



# From sequencer to polished reads for bacteria

September 2023

# General information

- 2 sessions - 9-13 GMT+2
  - 12th of September
  - 14th of September
- Sessions will be recorded!
- Q&A bottom to ask questions
- Please take your time to evaluate on EVA
- Two tech experts
  - Rodrigo

# Overall objectives

- Inspect the raw fastq files from the sequencer
- Understand how sequences can differ based on the preparation methodology
- Explain the differences between sequencing reads from Illumina and Nanopore
- Assess the quality of sequence data and trim low quality data
- Learn when to decontaminate and when to re-sequence
- Generate an overview of an entire dataset using different tools

# Course intro

## Day 1 Sep 12th

### Theory

- Course intro
- Illumina sequencing theory
- Nanopore sequencing theory
- Sequencing technique comparison

### Practical

- Practical 1 intro
- First look at dataset
- Contamination control
- Raw read QC

## Day 2 Sep 14th

### Interactive

- Session 1 recap
- Interpret raw read QC results
- Interpret contamination control results

### Practical

- Short-read trimming
- Long-read trimming
- Decontamination
- QC parameters for a whole dataset
- Course summary

# Presenters

## Kasper Thystrup Karstensen

### Education:

- Bachelor's in Medical Laboratory Technology
- Master's in Bioinformatics and Systems Biology

### Experience:

- Background in RNA-related work
- Employed at SSI since January 2022, specializing in antimicrobial resistance surveillance

## Astrid Rasmussen

### Education:

- Bachelor's in Biology
- Master's in Bioinformatics and Systems Biology

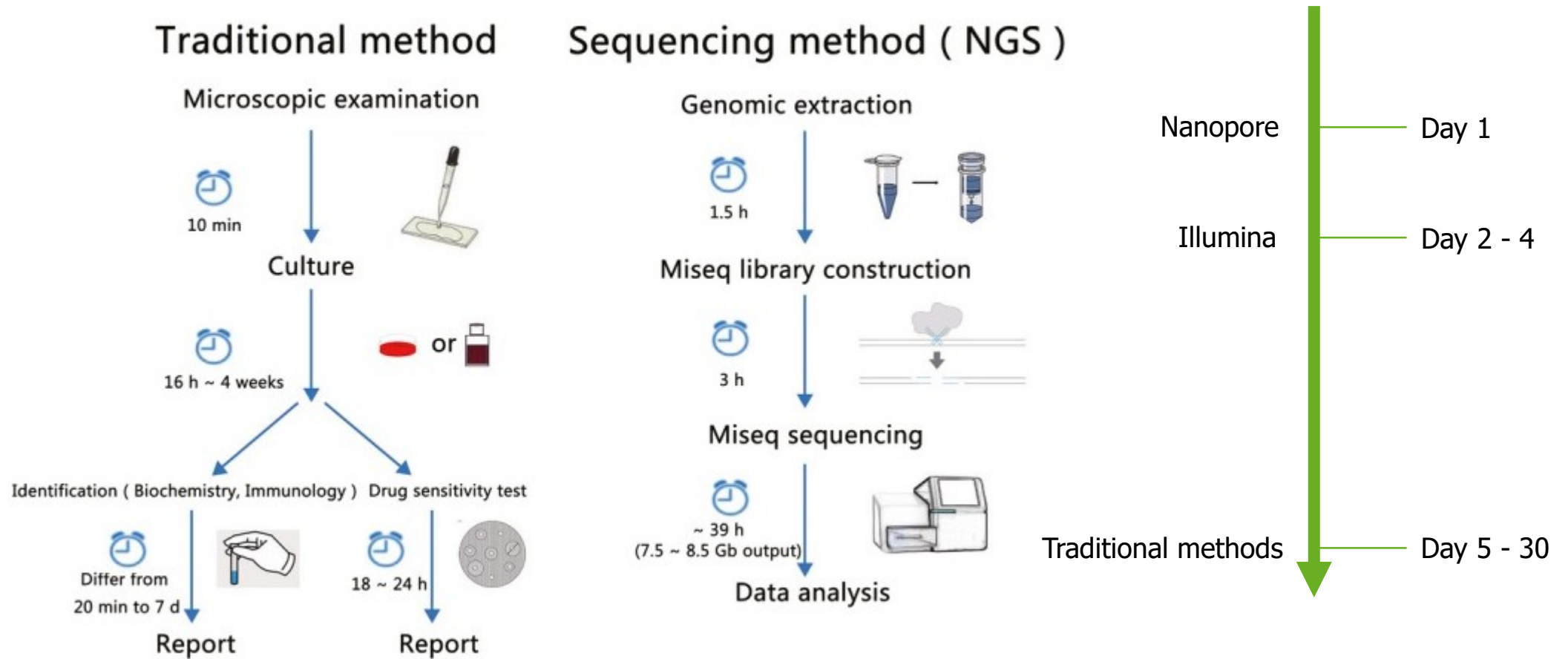
### Experience:

- Employed at SSI since January 2022, specializing in antimicrobial resistance surveillance and Nanopore sequencing

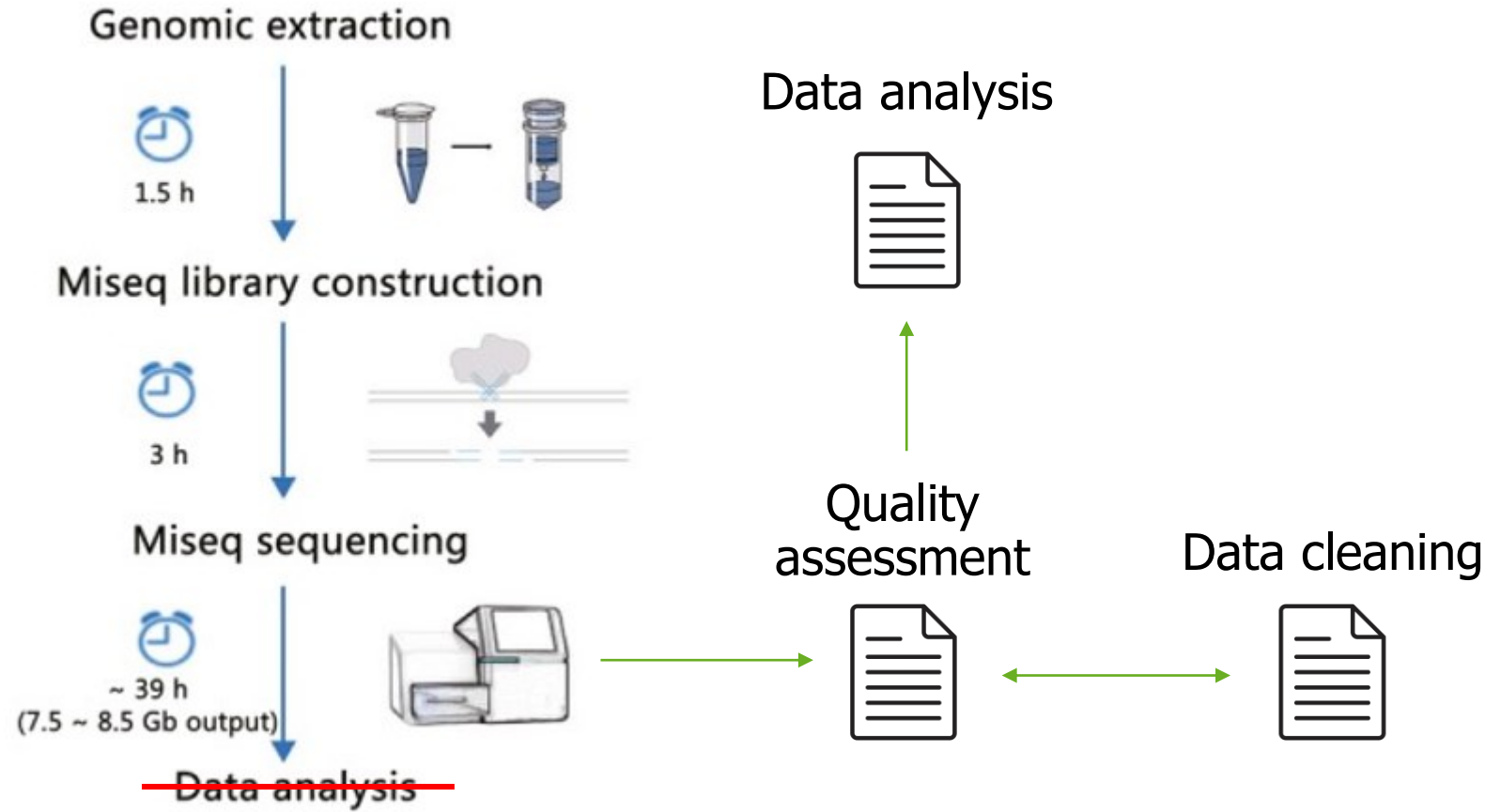
# How about you?

Please have a phone or another browser window at hand...

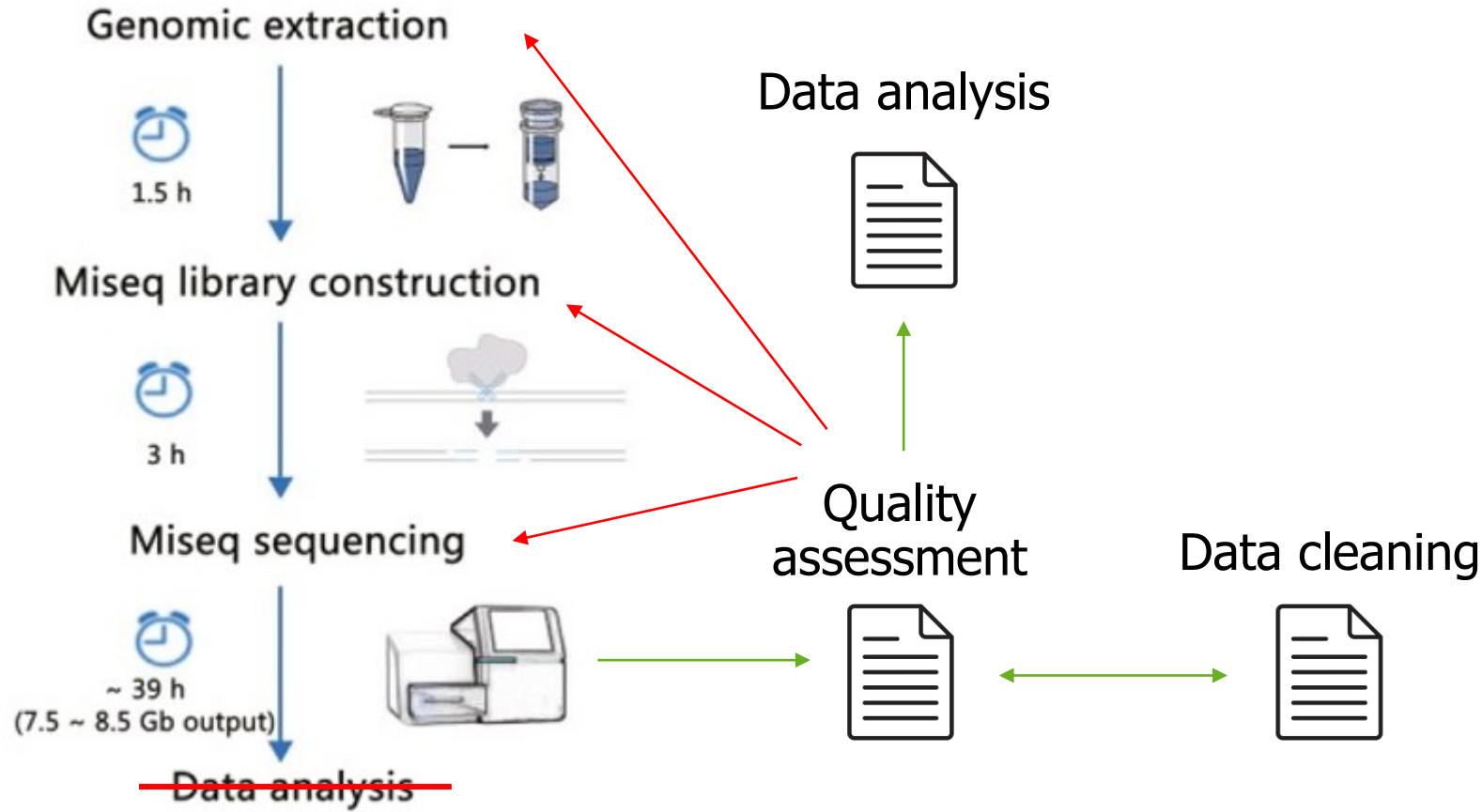
# Sequencing for bacterial typing – A perspective from a clinical setting



# Sequencing method ( NGS )



# Sequencing method ( NGS )



# Course page



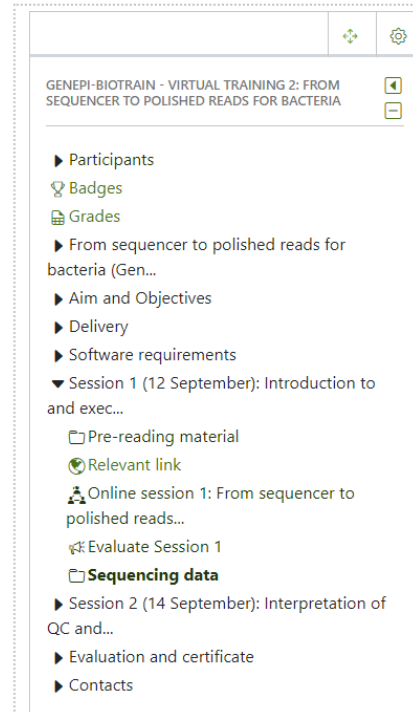
# Life demo

# Dataset

Find datasets on EVA

Session1 > Sequencing data

We assume that you have pre-installed FastQC and NanoPlot



GENEPI-BIOTRAIN - VIRTUAL TRAINING 2: FROM SEQUENCER TO POLISHED READS FOR BACTERIA

- ▶ Participants
- 🏆 Badges
- 📚 Grades
  - ▶ From sequencer to polished reads for bacteria (Gen...
  - ▶ Aim and Objectives
  - ▶ Delivery
  - ▶ Software requirements
  - ▼ Session 1 (12 September): Introduction to and exec...
    - 📁 Pre-reading material
    - 🔗 Relevant link
    - 👤 Online session 1: From sequencer to polished reads...
    - 📄 Evaluate Session 1
    - 📁 **Sequencing data**
  - ▶ Session 2 (14 September): Interpretation of QC and...
  - ▶ Evaluation and certificate
  - ▶ Contacts

## Sequencing data

Manually mark this activity when complete

I have completed this activity

- 📁 illumina
  - 📄 Ec001.illumina\_R1.fastq.gz
  - 📄 Ec001.illumina\_R2.fastq.gz
  - 📄 Ec002.illumina\_R1.fastq.gz
  - 📄 Ec002.illumina\_R2.fastq.gz
  - 📄 Ec003.illumina\_R1.fastq.gz
  - 📄 Ec003.illumina\_R2.fastq.gz
  - 📄 Ec004.illumina\_R1.fastq.gz
  - 📄 Ec004.illumina\_R2.fastq.gz
  - 📄 Ec005.illumina\_R1.fastq.gz
  - 📄 Ec005.illumina\_R2.fastq.gz
- 📁 nanopore
  - 📄 Ec001\_super.fastq.gz
  - 📄 Ec002\_super.fastq.gz

Download folder

Edit

From sequencer to polished reads for bacteria

# ILLUMINA SEQUENCING THEORY

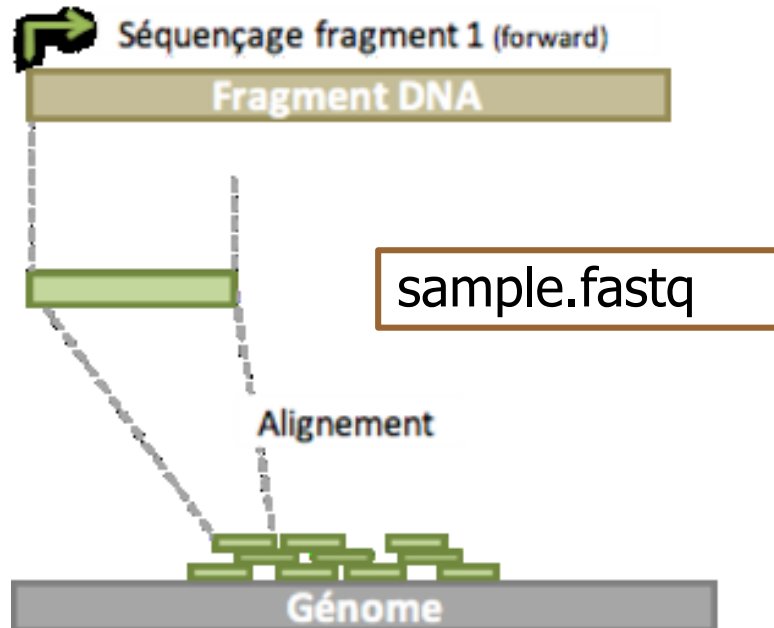
# ILO's

- A fundamental understanding of read type and read fragment sizes
- A crude understanding of tagmentation, size selection of fragments, as well as sequencing by synthesis
- Finally, knowledge of the output from the sequencer.

