

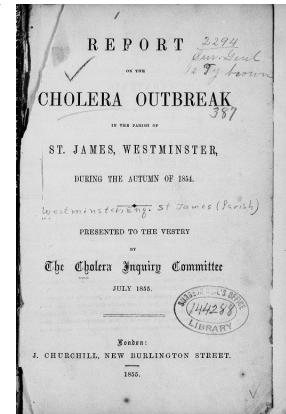
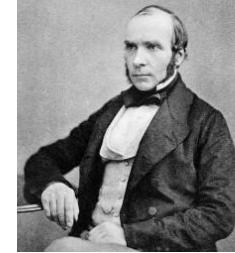
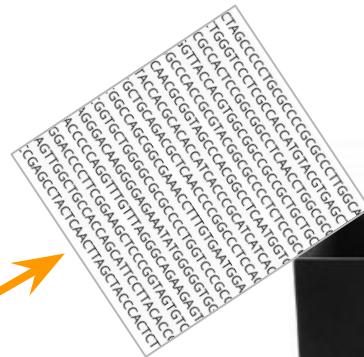
# Introduction to bioinformatics for microbial genomics

A/Prof Torsten Seemann



# Introduction

# Bioinformatics



# What's inside the black box?



# It's black boxes, all the way down.



Data / file formats



Software / analysis



Results / interpretation



Reporting / action

# We use real black boxes too!

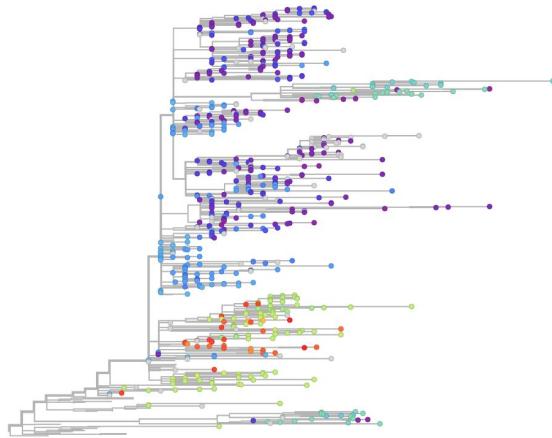


# Traditional public health and clinical microbiology

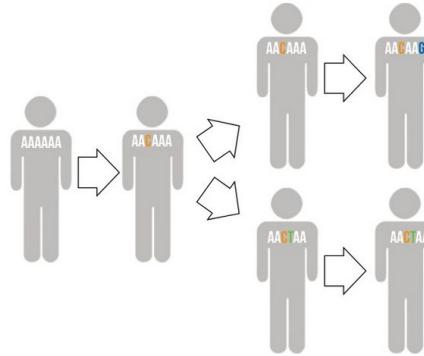
# The roles of a public health laboratory



Diagnostics



Surveillance



Outbreak response

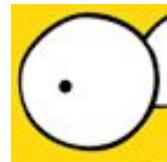
# Tracking a diverse range of pathogens



Pathogen: *Homerbacter simpsonii*

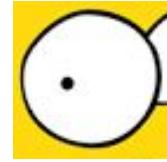
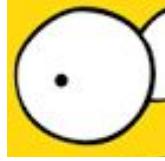


# Focus on a small informative section



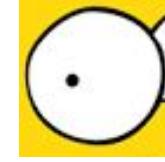
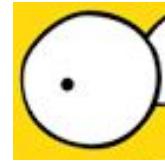
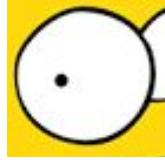
Genotyping via one or more marker genes.

# Another sample arrives at the lab



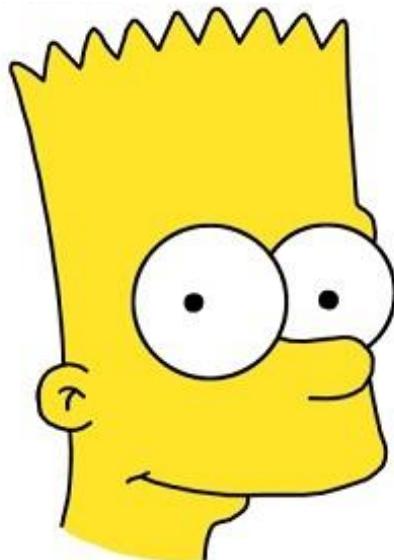
Looks related.

# Another sample from same city



Possible outbreak ?

D'oh !



# How can we improve this?

Of course you know the answer.

