

# **BIOL 343**

# **Applied Bioinformatics I**

## **FASTQ files**

**Dr. Nic Wheeler**

# Learning Objectives

You will be able to:

- 1.

# Stages of sequencing data analysis

Primary, secondary, and tertiary

## Primary

Occurs on the machine in real-time

Provides base calls and quality scores

Illumina Connected Analytics

## Secondary

*Demultiplexing and compression*  
(DRAGEN)

Alignment

Assembly

Specialized tools selected by the end user

## Tertiary

Questions of biological relevance

Differential expression, variant calling, etc

Specialized tools selected by the end user

# Stages of sequencing data analysis

Primary, secondary, and tertiary

## Primary

Occurs on the machine in real-time

Provides base calls and quality scores

Illumina Connected Analytics

Produces FASTQ files

## Secondary

*Demultiplexing and compression*  
(DRAGEN)

Alignment

Assembly

Specialized tools selected by the end user

## Tertiary

Questions of biological relevance

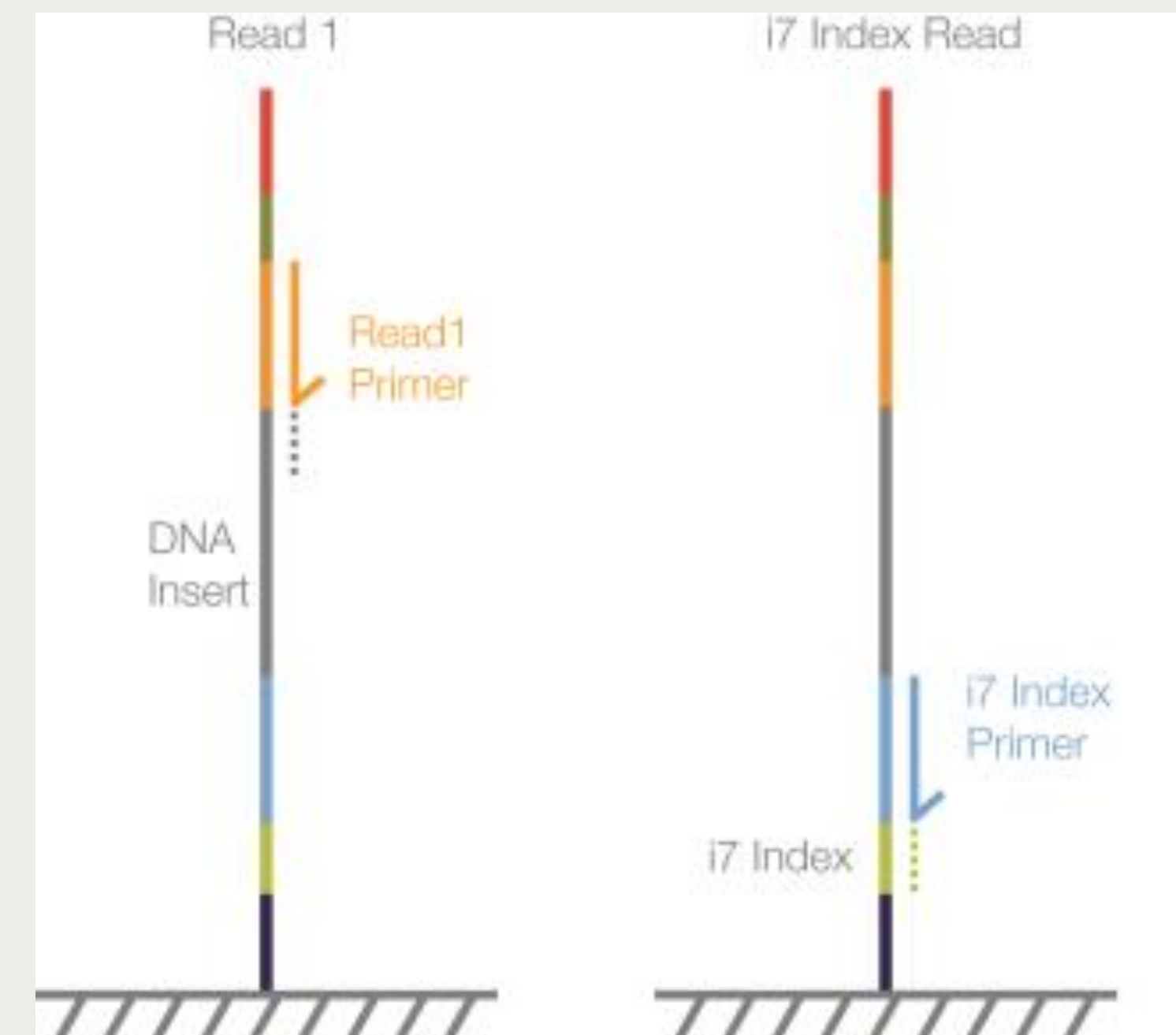
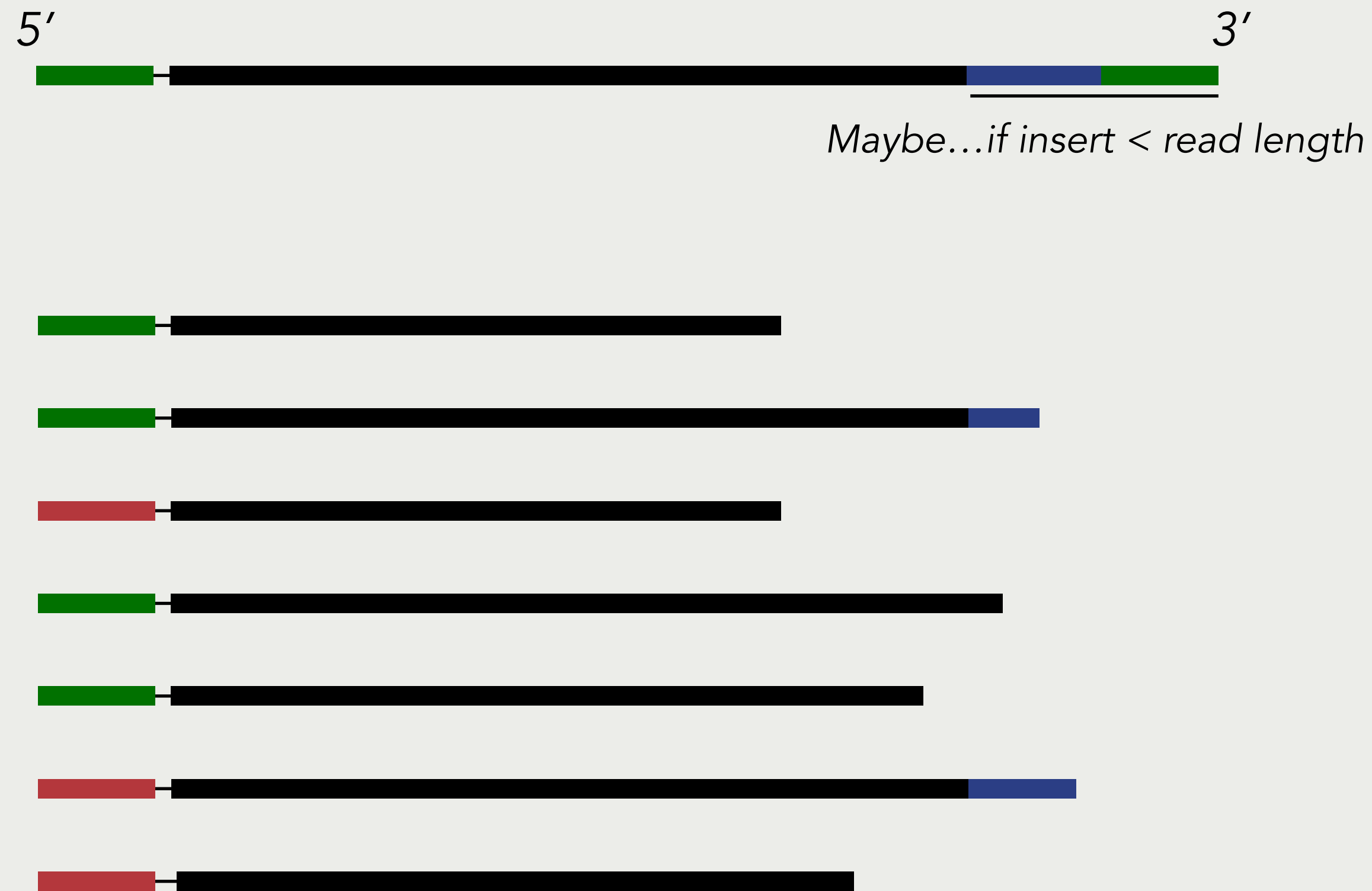
Differential expression, variant calling, etc

