# NGS - quality control, alignment, visualisation

Quality control + database retrieval

# Why Quality control?

1. How is the base quality?
2. What is the read length?
3. Are there adapters/barcodes in my sequences?
4. Are there overrepresented sequences?

# Dedicated software

- Manufacturers' software
- Illumina: fastQC
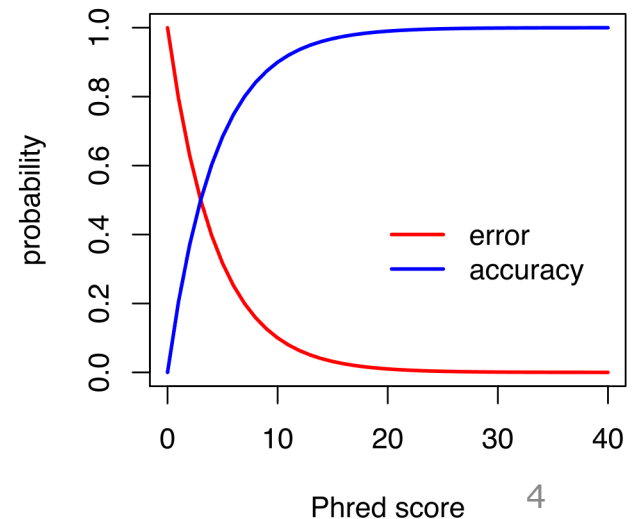- ONT: pycoQC
- ONT + PacBio: NanoPlot
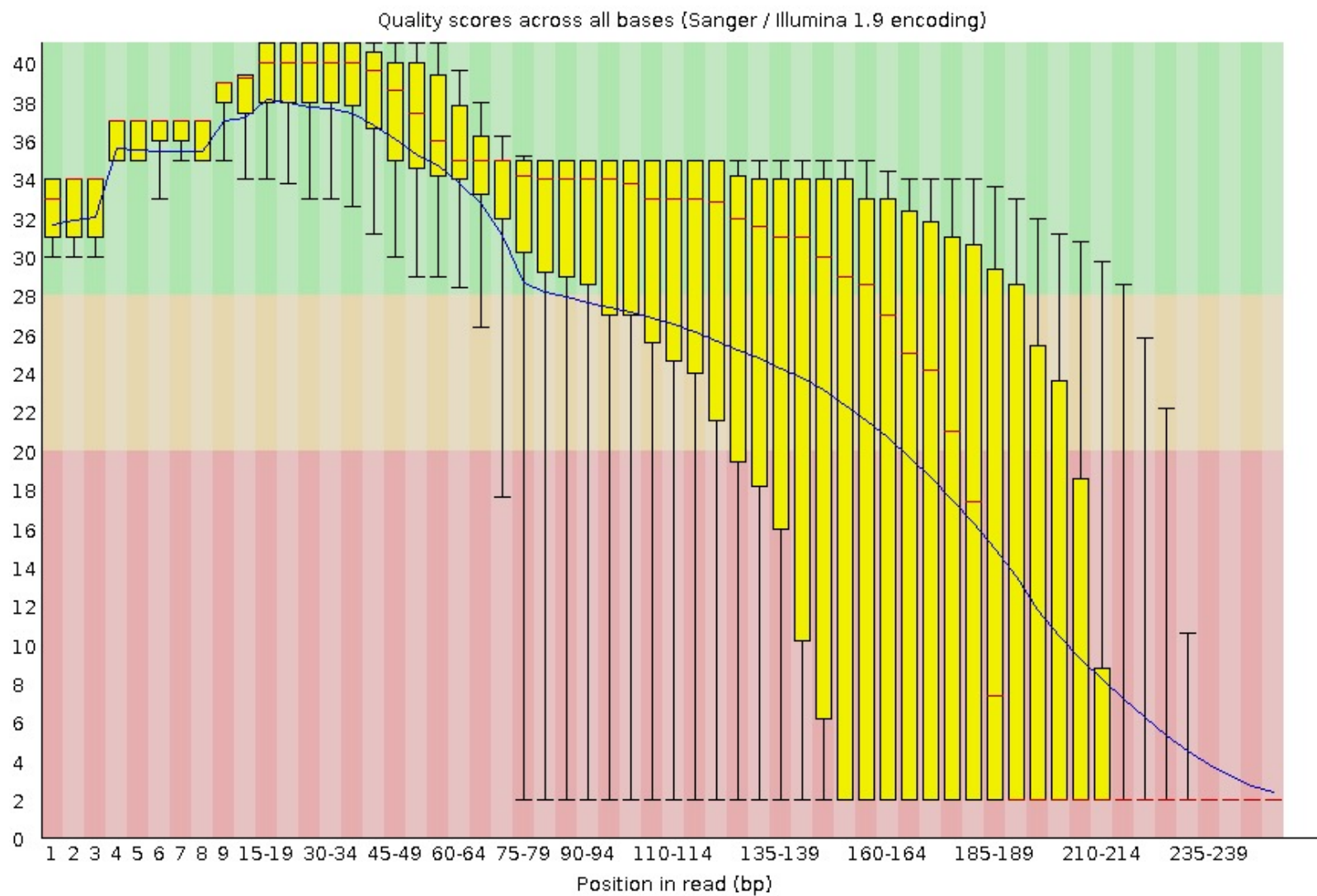
# fastq

