# Bio 296: Microbial Bioinformatics Introduction to DNA extraction and Sequencing Technologies
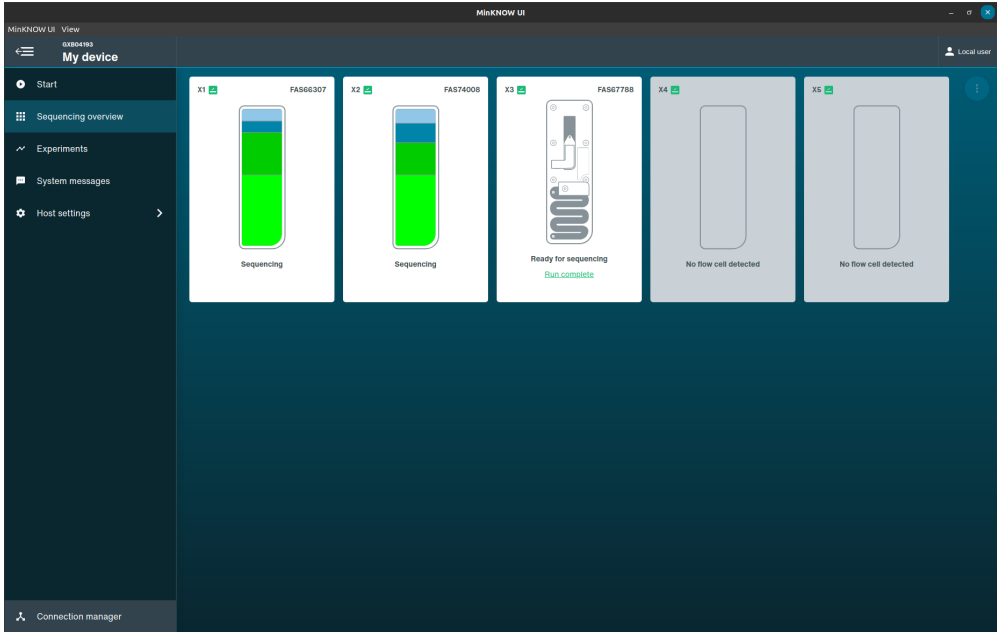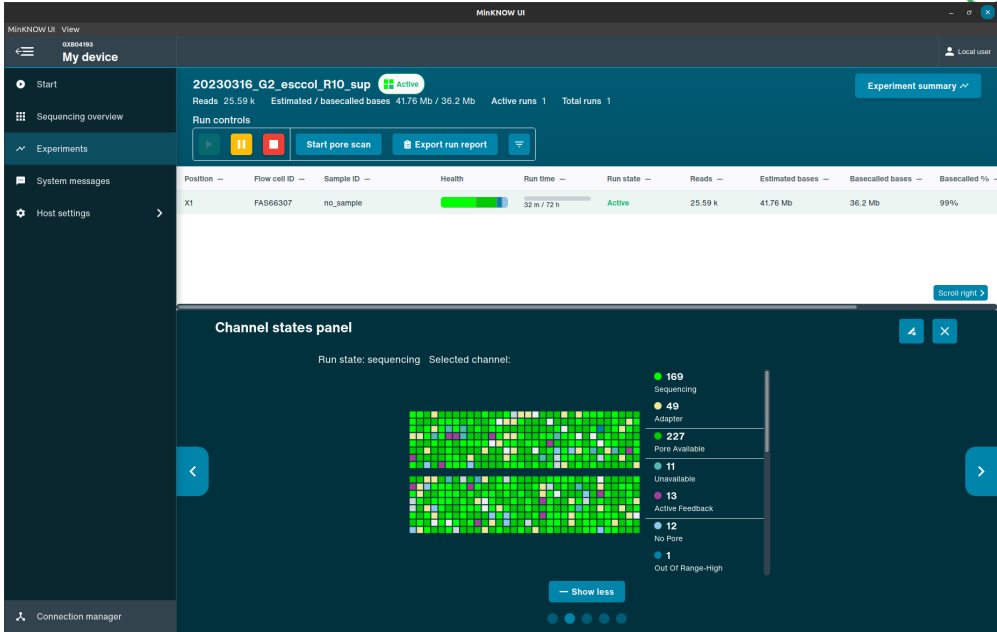
