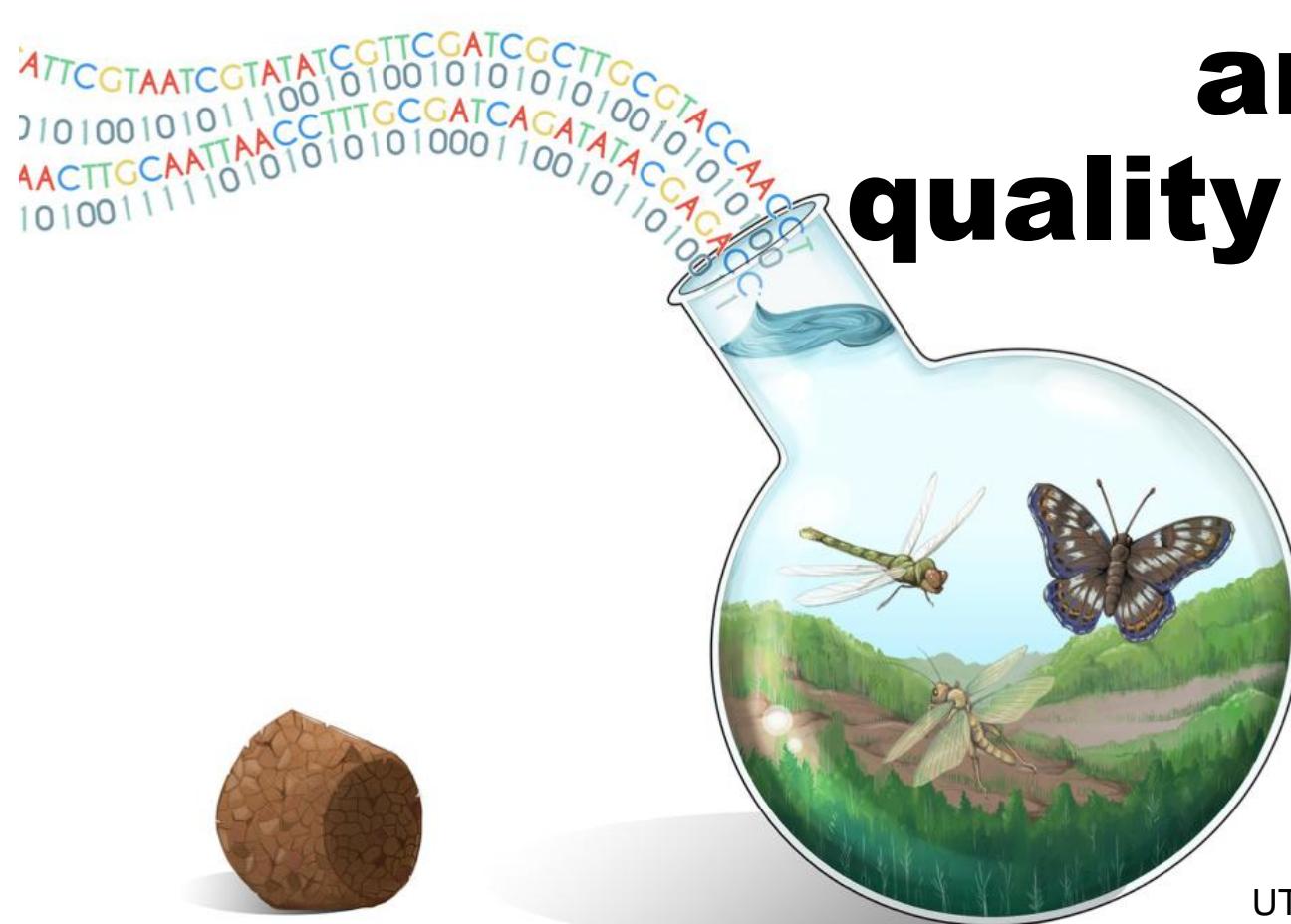




Sequencing data formats and quality control



Vladimir Mikryukov

UT International Summer University,
Tartu, August 01-05 2022

FASTA

The diagram illustrates the structure of FASTA sequence files. It shows three entries, each consisting of a header (prefixed with '>SeqID') and a sequence of DNA bases. The first two headers are circled in red, and blue brackets labeled 'Header' point to them. The third header is also circled in red, and a blue bracket labeled 'Sequence' points to the entire entry.

>SeqID1
GGGTGGACGGTTATCCACCATTGGCAGCGGGTATTGCTCA

>SeqID2
GGGTGGACGGTTATCCACCATTGGCAGCGGGTATTGCTCA

>SeqID3
GGGTGGACGGTTATCCACCATTGGCAGCGGGTATTGCTCA
GTGAGTCATCGAATCTTGAACGCACATTGCGCCCTCTGGT

>SeqID1 ; **tax=d:Eukaryota, p:Arthropoda, c:Entognatha**
GGGTGGACGGTTATCCACCATTGGCAGCGGGTATTGCTCA

>SeqID2 ; **tax=d:Eukaryota, p:Arthropoda, c:Insecta**
GGGTGGACGGTTATCCACCATTGGCAGCGGGTATTGCTCA

>SeqID3 ; **tax=d:Eukaryota, p:Arthropoda, c:Entognatha**
GGGTGGACGGTTATCCACCATTGGCAGCGGGTATTGCTCA

FASTQ

Label

@M03023:624:JRJBJ:18525:1633

Sequence

GCGAGTGGTGCAAAGCTTAAACTAAGGAAACGGAAGGGCACCACAGGA

+

AACGGGGFFGGG:::?:?FGACFGGG@?GGFGGGFGGGDGGGCEBFG8<F

Base = T

Quality = F

Phred score = 37

*Quality scores
(as ASCII characters)*

FASTQ

FASTQ (paired-end data)

* _R1 .fastq

```
@A00459:120:HHFYJDRXY:1:2101:1515:1016 1:N:0:AGTCAGGT+GTGGTTAC  
TNTACGGCGTTCTTCATCGATGGGAGAACCAAGAGATCCGTTCTCAAAGTTGTATTTAGTTCTGCCCG  
+  
F#FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF  
@A00459:120:HHFYJDRXY:1:2101:1642:1016 1:N:0:AGTCAGGT+GTGGTTAC  
ANGCGAACTAATCGACCAGCGGAGGGATCATTAATGAATAAACTCGGTGGATTGTTGCTGGCTCTAGGA  
+  
F#FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF:FFFF  
@A00459:120:HHFYJDRXY:1:2101:1325:1031 1:N:0:AGTCAGGT+GTGGTTAC  
TNTGCTCGTCTTCATCGATGGGAGAGCCAAGAGATCGGTGCTGAAAGTTGTATAGTTAATGACTG  
+  
F#FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF:FFFFFFFFFF:FFFF  
@A00459:120:HHFYJDRXY:1:2101:1362:1031 1:N:0:AGTCAGGT+GTGGTTAC  
ANCAGCTCAGATCGACCAGCGGAGGGATCATTATTGAATAAAACCTGGCGTGGTAGCTGGCTCTCG  
+  
F#FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF:FFFFFFFFFF:FFFF  
@A00459:120:HHFYJDRXY:1:2101:2374:1031 1:N:0:GGTCAGGT+GTGGTTAC  
ANGCATCGCACTGACCAGCGGAGGGATCATTAATGAATAAACTTGGTAGATTGTTGCTGGCTCTAGGA  
+  
F#FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF:FFFF
```