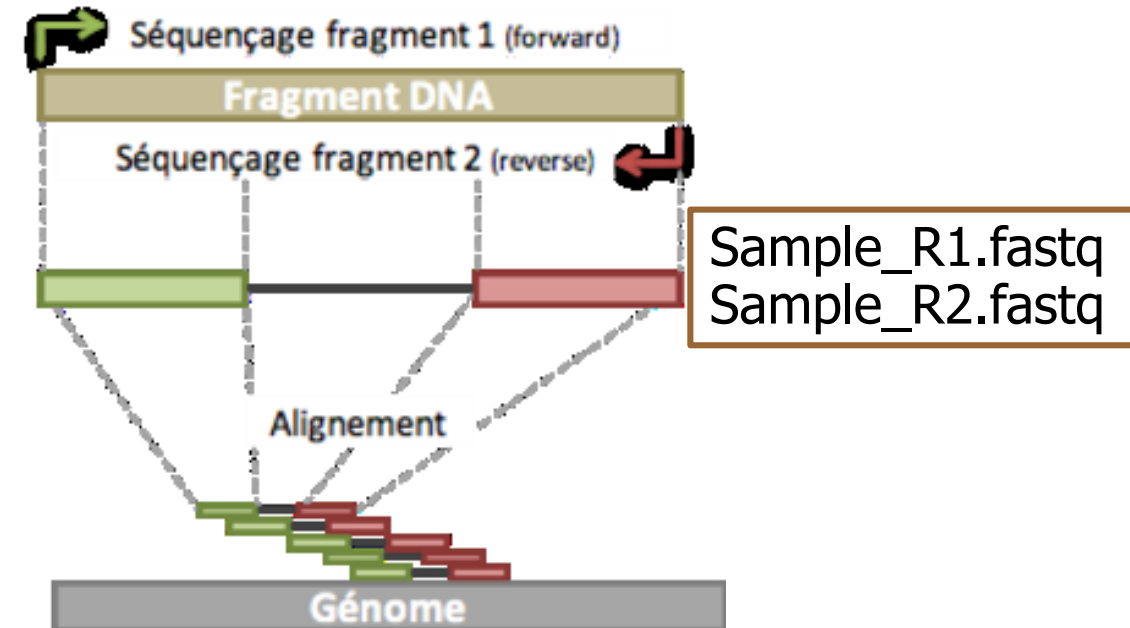
# Illumina read types

Single-end vs Paired-end reads

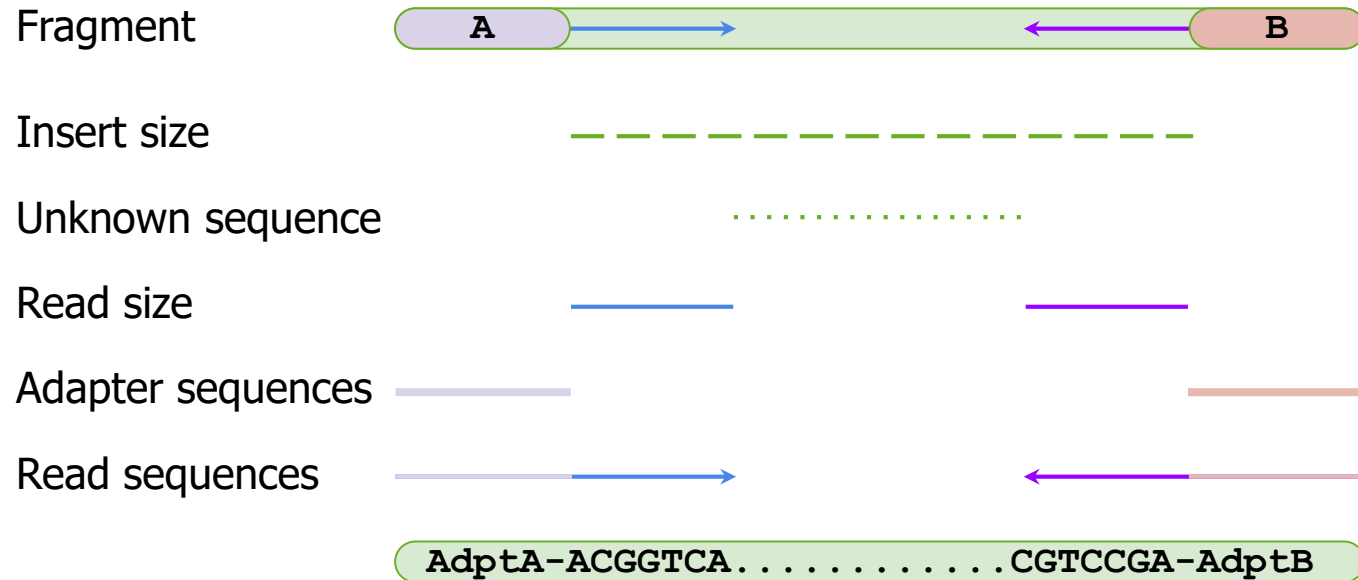
## Single-end



## Paired-end



# Read size



```
@sample1_Mate1_readX
AdptA-ACGGTCA
...
@sample1_Mate2_readX
AdptB-AGCCTGC
...
@sample1_Mate1_readY
```

# Scenarios

## Overlap



Insert size



Read sequences



Overlap



Proof reading



## Read-through



Data loss



Proof reading

## Repeat capture



Repeat capture



Unknown



# Illumina Nextera XT in summary

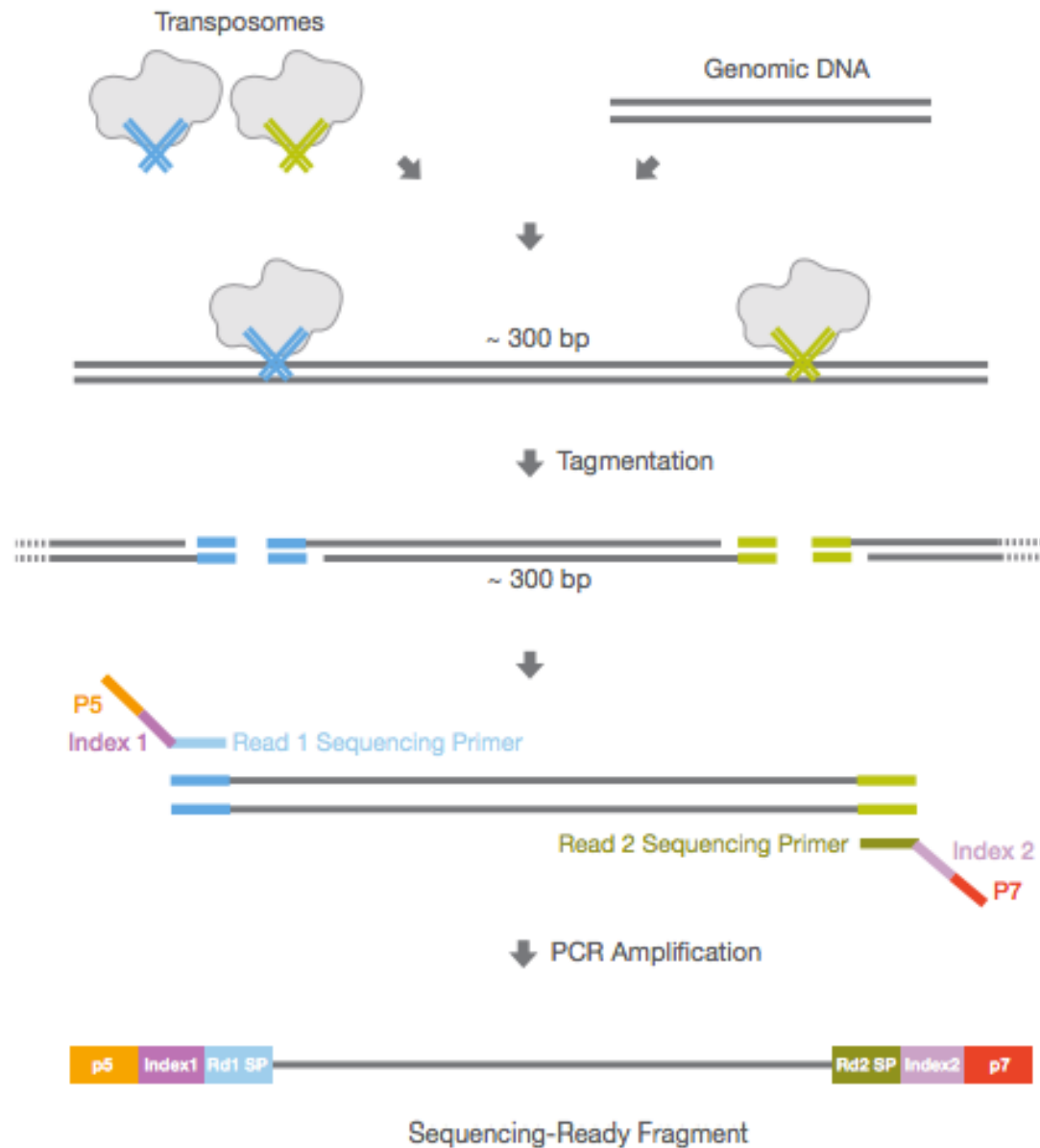
## 1) Sample prep

- gDNA extraction
- Pre-normalization

## 2) Library preparation

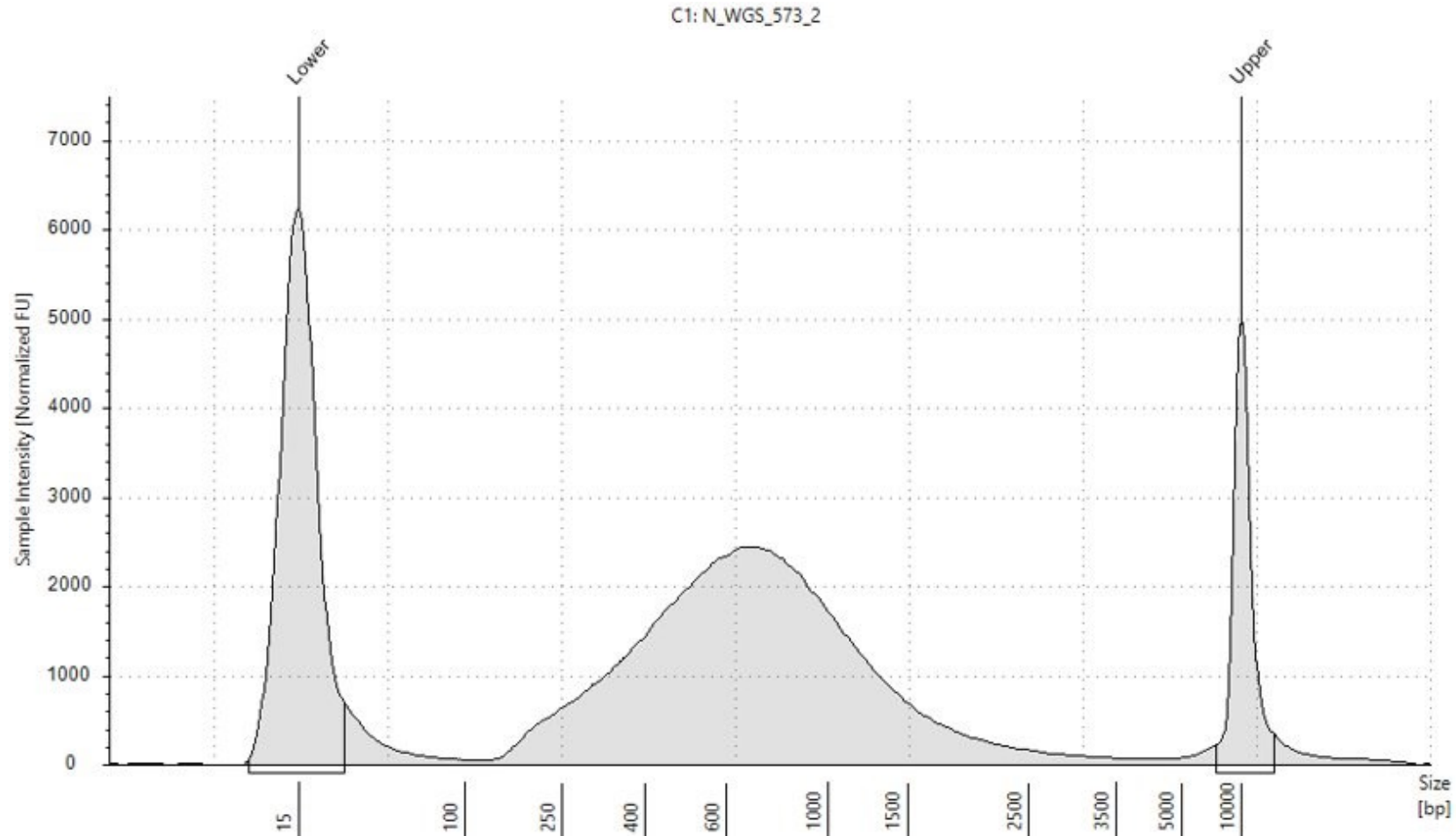
- Tagmentation
- Index PCR
- Normalization and pool