Tim Roloff
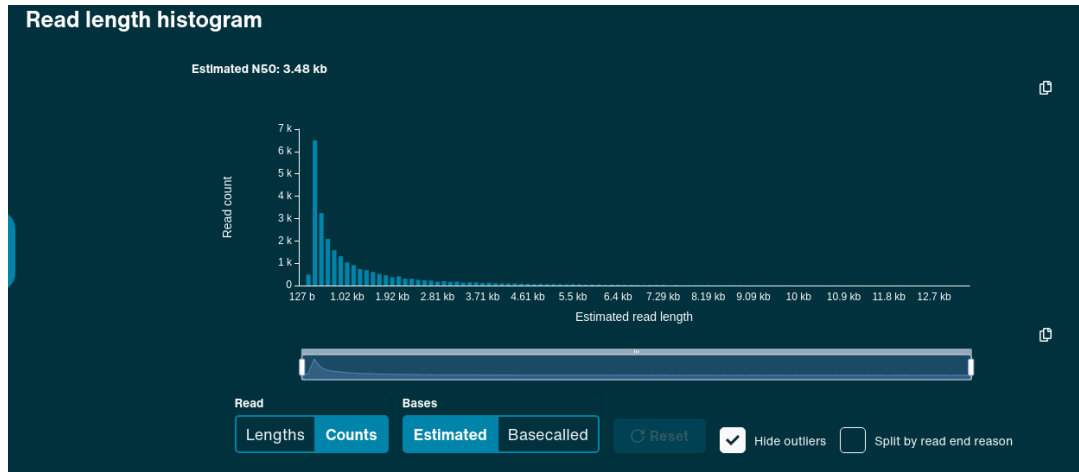
# Pores sequencing
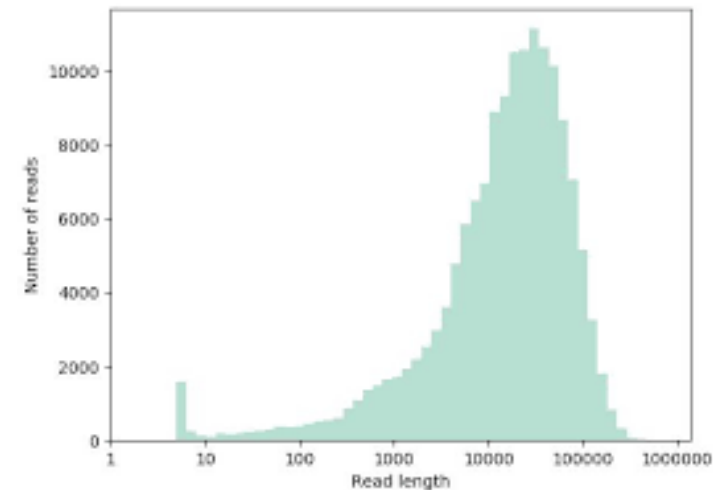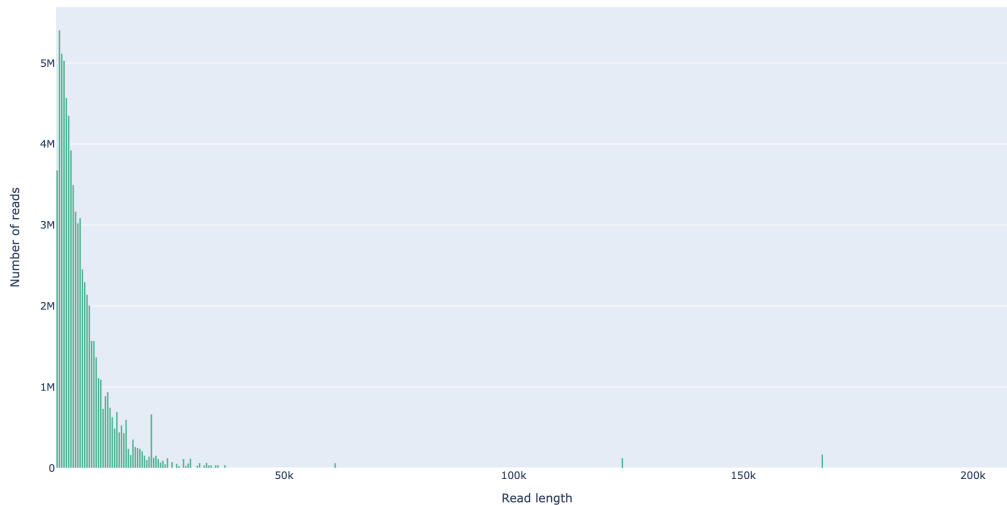
Green is good!