* _R2 .fastq

```
@A00459:120:HHFYJDRXY:1:2101:1515:1016 2:N:0:AGTCAGGT+GTGGTTAC  
ACACAGTCCTGACGACCAGCGGAGGGATCATTAGTCATACAACCGGGGAATCCACTCTGTGGGCCCAACC  
+  
FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF  
@A00459:120:HHFYJDRXY:1:2101:1642:1016 2:N:0:AGTCAGGT+GTGGTTAC  
TTCGCTACGTTCTTCATCGATGCGAGAGCCAAGAGATCGGTGTTGAAAGTTGTATAAGTTAAAGCCTA  
+  
FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF:FFFF  
@A00459:120:HHFYJDRXY:1:2101:1325:1031 2:N:0:AGTCAGGT+GTGGTTAC  
ACGTAACCACGTCGACCAGCGGAGGGATCATTATTGAATACGAATCGGTCTGATGCTGGCCCTCACCGG  
+  
:FFFFFFFFFFFFFFFFFFFFFFFFFFFFFF:FFFFFFFFFF:FFFFFFFFFF:FFFFFFFFFF:FFFF  
@A00459:120:HHFYJDRXY:1:2101:1362:1031 2:N:0:AGTCAGGT+GTGGTTAC  
TTCGCTACGTTCTTCATCGATGCGAGAGCCAAGAGATCGGTGCTGAAAGTTGTATAGTTAGGCAC  
+  
FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF:FFFFFFFFFF  
@A00459:120:HHFYJDRXY:1:2101:2374:1031 2:N:0:GGTCAGGT+GTGGTTAC  
TTTGCAGCGTCTTCATCGATGCGAGAGCCAAGAGATCGGTGTTGAAAGTTGTATTAAGTTATAAGGC  
+  
FFFFFFFFFFFFFFFFFFFFFFFFFFFFFF:FFFFFFFFFFFFFFFF:FFFFFFFFFF:FFFF:FF
```

Interleaved

{

```
@A00459:120:HHFYJDRXY:1:2101:1515:1016 1:N:0:AGTCAGGT+GTGGTTAC  
TNTACGGCGTTCTTCATCGATGGGAGAACCAAGAGATCCGTTCTCAAAGTTGTATTTAGTTCTGCCGA  
+  
F!FFFFFFFFFFFFFFFFFFFFFFFFFFFFFF:FFFFFFFFFFFFFFFFFFFFFFFFFFFF  
@A00459:120:HHFYJDRXY:1:2101:1515:1016 2:N:0:AGTCAGGT+GTGGTTAC  
ACACAGTCCTGACGACCAGCGGAGGGATCATTAGTCATACAACCGGGGAATCCACTCTGTGGGCCCAACC  
+  
FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF  
@A00459:120:HHFYJDRXY:1:2101:1642:1016 1:N:0:AGTCAGGT+GTGGTTAC  
ANGCGAACTAATCGACCAGCGGAGGGATCATTAATGAATAAACTCGGTGGATTGTTGCTGGCTCTAGGAG
```

File extensions

* .**fasta**, * .**fa**

* .**fastq**, * .**fq**

* .**fasta.gz**, * .**fa.gz**

* .**fastq.gz**, * .**fq.gz**

File extensions

* .fastq.gz

... > nnn > ⚡ Set01 1 R1.fastq.gz

1 U_S ? B_S N_L N_L n ? N_L ? ? + ? H ? - < ? ? [Q Q ? ? ?) % 1 ? ? ? ' E_O ? ? ? 7 ? ? ? ? D_E_L ? ? ? 6 ? ? B_S *] U M ? ? ? - ? ? C ? d b B_E_L ? q ? e | x ? 송 ? a ? ?
2 | (" j ? B_E_L ? U ? ? ? W { ? ? * r ? N_A_K ? ? ? ? I W ? ? ? H ? ? ? ? ' (H) u ? ? ? ? ? D_O_3 \ @ ? ? q T ? Q ? ? * ? ? ? 8 ? W ? u ? e ? ? S_U_B_F_S ? ? ? F { ? o ? ? ? ? : ? G_S ? \ T ? X _ ? ? ?
3 ? c ? } E_M ? ? ? 7 ? ? ? 1 ? | ? f ? @ ? ? U_S ? D_E_L E_S_D ? ? ? 0 { ? d ? E_O_* p { ? ? D_C_2 ? T ? B_S ? ? ? 1 , d ? K ? ? ? R_S ? J ? ? ? ?) ? N K ? o ? ? F_S { Q ? , > ? r P_S_S ? j ? ? V_T ? X ? C_A_N ?
4] ? i ? B_S D_E_L ? / J p ? ? ? c R M E ? I S 2 \ L * & V = ? \$? D_O_3 E_O / : F_S u ? \ ? N ? Z : E_Q d f y A ? E_S_C ? p ? K ? ? ? ? L ? ? ? ? * ? F ? E_N_Q ? Q ? ? ? = z + ? ? ? L | ? ? ? z
5 ? % ? ? ? C ? K X ? F_S ? Z 6 ? S_V_N ? J ? . ? } ? K ? \$? ? C_K ? ? ? (u ? ? > ? ? ? S_U_B ? U_S ? ? ? M [t 6 ? ? ? 1 b ^ E_T_B ? E 4 w ? ? E_T_B ? Q ? ^ ~ ? ? ? U_S ? l ? E_N_Q p ? C_A_N ? G_S ? F_S ? ? ?
6 ? P_D_E_S_U_B ? : ? - G G ? D_C_2 ? ? ? P ? D_O_4 (? B ? n q n t ? G_S ? ? ? E ? ? 4 e - ? ? [? ? \ N - ? m ? ? M ? N ? _ j ? y R ? G_S h L D % ? C_I 4 ? ? ? n ? ? k ? ? U_S u ? ? ? @ ? ? D_C_3 (?
7 ? ~ ? S B * ? ? ? t | P ? : e ? E_T_B ? ? ? ? ? # p ? G_S ? ? ? U_L ? E_T_X ? S E ? Y_T ? z 1 p } 1 ? m ? ? ? ? ? ? K V ? ? ? - J ? ? ? 8 ? ? . Q ? ? ? ? D_E_L ? C_A_N ? U L * ? v 3 ? ? S_U_B ? ? ? D_C_2 ? j ? G_S ? ? ?
8 ? ? ? ? ? ? 0 ? E_T_X ? f H ? [? j { S_V_N ? t ? U_S ? y \$? ? ? N_A_K ? ? ? C (? + ? H E ? ? ? ? Q ? D_E_L ? ? m ? ? c & ? D_E_L ? 7 9 ? - ? F_F ? x q ? 7 ? ? ? ? ? ^ N_U_L ? X ? ? ? L ? S_U_B ? S_T_X ? 9 ? T ; ? ? S_T_X ? U_S ? S_U_B ? C_I
9) C ? U ? ? ? K ? D_E_L ? ? D_C_2 ? [? Y E O ? F ? E_M ? K w & q U M @ ? ? ? l \ ^ ? g ? ? ? F : L = ? R_S ? 2 ? U_S ? E_S_C ? G ? S_U_B ? D_C_3 ? i ? B_S . J 7 ? N_A_K ?
10 ? ? ? l ? G_S ? K ? F_S ? / ? b c ? E_S_C ? ? ? ? ? ? N ? ? ? L - ? Z ? ? ? ? N_A_K : ? ? U ^ r ? _ ? 6 ? ? ? ? (? V S z % ? ? @ ? ? ? # ? ? ;) ? ? N_A_K ? E_N_Q ? " & ? ? ? ? [
11 ? ? ? ? + ? S_V_N ? + ? ? ? B ? S ? \ 9 % ? ? A_K ? j ? ? ? ? ? G ?) F_O_T ? n ^ ? K n b ? ? ? G_S = ? ? % ? ? D_C_2 ? ! ? \$? ~ ? y ? I w ? 1 ? E_O_T ? ? ? ? ' @ ? ? ? F ? * g P = @ F_F ? S_V_N ? 4 3 ?
12 ? Y ? K ? S_U_B ? ? a ? A ~ ? % w ? D_C_3 ? \ ? P ? E_S_C ? ? ? | 4 ? B_S * ? ? & C_A_N " ? ? 跌 ? ? ? ? i g @ e ! y ? E_N_Q ? ? ? S_I ? R ? ? & A_K ? E_O_T ? ? ? X j ? ? } & & ? E_O_T ? S_I ? ? ? ? ? , + ? N
13 D_E_T ? V_T ? " L 5 ? , t ? ? B e ? ? ? 7 Q ? ? x ? 7 C d ? F ? ? ? ? ^ ? ? x ? D_E_L ? R 9 @ ? ? # ? ? ? ? D_E_L ? ? ? ? + (I ? ? ? ? ? ? l C ? ? H z ? D_S (0 (D_C_3 ? 2 ? d ? ? ?
14 N_A_K ? \$ o ? D_C_3 ? m ? D_O_2 ? I ? Z ? } e ? & ? ? ? ? Q ? ? ? F_S % ? - { ? c ?
15 ? ? ? K + ? c ? ? ? ? ?] ? ? ? ? _ d p ? Z ? D_E_L ? ? ? ? N_U_L ? ? ? ? \ ? ? Y ? ? ? ? r \ D_E ? 9 ? ? n ? v ? D_C_1 ? N ? / D_E ? ? ? E ? D_C_4 ? } E_N_Q ? 8 ? B_S ? ? ? ? G ? E_S_C ?
16 ? ? ? = , 4 > ? ? ? ? ? U_S ? P r ? ? ? ? " ? F_S & S_I ? ? b ? E_O_T ? W ? ? E_T_X ? 0 4) B_E_L ? ? ? ? ? & ? y G ? E_M ? A ? S_I ? V_T ? S_O_H ? ? ? P ? x * ? * G_S ? S_Y_N ? ? ? ? D_C_2 ? D_E_L ? S_S ? ? ? ? ? S_I ? ? ? ? ? I ? ? E_M ?
17 C_A_N ? ? ? J ? ; H B ? \$? ? e ^ * ~ V ? o ? ? ? ? " J ? ? ? ? N ? ? b V ? B_E_L ? d c ? ? ? ^ ? ? ? X + d 1 ? \ v \ ? ? ? ? w ? ? ? ?] q ? ? (\ ? ? ? ? J ? ? ? ? w ? A_C_K ? ? ? ? ? ? ? ? ?
18 ? ? U ? ? ? V_T ? 0 1 ? ? ? ? Y ?