# Whole genome sequencing

- :: Use 100% of the genome
- :: Single nucleotide resolution
  - : Infer source and transmission of infection
- :: Full complement of genes
  - : Track mobile elements and antimicrobial resistance



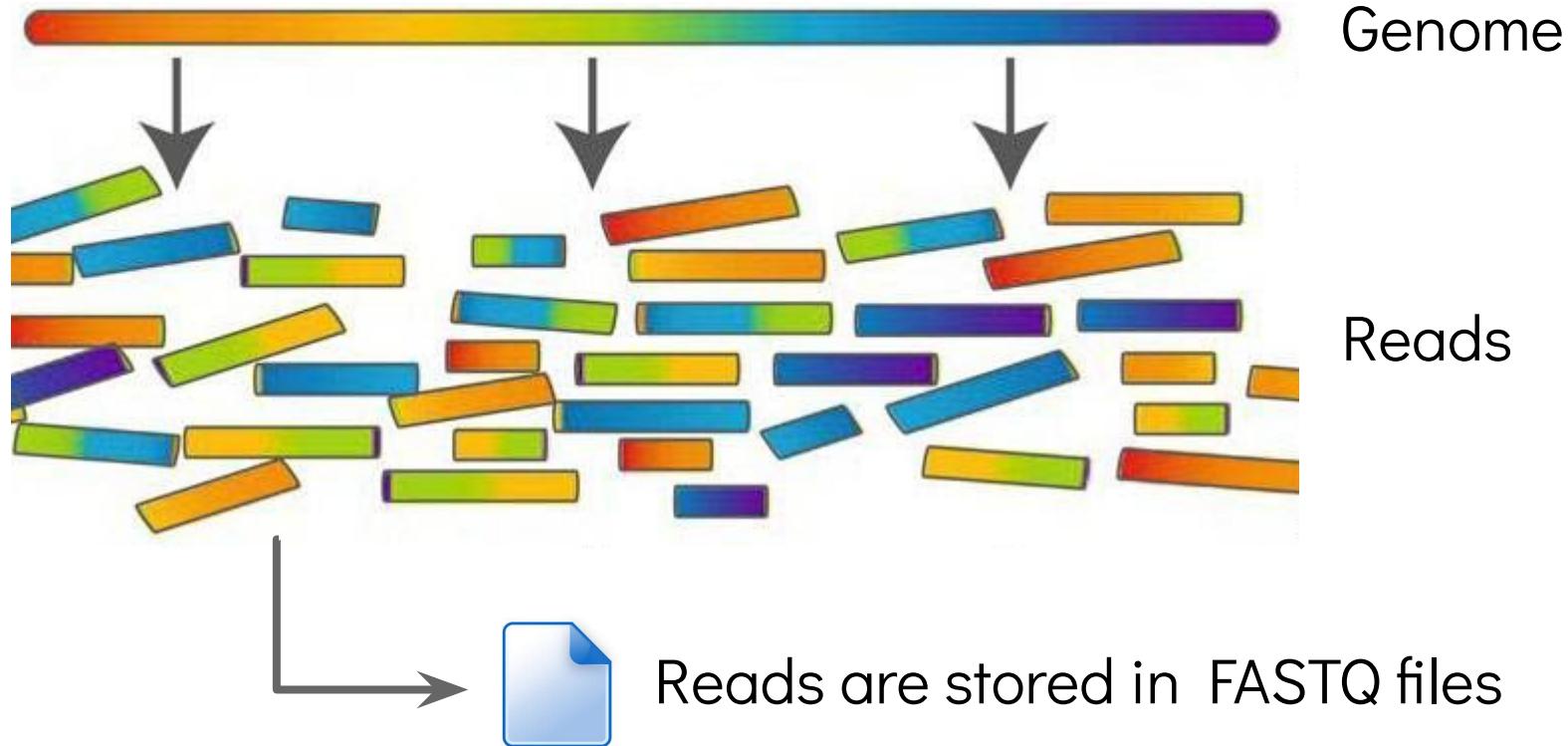
# Routine sequencing of isolates



Can sequence 100s of genomes per week

# Whole genome sequencing

# The currency of genomics



# Types of sequence reads



100 - 300 bp