For a good run <mark>Sequencing</mark> > <mark>Pores available</mark>
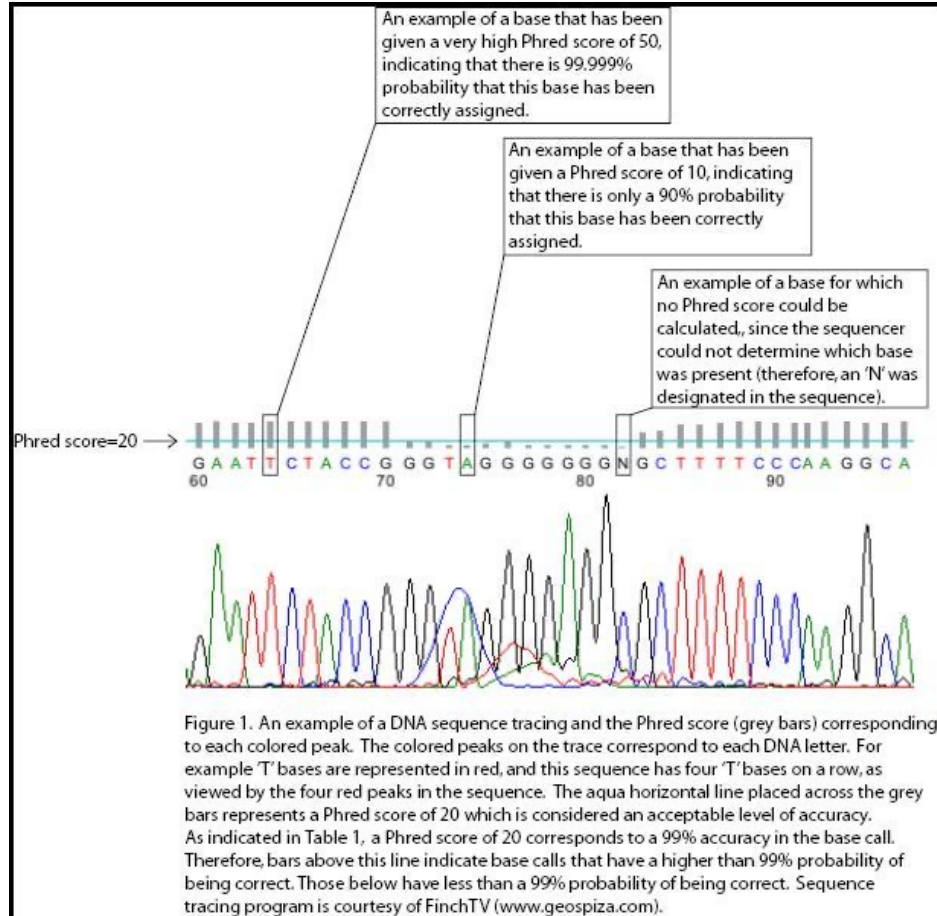
# Read length distribution



