

Long-read sequence analysis

QC & alignment

fastq

reads.fastq

```
@D00283R:66:CC611ANXX:4:2311:2596:2330 1:N:0:TCCGGAG
ACTCTACGCTCAATAAAGATTTCTGATACGGCTCCTGAAATGCAGAATGAGT
+
B/<<<B<FFFFFFFFFBBFFFBFFFBFFFF/FFFFFFFF/BFFFBFFF
```

title, starts with @

nucleotide sequence

optional description

base quality

fastq

fasta + basequality (fasta + q = fastq)

$$BASEQ = -10\log_{10} \Pr\{base\ is\ wrong\}$$

$$-10\log_{10} (0.01) = 20$$

$$-10\log_{10} (0.1) = 10$$

$$-10\log_{10} (0.5) = 3$$

Quiz question

What kind quality characteristics are important to long-read sequencing reads but less important for Illumina sequencing?

- A. Base quality
- B. Read length
- C. GC content
- D. Adapter content

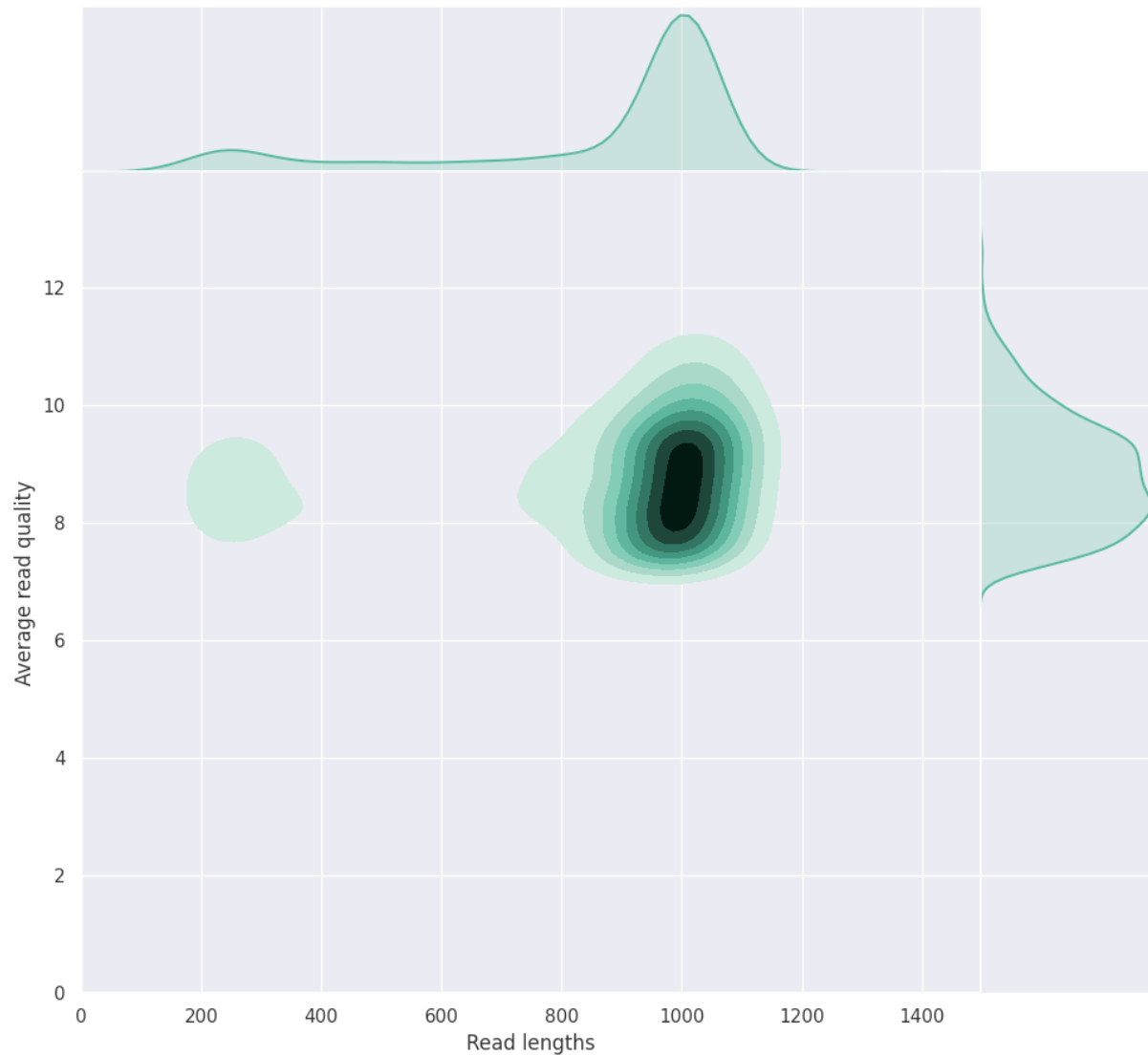
Read quality control

- Number of reads
- Read length (mean and spread)
- Base quality
- Overrepresented sequences
- GC content
- Demultiplexing statistics
- Run duration/location dependency
- Others?

Read quality software

- Software of manufacturer
- NanoPlot (<https://github.com/wdecoster/NanoPlot>)
 - Takes many input formats
 - Basic statistics (fastq based)
- PycoQC (<https://github.com/a-slide/pycoQC>)
 - Specific for ONT
 - Requires so-called sequencing_summary file
- FastQC (<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>)
 - Works also for long reads
 - Familiar output for many of us