5,000 - 45,000 bp



5,000 - 900,000+ bp

# FASTQ files

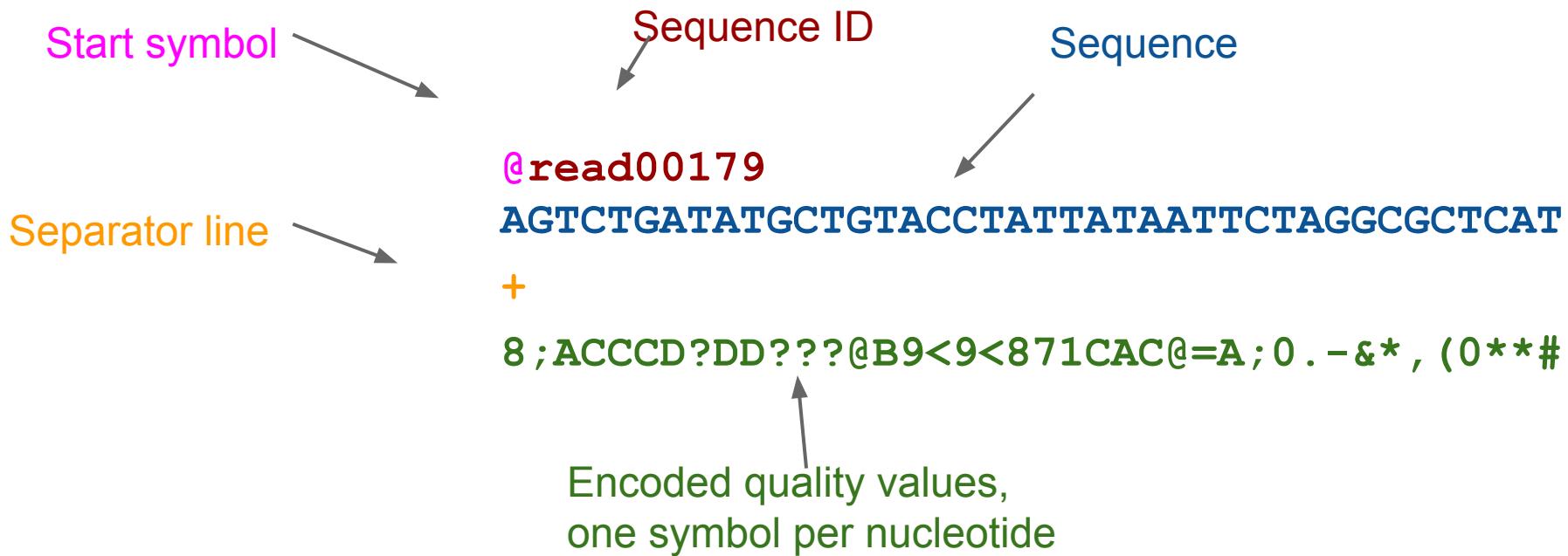
# The FASTQ format



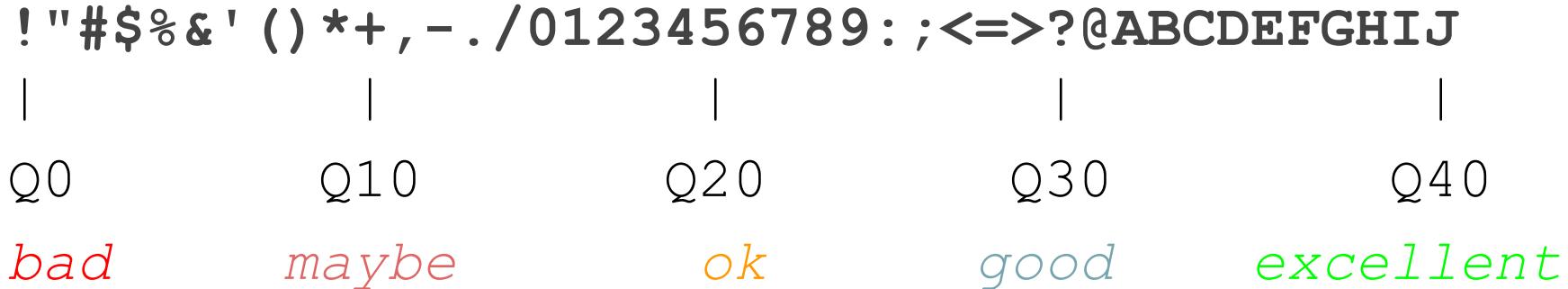
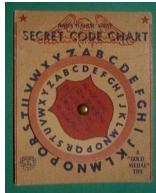
A single sequence read looks like this:

```
@read00179
AGTCTGATATGCTGTACCTATTATAATTCTAGGCGCTCAT
+
8 ;ACCCD?DD???@B9<9<871CAC@=A ;0 . - &* , (0***#
```