Specialized tools selected by the end user

# Stages of sequencing data analysis

## Primary analysis produces FASTQ data and demultiplexes it

Read 1



# Stages of sequencing data analysis

## Primary analysis produces FASTQ data and demultiplexes it

*Read 1*

FASTQ

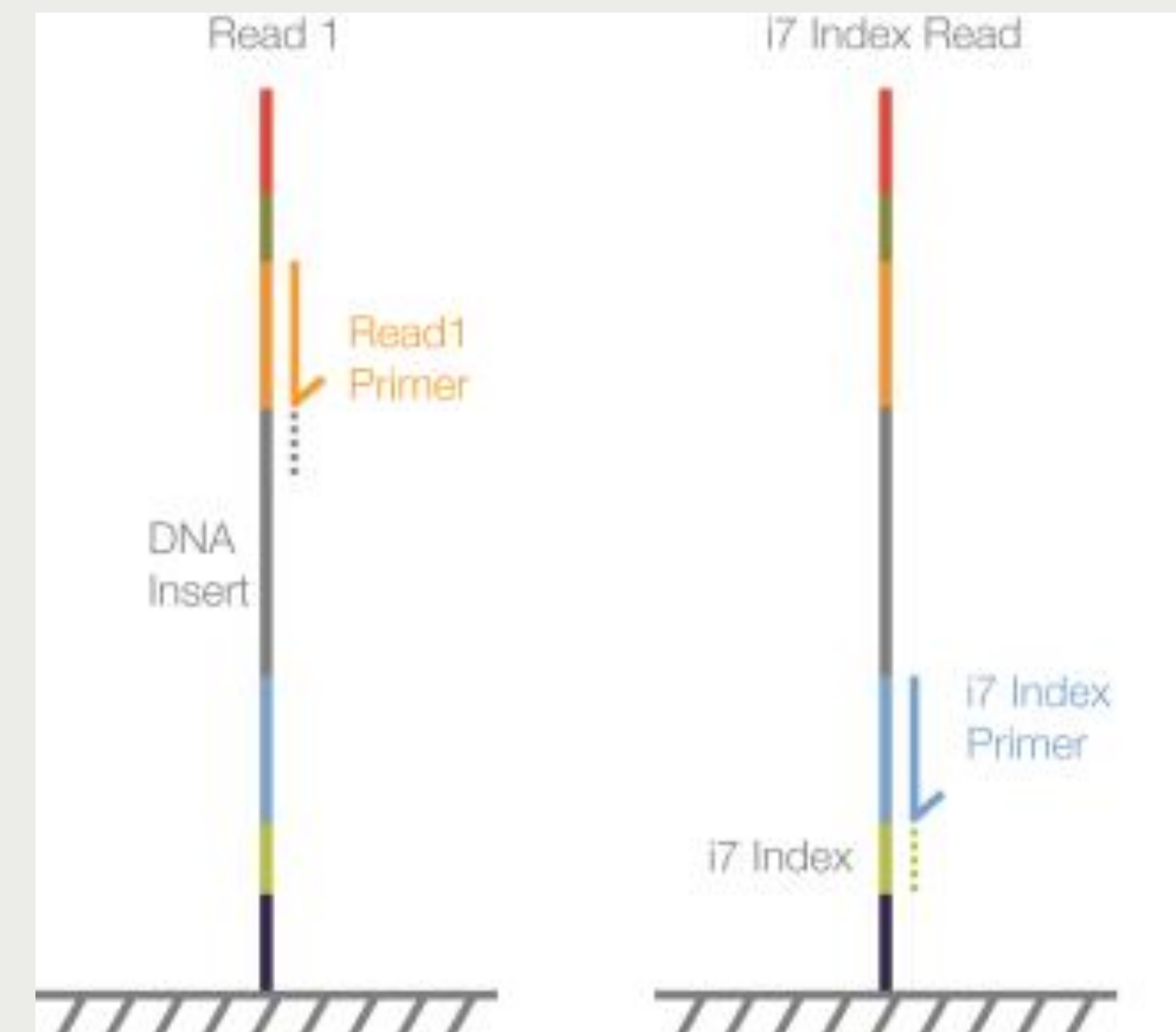


*Client 2*

FASTQ



*Client 1*



# Types of sequence storage files

## FASTA and FASTQ

### FASTA

- Used for storing generic primary sequence data
- Nucleotide or amino acids (ASCII format)
- Usually KB to MB in size
- Header defined by >
  - Sequence name and attributes
  - Spaces discouraged
- Sequence in the line immediately after
  - Can have new line characters (\n)

```
>Smp_104210.1 RH0 cdna:protein_coding
TACTTCATTTATCATTCTGGTAAGTAATGAGTTAACTAAA
TCTGTTCAATTTGTTTCACTAAATTTTAAATCAGAAATTTC
TTTTTTTTTAACACTATTTCTAAACTGTTAAATGCACATTT
ATTTTTC AATTCGTTTAATATCTAGTAGAGTAATCAGTC
TATGTTATTTTAATGAGAATCCTCATTAAAATACATTTCA
GATACTTGTTGAGTTCAATTGAAAAACATTCTCAGAAGGG
GTTTTGTGGAGATTTTCAGTATTTTCATAGTTGAAATCATG
AGTCATTTGAAGCTAAACCCCCATGGAAAACCTAGGAGCA
ATGGACGGCCGTCTCGTTGTATTGTGAGACTCCTCAGCAG
TACCCATCCACGATCCCGCCTCGTGAGATTCGAACCCAGG
ATCTACCAGTCTCGCGCCAGAGCGCTTAACCACTAGATAT
