# Metagenomics 101

# Session 2: Raw data & QC

Anna Heintz-Buschart
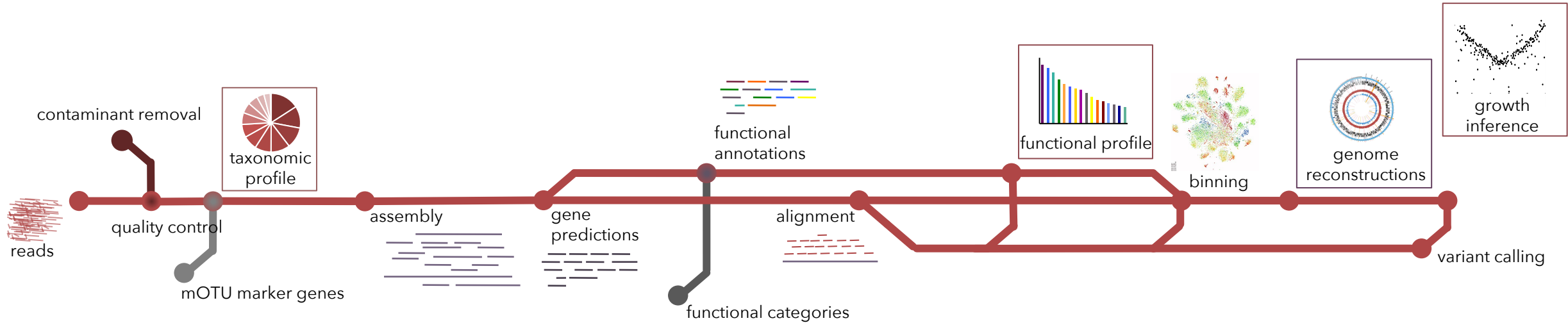
February 2022
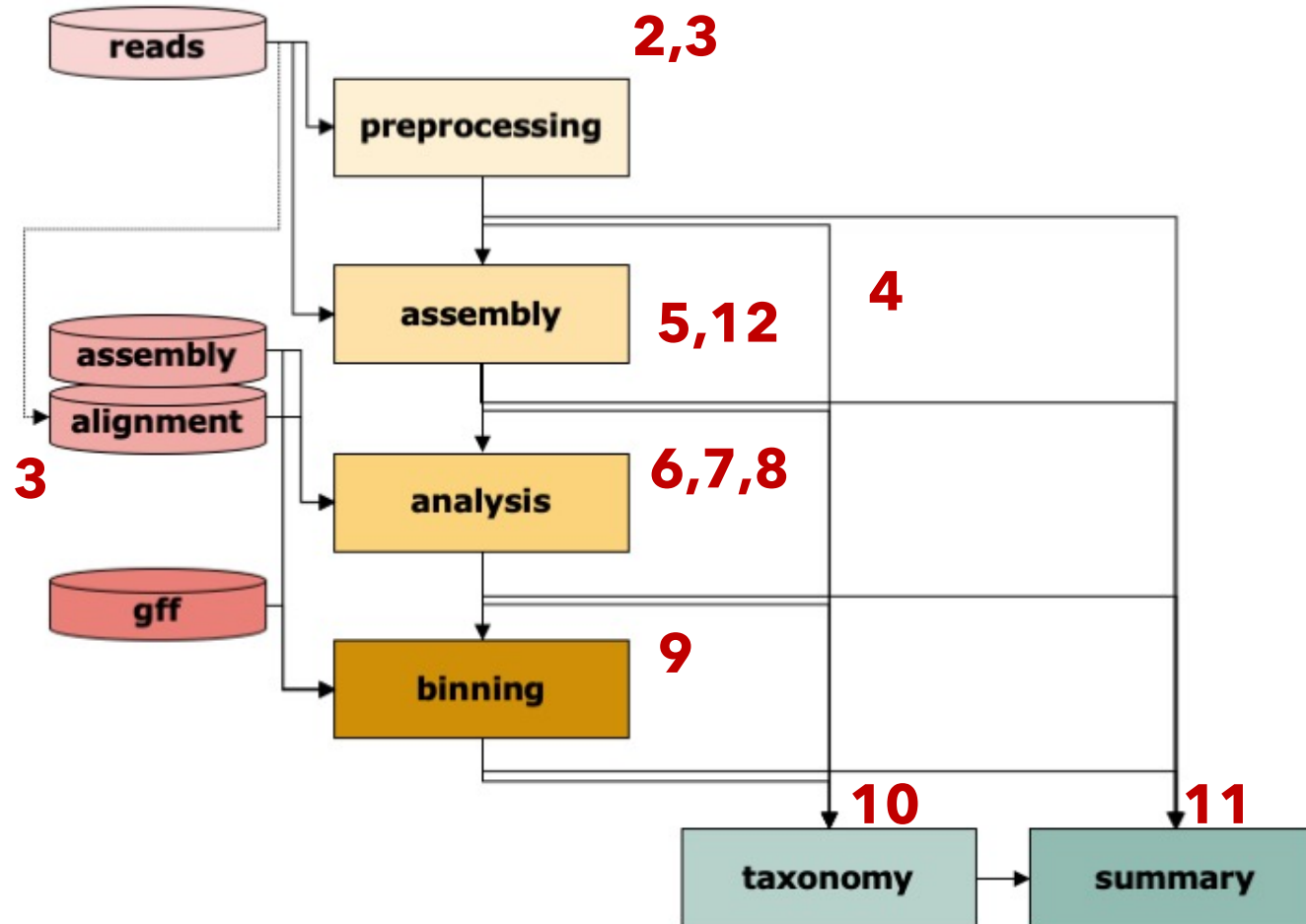
UNIVERSITY OF AMSTERDAM
Life Sciences