## 3) Sequencing



# Library size selection

## Bead-based size selection



# Library size selection

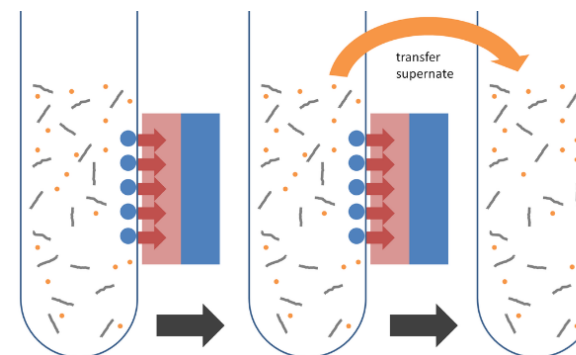
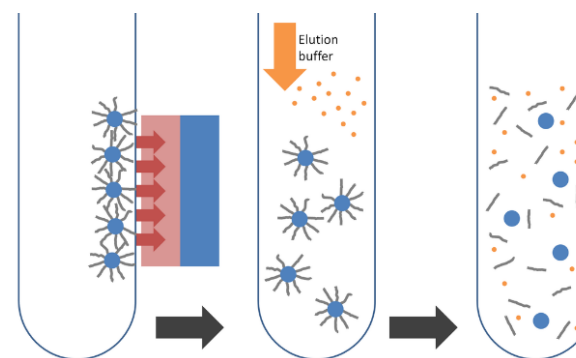
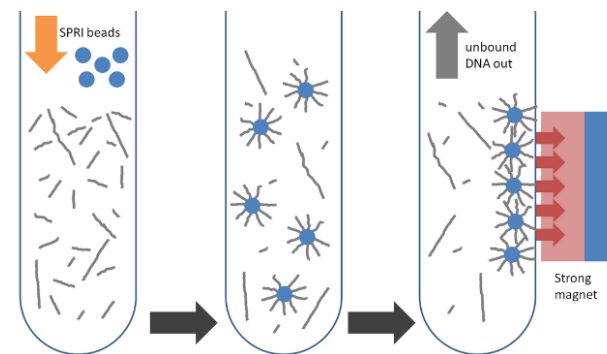
## Bead-based size selection

### INDEX PCR AND CLEANUP

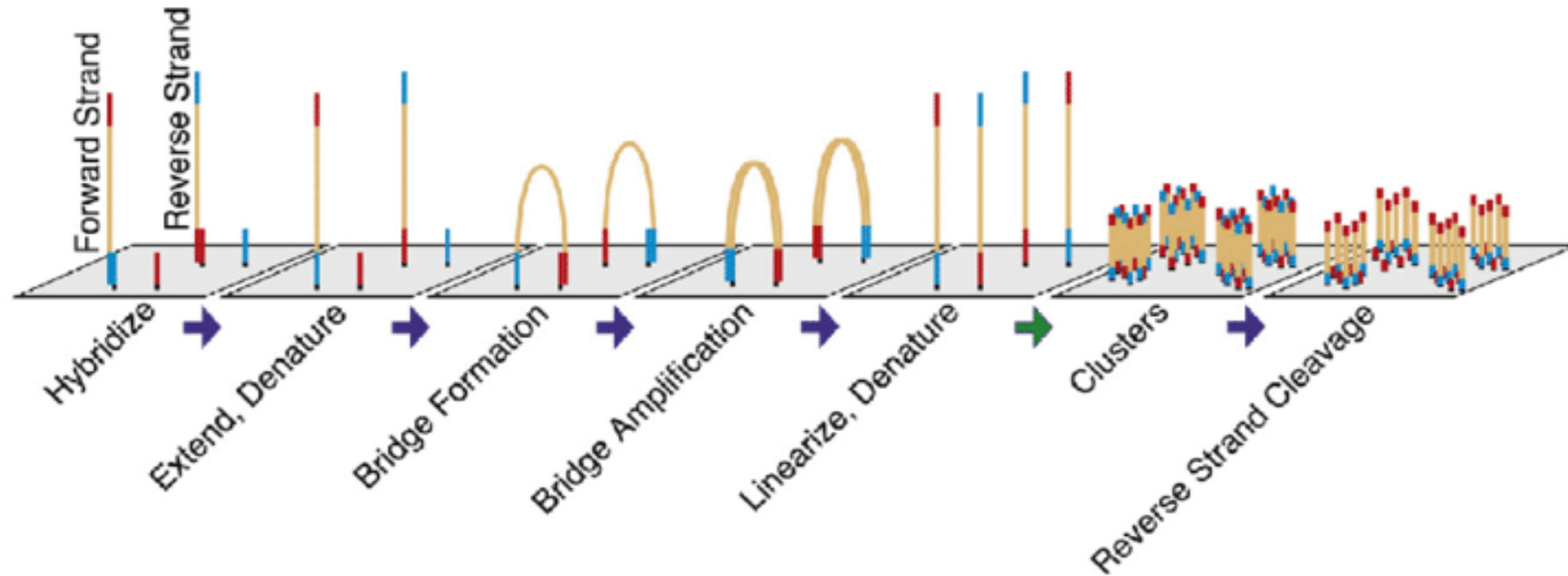
**BEFORE**  
TAGGED AND FRAGMENTED  
DNA (DIFFERENT SIZES)



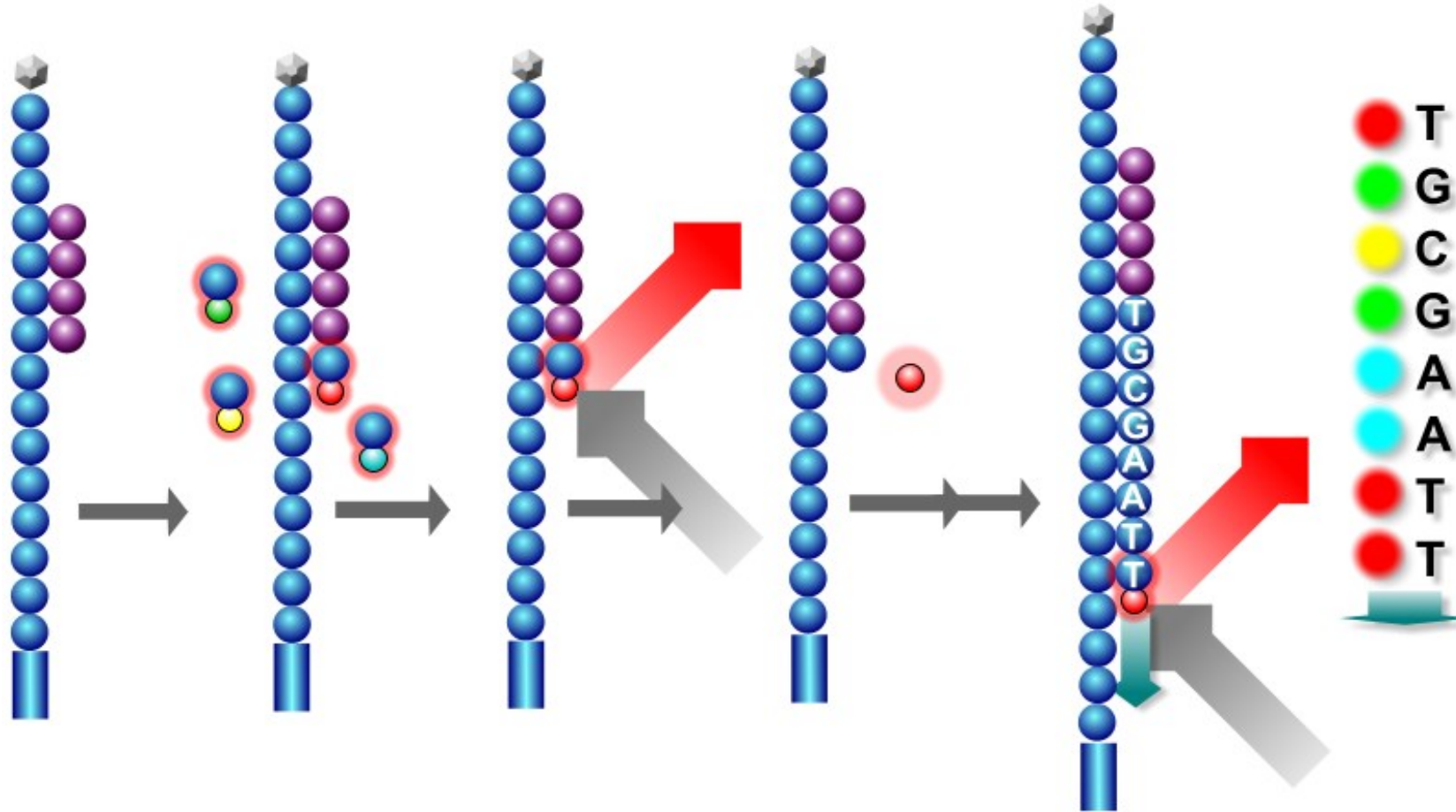
**AFTER**  
UNIFORMLY-SIZED LIBRARIES



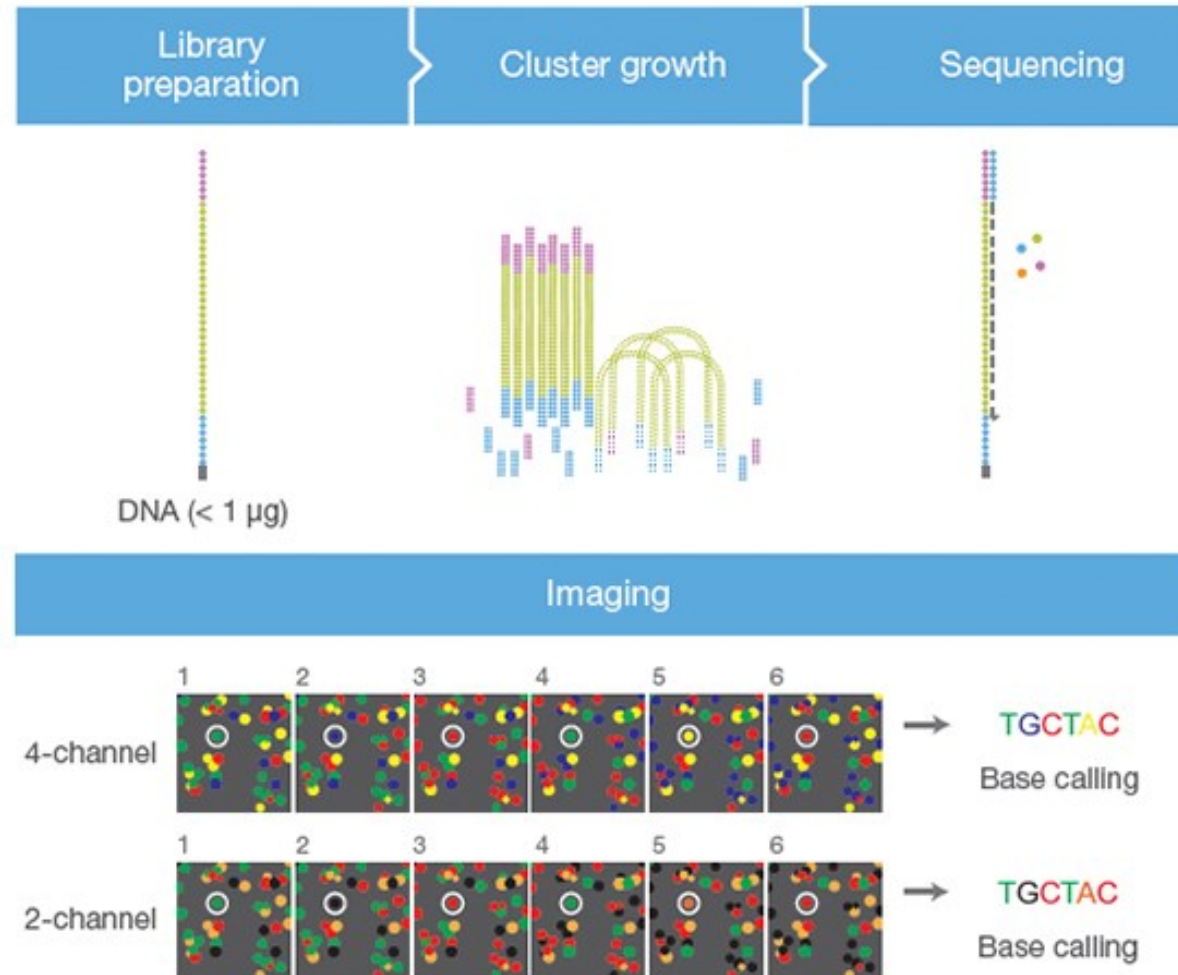
# Illumina cluster generation



# Sequencing by Synthesis



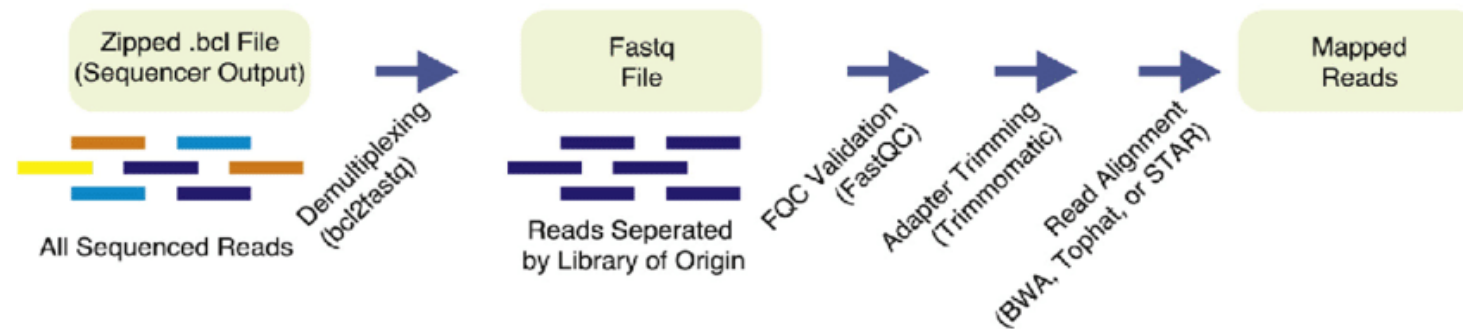
# Illumina cluster generation



# Illumina sequence data files

Illumina sequencer generates .bcl

Translated to fastq file format on the machine using bcl2fastq



# QUESTIONS

# Break!



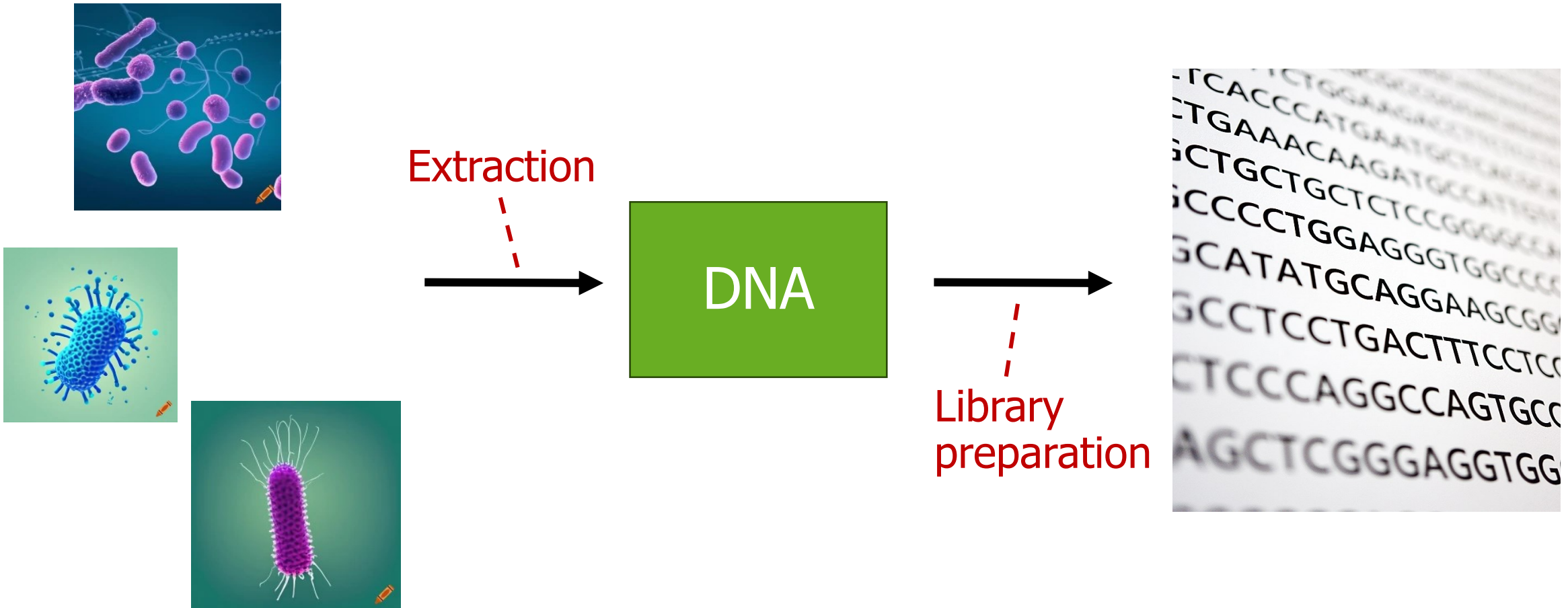
From sequencer to polished reads for bacteria