# FASTQ components



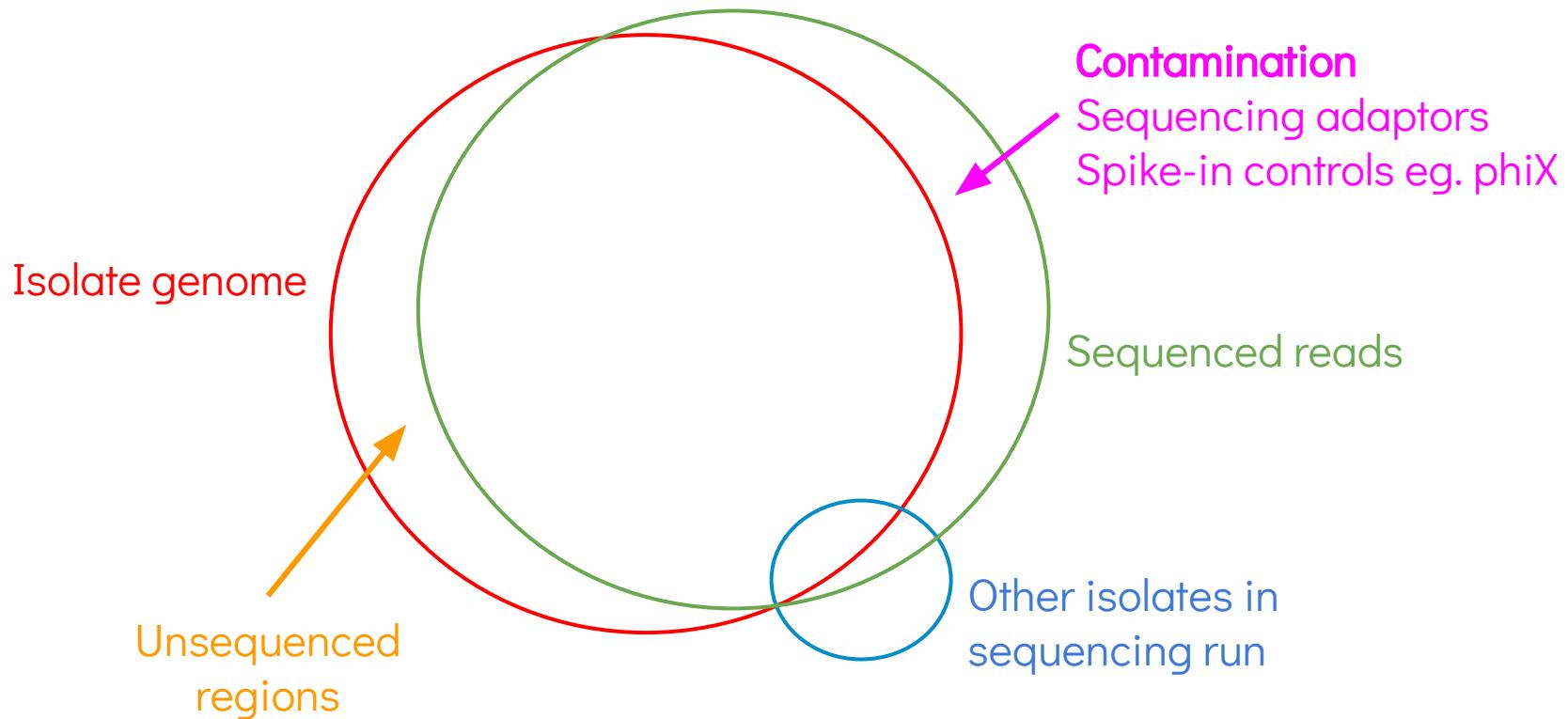
# FASTQ quality encoding



- Use characters to represent numbers
  - numbers measure quality
  - a reliability estimate for each base
  - estimated from the physical measuring process
  - formalized during the Human Genome Project