# Omics paradigm

# Metagenomics (+ other omics) pipeline

# Metagenomics (+ other omics) pipeline

# What does sequencing data look like?



Line 1: Name

Line 2: Sequence

Line 3: anything

Line 4: Quality at each position

.
.
.
.
.
.
.
.
.

as many as we have reads

(forward- & reverse files)

# Where do errors come from?

UNIVERSITY OF AMSTERDAM
Life Sciences

# Sequencing

short read sequencing

long read sequencing

# Short read sequencing

UNIVERSITY OF AMSTERDAM
Life Sciences

# Short read sequencing

Patterson, J., Carpenter, E.J., Zhu, Z. et al. BMC Genomics 20, 604 (2019)

UNIVERSITY OF AMSTERDAM
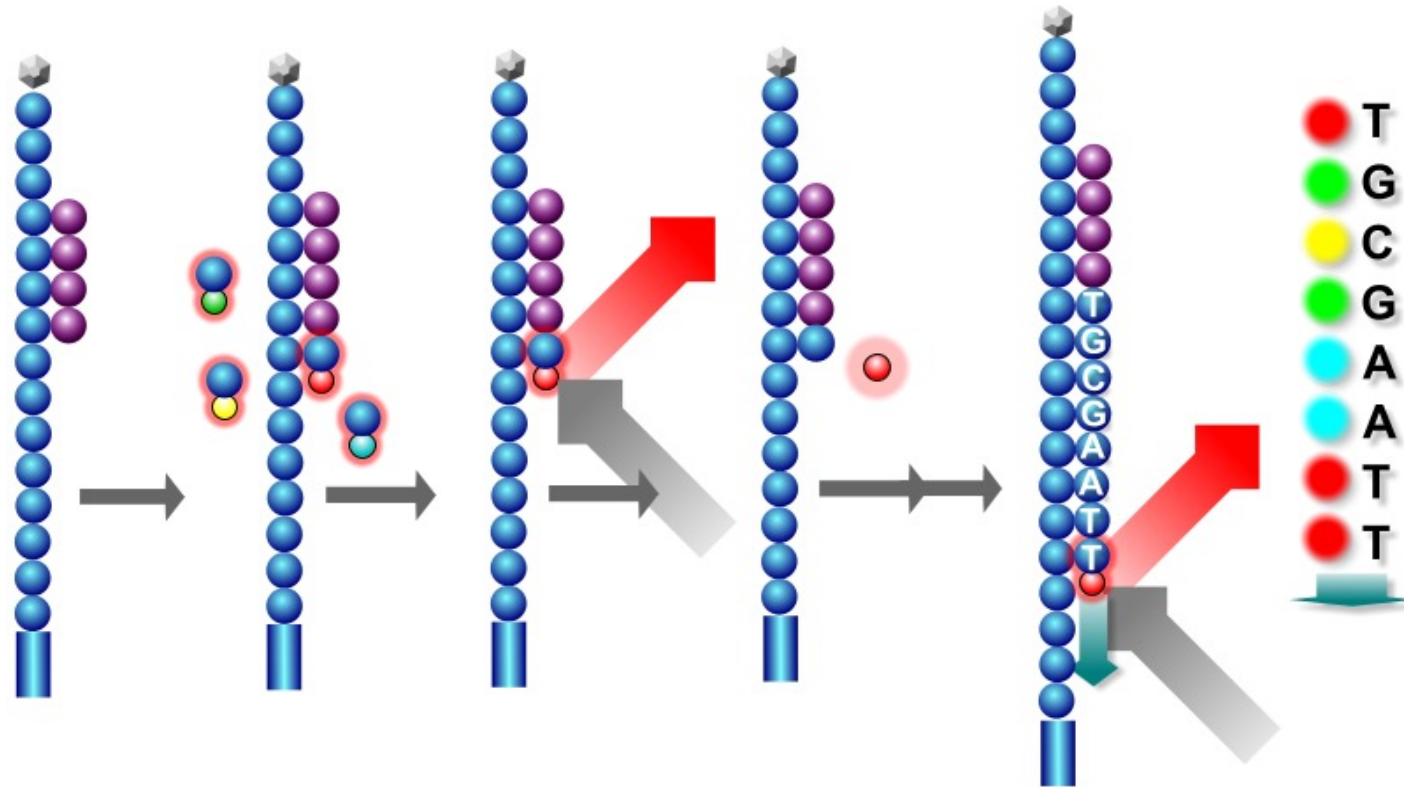Life Sciences

# Short read sequencing



A. Two representative DNA fragments from two unique samples, each attached to a specific barcode sequence that identifies the sample from which it originated.

B. Libraries for each sample are pooled and sequenced in parallel. Each new read contains both the fragment sequence and its sample-identifying barcode.