File extensions

* .**fasta**, * .**fa** }

* .**fastq**, * .**fq** } Uncompressed data

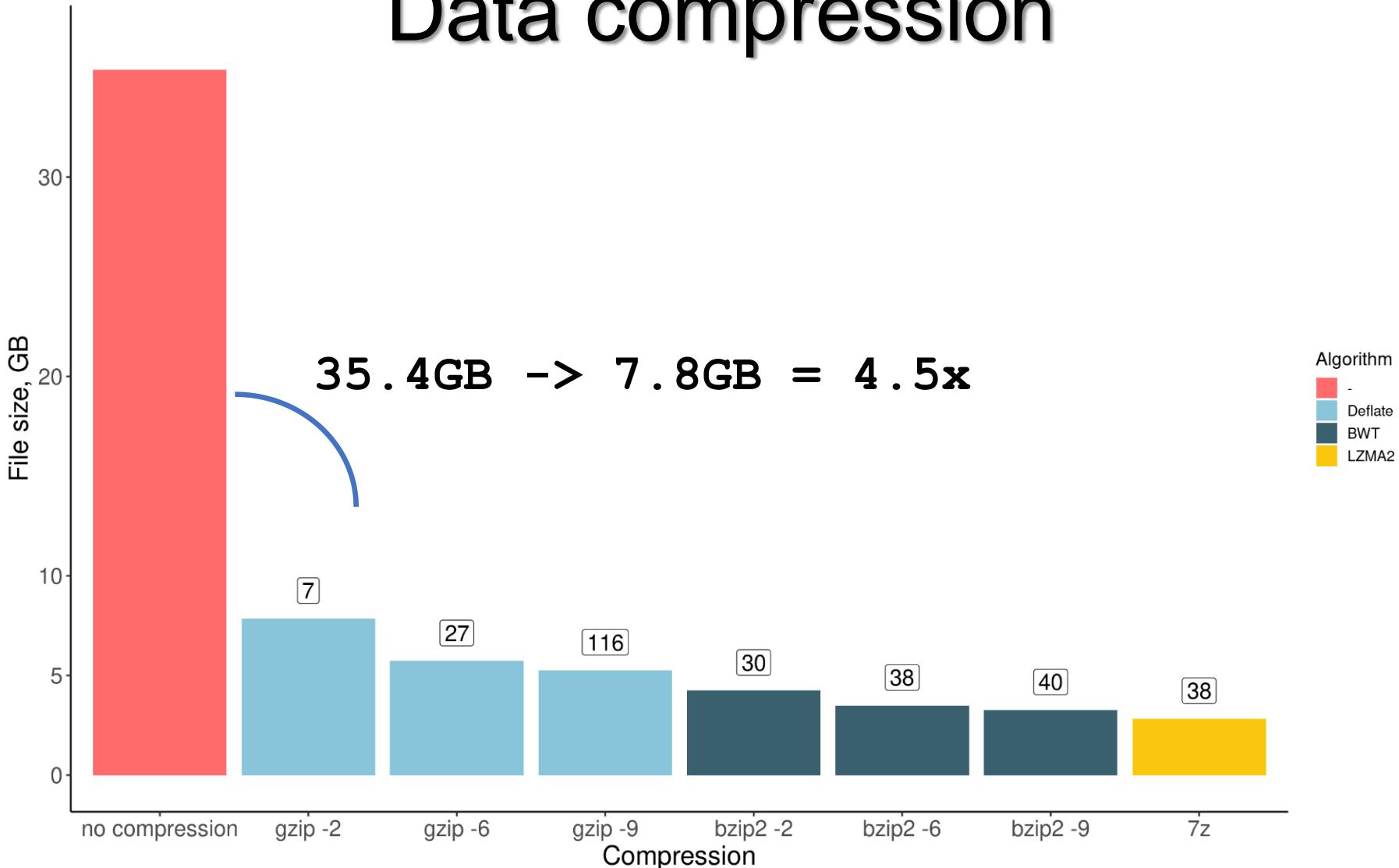
* .**fasta.gz**, * .**fa.gz** }

* .**fastq.gz**, * .**fq.gz** } gzip-compressed

* .**fasta.bz2**, * .**fa.bz2** }

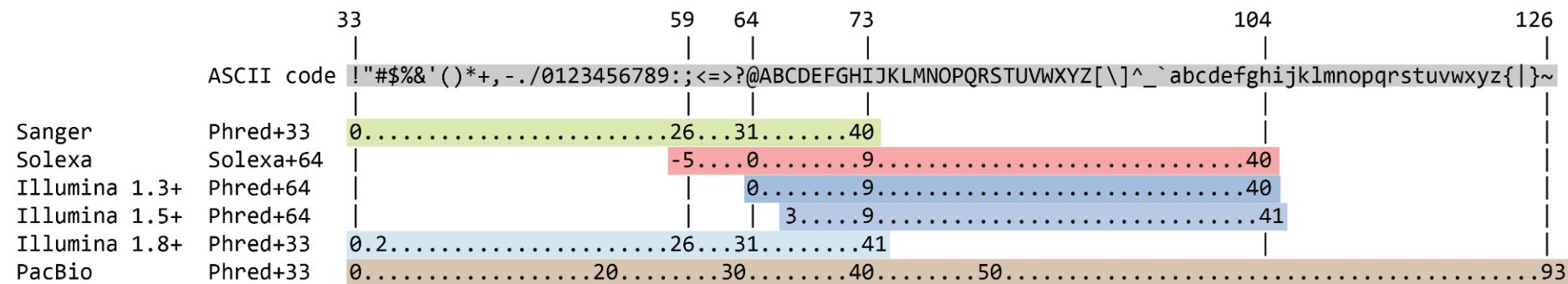
* .**fastq.bz2**, * .**fq.bz2** } bzip2-compressed

Data compression



Input = NovaSeq data, 67.9M sequences

Quality scores



Phred quality scores

$$Q = -10 * \log_{10} P \quad P = 10^{-Q/10}$$

ASCII	Q	P	Chance that the base is wrong (Error)	Accuracy (1 - Error)
\$	3	0.50	50%	50%
+	10	0.10	10%	90%
.	13	0.05	5%	95%
5	20	0.01	1%	99%
:	25	0.00316	0.316%	99.7%
?	30	0.00100	0.1%	99.9%
I	40	0.00010	0.01%	99.99%
S	50	0.00001	0.001%	99.999%
]	60	0.000001	0.0001%	99.9999%
~	93	0.00000001	0.0000001%	99.9999999%

} PacBio
HiFi
Reads

FASTQE



Bin Emoji

N ✌️

2-9 💀

10-19 💩

20-24 !