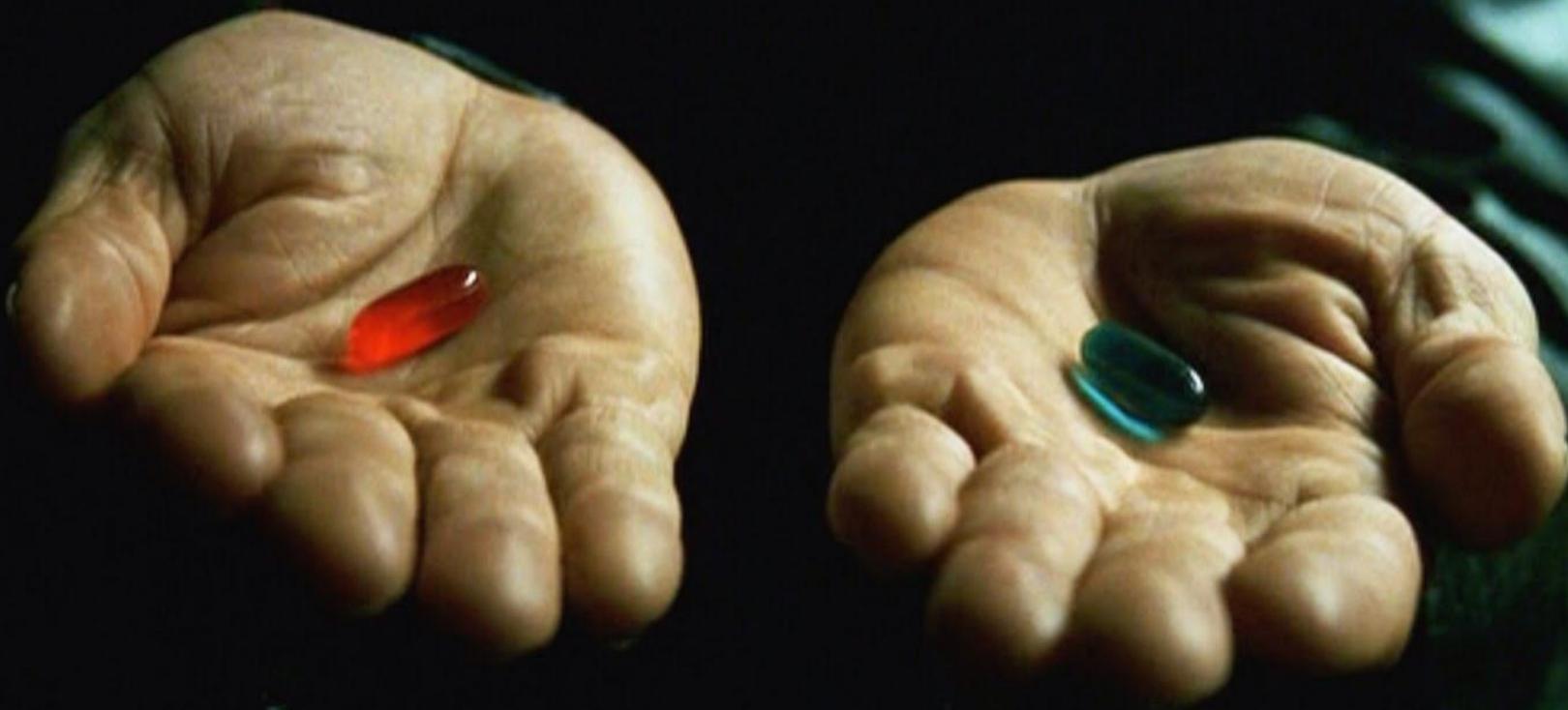
# Quality control

# What data do we really have?



**Got my reads, now what?**

# How to use raw sequencing reads





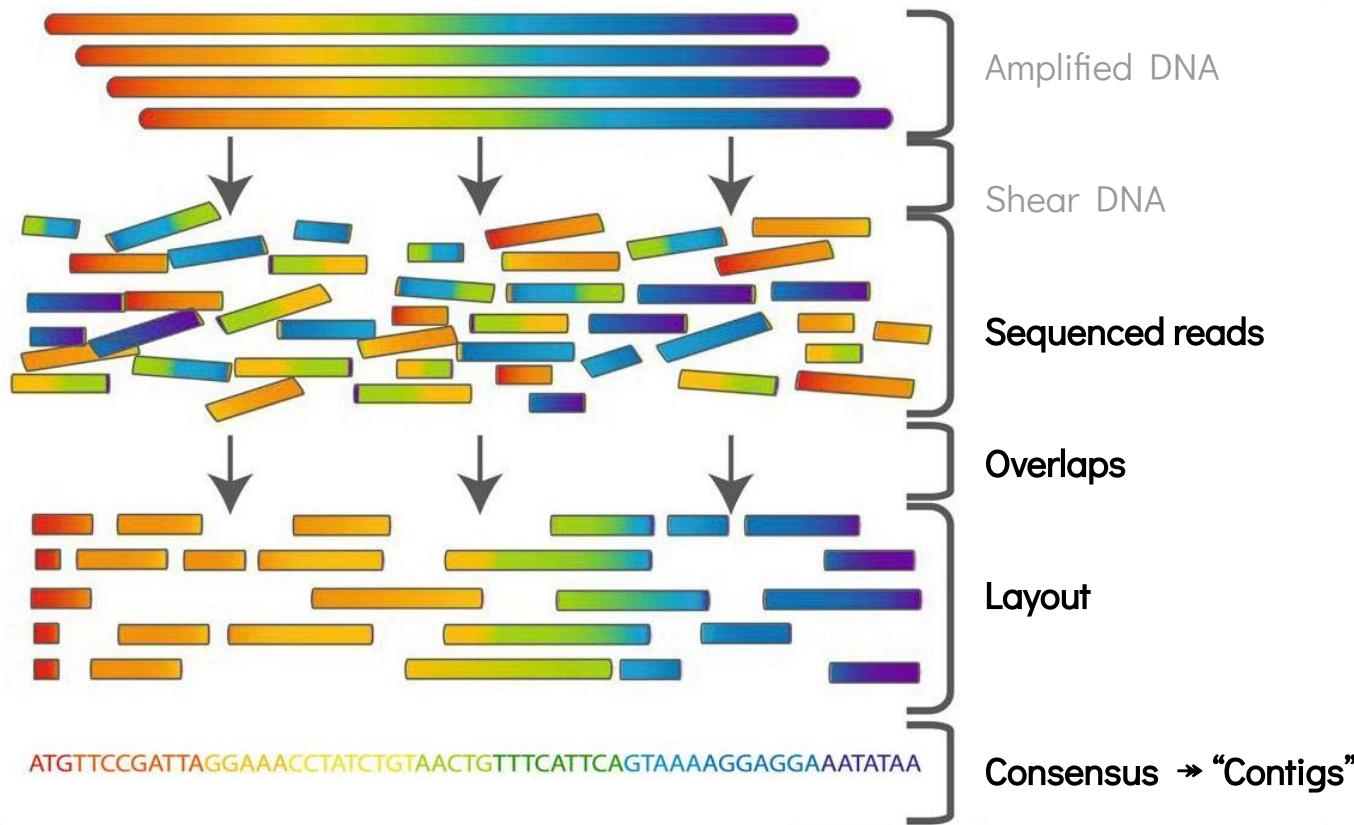
Assemble



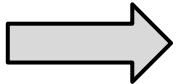
Align