C. Barcode sequences are used to de-multiplex, or differentiate reads from each sample.

D. Each set of reads is aligned to the reference sequence.

UNIVERSITY OF AMSTERDAM
Life Sciences

https://www.illumina.com/techniques/sequencing/ngs-library-prep/multiplexing.html

# Short read sequencing

# Short read sequencing

|  | MiSeq Series ⊕ | NextSeq 550 Series ⊕ | NextSeq 1000 & 2000 | NovaSeq 6000 |
|---|---|---|---|---|
| Run Time | 4–55 hours | 12–30 hours | 11-48 hours | ~13 - 38 hours (dual SP flow cells)<br>~13–25 hours (dual S1 flow cells)<br>~16–36 hours (dual S2 flow cells)<br>~44 hours (dual S4 flow cells) |
| Maximum Output | 15 Gb | 120 Gb | 360 Gb* | 6000 Gb |
| Maximum Reads Per Run | 25 million † | 400 million | 1.2 billion* | 20 billion |
| Maximum Read Length | 2 × 300 bp | 2 × 150 bp | 2 × 150 bp | 2 × 250** |

# Short read sequencing

UNIVERSITY OF AMSTERDAM
Life Sciences

# Short read sequencing

UNIVERSITY OF AMSTERDAM
Life Sciences

14

# Long read sequencing: Pacbio SMRT



Eid, J., et al. (2009) Science, 323(5910), 133–138.