25-29 😊

30-34 😂

35-39 😎

≥ 40 😍

@GOOD_SEQUENCE

N T G G C C C C G G G T C G A C G T G G C

+

J J J H H G G G F F F E E C B = ; ; : : 9

😍😍😍😎😎😎😎😎😎😎😎😎😎😎😎😎😎😎😎😎😎😎😎😎😎😎⚠️

@BAD_SEQUENCE

N T G G C C C C G G G T C G A C G T G G C

+

! 9 9 9 3 3 2 2 2 1 1 . . - - # \$ % & % (

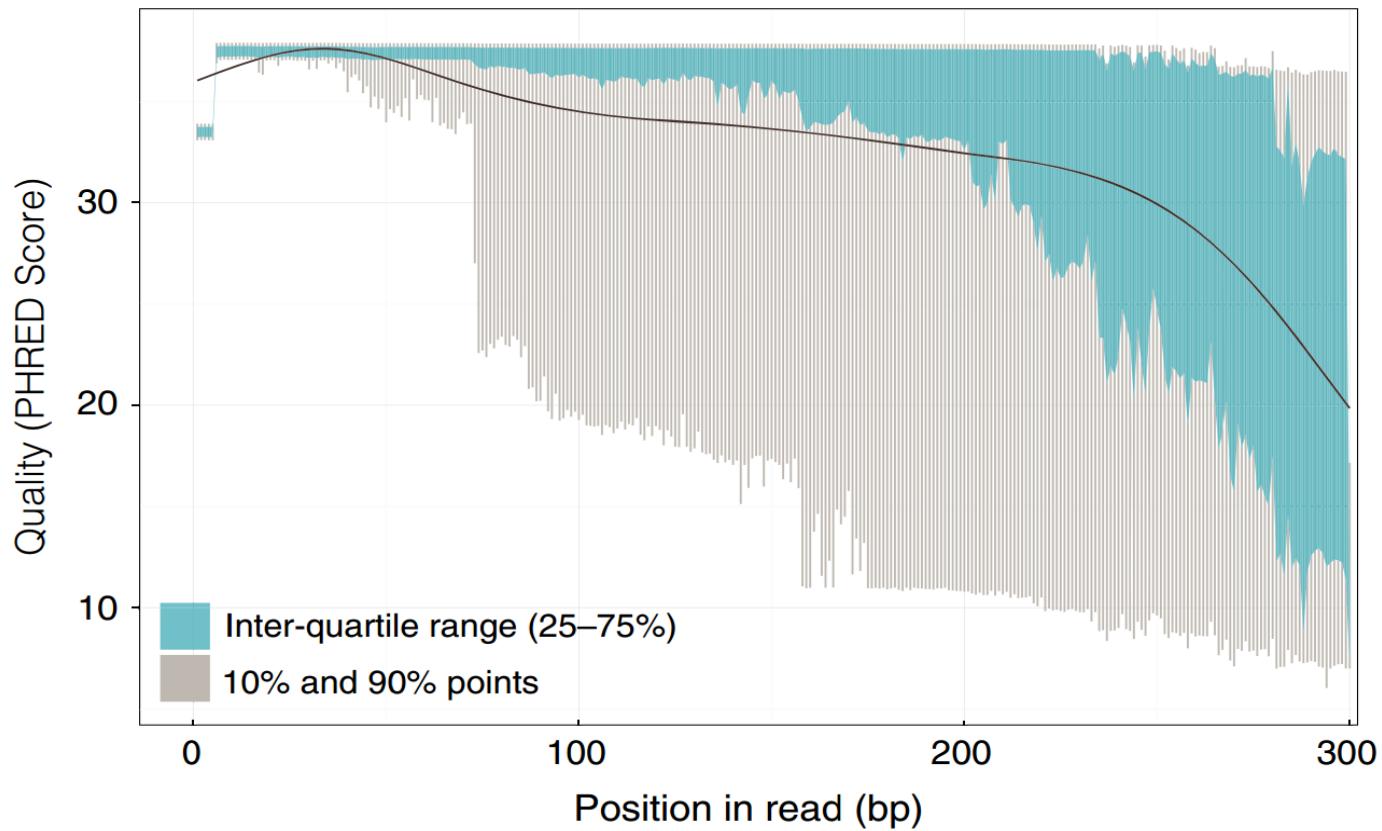
🚫⚠️⚠️⚠️💩💩💩💩💩💩💩💩💩💩💩💩💩💩💩💩💩💩💩💩

Quality Control

99.9%

99%

90%



Binned quality

- Illumina 2-Channel SBS Technology
(NovaSeq / NextSeq / MiniSeq / iSeq)

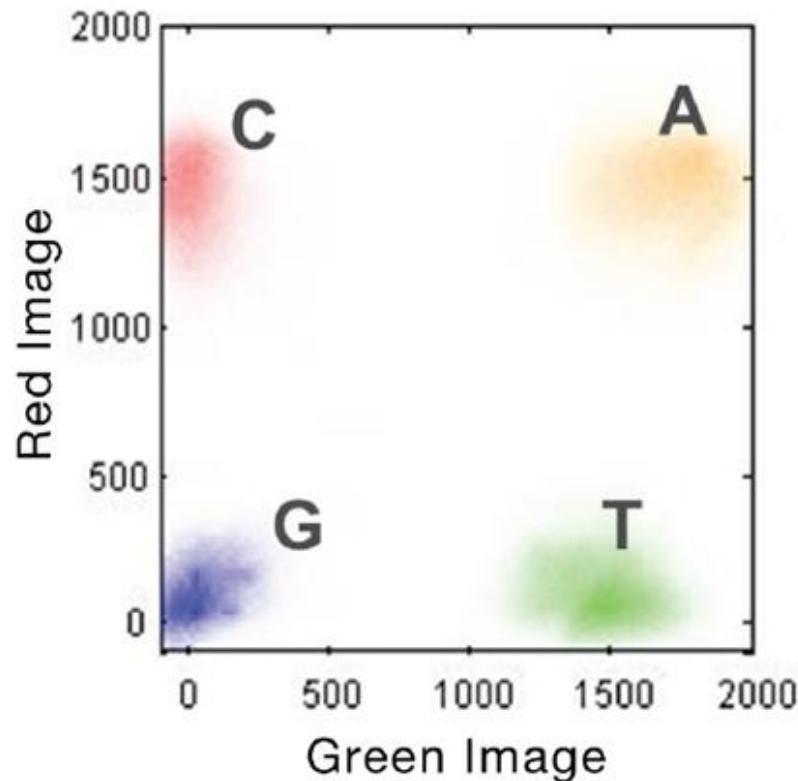
2, 11, 25, and 37

2, 12, 23, and 37

Poly-G tails

@A00551:123:H2MKG

...TCGCAGGTAGCGCGCCGCCGGGGGGGGGG
+
...:FF::F:FF,:F:, :, :, FFFFFFFFFF



Illumina 2-channel technology

Importance of quality control

Garbage In, Garbage Out



Your analysis is as good as your data!

Quality filtering

- Number of ambiguous (N) nucleotides $(10 + 20) / 2 \neq 15$
- Average quality score $(10 + 20) / 2 = 12.6$
- Maximum number of expected errors (MaxEE)
- Maximum expected errors as a percentage of read length (MEEP)
- Number / percentage of unqualified bases