# Nanopore sequencing theory

# ILOs

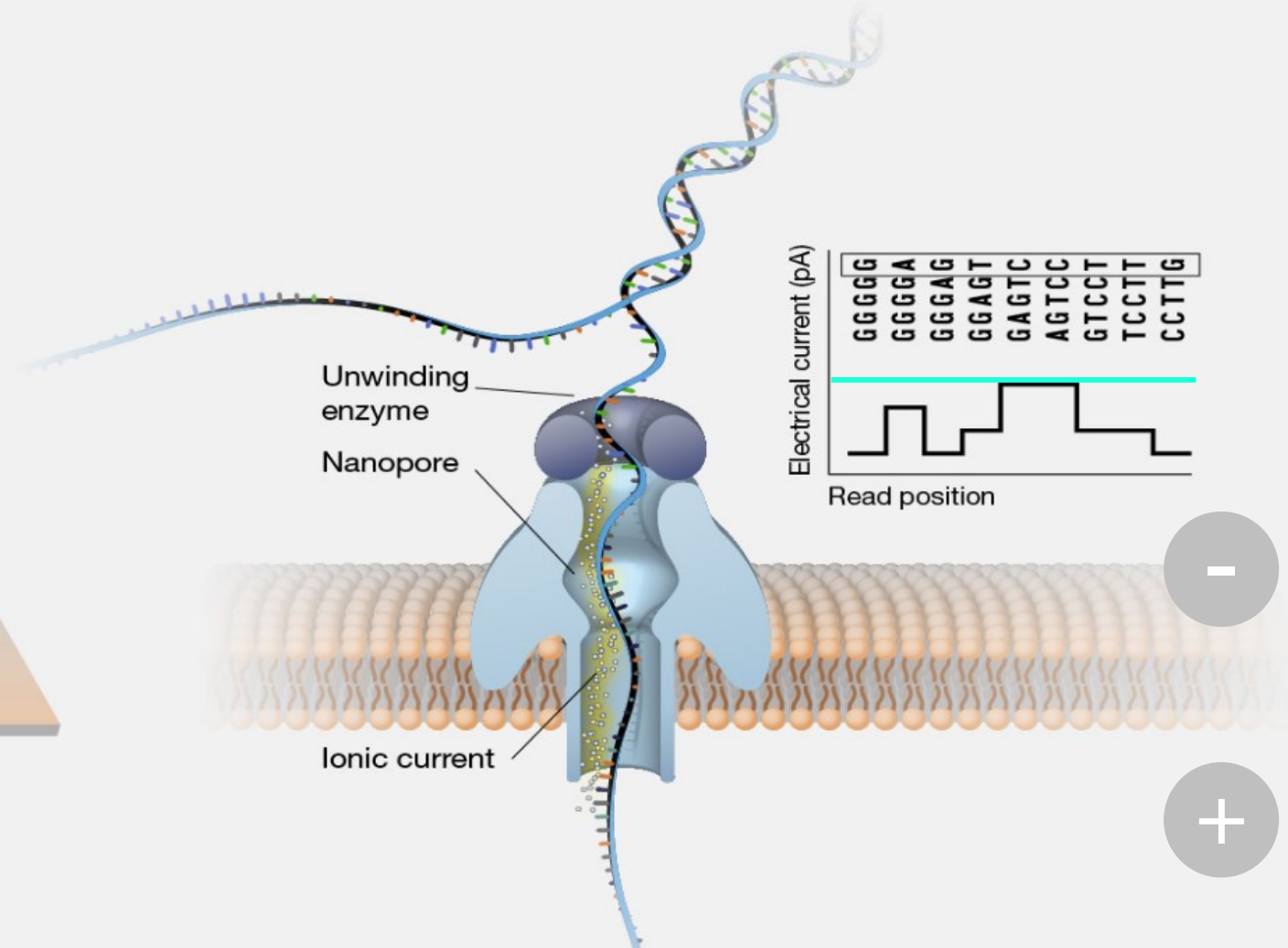
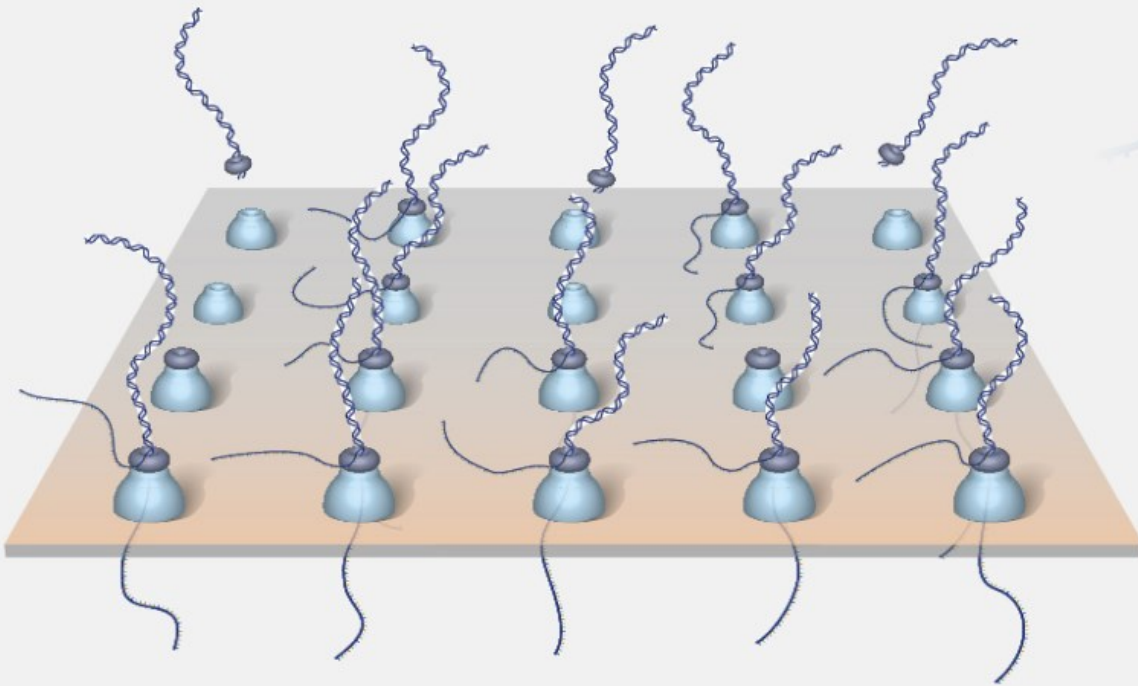
- Understand the fundamentals of Nanopore sequencing.
- Explain how library preparation affects Nanopore sequencing data.
- Distinguish between Nice-to-Know and Need-to-Know information.

# Sequencing fundamentals



# Nanopore sequencing

## Nanopore DNA sequencing



# Nanopore sequencing platform overview



Configuration	Platform				Techniques		Tech specifications	
Number of flow cells per device	1	1	1	5	2	2	24	48
Maximum number of channels per flow cell	512	512	512	512	2,675	2,675	2,675	2,675
Run time	72 Hours	72 Hours	72 Hours	72 Hours	72 Hours	72 Hours	72 Hours	72 Hours
Device TMO <sup>†</sup>	50 Gb	50 Gb	50 Gb	250 Gb	580 Gb	580 Gb	~7 Tb	~14 Tb
Maximum number of flow cells per year*	104	104	104	520	208	208	2,596	4,992
Offer sequencing as a service	No	No	No	Yes	Yes	Yes	Yes	Yes

# Nanopore sequencing platform overview



MinION and Flongle Flow Cell compatible

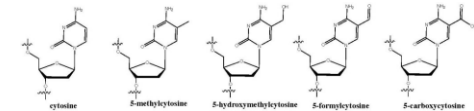
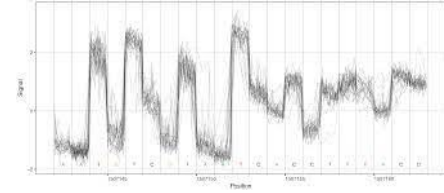
PromethION Flow Cell compatible

Configuration	MinION and Flongle Flow Cell compatible				PromethION Flow Cell compatible			
	Platform		Platform		Platform		Tech specifications	
Number of flow cells per device	1	1	1	1	2	2	24	48
Maximum number of channels per flow cell	512	512	512	512	375	2,675	2,675	2,675
Run time	72 Hours	72 Hours	72 Hours	72 Hours	72 Hours	72 Hours	72 Hours	72 Hours
Device TMO <sup>†</sup>	50 Gb	50 Gb	50 Gb	250 Gb	580 Gb	580 Gb	~7 Tb	~14 Tb
Maximum number of flow cells per year*	104	104	104	520	208	208	2,596	4,992
Offer sequencing as a service	No	No	No	Yes	Yes	Yes	Yes	Yes



# Nanopore sequencing data

- The raw electrical signal is saved in **pod5 files** and is translated, or base called, into nucleotide bases which are stored in fastq files.
- Each DNA strand will be represented as a read and the **length is determined by input DNA length**.
- With enough computer power, squiggles are base called in **real time** and you can start to analyse while still sequencing.
- **Base modifications** (methylation, glycosylation etc.) will appear as abnormal squiggles. Very interesting but difficult to basecall correctly.
- Relatively **high error-rate** specifically in homopolymeric regions and because of base modifications.

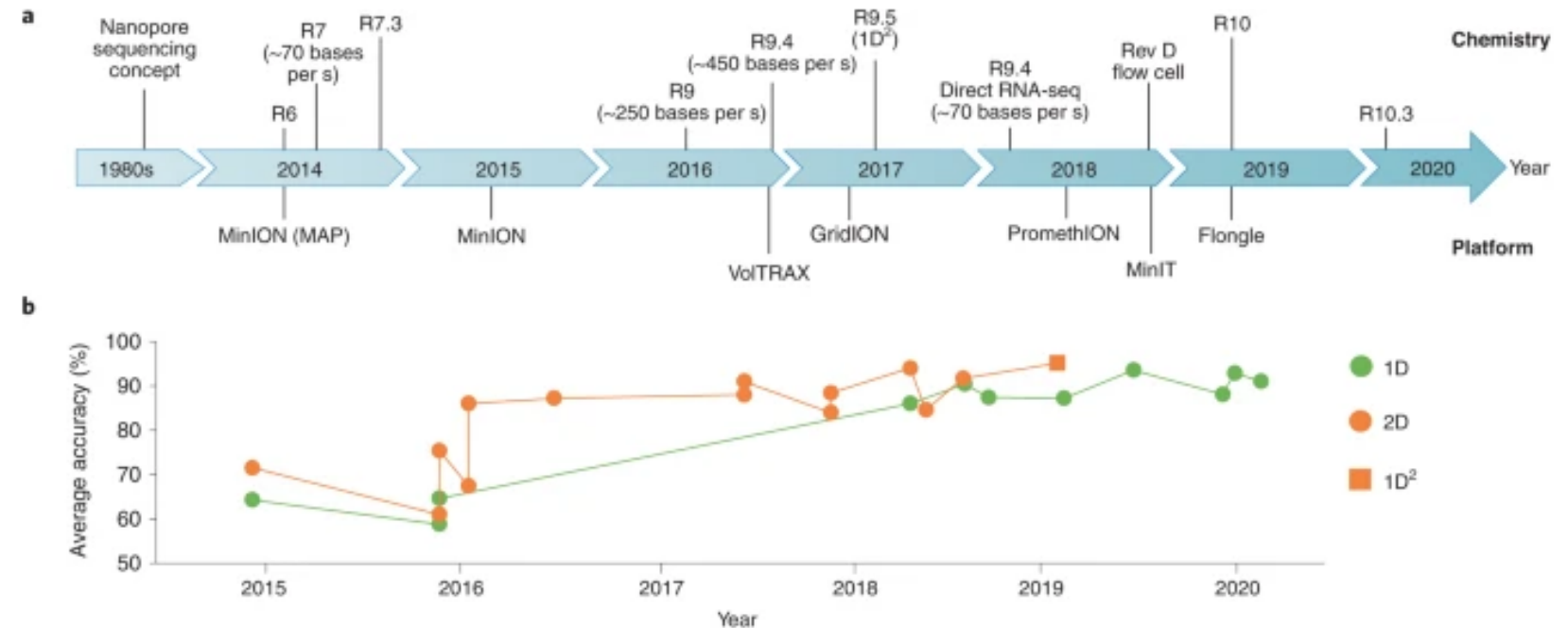


# A bit of history...

Nanopore sequencing was made available to the public in 2014.

- Since then the chemistry has changed many times..

**Fig. 2: ONT sequencing data improvement over time.**



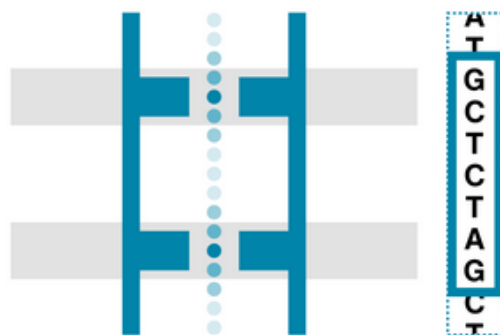
# A bit of history...

Nanopore sequencing was made available to the public in 2014.