- Shows the distribution of read lengths

- For Illumina data typically 1 bar as all reads are of the same length

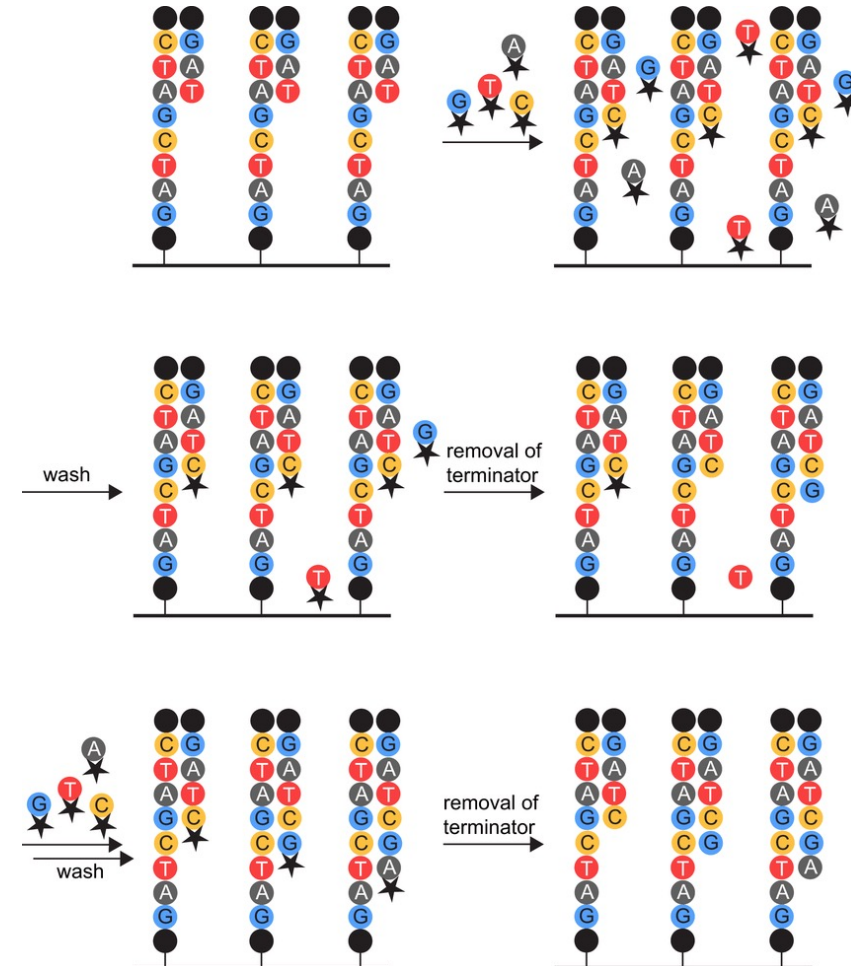- For Nanopore data heavy tailed

# Phred score

Origin of Phred score in Sanger sequencing data



Figure 1. An example of a DNA sequence tracing and the Phred score (grey bars) corresponding to each colored peak. The colored peaks on the trace correspond to each DNA letter. For example 'T' bases are represented in red, and this sequence has four 'T' bases on a row, as viewed by the four red peaks in the sequence. The aqua horizontal line placed across the grey bars represents a Phred score of 20 which is considered an acceptable level of accuracy. As indicated in Table 1, a Phred score of 20 corresponds to a 99% accuracy in the base call. Therefore, bars above this line indicate base calls that have a higher than 99% probability of being correct. Those below have less than a 99% probability of being correct. Sequence tracing program is courtesy of FinchTV (www.geospiza.com).