TCTTTTATTATGTTAGGAAGTAATAAAAGTTTTTCTTGAG
```



# Types of sequence storage files

## FASTA and FASTQ

# FASTQ

- Used for storing sequencing data and quality
- ASCII format
- Usually MB to GB in size
- Four lines per sequence:
  1. Header (@ instead of >)
  2. Sequence
  3. + (sometimes the seq id)
  4. Phred quality scores
- Regex for the block: @<seqname>\n<seq>\n+[<

```
@SRR26691082.1 NB501229:521:H5HFKBGXG:1:11101:4638:1052_TAACAN length=69  
TATAAGCACAACTGTCTATCACAAGAAAACCAACGAGATATTGCTTGTAACCTTTCTATTTAGAAACG  
+SRR26691082.1 NB501229:521:H5HFKBGXG:1:11101:4638:1052_TAACAN length=69  
EEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEAEEEEEEEEEEE/EEEEEEEEEEEEEEEEEEEE/EEAEAA  
@SRR26691082.2 NB501229:521:H5HFKBGXG:1:11101:14456:1055_ATTATN length=69  
TATACAAACTAAGACAGATATTTTTTTTAGATTTAGCTATCTTCAACTTTCATTCCATTATTGTAATTA  
+SRR26691082.2 NB501229:521:H5HFKBGXG:1:11101:14456:1055_ATTATN length=69  
EEEEEEEEEEEEEEEEEEEE6EEEEEEEEEEEEEEEEEEEEAEEEEEEE/EEEEEEEEEEEEEEEEEEEEEA  
@SRR26691082.3 NB501229:521:H5HFKBGXG:1:11101:20997:1055_ATCGCN length=69  
TATATGATAATTGGTGGTTCCAGACGTGCAGCTTGGAAGCGCCTCAATACTGTAGAGATGCGACGAAAG
```