Read lengths vs Average read quality plot

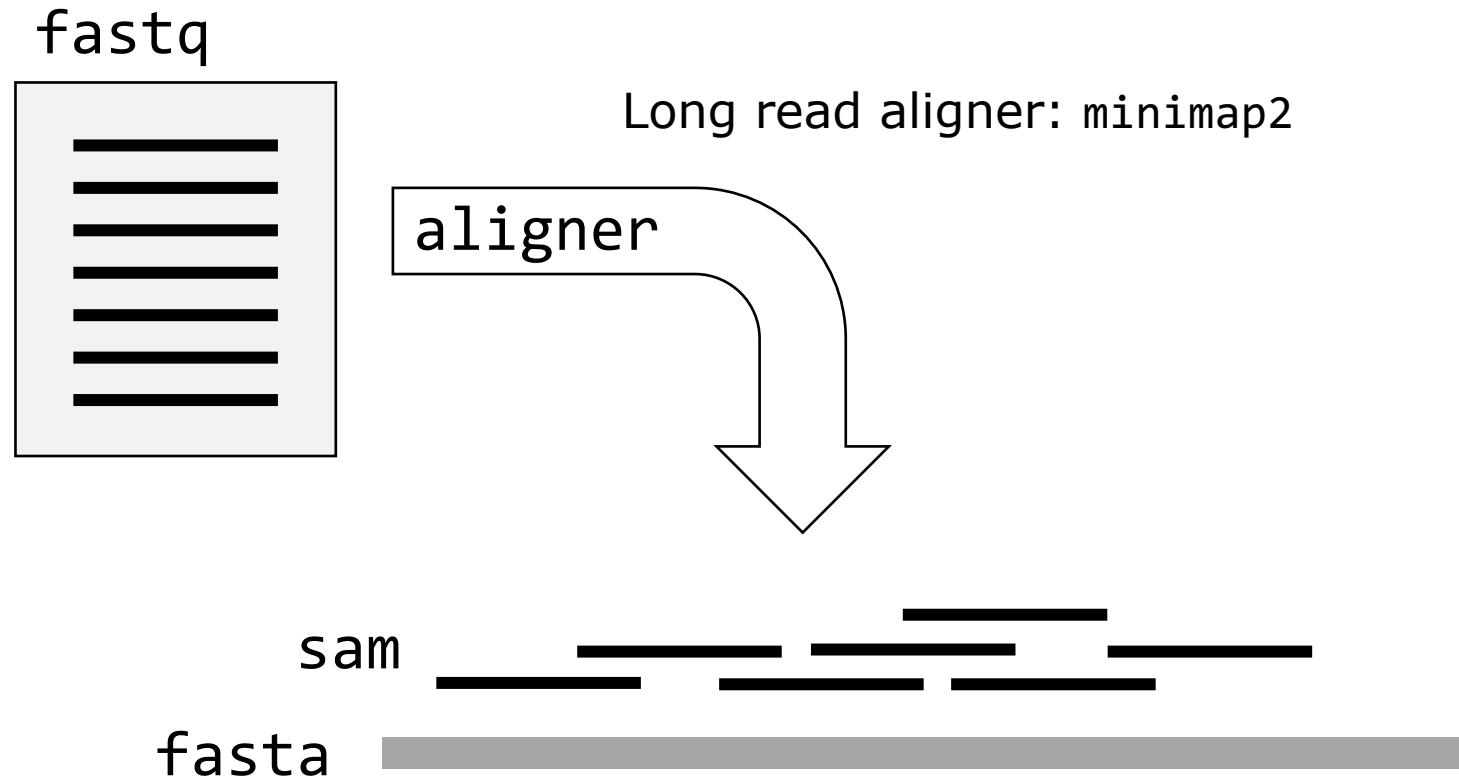