– Pre-/Post-phasing in Illumina data

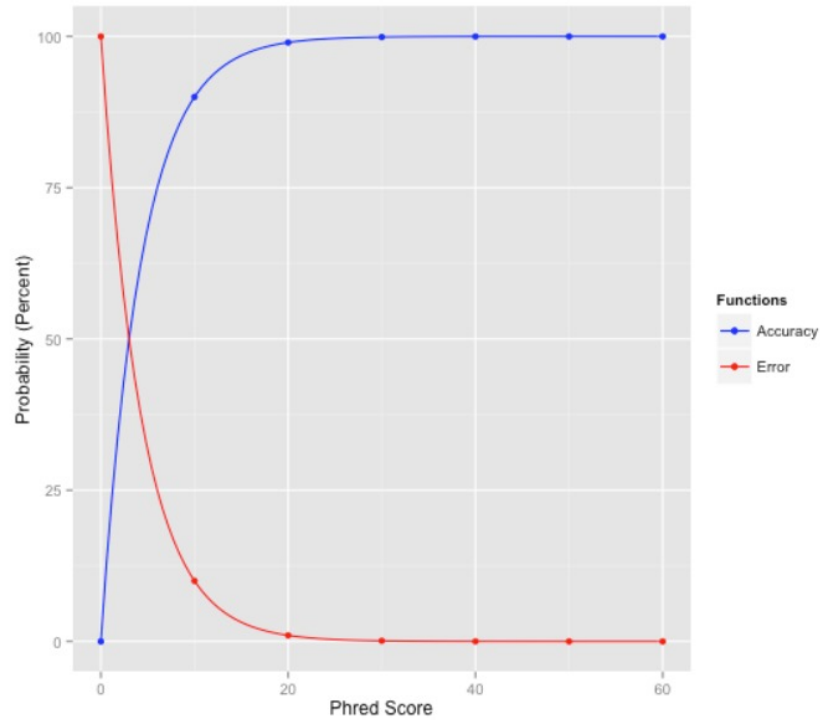# Phred score

– Q-score or Phred Quality Score

– Q30 considered gold standard for Illumina sequencing

  – 1 error in 1000 sequenced bases

Phred score is inversly correlated to accuracy

$$E = 10^{-\left(\frac{Q}{10}\right)}$$

$$Q = -10 \log E$$



| Phred Quality Score | Error | Accuracy (1 - Error) |
|---|---|---|
| 10 | 1/10 = 10% | 90% |
| 20 | 1/100 = 1% | 99% |
| 30 | 1/1000 = 0.1% | 99.9% |
| 40 | 1/10000 = 0.01% | 99.99% |
| 50 | 1/100000 = 0.001% | 99.999% |
| 60 | 1/1000000 = 0.0001% | 99.9999% |

# Q score for Nanopore data

– For Nanopore data there are no intensities to calculate Q-scores from

– Nanopore per base Q-scores are calculated based on then output of the neural network used by the base caller (e.g. guppy)

– Nanopore single read accuracies can be calculated by aligning reads to a reference sequence

$$\frac{N_{matches}}{N_{matches} + N_{mismatches} + N_{deletions} + N_{insertions}}$$