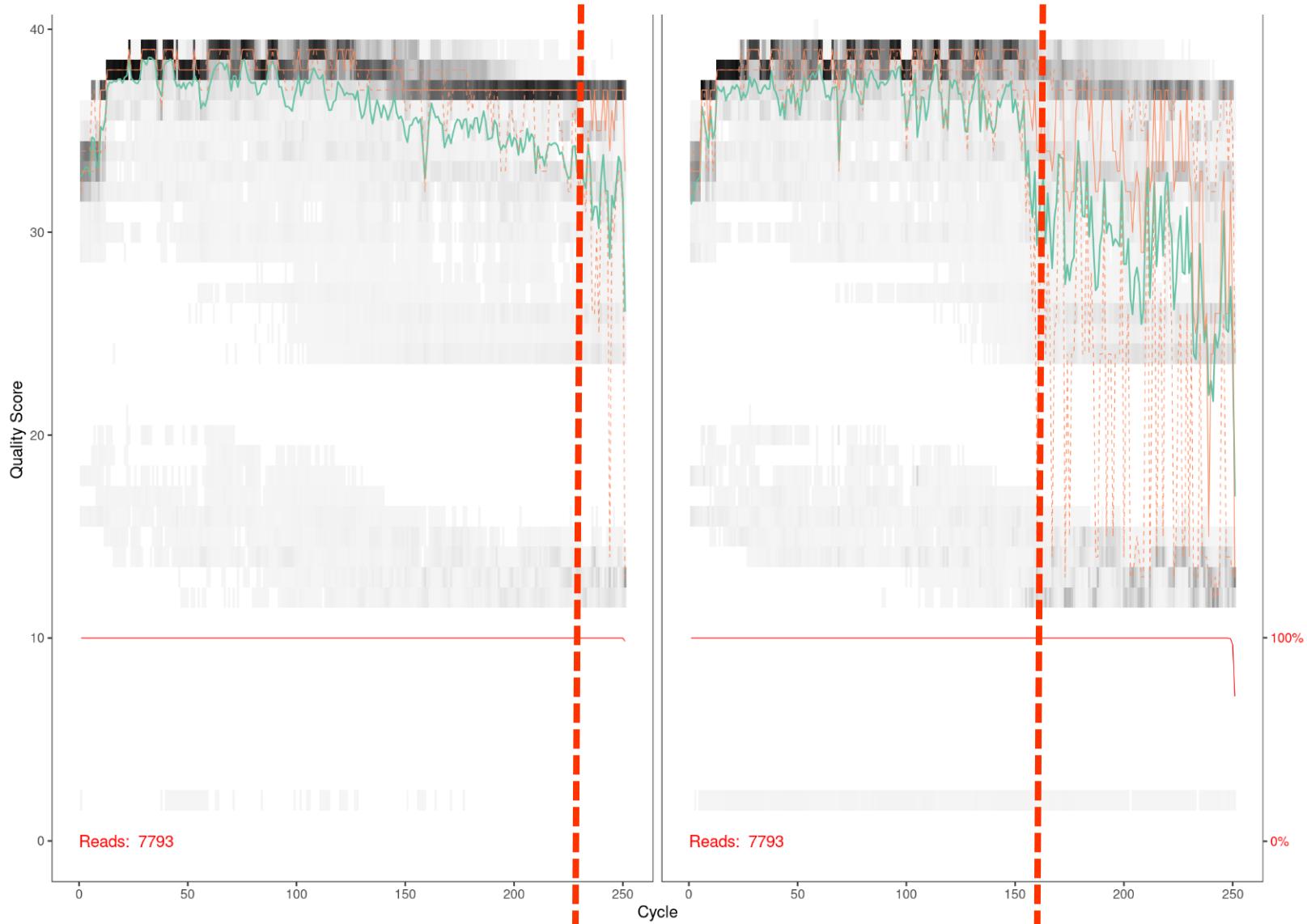
$$EE = \sum_i P_i$$

Edgar & Flyvbjerg (2015)
DOI:10.1093/bioinformatics/btv401

Trimming

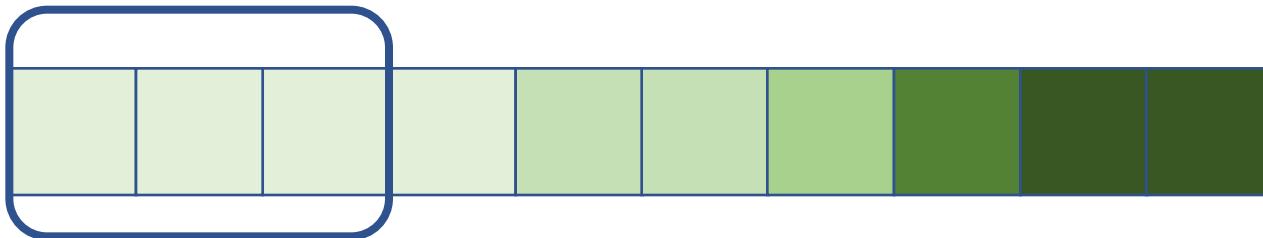
R1

R2

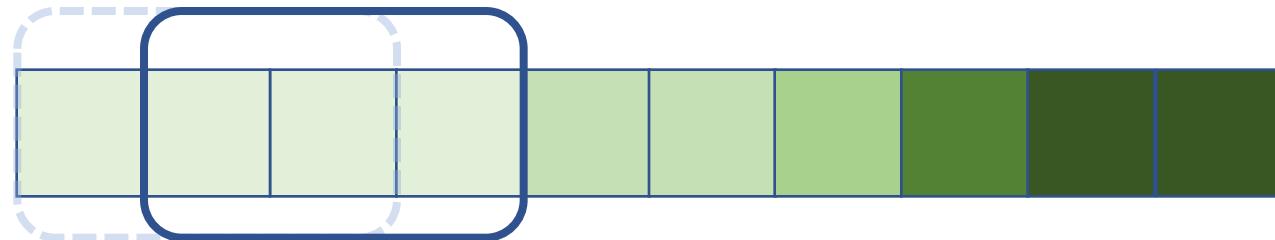


Sliding window trimming

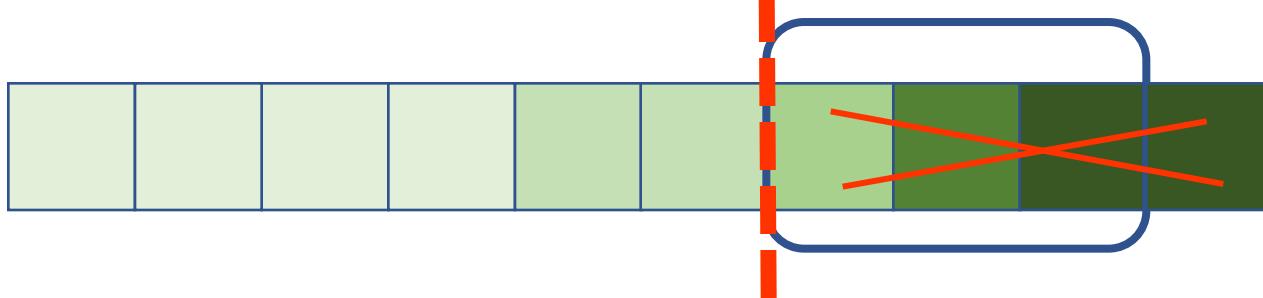
Window (size = 3)



Slide one bp forward, estimate average quality in the window

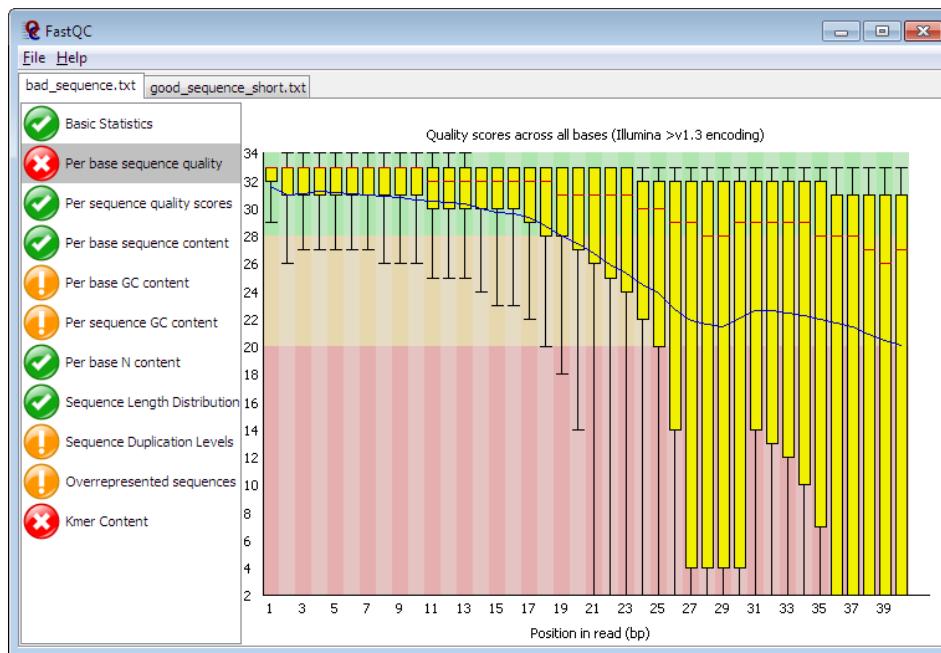


Cut the leftmost position in the window where the quality drops below the threshold