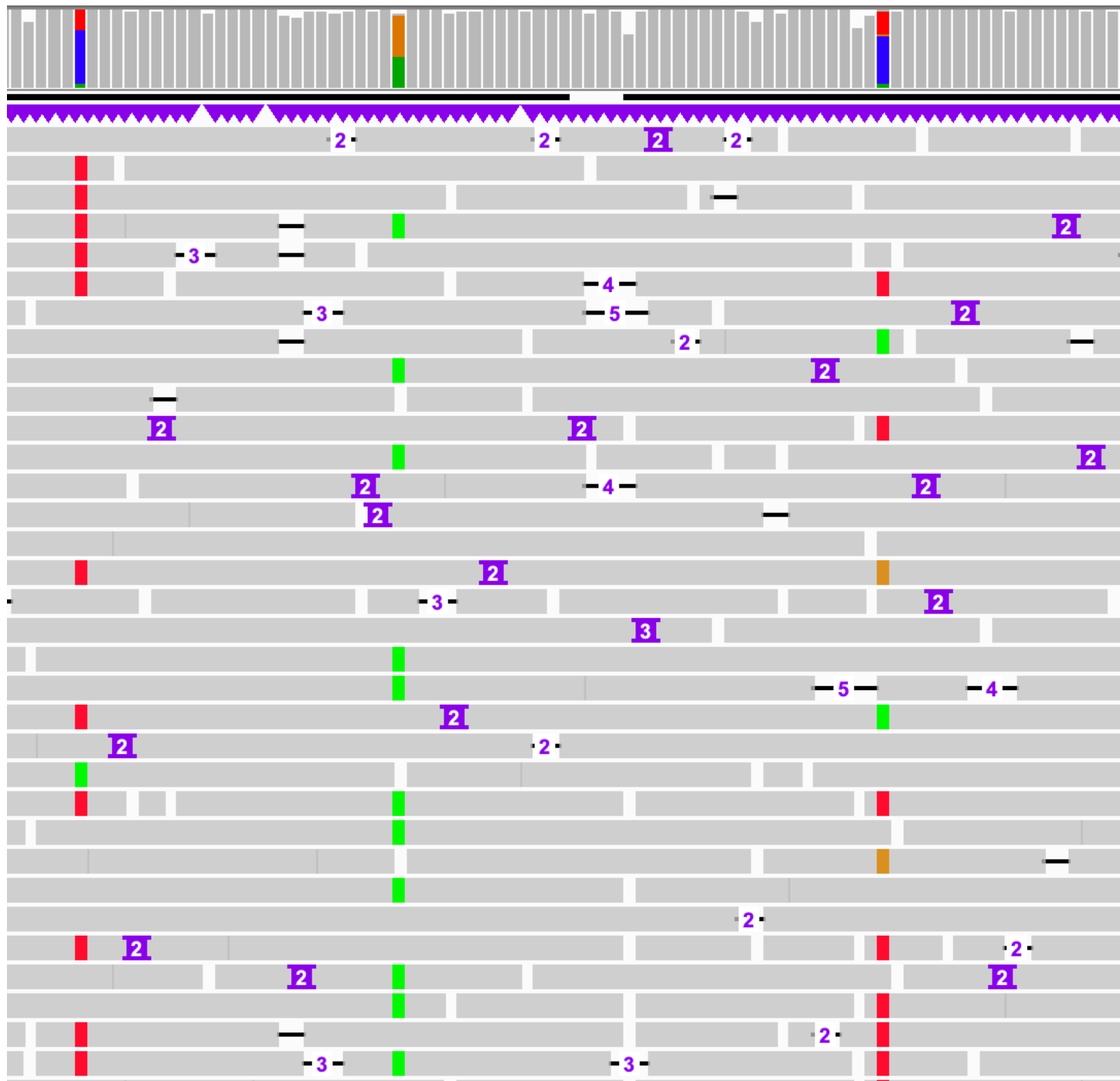
output of NanoPlot (<https://github.com/wdecoster/NanoPlot>)