- Since then the chemistry has changed many times..

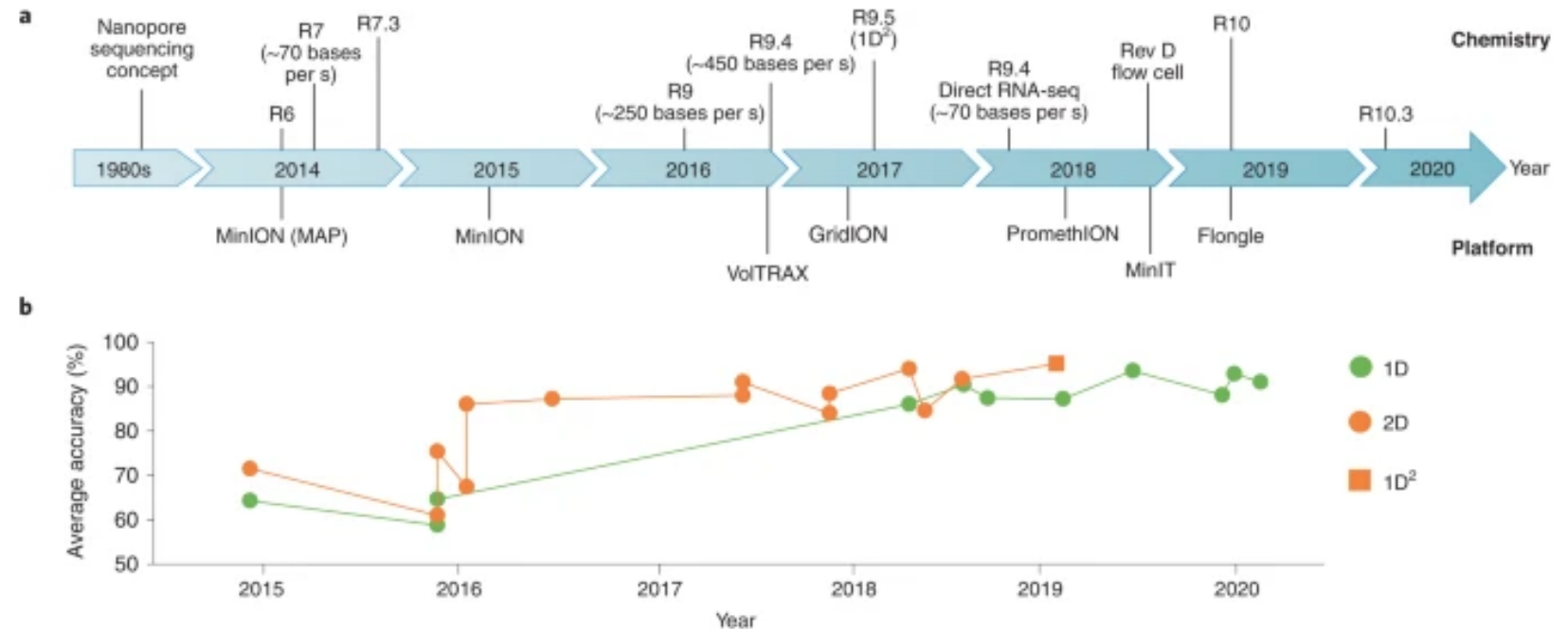


R9



R10

**Fig. 2: ONT sequencing data improvement over time.**

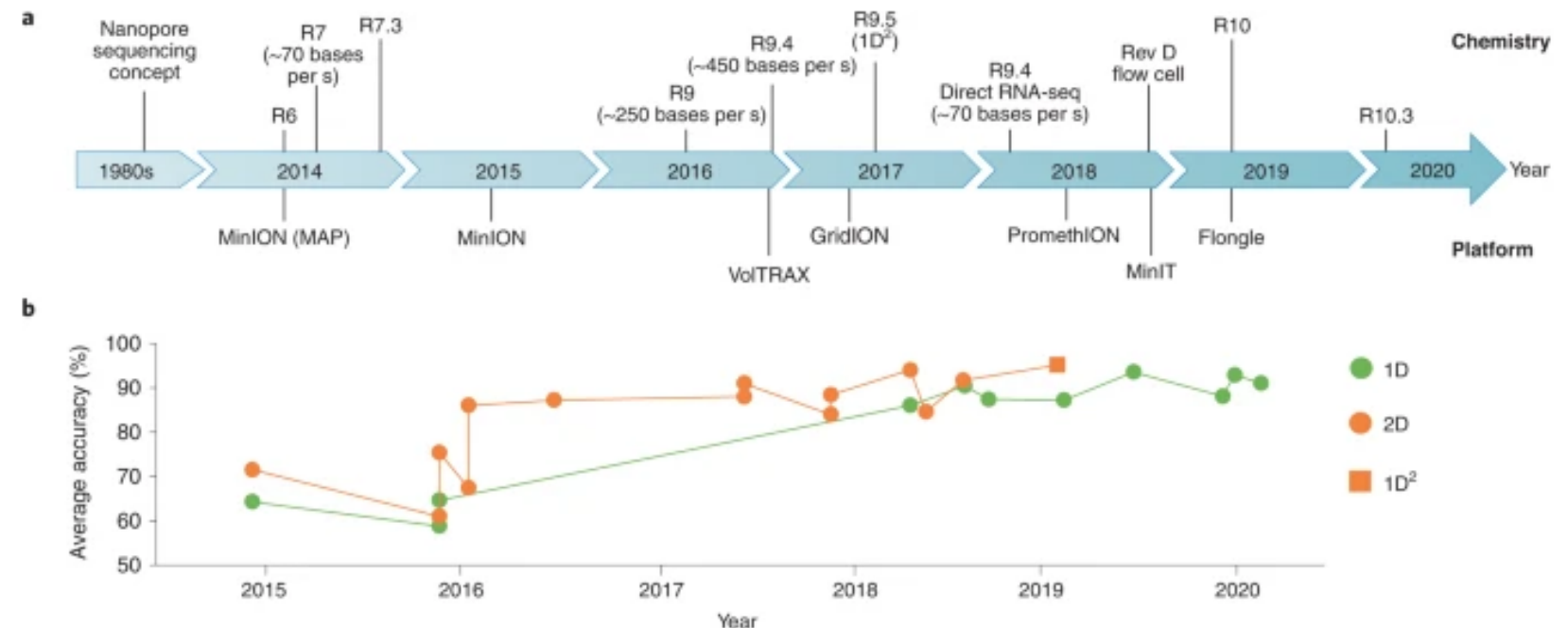


# A bit of history...

Nanopore sequencing was made available to the public in 2014.

- Since then the chemistry has changed many times..

**Fig. 2: ONT sequencing data improvement over time.**







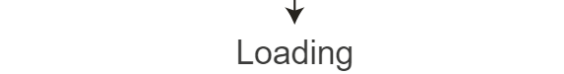
..and so has the software

- 2017 == Albacore
- 2019 == Guppy
- 2023 == Dorado

# Nanopore library prep

	Output optimised	Speed optimised	Ultra-long reads optimised	Low input optimised	Targeted sequencing	
	<a href="#">Ligation Sequencing Kit</a>	<a href="#">Rapid Sequencing Kit</a>	<a href="#">Ultra-Long DNA Sequencing Kit</a>	<a href="#">Rapid PCR Barcoding Kit</a>	<a href="#">16S Barcoding Kit</a>	<a href="#">Cas9 Sequencing Kit</a>
Preparation time	60 minutes	10 minutes	200* mins +1xO/N incubation	15 mins + PCR	10 mins + PCR	110 minutes
Input recommendation	1000 ng gDNA 100-200 fmol for amplicons	50 - 100 ng	6M Cells	1 - 5 ng	10 ng	1 - 10 µg
Fragmentation	Optional	Transposase-based	Transposase-based	Transposase-based	-	Cas9-dependent cleavage
Amplification	No	No	No	Yes	Yes	No
Barcode options	<a href="#">Native Barcoding Kit 24</a> <a href="#">Native Barcoding Kit 96</a>	<a href="#">Rapid Barcoding Kit 24</a> <a href="#">Rapid Barcoding Kit 96</a>	-	12 plex	24 plex	<i>In development</i>
Typical output	●●●	●●○	●●○	●●○	●●○	●○○
Adaptive sampling	✓	✓	✓	✓	-	<i>In development</i>
Methylation included	✓	✓	✓	-	-	✓

# Nanopore library prep

	Output optimised	Speed optimised	Ultra-long reads optimised	Low input optimised	Targeted sequencing		
	<b>Ligation Sequencing Kit</b>	<b>Rapid Sequencing Kit</b>	<b>Ultra-Long DNA Sequencing Kit</b>	<b>Rapid PCR Barcoding Kit</b>	<b>16S Barcoding Kit</b>	<b>Cas9 Sequencing Kit</b>	
Preparation time	60 minutes	10 minutes	200* mins +1xO/N	15 mins + PCR	10 mins + PCR	110 minutes	
Input recom	Double stranded DNA e.g. gDNA, locus-specific amplicons, cDNA amplicons 			1 - 5 ng	10 ng	1 - 10 µg	
Fragmentati	Optional fragmentation or size selection 			ased	Transposase-based	-	Cas9-dependent cleavage
Amplificati	End-prep and nick repair 			Yes	Yes	No	
Barcode opt	Ligation of sequencing adapters 			12 plex	24 plex	<i>In development</i>	
Typical outp	Loading 			●●○	●●○	●○○	
Adaptive sar				✓	-	<i>In development</i>	
Methylation included	✓	✓	✓	-	-	✓	

60 min

# Nanopore library prep

	Output optimised	Speed optimised	Ultra-long reads optimised	Low input optimised	Targeted sequencing	
	<a href="#">Ligation Sequencing Kit</a>	<b><a href="#">Rapid Sequencing Kit</a></b>	<a href="#">Ultra-Long DNA Sequencing Kit</a>	<a href="#">Rapid PCR Barcoding Kit</a>	<a href="#">16S Barcoding Kit</a>	<a href="#">Cas9 Sequencing Kit</a>
Preparation time	60 minutes	10 minutes	200* mins + 1xO/N	15 mins + PCR	10 mins + PCR	110 minutes
Input recom	Double stranded DNA e.g. gDNA, locus-specific amplicons, cDNA amplicons			Barcoded transposome complex		
Fragmentati	Optional fragmentation or size selection		ased	gDNA		
Amplificati	End-prep and nick repair			Cleavage of gDNA by transposase and addition of barcodes		
Barcode opt	Ligation of sequencing adapters			Pooling of barcoded libraries, SPRI clean-up		
Typical outp				Attachment of sequencing adapters		
Adaptive sar				Loading		
Methylation included	✓	✓	✓	-	-	✓

60 min

< 60 min

# Nanopore library prep

## Choice of library kit and sequencing platform will affect:

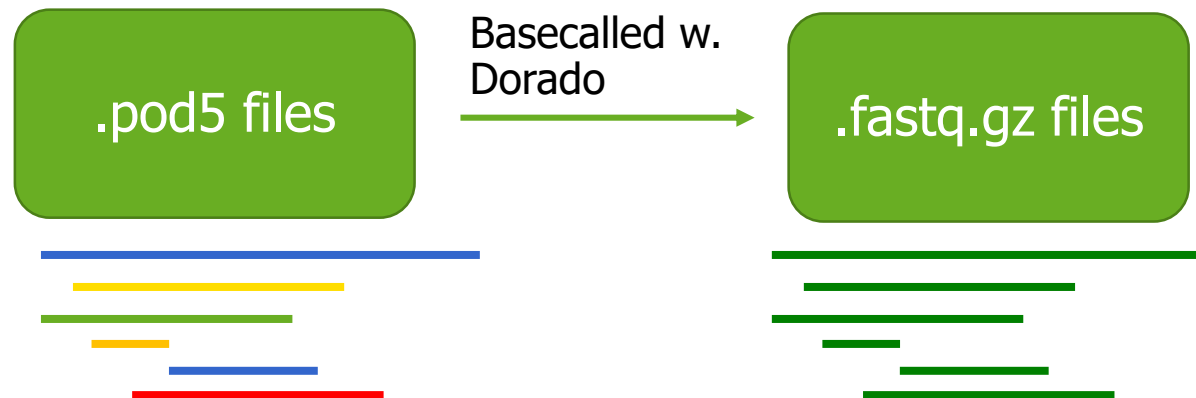
- Amount of data produced
- Length of reads
- Coverage
- Overall quality

# Nanopore sequence data files

Nanopore sequencer generates demultiplexed .pod5 files

Translated to fastq file format with dorado

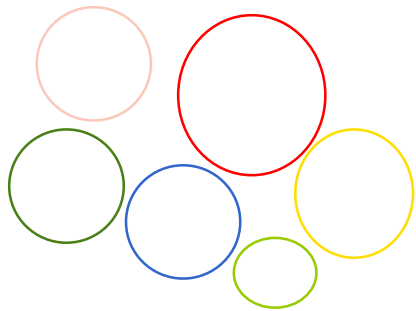
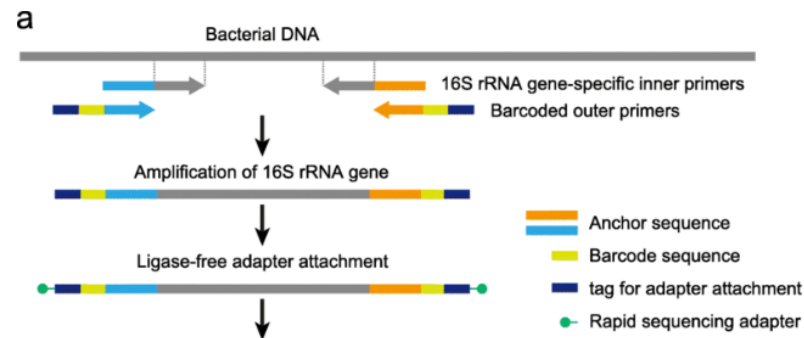
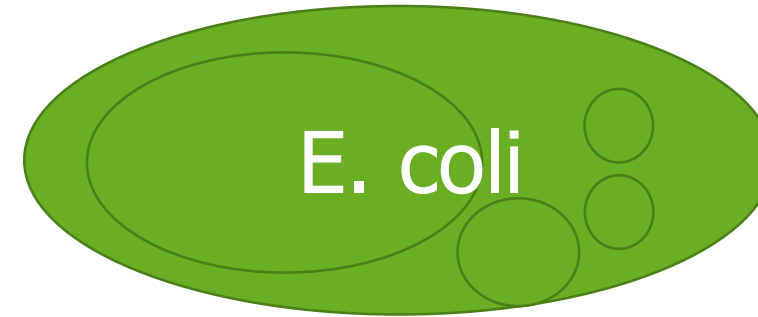
- Enable barcode trimming in MinKNOW



# Need-to-Know

## What was sequenced?

- Single bacterial isolate
- 16S
- Metagenomic sample



# Nice-to-Know

- Which flow cell version was used?
- Which library was used?
- Which base caller was used?
- Has it been quality filtered, if so on which q-score?
- Have adapters/barcodes been trimmed?