UNIVERSITY OF AMSTERDAM
Life Sciences

# Long read sequencing: Pacbio SMRT



Start with high-quality double stranded DNA

Ligate SMRTbell adapters and size select

Anneal primers and bind DNA polymerase

Circularized DNA is sequenced in repeated passes

The polymerase reads are trimmed of adapters to yield subreads

Consensus is called from subreads

**HiFi READ**
>99.9% accuracy

UNIVERSITY OF AMSTERDAM
Life Sciences

# Long read sequencing: ONT



Sample preparation

Inhomogeneous translocation speed
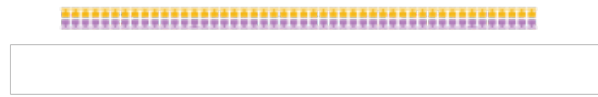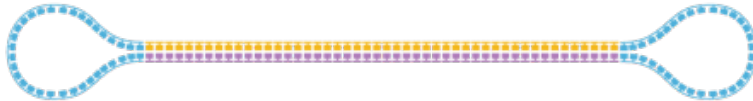
Low signal-to-noise

Multiple nucleotides simultaneously

Updates to MinION chemistry

Nanopore sequencing

current (pA)

Time (ms)

Signal detection

Loss of information by event-calling

Faulty segmentation

Difficulty predicting homopolymer length

Unrepresentative training data

Ignoring DNA modifications

HMM- vs RNN-based

Raw base calling

Optimized training data

Modeling strand progression

A T T C A G C

Raw signal processing(optional) and base calling

Consensus calling

Polishing from raw data

Post-sequencing processing

Rang, F.J., Kloosterman, W.P. & de Ridder, J. Genome Biol 19, 90 (2018)
Xu, L., Seki, M. J Hum Genet 65, 25–33 (2020)

UNIVERSITY OF AMSTERDAM
Life Sciences

# Quality?

```
GGGGGGGG9BGGGGGGFFGGGGG#######::DGGGGGFGGGGGGGGGGGGGGGG#9CFGGGGGGGG#:DFG#:#######:###
GF#######/2/##2############0-0###)2###*0/#01<DFF7#)07F:FEF:FFFFFFFFFF<BFFBF?<7?FFFFF
```

Phred+33 score:

```
 !"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHI
  |                            |    |      |
 0.2.........................26...31........
```

# Quality scoring

An important technical aspect of our work is the use of log-transformed error probabilities rather than untransformed ones, which facilitates working with error rates in the range of most importance (very close to 0). Specifically, we define the quality value $q$ assigned to a base-call to be

$$q = -10 \times \log_{10}(p)$$

where $p$ is the estimated error probability for that base-call. Thus a base-call having a probability of $1/1000$ of being incorrect is assigned a quality value of 30. Note that high quality values correspond to low error probabilities, and conversely.

UNIVERSITY OF AMSTERDAM
Life Sciences

# Quality scoring



```
!"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHI
 |                             |   |      ||
0.2..........................26...31.......
```

| error | | 1 | 0.63 | 0.1 | 0.01 | 0.0025 | 0.00079 | 0.0001 |
| probability | | | | | | | | |

$$q = -10 \times \log_{10}(p)$$

# Quality scoring

```
SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS.................................................
..................................XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX.............
.............................................IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII...............
..............................................JJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJ...............
..LLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLL.................................................
!"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz{|}~
|                                 |    |       |                                         |        |
33                                59   64      73                                        104      126
0...............................26...31.......40
                       -5....0.........9...............................40
                            0.........9...............................40  ←
                            3.....9...............................40
0.2.....................26...31........41                                  ←
```