– For more details and formulas see https://labs.epi2me.io/quality-scores/

# Chastity filter

Illumina only

% PF = clusters passing filter

$$\frac{Ia}{Ia + Ib} > 0.6$$

Illumina run statistics show % PF

Different machines / sequencing
kits have different % PF cutoffs

Ia = Intensity of brightest base

Ib = Intensity of second brightest base

# FASTA file format

Start            Unique Sequence Header

```
1 - >BY999847.1 BY999847 Moon Jellyfish cDNA library Aurelia aurita cDNA
      clone Aa_plW_142145_H14, mRNA sequence
2 - AAAATACCGCATGATTGTTCGTTTCACAAACAAAGATATAGCTTGCCAGATAGCGTATGCCAGATTGCAA
3 - GGAGATGTGATCATTTGTGCAGCTTATGCTCATGAACTCCCAAGATATGGTGTCAAGGTCGGGTTGACCA
4 - ACTATGCAGCTGCTTATTGCACTGGCCTCTTGCTCGCAAGAAGGCTCCTTTCAAAATTGAAATTGGCTGA
5 - CACTTACAAAGGTTGTGAAGAAGTGAATGGTGATGAATACCTTGTGGAAGGAGAGGAGGGACAGCCTGGA
6 - CCTTTCCGTTGTTACCTTGATATTGGCCTTGCCAGAACCTCAACTGGTGCCAAGATCTTTGGTGCATTGA
7 - AAGGTGCAGTTGATGGTGGACTTGACATCCCACACAGCAACACGAGATTCCCTGGTTATGACAATGAAGC
8 - AAAGGAATTTGACCCAGAGGTGCACAGACAACACA...
...
```

Sequence (nucleotide or protein)

File Suffix: sequence(s).fa, sequence(s).fasta

Special cases: sequences.mfa (multiple - aligned - sequences)
                  sequences.afa (aligned sequences)

https://www.gdc-docs.ethz.ch/GeneticDiversityAnalysis/GDA20/handouts/03_GDA20_NGS_QC.pdf

# FASTQ file format

Sequence (Read) Header

Start      Nucleotide Sequence (A, C, G, T, N)      +Sequence ID

```
1 -  @HWI-ST486:166:C06K9ACXX:7:1101:1443:1995  1:N:0:ACAGTG
2 -  GCCCAGCGTGGGCGAGCCGCACGGCACCATCCTCTGGCACACCCTCTCCTC
3 -  +
4 -  BCCFFFFDFHHHFJJJJJJJIIJJJJJIJIIJJGIHHHHHFFFDDDEDDDC
     @HWI-ST486:166:C06K9ACXX:7:1101:2519:1936  1:N:0:ACAGTG
     NTGCACACGAATTCGCCGTTGTGGCAGCTCGAGTACCTGTGGTTCGCCGAG
     +
     #1=DDDFFHHHFFIGIIJJJJJIIJIIJ;?*?.?/9*88BBBF.;7@@E###
```

$N_{rows} = 4$

ASCII encoded quality scores per base

File Suffix: reads.fq, reads.fastq

Special cases: read_R[12].fq (> paired reads)
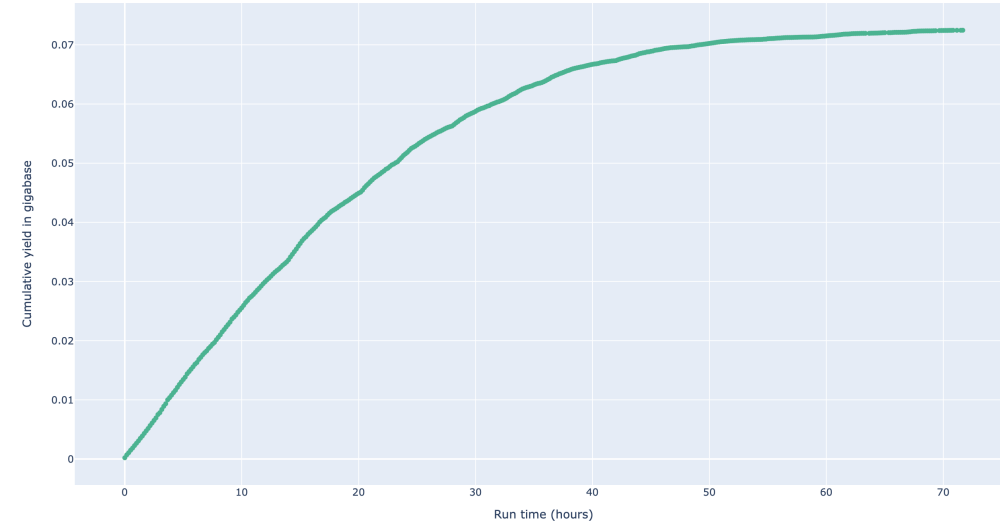
read_I[12].fq (> index)