# QUESTIONS



# Summary

# Summary of the two methods

## Illumina

- Library
  - + 1 ng/μl input required
  - Amplification bias introduced
- Reads
  - + High accuracy > 99%
  - Very short sequences
- Sequencing
  - + High yield: Gigabases
  - 4h to 2.5 days

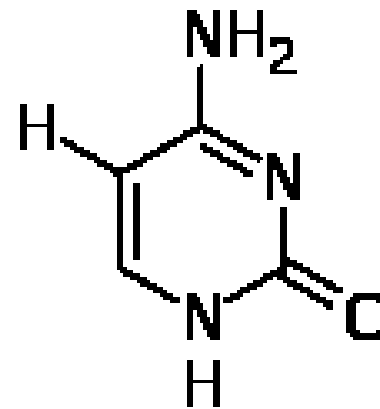
## Nanopore

- Library
  - + Bias-free
  - High DNA purity required
- Reads
  - + Very long reads (> 2Mb)
  - Lower accuracy 92 – 99%
- Sequencing
  - + Real-time
  - Fixed max. Yield / flowcell

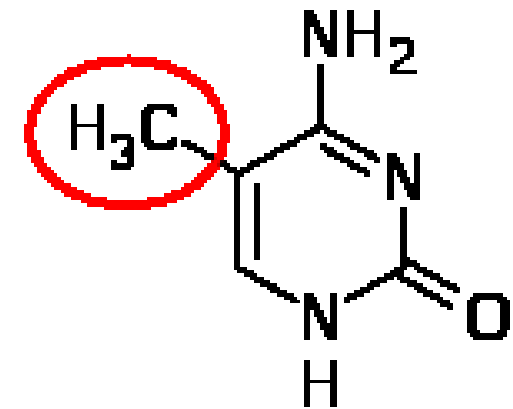
# Nanopore sequencing

## Important differences from NGS

- Read lengths are determined by input DNA length
- Relatively high error-rate
- Base modifications



cytosine



methylated  
cytosine

# Which one to chose

It depends...

We use both Illumina and Nanopore data..

We even combine their data.

# Break!

**And now to something completely different...**

# Intro to Session 1 Practical

- Get to know your dataset
- Intro to contamination control
- Perform quality assessment on Illumina data
- Perform quality assessment on Nanopore data

From sequencer to polished reads for bacteria

# Get to know your dataset

September 2023, Statens Serum Institut

# Whole genome sequencing data

GENEPI-BIOTRAIN - VIRTUAL TRAINING 2: FROM SEQUENCER TO POLISHED READS FOR BACTERIA

- ▶ Participants
- 🔗 Badges
- 📅 Grades
  - ▶ From sequencer to polished reads for bacteria (Gen...
  - ▶ Aim and Objectives
  - ▶ Delivery
  - ▶ Software requirements
  - ▼ Session 1 (12 September): Introduction to and exec...
    - 📁 Pre-reading material
    - 🌐 Relevant link
    - 👤 Online session 1: From sequencer to polished reads...
    - 🗣️ Evaluate Session 1
    - 📁 **Sequencing data**
  - ▶ Session 2 (14 September): Interpretation of QC and...
    - ▶ Evaluation and certificate
    - ▶ Contacts

## Sequencing data

Manually mark this activity when complete  I have completed this activity

- illumina
  - Ec001.illumina\_R1.fastq.gz
  - Ec001.illumina\_R2.fastq.gz
  - Ec002.illumina\_R1.fastq.gz
  - Ec002.illumina\_R2.fastq.gz
  - Ec003.illumina\_R1.fastq.gz
  - Ec003.illumina\_R2.fastq.gz
  - Ec004.illumina\_R1.fastq.gz
  - Ec004.illumina\_R2.fastq.gz
  - Ec005.illumina\_R1.fastq.gz
  - Ec005.illumina\_R2.fastq.gz
- nanopore
  - Ec001\_super.fastq.gz
  - Ec002\_super.fastq.gz

Download folder

Edit

5 Illumina sequenced isolates

2 Nanopore sequenced isolates

# Whole genome sequencing data

GENEPI-BIOTRAIN - VIRTUAL TRAINING 2: FROM SEQUENCER TO POLISHED READS FOR BACTERIA

- ▶ Participants
- 🔗 Badges
- 📄 Grades
  - ▶ From sequencer to polished reads for bacteria (Gen...
  - ▶ Aim and Objectives
  - ▶ Delivery
  - ▶ Software requirements
  - ▼ Session 1 (12 September): Introduction to and exec...
    - 📁 Pre-reading material
    - 🌐 Relevant link
    - 👤 Online session 1: From sequencer to polished reads...
    - 📊 Evaluate Session 1
    - 📁 **Sequencing data**
  - ▶ Session 2 (14 September): Interpretation of QC and...
  - ▶ Evaluation and certificate
  - ▶ Contacts

## Sequencing data

Manually mark this activity when complete  I have completed this activity

- illumina
  - Ec001.illumina\_R1.fastq.gz
  - Ec001.illumina\_R2.fastq.gz
  - Ec002.illumina\_R1.fastq.gz
  - Ec002.illumina\_R2.fastq.gz
  - Ec003.illumina\_R1.fastq.gz
  - Ec003.illumina\_R2.fastq.gz
  - Ec004.illumina\_R1.fastq.gz
  - Ec004.illumina\_R2.fastq.gz
  - Ec005.illumina\_R1.fastq.gz
  - Ec005.illumina\_R2.fastq.gz
- nanopore
  - Ec001\_super.fastq.gz
  - Ec002\_super.fastq.gz

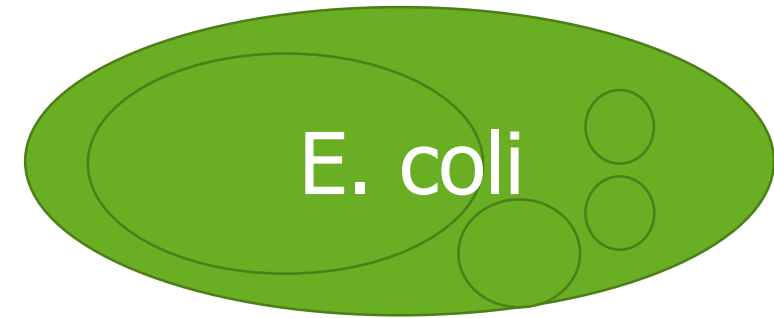
Download folder

Edit

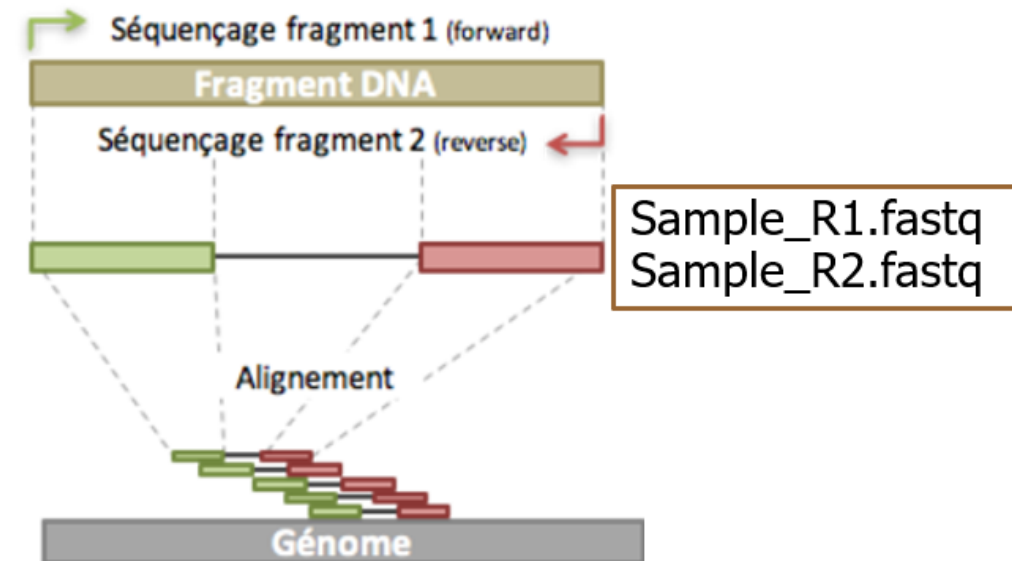
# What to expect

Whole genome sequencing of E. coli

Illumina 2 x 250 bp paired end reads



## Paired-end





# The FASTQ format

```
Read name: @Sequence_ID
Read Seq:  A T C G A T C G
Separator: +
Phred-ASCII:F D G F J J S K
Phred:      36 34 38 36 40 40 41 39
```

} fastq

# Phred quality scores

Phred quality scores are logarithmically linked to error probabilities

Phred Quality Score	Probability of incorrect base call	Base call accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1000	99.9%
40	1 in 10,000	99.99%
50	1 in 100,000	99.999%
60	1 in 1,000,000	99.9999%

$$P = 10^{-Q/10}$$

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

P = the probability that a base is incorrect

Q = Phred quality score

# Number of reads and bases in fastq files

*Read count*

```
echo $(zcat sample.fastq.gz|wc -l)/4|bc
```

A pipe



*Character count*

```
zcat sample.fastq.gz | paste - - - - | cut -f 2 | tr -d '\n' | wc -c
```

# Number of reads and bases in fastq files

*Read count*

```
echo $(zcat sample.fastq.gz|wc -l)/4|bc
```

A pipe



*Character count*

```
zcat sample.fastq.gz | paste - - - - | cut -f 2 | tr -d '\n' | wc -c
```



@seqID1	ATGCTGAAATCA	+	HVIKILIFDPFBI
@seqID2	GCTAGGTACCAT	+	2222;;;<=?CCH
@seqID3	ATGCTGAAATCA	+	(+.9>;;76667<
@seqID4	GCTAGGTACCAT	+	@(&&&&&&%'
@seqID5	ATGCTGAAATCA	+	DEGHKHJJKKEC
@seqID6	GCTAGGTACCAT	+	LEDCCEIB@@@



From sequencer to polished reads for bacteria

# Contamination

September 2023, Statens Serum Institut

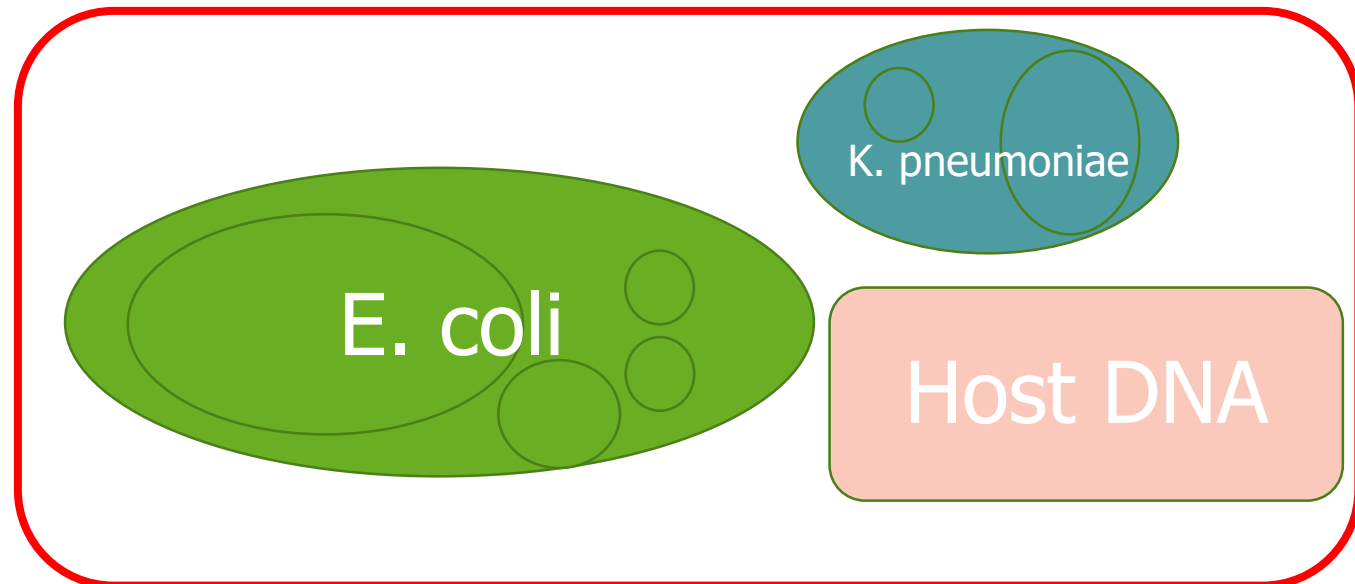
# Contamination in the context of sequencing data



- Single isolate sequencing: any sequence that does not belong to the organism sequenced.
- Can be interspecies contamination: sequences from another species or intraspecies contamination: sequences from the same species.

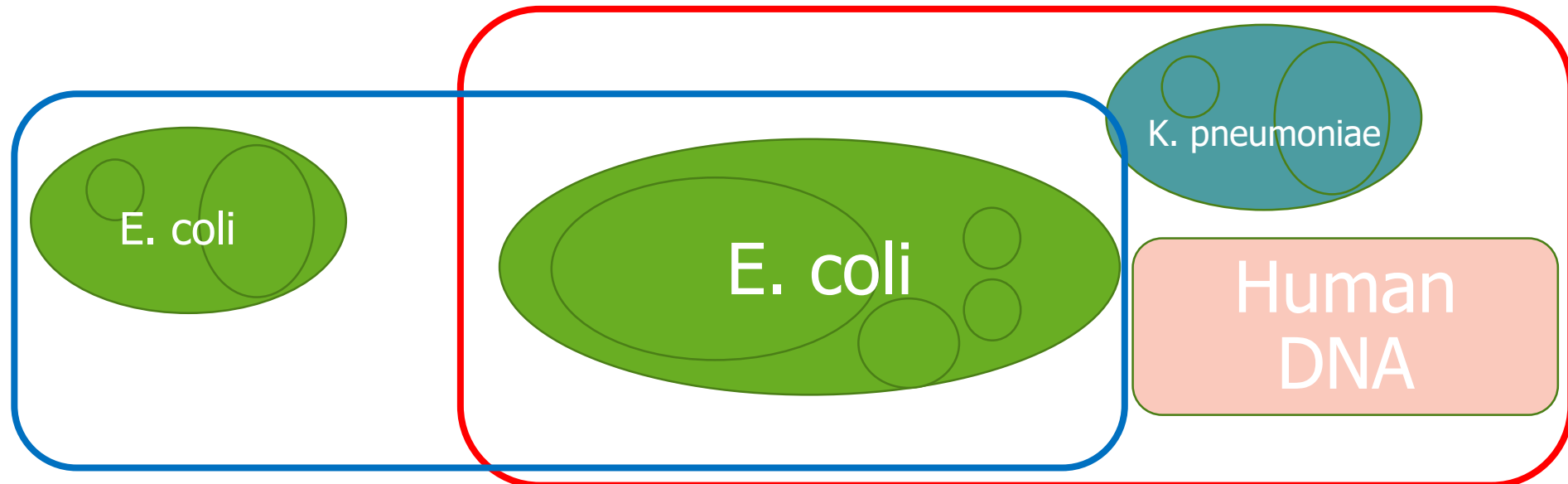
# Contamination in the context of sequencing data

- Single isolate sequencing: any sequence that does not belong to the organism sequenced.
- Can be **interspecies** contamination: sequences from another species or intraspecies contamination: sequences from the same species.