Quality check

FastQC



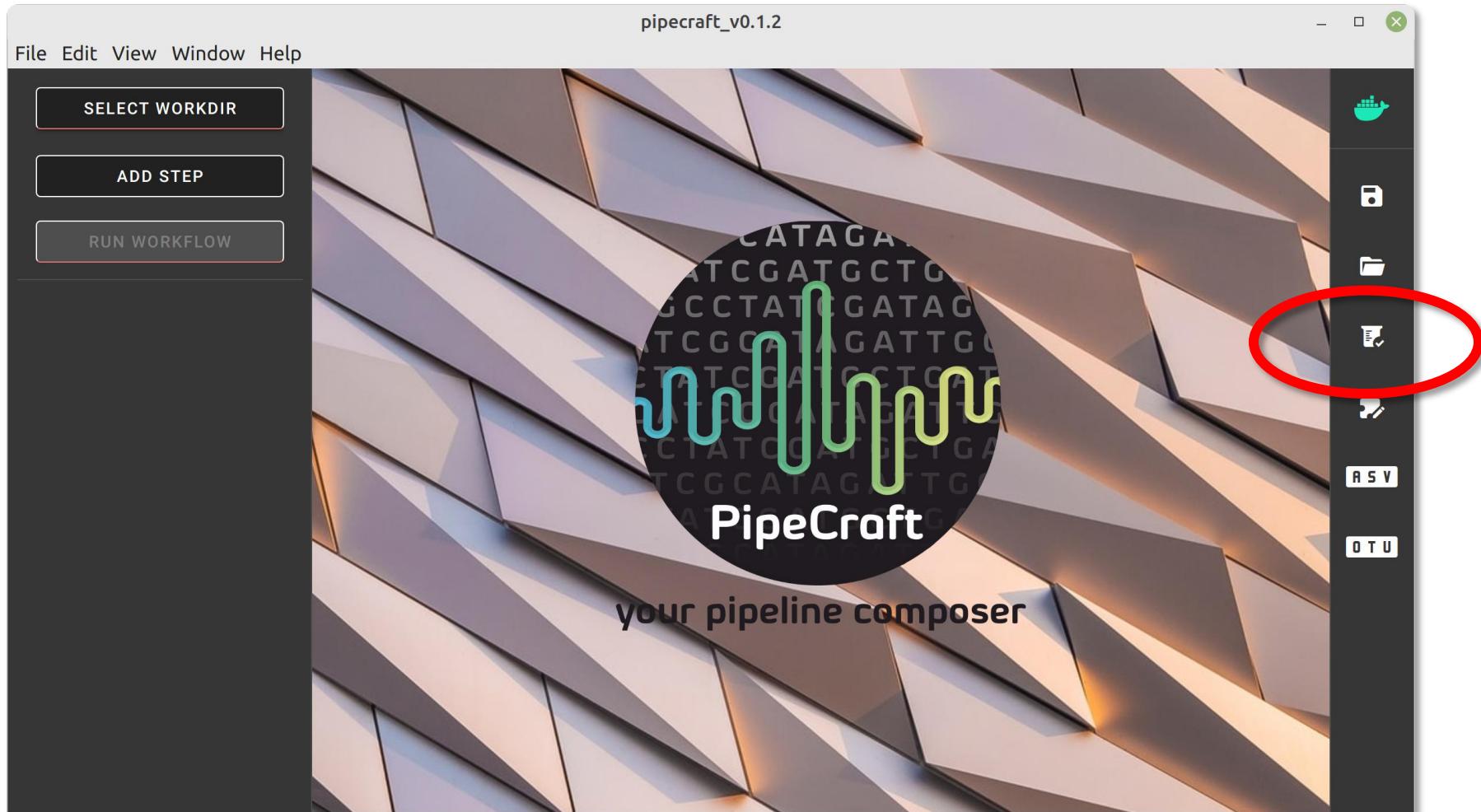
Andrews et al. (2010)

MultiQC



Ewels et al. (2016)
DOI:10.1093/bioinformatics/btw354

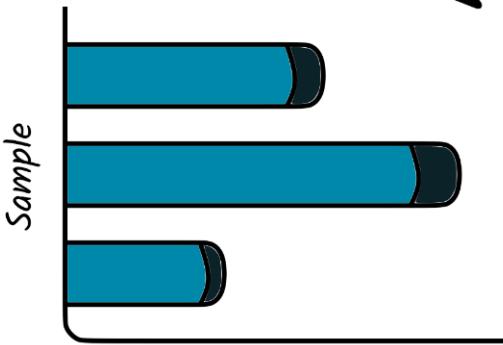
Quality check in PipeCraft2



FASTQC - Sequence count

ZandraSelina

Not good, if equal number of reads were requested! Maybe uneven sequencing pool or failed demultiplexing?



Not good! Many duplicates - is the library over-amplified?



Beautiful! Go reward yourself with a pizza!



CC-BY 4.0

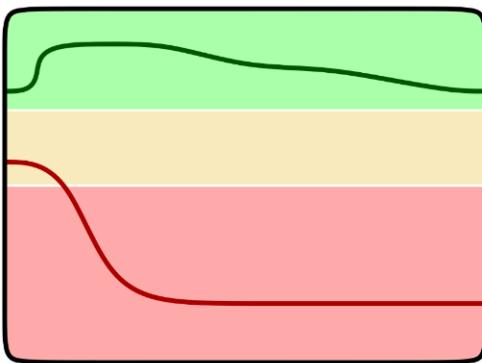
FASTQC - Sequence quality

ZandraSelina

Not good! Maybe one pair failed during sequencing?

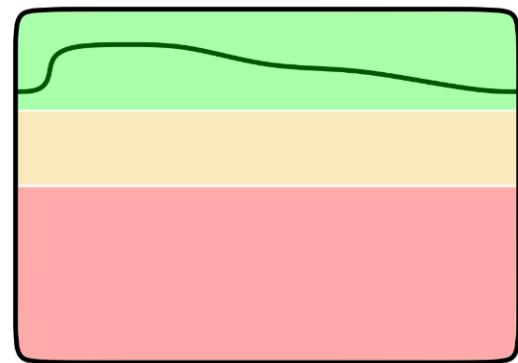
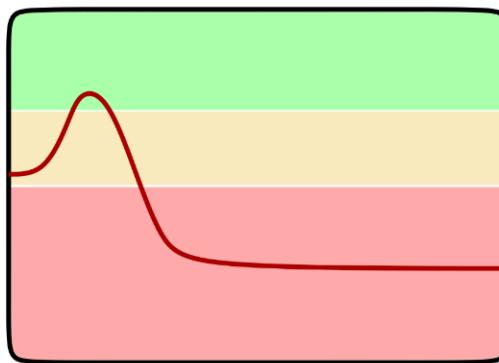
Not good! Was there a cycle chemistry issue during the run?

Beautiful! Go reward yourself with some cookies!



R1

R2

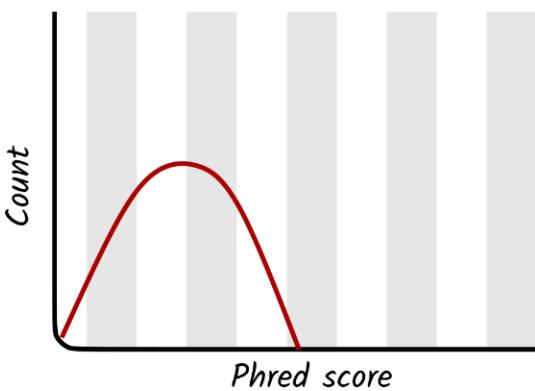


CC BY 4.0

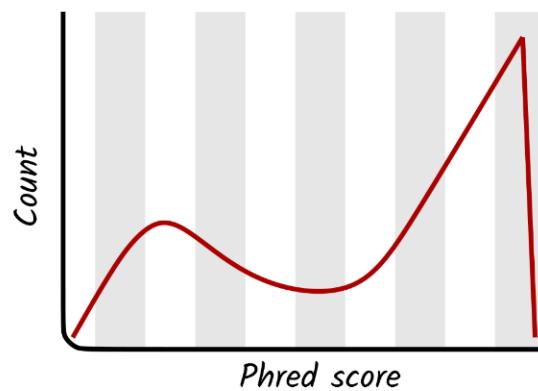
FASTQC - Per sequence quality score

ZandraSelina

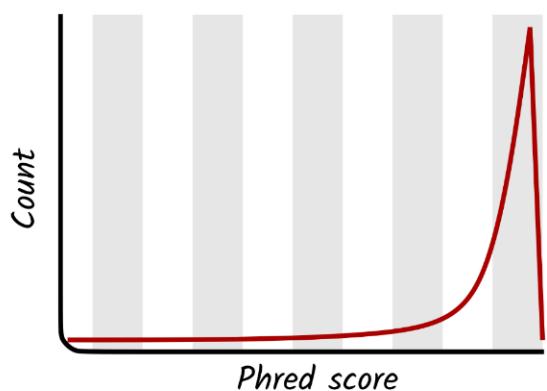
Not good! All reads have low quality, was it a bad sequencing run?