# Genome assembly

# *De novo* genome assembly



# *De novo* genome assembly



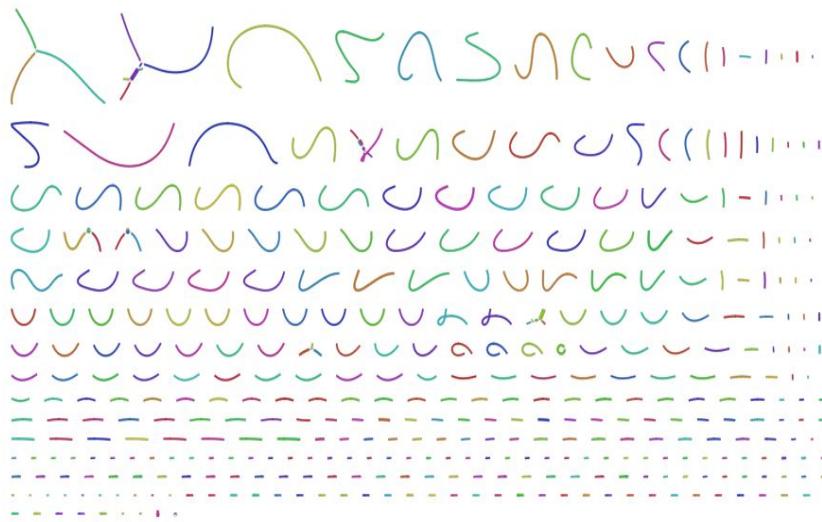
“From scratch”

# *De novo* genome assembly

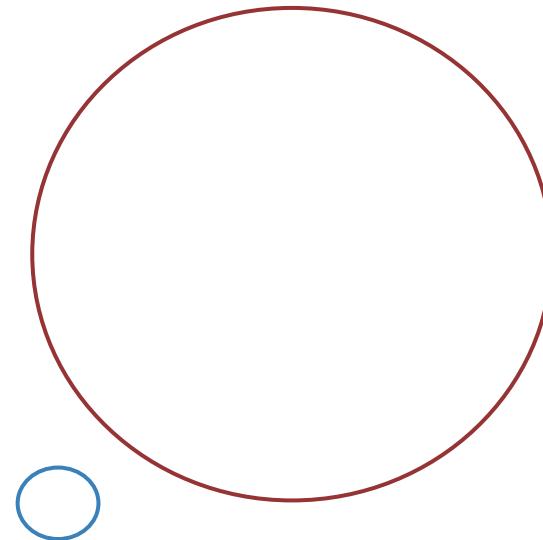


Usually sequencing a *population* of cells

# The effect of read length



150 bp - Illumina (short)



8,000 bp - Nanopore (long)