# Types of sequence storage files

## FASTA and FASTQ

# FASTQ Header

- @SRR26691082.1 - SRR ID and read #
- NB501229 - instrument name
- 521 - Run ID
- H5HFKBGXG - flow cell ID
- 1 - flow cell lane
- 11101 - flow cell tile #
- 4638 - x coordinate in the tile
- 1052 - y coordinate in the tile
- TAACAN - index sequence

```
@SRR26691082.1 NB501229:521:H5HFKBGXG:1:11101:4638:1052_TAACAN length=69  
TATAAGCACAACTGTCTATCACAAGAAAACCAACGAGATATTGGCTTGTAACCTTTCTATTTAGAAACG  
+SRR26691082.1 NB501229:521:H5HFKBGXG:1:11101:4638:1052_TAACAN length=69  
EEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEAEEEEEEEEEEE/EEEEEEEEEEEEEEEEEEEE/EEAEAA  
@SRR26691082.2 NB501229:521:H5HFKBGXG:1:11101:14456:1055_ATTATN length=69  
TATACAAACTAAGACAGATATTTTTTTTAGATTTAGCTATCTTCAACTTTCATTCCATTATTGTAATTA  
+SRR26691082.2 NB501229:521:H5HFKBGXG:1:11101:14456:1055_ATTATN length=69  
EEEEEEEEEEEEEEEEEEEE6EEEEEEEEEEEEEEEEEEEEEEEEEEEE/EEEEEEEEEEEEEEEEEEEEAAAA  
@SRR26691082.3 NB501229:521:H5HFKBGXG:1:11101:20997:1055_ATCGCN length=69  
TATATGATAATTGGTGGTTCCAGACGTGCAGCTTGGAAGCGCCTCAATACTGTAGAGATGCGACGAAAG
```

# Types of sequence storage files

## FASTA and FASTQ

### Phred quality

- String of the same length as <seq>
- Each character represents the Phred quality of the corresponding nt
- Represents the likelihood that the base call was correct

$$\$Q = -10\log_{10}(\$e)$$

Where \$e is the error probability

- If the quality of a base call is 30, the probability that it is wrong is 0.001. In other words, given 1000 base calls with Q=30, one of them is wrong in average.
- Minimum of 0 (!) and maximum of 42 (K)

```
@SRR26691082.1 NB501229:521:H5HFKBGXG:1:11101:4638:1052_TAACAN length=69
TATAAGCACAATCTGCTATCACAAGAAAACCAACGAGATATTTGCTTGTAACCTCTTCTATTTAGAAACG
+SRR26691082.1 NB501229:521:H5HFKBGXG:1:11101:4638:1052_TAACAN length=69
EEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEE
@SRR26691082.2 NB501229:521:H5HFKBGXG:1:11101:14456:1055_ATTATN length=69
TATACAACTAAGACAGATATTTTTTTTAGATTTAGCTATCTTCACTTTCATTCCATTATTGTAATTA
+SRR26691082.2 NB501229:521:H5HFKBGXG:1:11101:14456:1055_ATTATN length=69
EEEEEEEEEEEEEEEEEEEE6EEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEE
@SRR26691082.3 NB501229:521:H5HFKBGXG:1:11101:20997:1055_ATCGCN length=69
TATATGATAATTGGTGGTTCCAGACGTGCAGCTTGAAGCGCCTCAATACTGTAGAGATGCGACGAAAG
```

### Phred+33 encoding

- \$Q + 33, take the corresponding ASCII character
- Use [this chart](#) to get \$Q for E, /, and A