fasta + basequality (fasta + q = fastq)

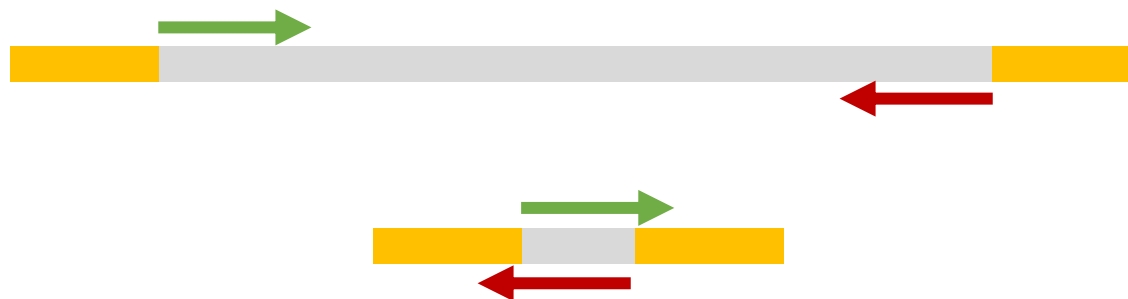$$BASEQ = -10log_{10} \Pr\{base\ is\ wrong\}$$

$$-10log_{10}(0.01) = 20$$
$$-10log_{10}(0.1) = 10$$
$$-10log_{10}(0.5) = 3$$

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

Position in read (bp)

5

# Trimming

- Find and remove:
  - Regions or reads with low base quality
  - Adapter sequences
- Software: `cutadapt` (or `trimmomatic`, `trim_galore`, `bbduk` ..)

# Databases



**INSDC**: **I**nternational **N**ucleotide **S**equence **D**atabase **C**ollaboration  8

# BioProject (Former DRA Study)

**BioProject** <span style="color:red">PRJD</span>
- Project description
- Grants
- Publications

# BioSample (Former DRA Sample)

**BioSample** <span style="color:red">SAMD</span>

**BioSample** <span style="color:red">SAMD</span>

**BioSample** <span style="color:red">SAMD</span>
- Sample description
- Taxonomy ID

# Sequence Read Archive

**Experiment** <span style="color:red">DRX</span>
- Library layout
- Sequencing platform

**Run** <span style="color:red">DRR</span>

**Run** <span style="color:red">DRR</span>

**Run** <span style="color:red">DRR</span>
- Data files

Sequence data files (fastq, BAM)

DDBJ

<span style="color:red">Prefix of accession number</span>

9

# Command line tools

- Retrieve raw data: SRA-tools
  - `prefetch`
  - `fastq-dump`

- Retrieve sequences: Entrez Direct
  - `esearch`
  - `efetch`