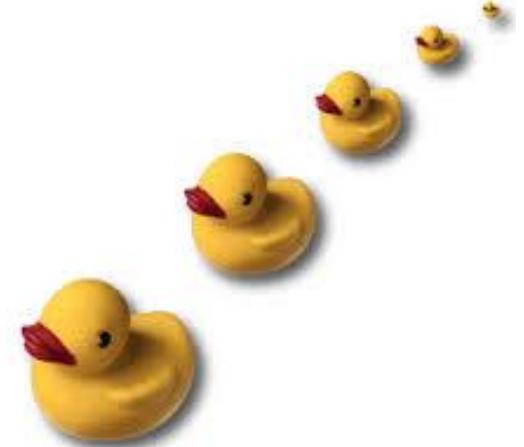
# Align to reference

# Align to reference

Seven short 4bp reads

AGTC TTAC GGGA CTTT

TAGG TTTA ATAG



Aligned to 31bp reference

**AGTCTTTATTATAGGGAGCCATAGCTTTACA**

AGTC                TAGG                ATAG                TTAC

TTTA                GGGA                CTTT

# Ambiguous alignment

Eight short 4bp reads

AGTC TTAC GGGA CTTT

TAGG TTTA ATAG **TTAT**



Aligned to 31bp reference

AGTCTTATTATAGGGAGCCATAGCTTACA

AGTC TAGG ATAG TTAC

TTTA

GGGA

CTTT

**TTAT**

**TTAT**

*D'oh!*

# Using reads, assemblies, and alignments

# Finding differences

SNP      Deletion      Reference

**AGTCTGATTAGCTTAGCTTGATAGCGCTATATTAT**

AGTCTGATTAGCTTAGAT

ATTAGCTTAGATTGTAG

CTTAGATTGTAGC-C

TGATTAGCTTAGATTGTAGC-CTATAT

TAGCTTAGATTGTAGC-CTATATT

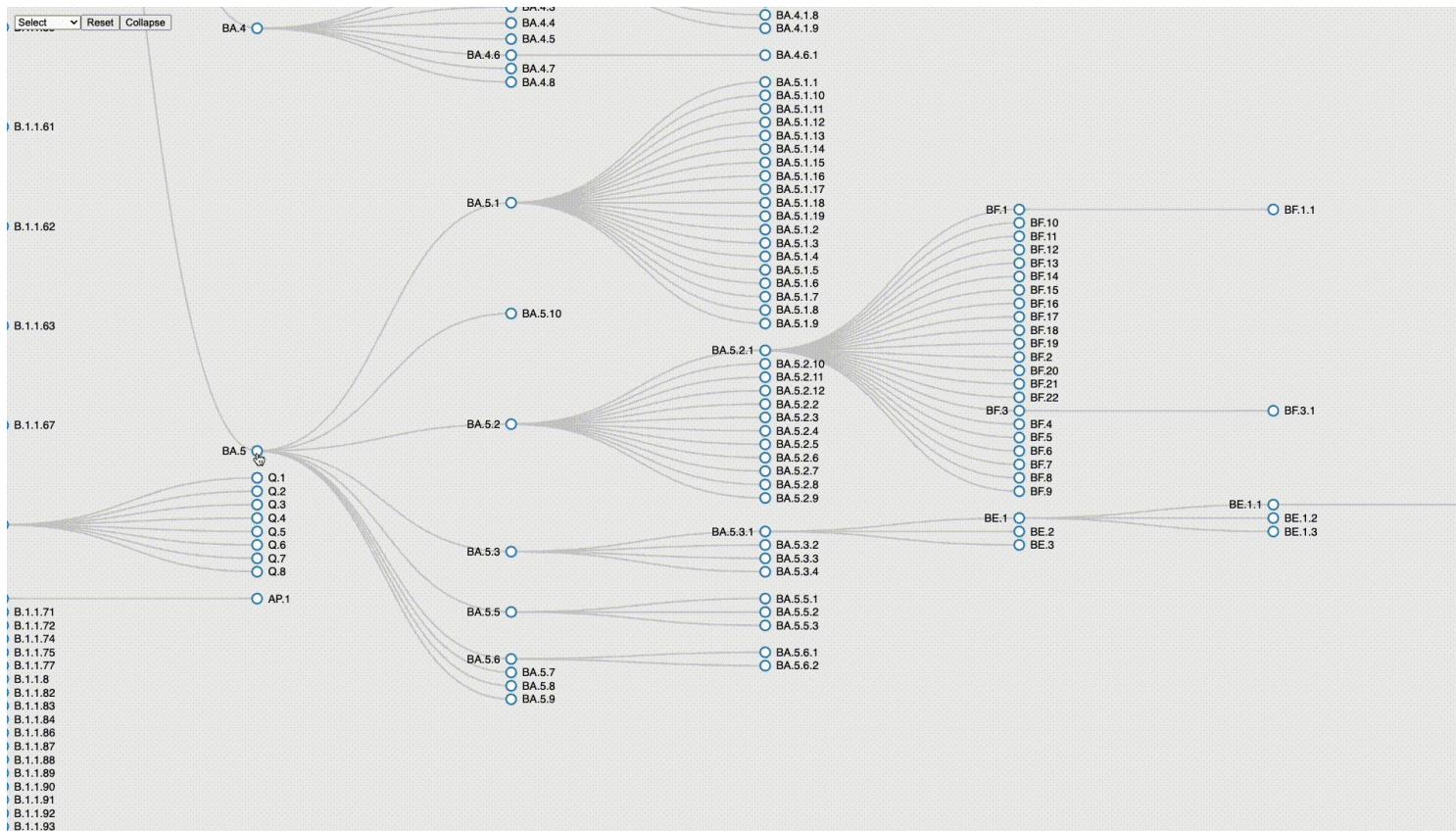
