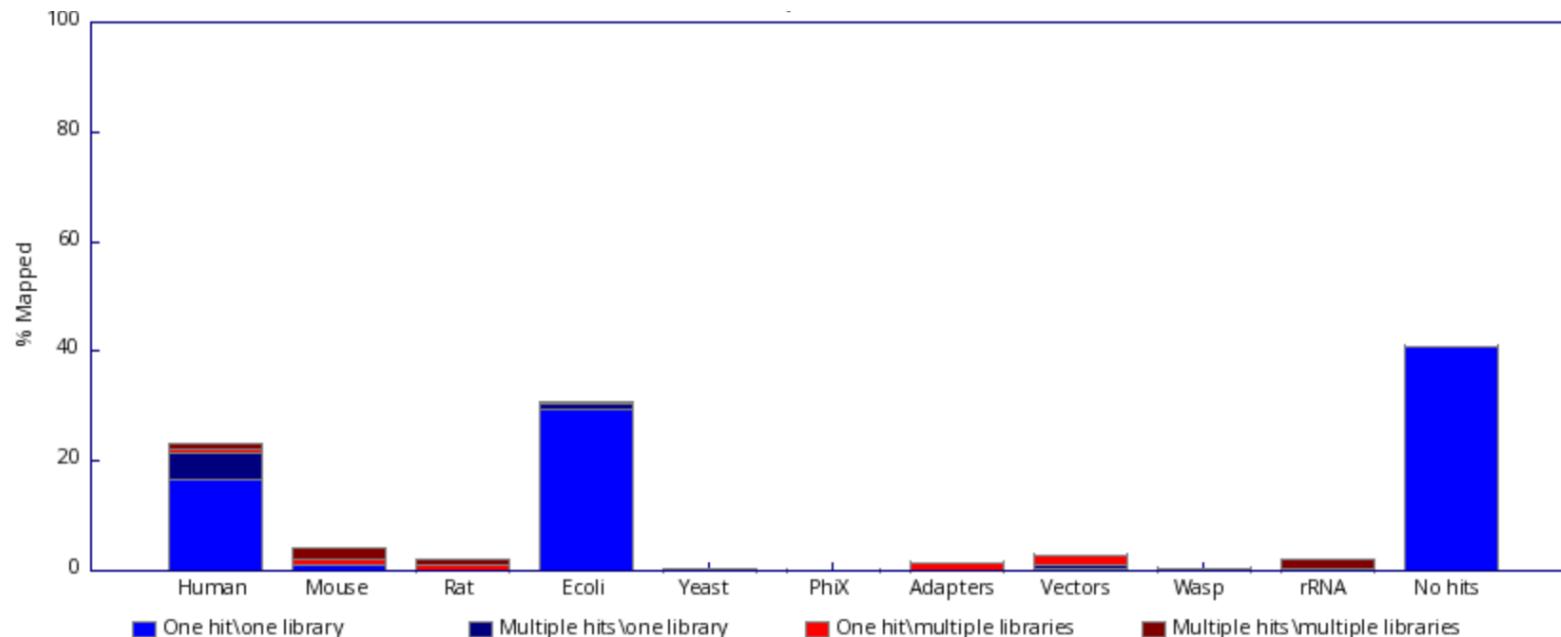
Fastqscreen – check for sample contamination

- Compare sequencing reads to databases of known sequences
 - Report top matches
 - In clonal sample, uniquely mapped reads should come from only a single organism



Frequently sequenced organisms rRNA genes (Silva) Frequent contamination

Contamination Check

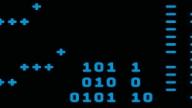




10

01

101



01

01

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

10

Sequencing data pre-processing tasks

- Trimming: remove bases from read end(s)
 - Adaptor sequence
 - Low quality bases
- Filtering: remove reads
 - Low quality reads
 - Contaminating sequences
 - Low complexity reads (repeats)
 - Short (<20bp) reads – they slow down mapping software

10
01
0101010 01
101 10
10 10
010 01
01 10
.. f g c z ..
01 10
01 10

Tools for pre-processing sequencing data

- PRINSEQ
 - <http://prinseq.sourceforge.net/>
 - Quality/hard trimming, quality filtering, reformat, ...
- Trimmomatic
 - <http://www.usadellab.org/cms/?page=trimmomatic>
 - Adaptor trimming, quality trimming &filtering, ...
- fastp*
 - <https://github.com/OpenGene/fastp>
 - Adaptor trimming, quality trimming, deduplication, parallel processing
- FlexBar (FAR)
 - <http://sourceforge.net/projects/theflexibleadap/>
 - Flexible barcode detection and adapter removal

Summary

- Always generate quality plots for all data sets
 - Interpretation of the plots need knowledge about the samples and libraries
 - Trim and/or filter data if needed
 - Always trim and filter away low quality data for variant analysis

10
01
101010
101
1001
10
0101
10
0101
10
0101
10
0101
10
01

Questions ?!



10

01

101

010

01

101

10

0

01

0

1

01

1

0

1

01

1

0

1

01

1

0

1

01

1

0

1

01

1

0

1

01

1

0

1

01

1

0

1

01

1

Raw QC Clicker Questions