Not good! You have a bimodal peak, maybe an artefact from location on flowcell?



Beautiful! Go reward yourself with a good book!



CC BY 4.0

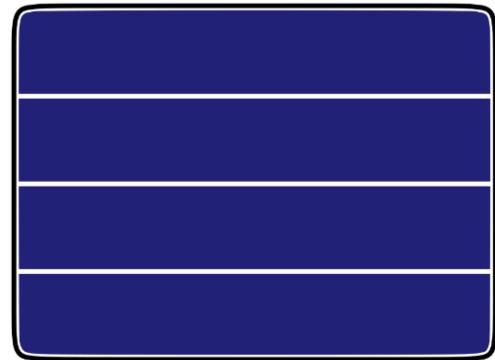
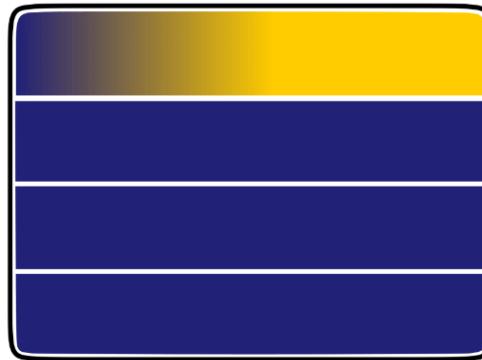
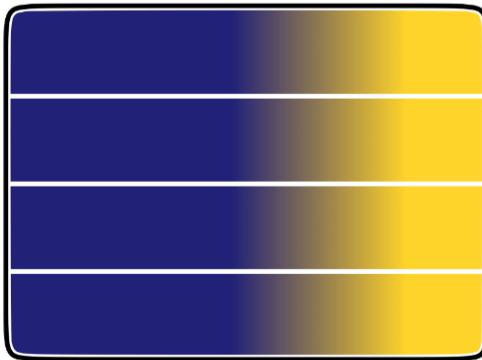
FASTQC - Per base sequence content

ZandraSelina

Not good! Did the sequencing chemistry run out? Or do you have very short reads?

Not good! Did one of the lanes fail?

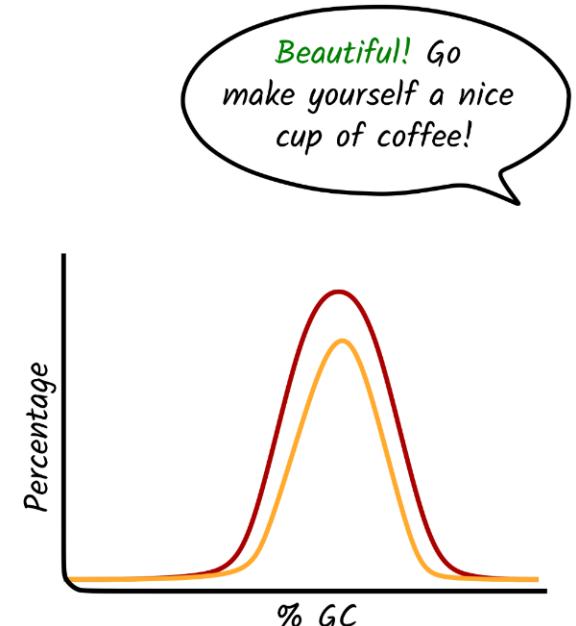
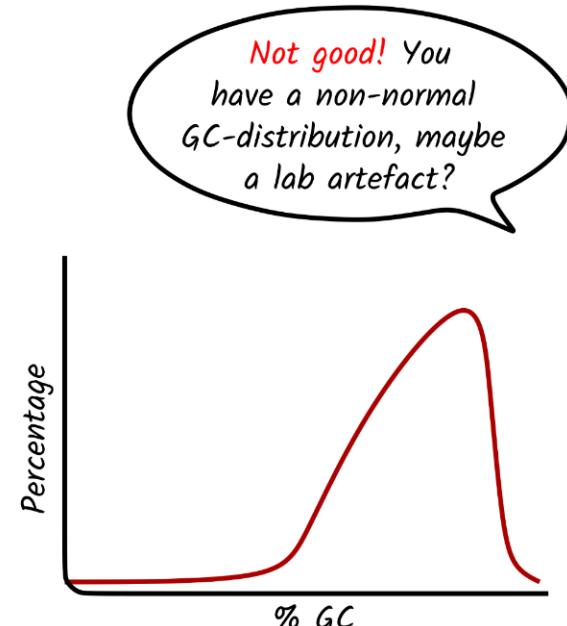
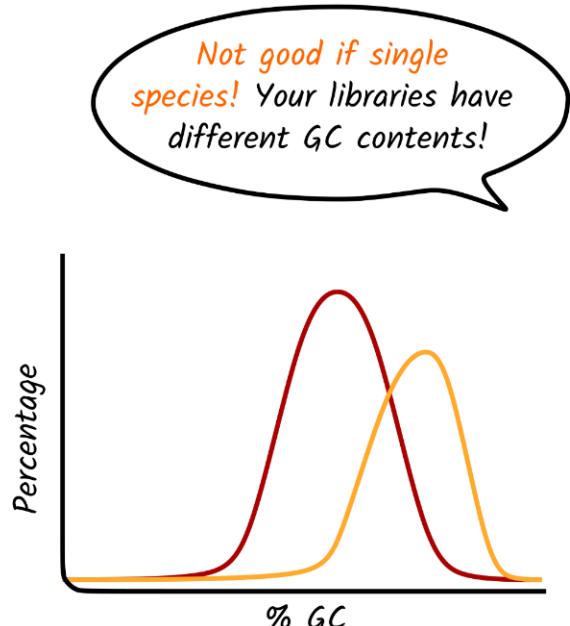
Beautiful! You deserve to go take a nap!



CC BY 4.0

FASTQC - Per sequence GC-content

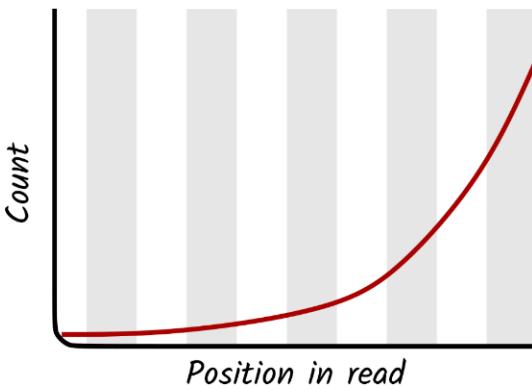
ZandraSelina



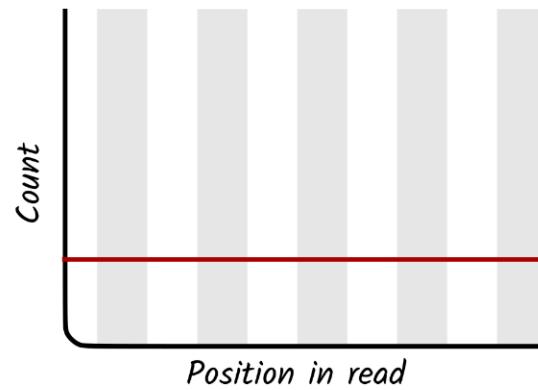
CC BY 4.0

FASTQC - Per base N content

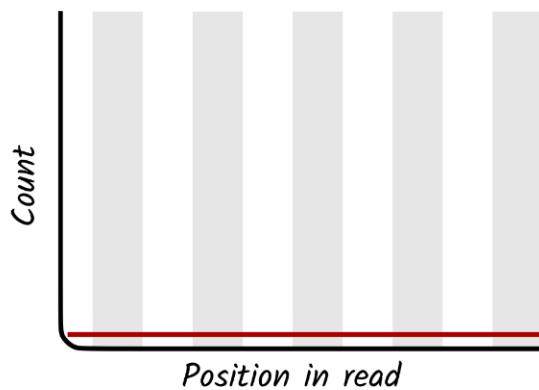
Not good! Did you run out of sequencing chemistry?