Quality trimming

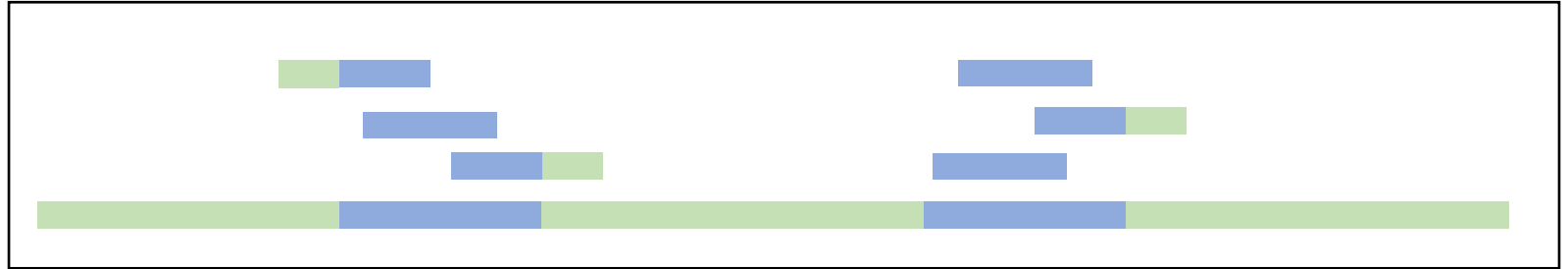
- Removal of:
 - Low quality sequences
 - Adapters/barcodes
- Oxford nanopore: On-instrument (guppy)
- PacBio:
 - On-instrument
 - During CCS generation (pbccs)

Read alignment (phred)





Mapping quality



$$MAPQ = -10\log_{10} \Pr\{\text{mapping position is wrong}\}$$

$$-10\log_{10} (0.01) = 20$$

$$-10\log_{10} (0.5) = 3$$

sam header

```
@HD      VN:1.0  SO:coordinate
@SQ      SN:U00096.3      LN:4641652
@PG      ID:bowtie2      PN:bowtie2      VN:2.4.1      CL: bowtie2-
align-s --wrapper basic-0 -x ref.fasta -1 reads_1.fastq -2
reads_2.fastq"
```

SAM column	example
read name	SRR519926.5
flag	89
reference	chr20
start position	61
mapping quality	42
CIGAR string	150M
reference name mate is mapped	=
start position mate	476
fragment length	515
sequence	CATCACCATTCCCAC
base quality	@>4:4C@89+&9CC@
optional	AS:i:-2
optional	XN:i:0

samtools

- Convert .sam files into (a.o.)
 - .fastq
 - .bam (compressed .sam)
- Subset based on:
 - flag
 - region
- Ordering
- Mark alignment duplicates
- And many other things

Long-reads & fastq

- fastq format is limited to:
 - base
 - base-quality
- Long-read technologies -> need to store more information:
 - PacBio: (unaligned) bam
 - ONT: fast5