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)
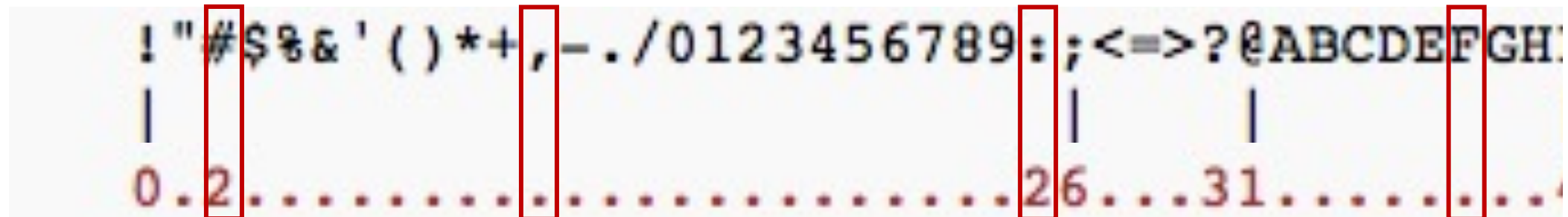
# Quality scoring

novaseq data:

```
@SRR15010442.1
CCTGTTTGCTCCCCACGCTTTCGCGCCTCAGCGGCAGTTACAGACCAAAAAGCCGCCTTCGCCACTGGTGTTC
CTCCACATCTCTACGCATTTCACCGCTACACGTGGAATTCTACCCCCC
+
F:FFFF,FF:FFFF:FFFFFFFFFFFFFFF:FFF,FFFFF,FFFFFFFF:FFFFFFFFFFFFFF:FFFFFF:FF
:FFFFFFFFFF,FF:F:FFFFFFFF:F:F:FFFFFF,F,FF,FFFFFF
```