Not good! Maybe a failed lane? Or are host reads masked, if it is public data?

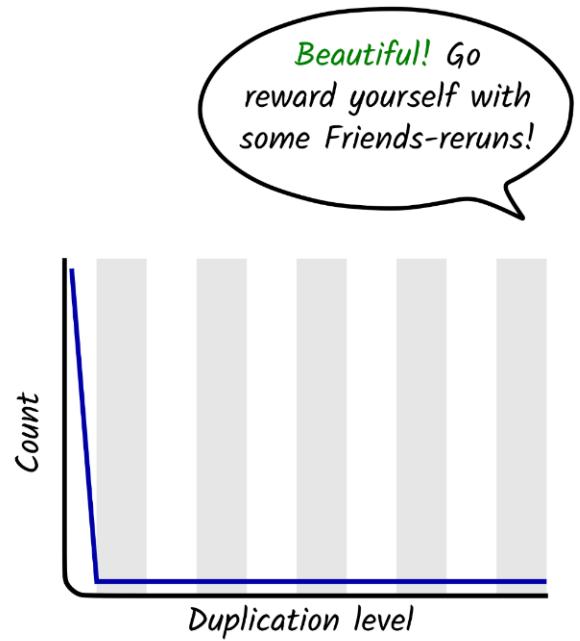
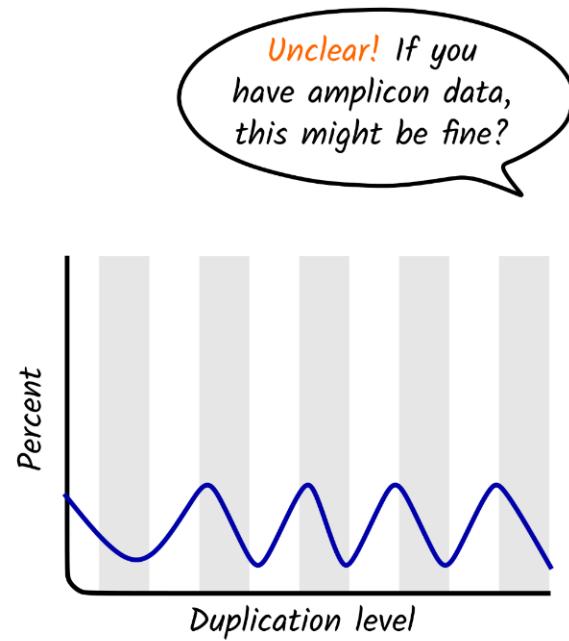
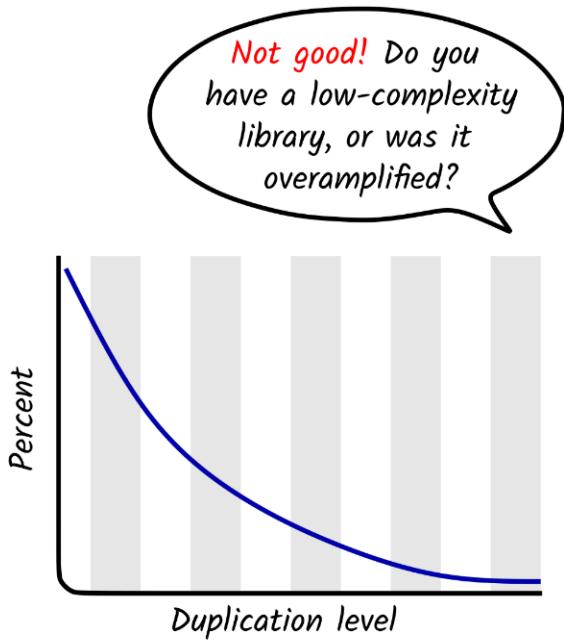


Beautiful! You deserve a pumpkin spice latte!

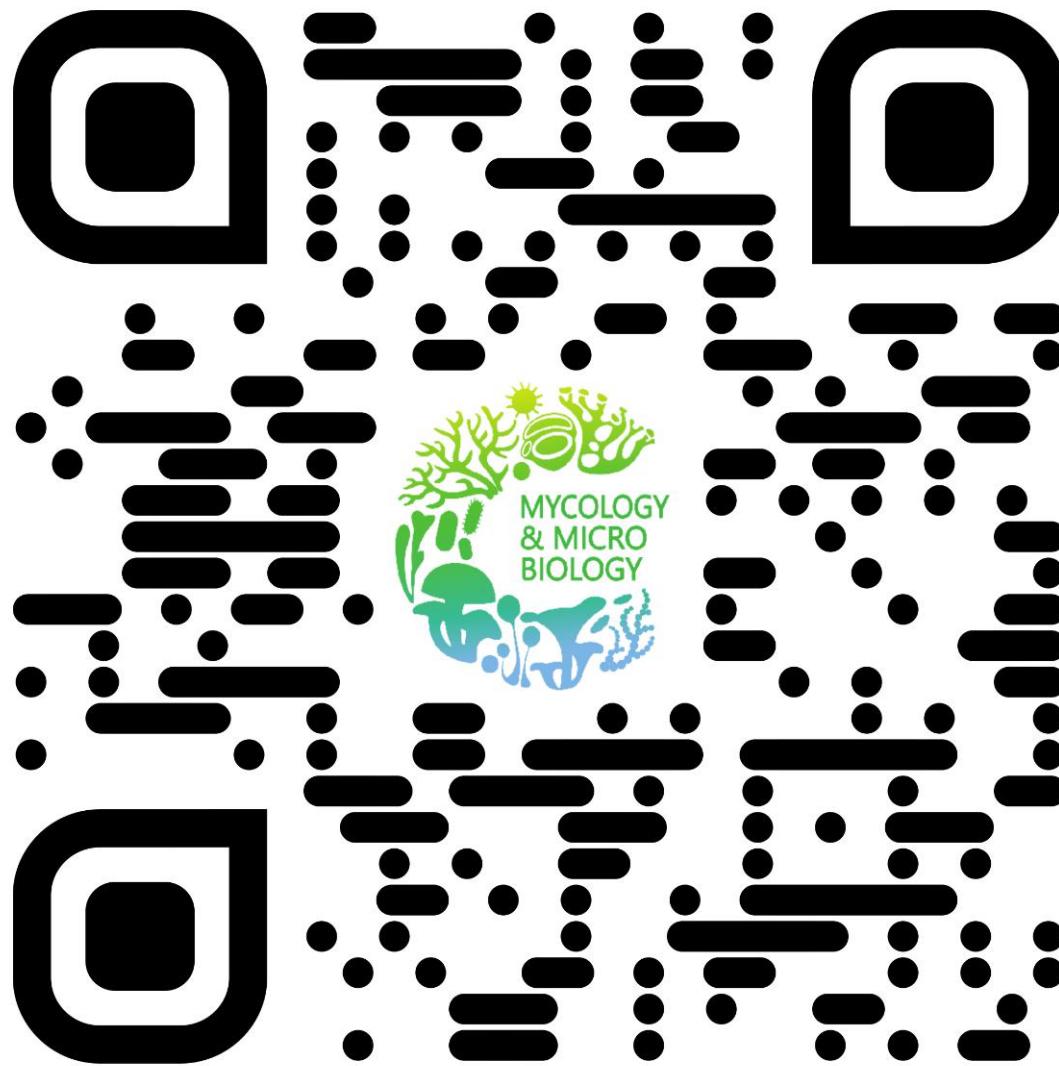


FASTQC - Sequence duplication levels

ZandraSelina



CC BY 4.0



<https://mycology-microbiology-center.github.io/Metabarcoding2022/>

<https://t.ly/dvZS>

main ▾

Metabarcoding2022 / data /



vmikk qc

..

DB

MiSeq_Fungi

QC

Set01

Set02

Set03

README.md

main ▾

Metabarcoding2022 / data / QC /



vmikk add data for QC

..

MiSeq_R1.fq.gz

MiSeq_R2.fq.gz

NovaSeq_R1.fq.gz

NovaSeq_R2.fq.gz

PacBio.fq.gz

<https://github.com/Mycology-Microbiology-Center/Metabarcoding2022/tree/main/data>