# Contamination in the context of sequencing data

- Single isolate sequencing: any sequence that does not belong to the organism sequenced.
- Can be **interspecies** contamination: sequences from another species or **intraspecies** contamination: sequences from the same species.



# Possible sources of contamination in sequencing data

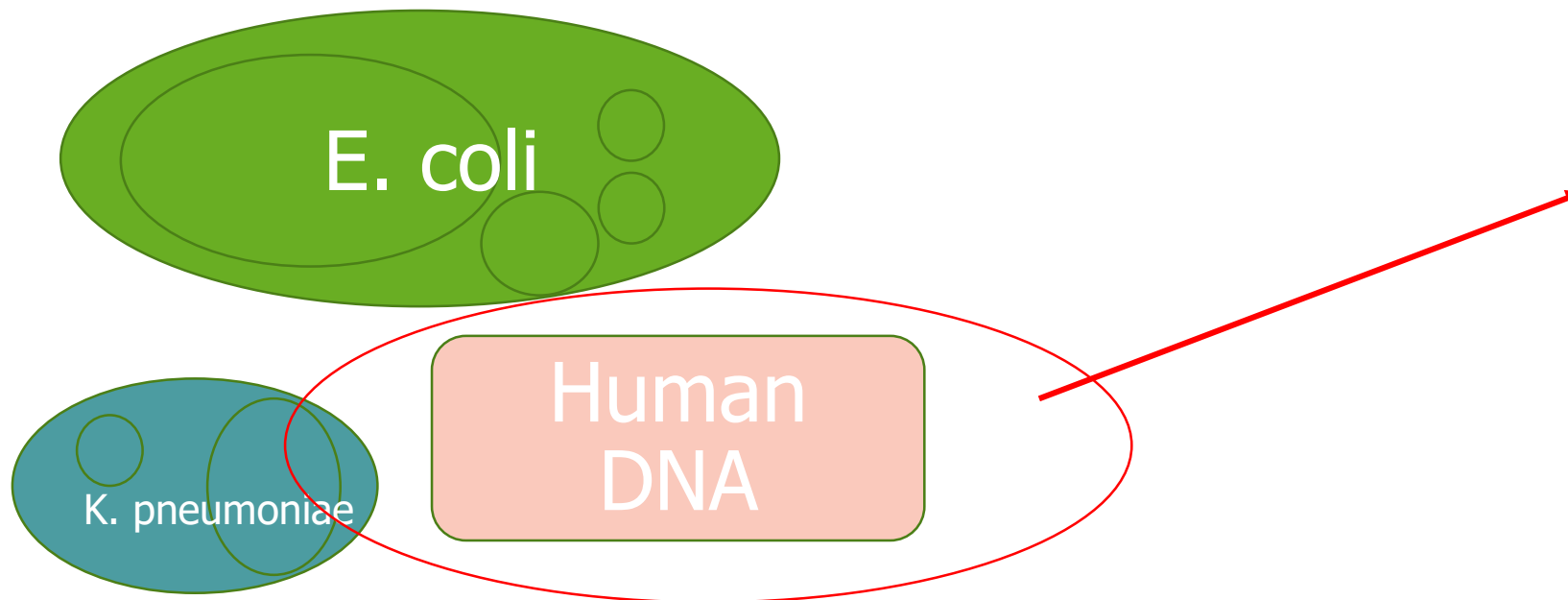
- Bacterial contamination in laboratory reagents and sequencing kits
- Cross-contamination across samples
- Host contamination
- Carryover contamination from previous sequencing runs

# When can you do something about contaminated sequencing data

- You can always remove unrelated sequences if you know what your target sequences are.
  - You rarely know exactly what you are looking for.

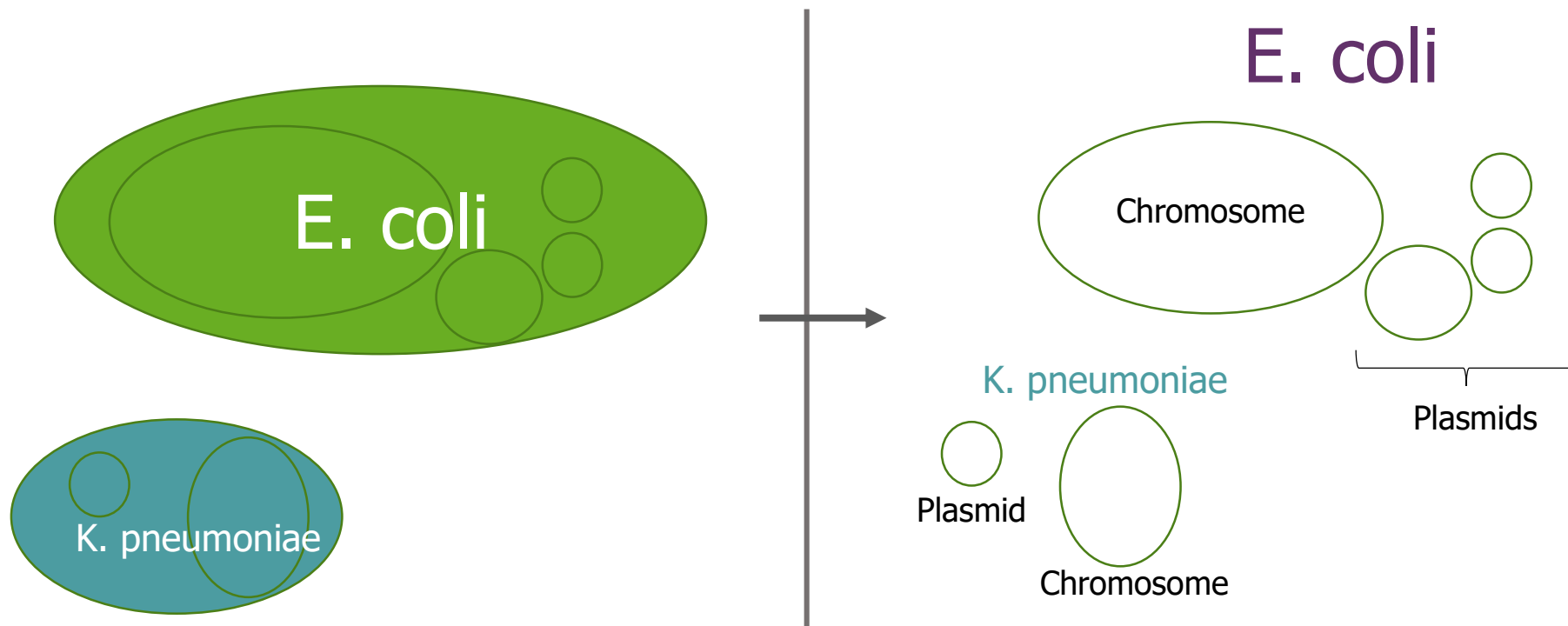
# When can you do something about contaminated sequencing data?

- You can always remove unrelated sequences if you know what your target sequences are.
  - You rarely know exactly what you are looking for but you might know exactly what you are not looking for.



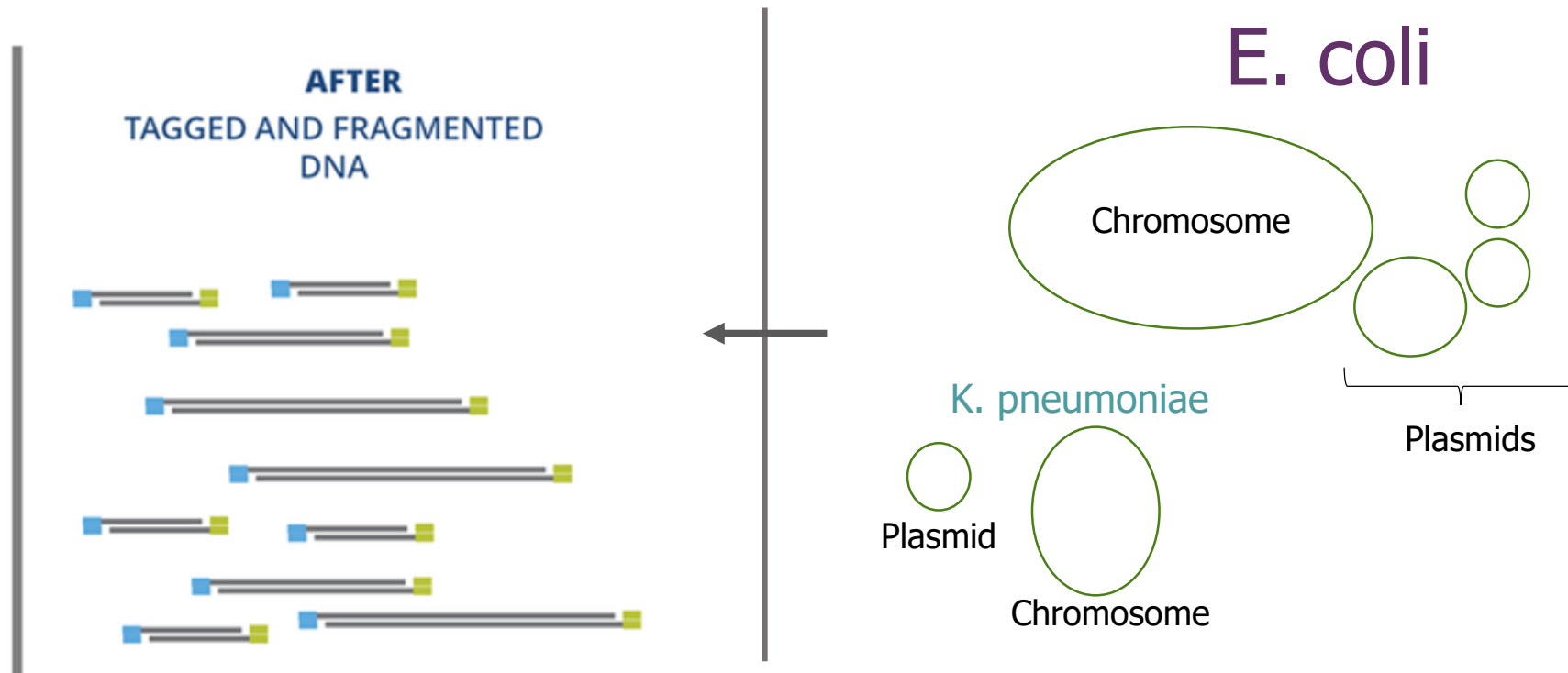
# When can you do something about contaminated sequencing data?

- You can always remove **unrelated** sequences if you know what your target sequences are.



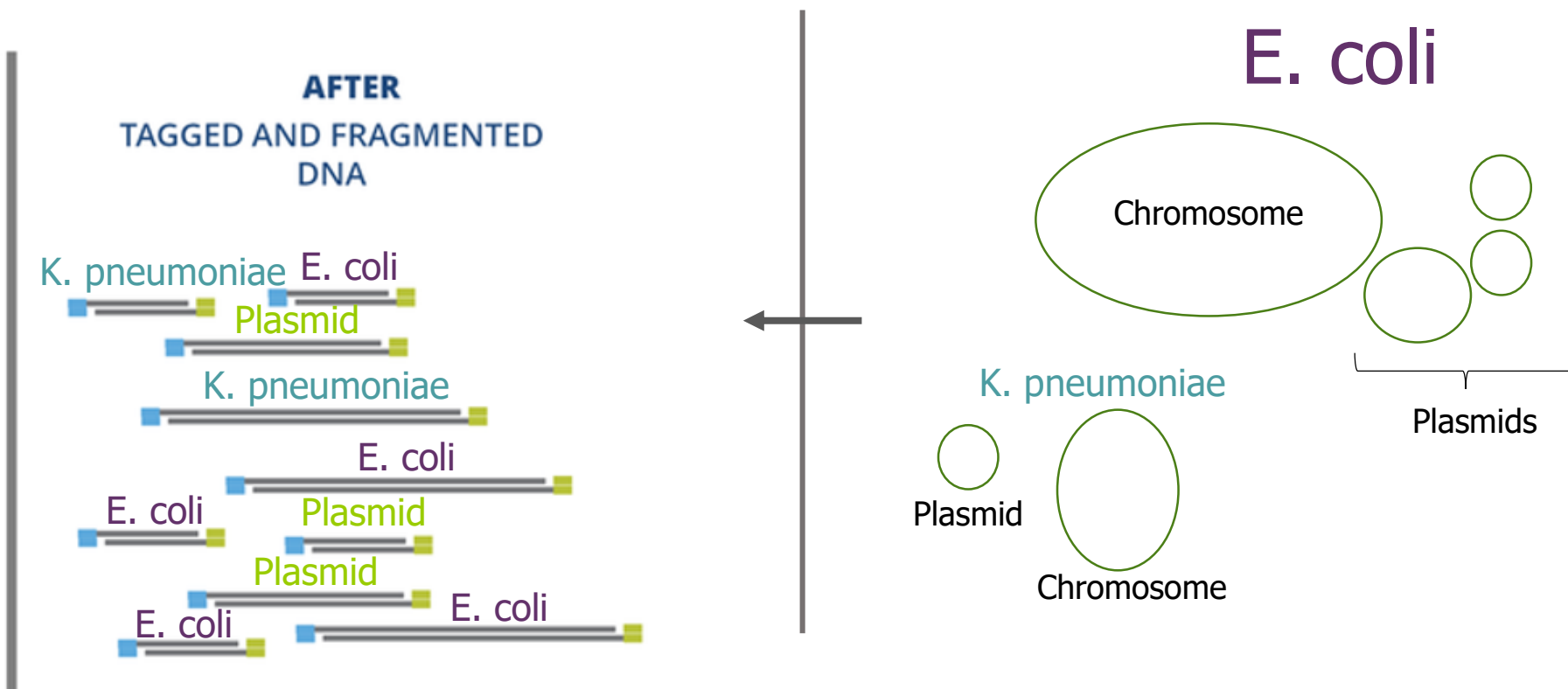
# When can you do something about contaminated sequencing data?

- You can always remove **unrelated** sequences if you know what your target sequences are.