# Quality scoring

novaseq data:

```
@SRR15010442.1
CCTGTTTGCTCCCCACGCTTTCGCGCCTCAGCGGCAGTTACAGACCAAAAAGCCGCCTTCGCCACTGGTGTTC
CTCCACATCTCTACGCATTTCACCGCTACACGTGGAATTCTACCCCCC
+
F:FFFF,FF:FFFF:FFFFFFFFFFFFFF:FFF,FFFFF,FFFFFFFF:FFFFFFFFFFFFFF:FFFFFF:FF
:FFFFFFFFFF,FF:F:FFFFFFFF:F:F:FFFFFF,F,FF,FFFFFF
```



```
error         0.63              0.079                      0.0032            0.0002
probability
```
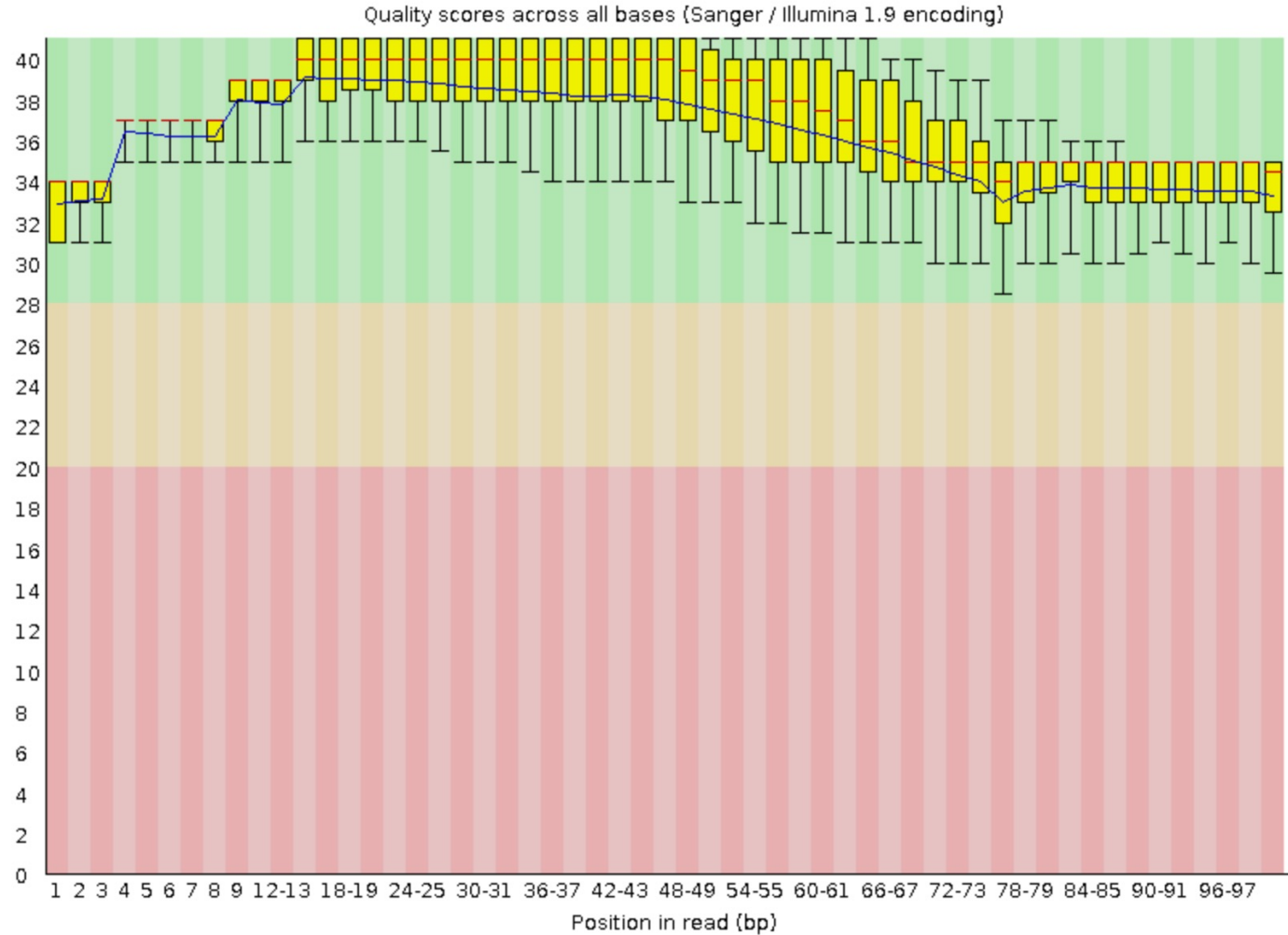
# Quality reports

summary of 1 dataset:



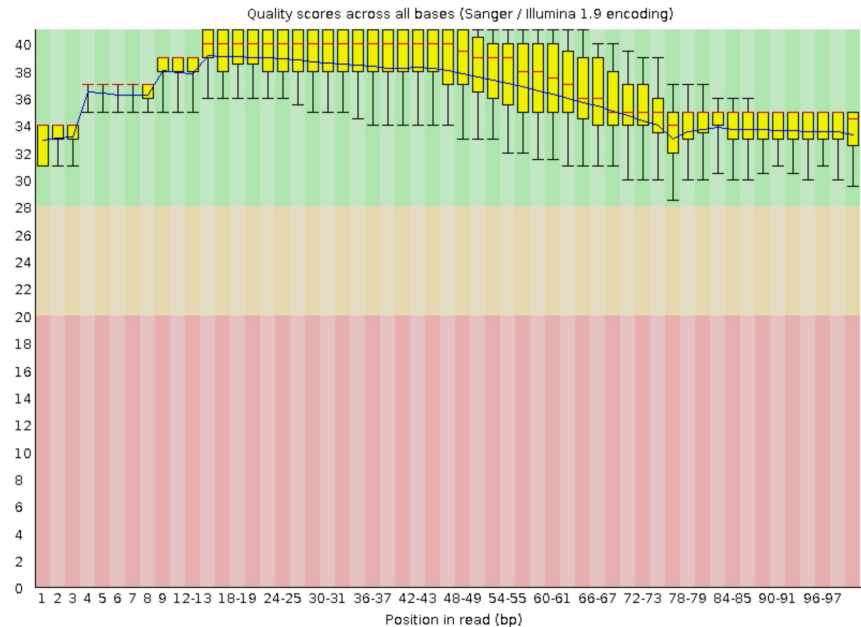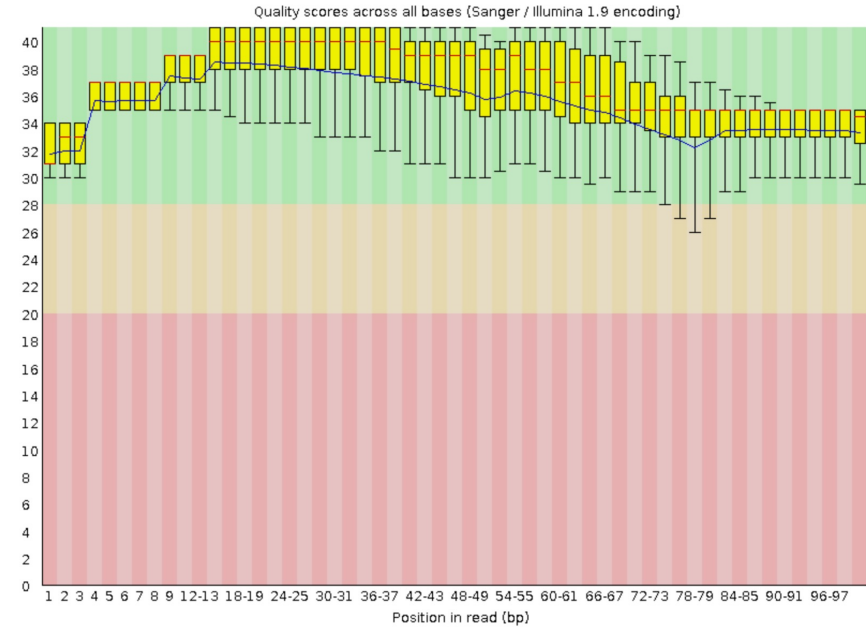Quality scores across all bases (Sanger / Illumina 1.9 encoding)

# Quality reports

summary of 1 dataset:

forward reads:

reverse reads:

# Quality reports

summary of 1 dataset:

UNIVERSITY OF AMSTERDAM
Life Sciences

# Quality reports

summary of 1 dataset:

UNIVERSITY OF AMSTERDAM
Life Sciences
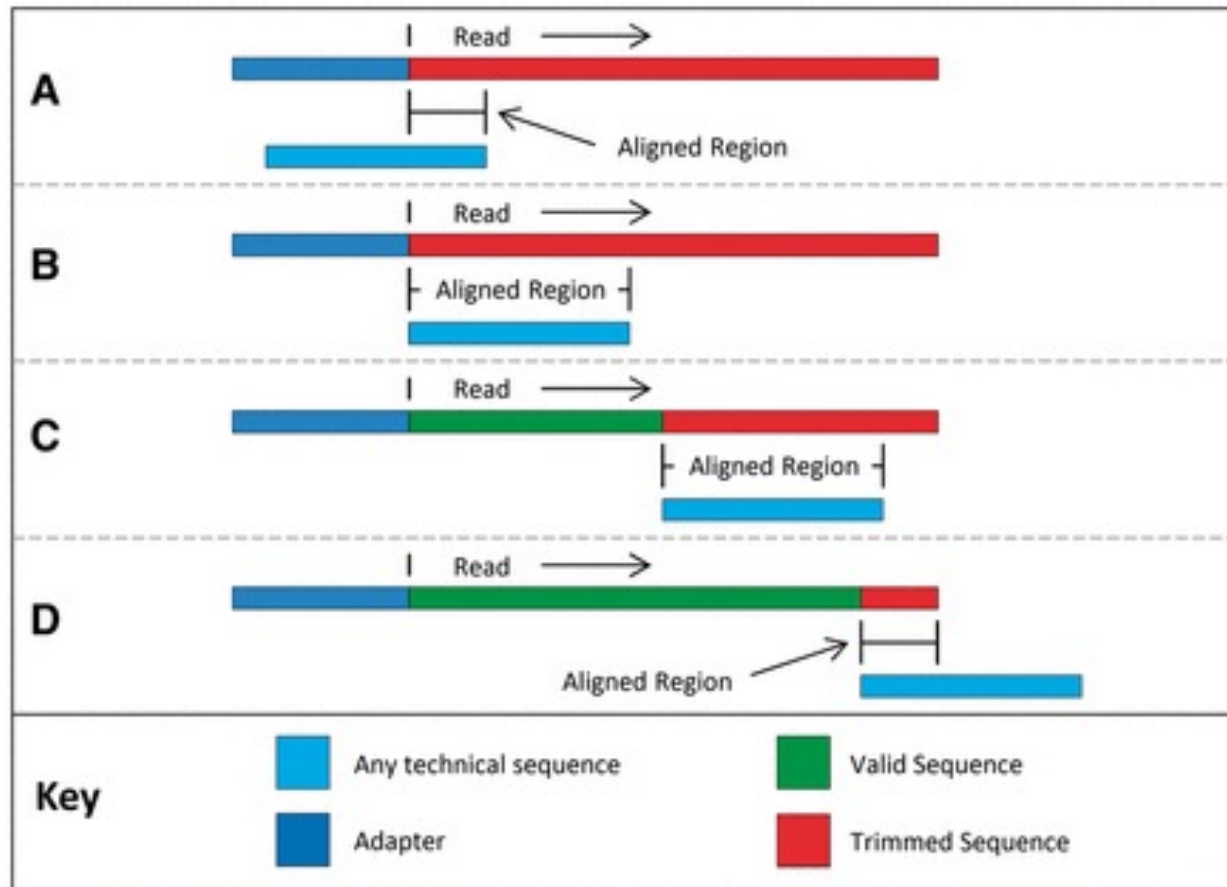
# Data preprocessing: filtering & trimming

- remove adapter sequences

- remove low-quality ends

- remove dark-cycle poly-G ends

UNIVERSITY OF AMSTERDAM
Life Sciences

# Data preprocessing : filtering & trimming

UNIVERSITY OF AMSTERDAM
Life Sciences

# Data preprocessing – remove contaminants!

- remove uninformative sequences:

- phiX spike-in

- host genome

- for rRNA-depleted RNAseq: remove rRNA

UNIVERSITY OF AMSTERDAM
Life Sciences

# Data preprocessing – remove contaminants!

SIGS **Standards in Genomic Sciences**

**COMMENTARY**                                          **Open Access**

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

Supratim Mukherjee[1*], Marcel Huntemann[1], Natalia Ivanova[1], Nikos C Kyrpides[1,2] and Amrita Pati[1]

## Contamination in Reference Sequence Databases: Time for Divide-and-Rule Tactics

*Valérian Lupo[1,2], Mick Van Vlierberghe[1], Hervé Vanderschuren[3], Frédéric Kerff[2], Denis Baurain[1*] and Luc Cornet[1,3]*

**Genome Biology**

PLOS COMPUTATIONAL BIOLOGY

RESEARCH ARTICLE
## Removing contaminants from databases of draft genomes

**Jennifer Lu[1,2*], Steven L. Salzberg[1,2,3]**

**METHOD**                                          **Open Access**

## Terminating contamination: large-scale search identifies more than 2,000,000 contaminated entries in GenBank

Martin Steinegger[1,2,3*] and Steven L. Salzberg[2,4,5]

UNIVERSITY OF AMSTERDAM
Life Sciences

# Thanks for your attention!

a.u.s.heintzbuschart@uva.nl

SP C2.205

github.com/a-h-b

twitter.com/_a_h_b_

UNIVERSITY OF AMSTERDAM
Life Sciences