# Types of sequence storage files

## FASTA and FASTQ

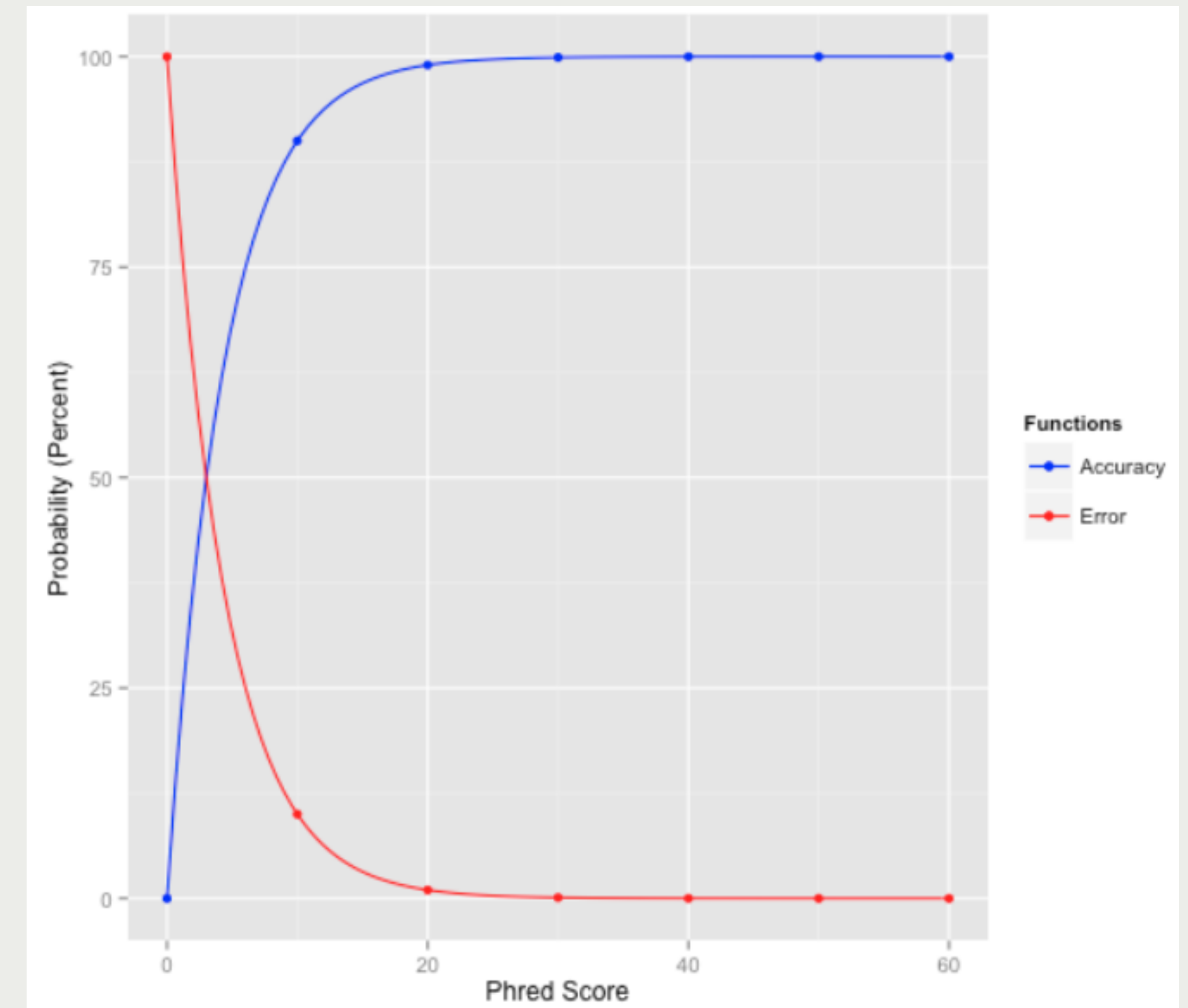
### Phred quality

- String of the same length as <seq>
- Each character represents the Phred quality of the corresponding nt
- Represents the likelihood that the base call was correct

$$Q = -10\log_{10}(e)$$

Where  $e$  is the error probability

- If the quality of a base call is 30, the probability that it is wrong is 0.001. In other words, given 1000 base calls with  $Q=30$ , one of them is wrong in average.
- Minimum of 0 (!) and maximum of 42 (K)



# FASTQ databases

## NCBI SRA & ENA

### NCBI SRA

- <https://www.ncbi.nlm.nih.gov/sra>
- Sequence Read Archive
- Free database of FASTQ (or BAM) files
- Most public sequencing data can be found on SRA
  - Europeans post to ENA
- Dedicated software: SRA Toolkit (conda install bioconda::sra-tools)
  - Already installed in bio1343 environment