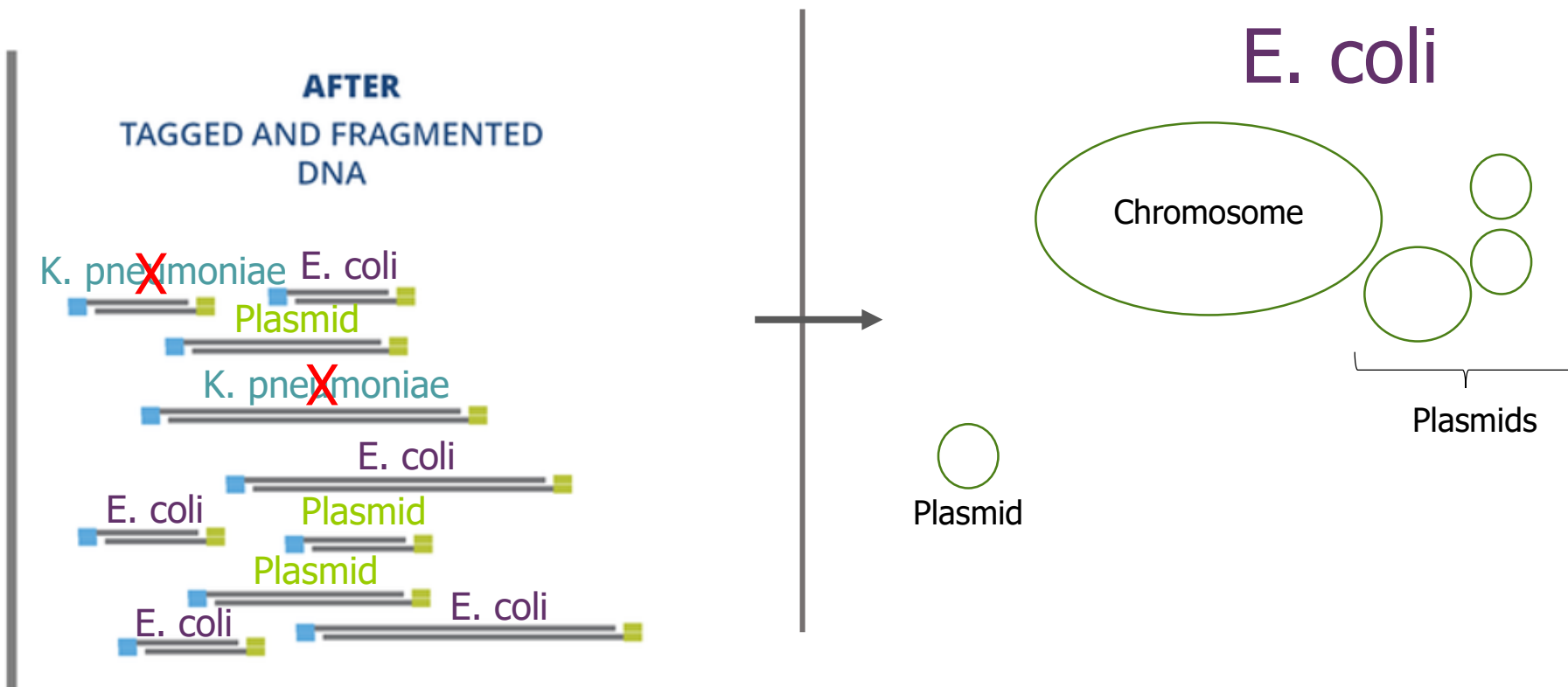
# When can you do something about contaminated sequencing data?

- You can always remove **unrelated** sequences if you know what your target sequences are.



# When can you do something about contaminated sequencing data?

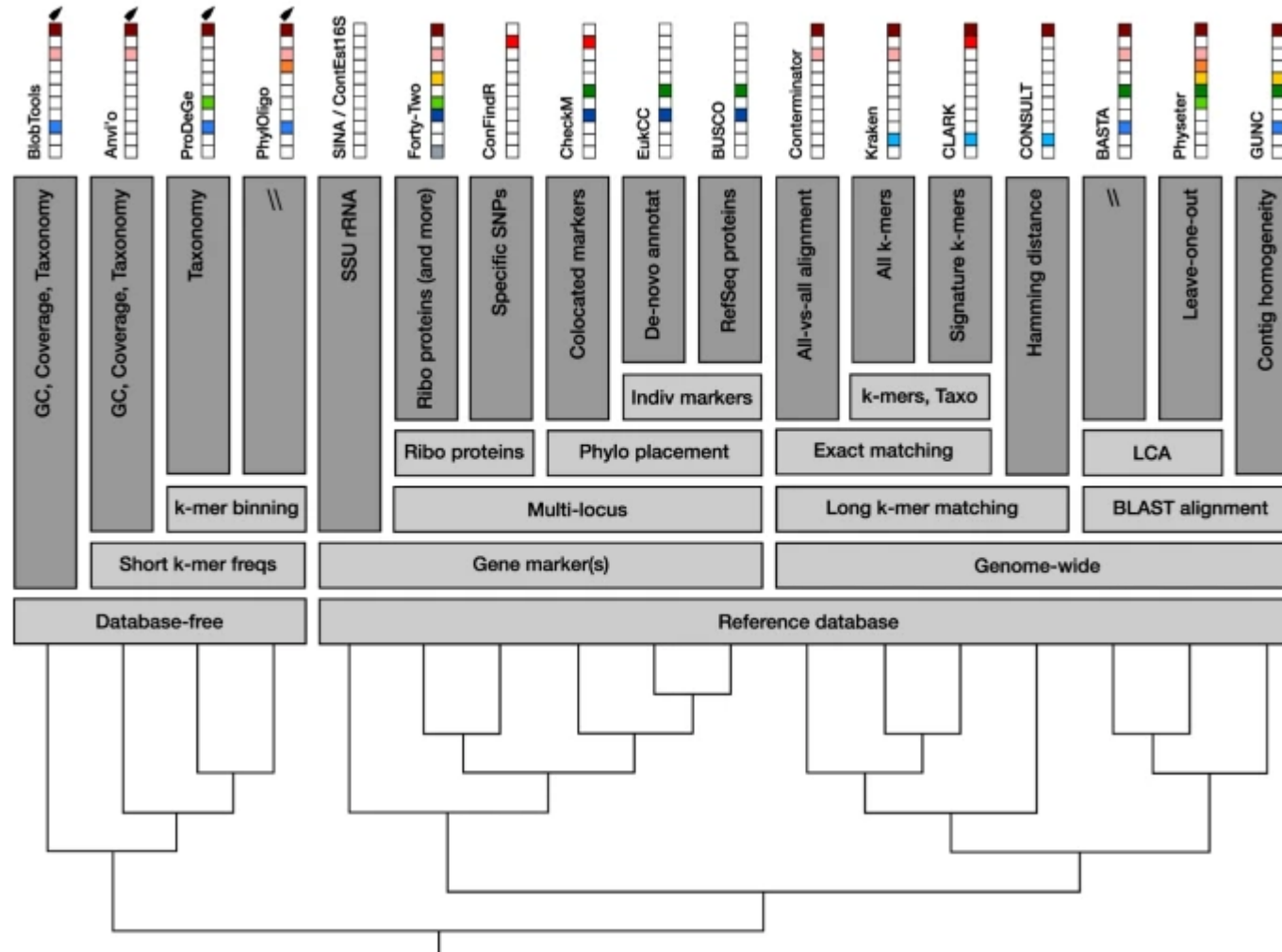
- You can always remove **unrelated** sequences if you know what your target sequences are.



# When can you do something about contaminated sequencing data?

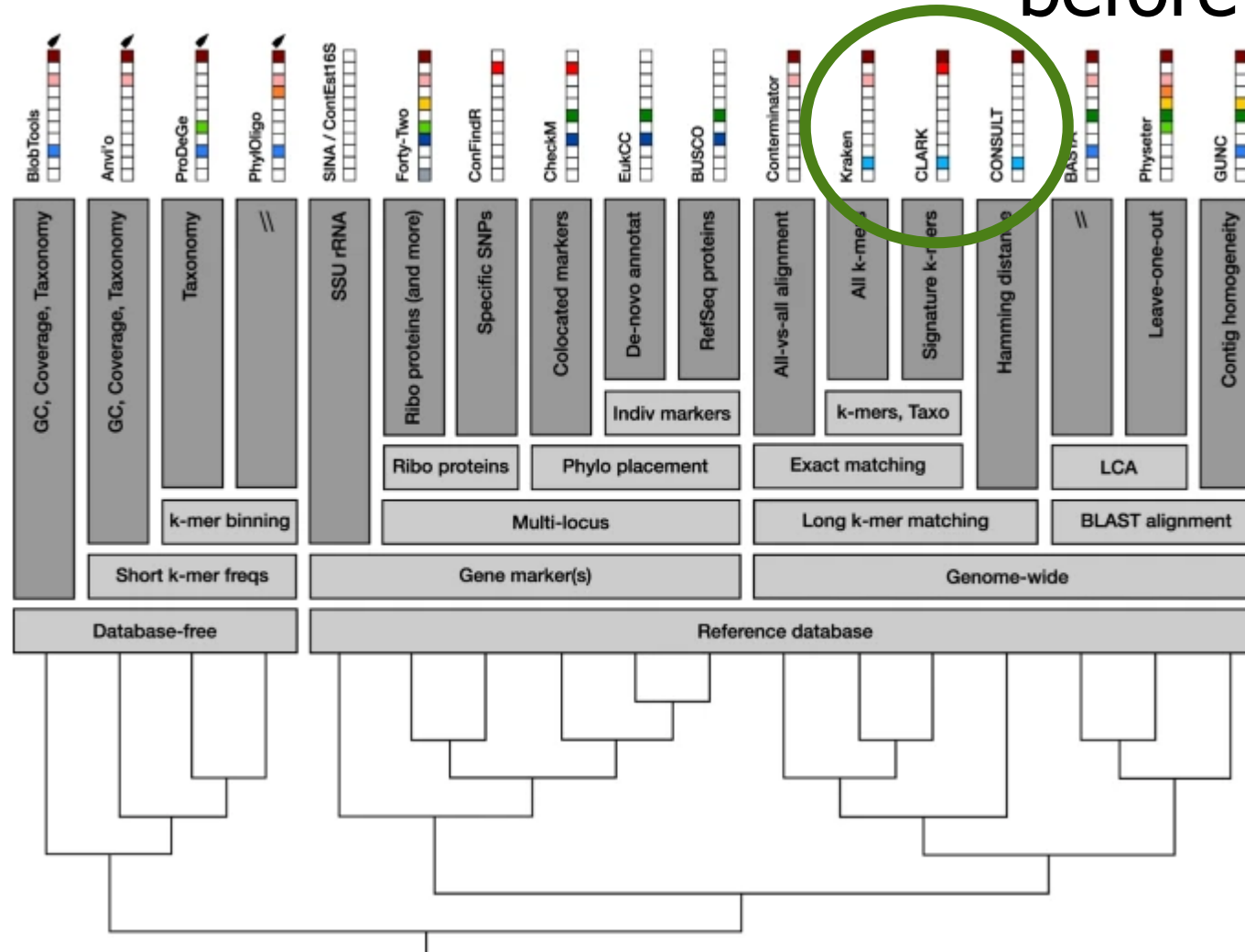
- You can always remove **unrelated** sequences if you know what your target sequences are.
- Consider which questions you want to answer with your sequencing data.
- When in doubt, re-sequence.

# Contamination control tools



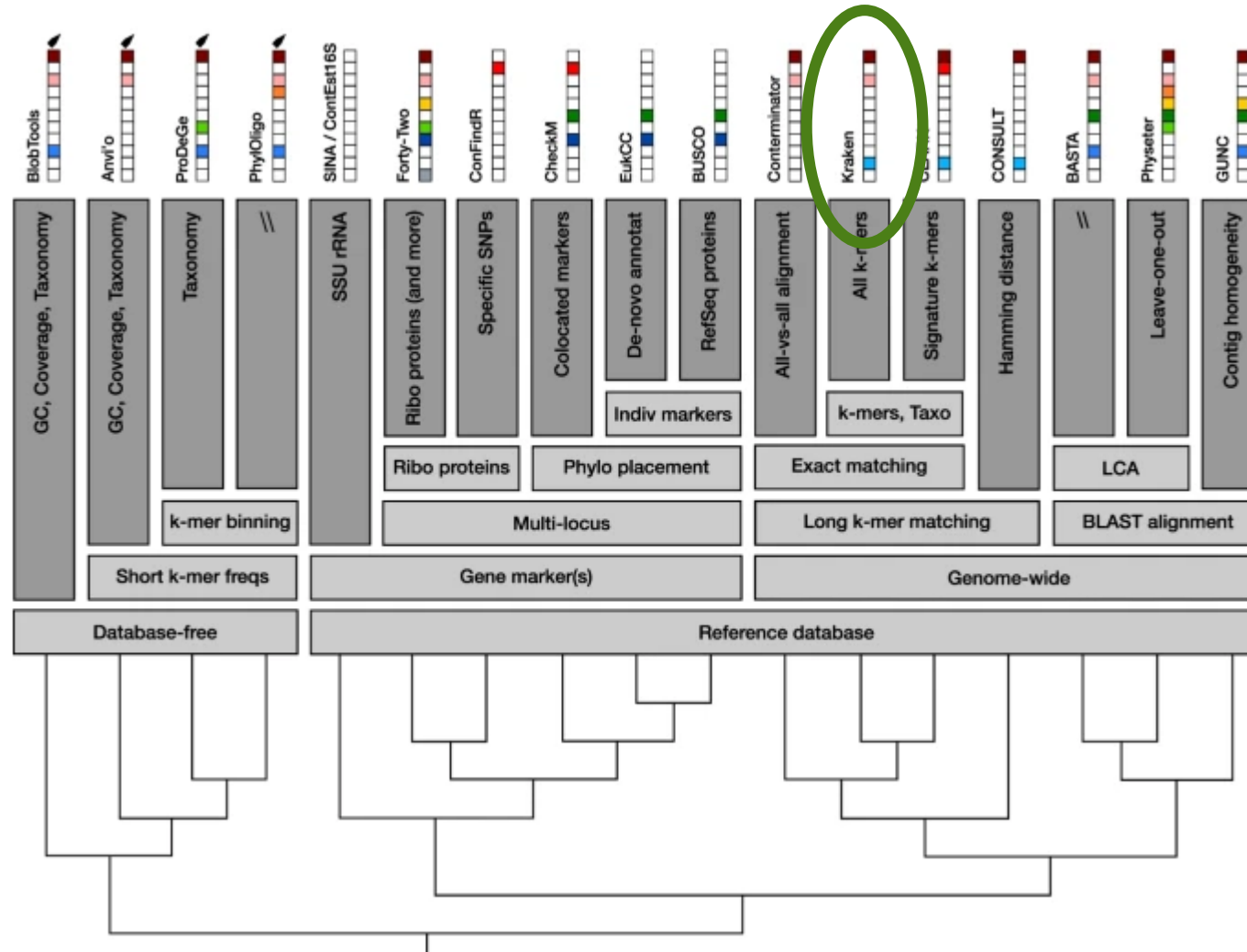
# Contamination control tools

Can estimate contamination before assembly



# Contamination control tools

## Kraken



# Contamination control tools

Kraken, a taxonomic sequence identifier.

- Exact alignments of k-mers and a novel classification algorithm.
- (<https://github.com/DerrickWood/kraken>)
- made for short reads.

**Kraken 1 | KrakenUniq | Kraken 2 | Kraken2Uniq | MiniKraken**

```
kraken2 --output [output/name] --db [db] --report sample_report --gzip-compressed [input.fastq.gz]
```

# Contamination control tools

## KMA

- A mapping method for both short and long reads
- <https://bitbucket.org/genomicepidemiology/kma/src/master/>

#download fasta files for database

#index database

```
kma index -i [templates.fsa.gz] -o [database/name]
```

#run kma

```
kma -i [someLongReads.fq.gz] -o [output/name] -t_db [database/name] -bcNano -bc 0.7
```

## Part 1

# Illumina quality assessment

# Practicals

Head to the course page on EVA

Navigate to session one and click into the Part 1 – Illumina Quality Assessment page

GL HF!

## Part 2

# Nanopore quality assessment

# Practicals

Head to the course page on EVA

Navigate to session one and click into the Part 2 –Nanopore quality control page

Enjoy!

# Acknowledgements

The creation of this training material was commissioned by ECDC to Statens Serum Institut (SSI) with the direct involvement of Sharmin Baig, Søren Hallstrøm, Kasper Thystrup Karstensen & Astrid Rasmussen