# Amount of data generated

Data accumulates over time

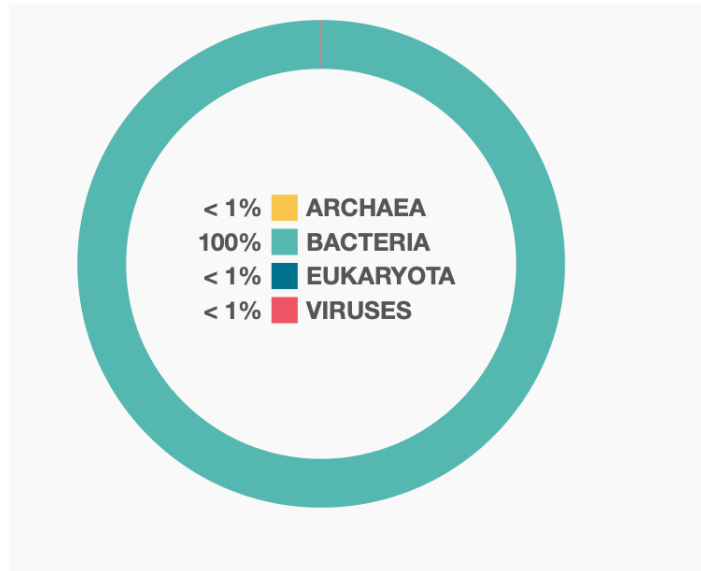Files of 4000 sequences saved

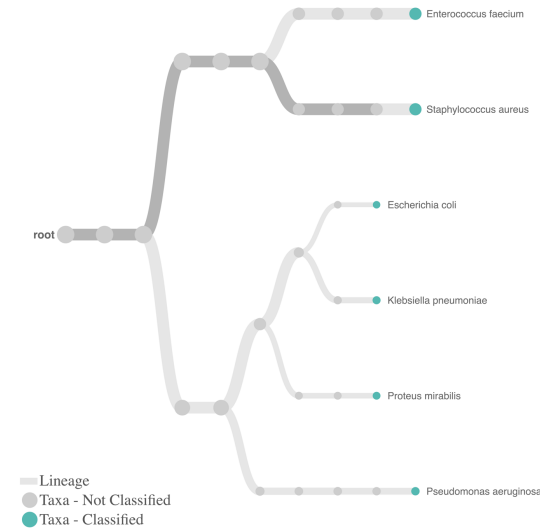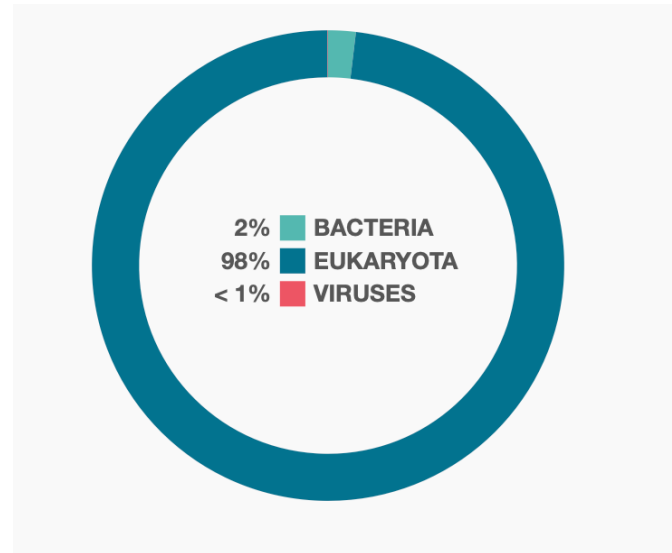Max run time 72 hours

Typical output > 25 GB / flow cell

# Species detected

Pure Culture

Tissue sample



Various approaches to identify species from sequencing data
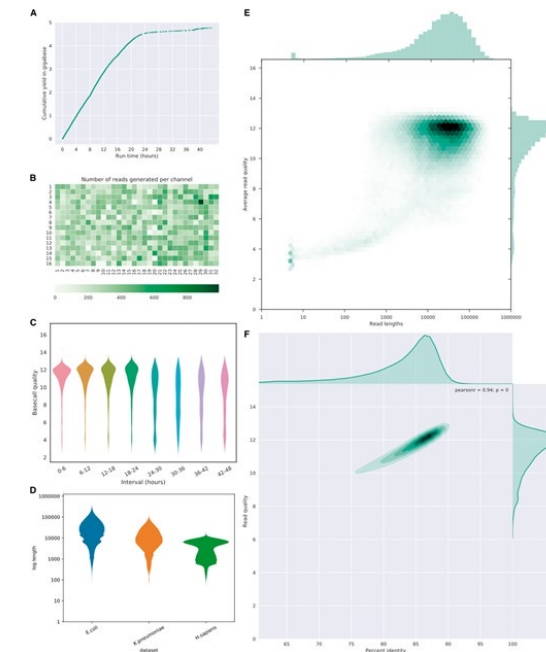rMLST
WIMP - what is in my pot
FASTQ screen – screen against a fiven set of databases
KRAKEN2 – use Kmer approach to
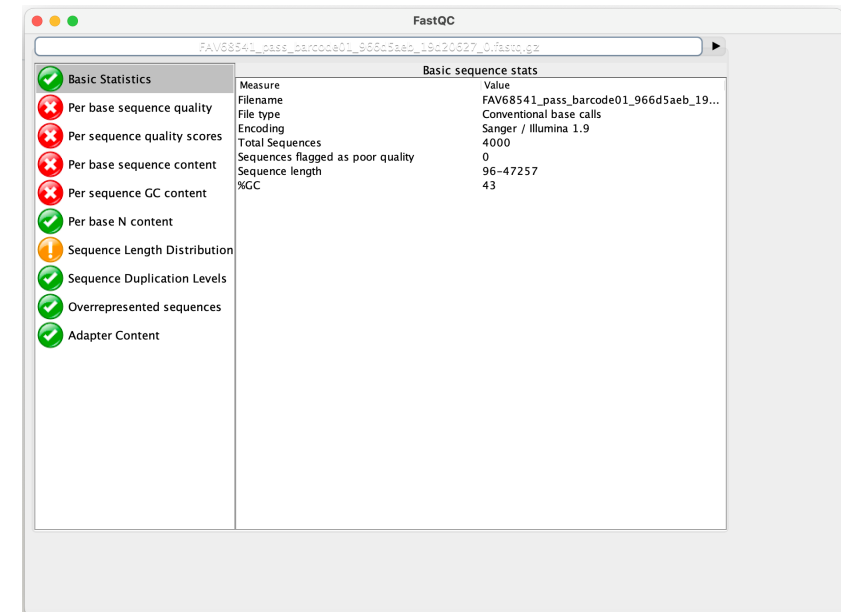
# Software for QC

Nanoplot - tool for quality control analysis for nanopore reads.

- Input - fastq or sequencing summary file.

- Output - HTML report, summary table and plot images

- https://github.com/wdecoster/NanoPlot



FastQC – tool for quality control analysis for fastq files

- Input – fastq files

- Output – HTML report with different graphs

- https://www.bioinformatics.babraham.ac.uk/projects/download.html#fastqc

# Software from ONT for data analysis

Epi2me – online platform for on the fly data analysis with various modules

EPI2MELabs downloadable version for local analysis



https://epi2me.nanoporetech.com/report-387749

https://labs.epi2me.io/downloads/

Useful worklfows:

FASTQ – WIMP

- QC of reads

- What is in my pot (WIMP): bacterial classification using Centrifuge classification engine by Johns Hopkins University

FASTQ – Antimicrobial Resistance

- QC of reads

- What is in my pot (WIMP): bacterial classification using Centrifuge classification engine by Johns Hopkins University

- Antimicrobial resistance: search against CARD database

# Some examples



– What has happened here?