TAGATTGTAGC-CTATATTA

TAGATTGTAGC-CTATATTAT

Reads

The diagram illustrates the process of finding differences between a reference genome and multiple sequencing reads. The reference genome is shown at the top in bold black text. Below it, several sequencing reads are listed, each starting with a portion of the reference genome followed by variations. Vertical dashed lines are used to highlight specific differences: one line marks the position of an SNP (a change from 'G' to 'C' at position 10), and another line marks a deletion (the loss of the 'CT' dinucleotide at position 6). The word "Reference" is positioned to the right of the top sequence, and "Reads" is positioned to the right of the bottom sequences.

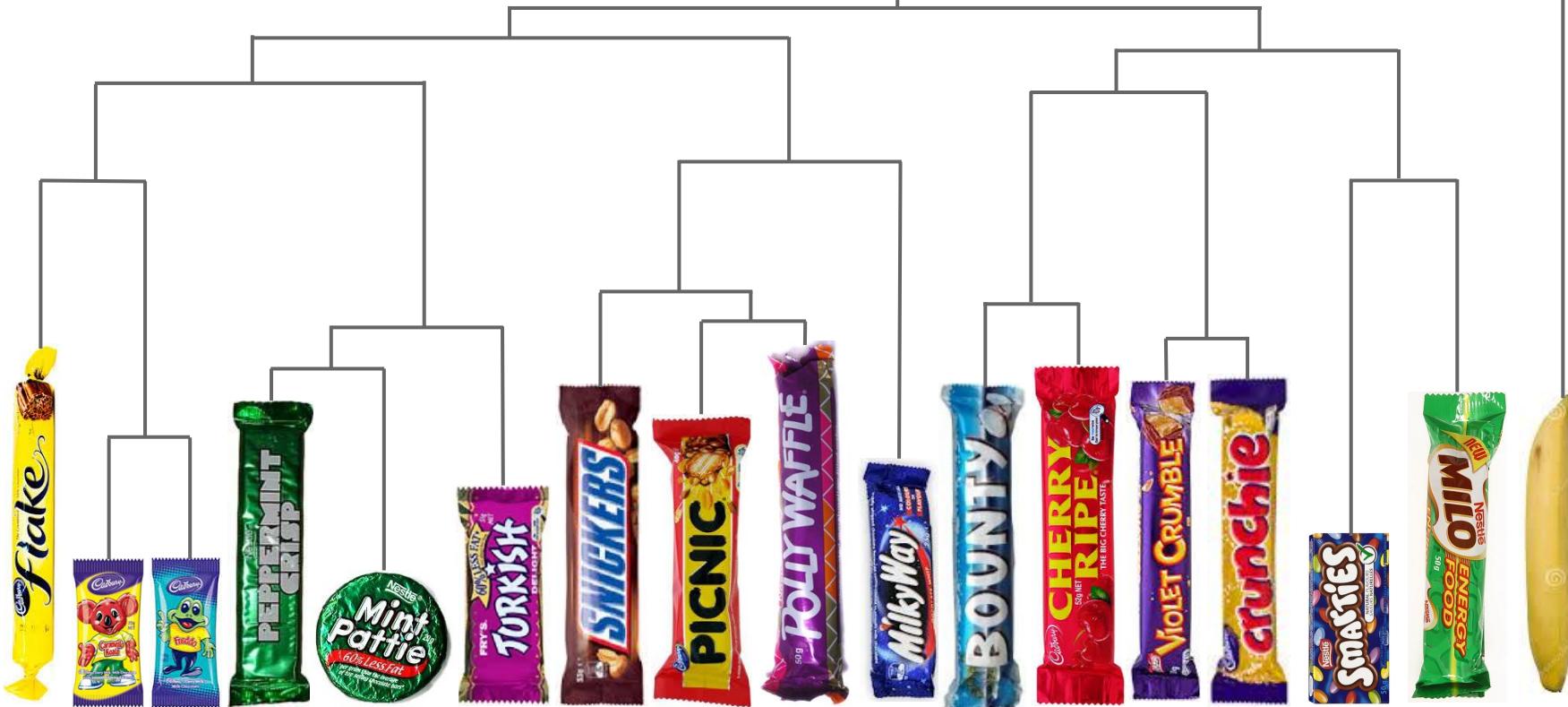
# Genotyping / lineages



# Antimicrobial resistance

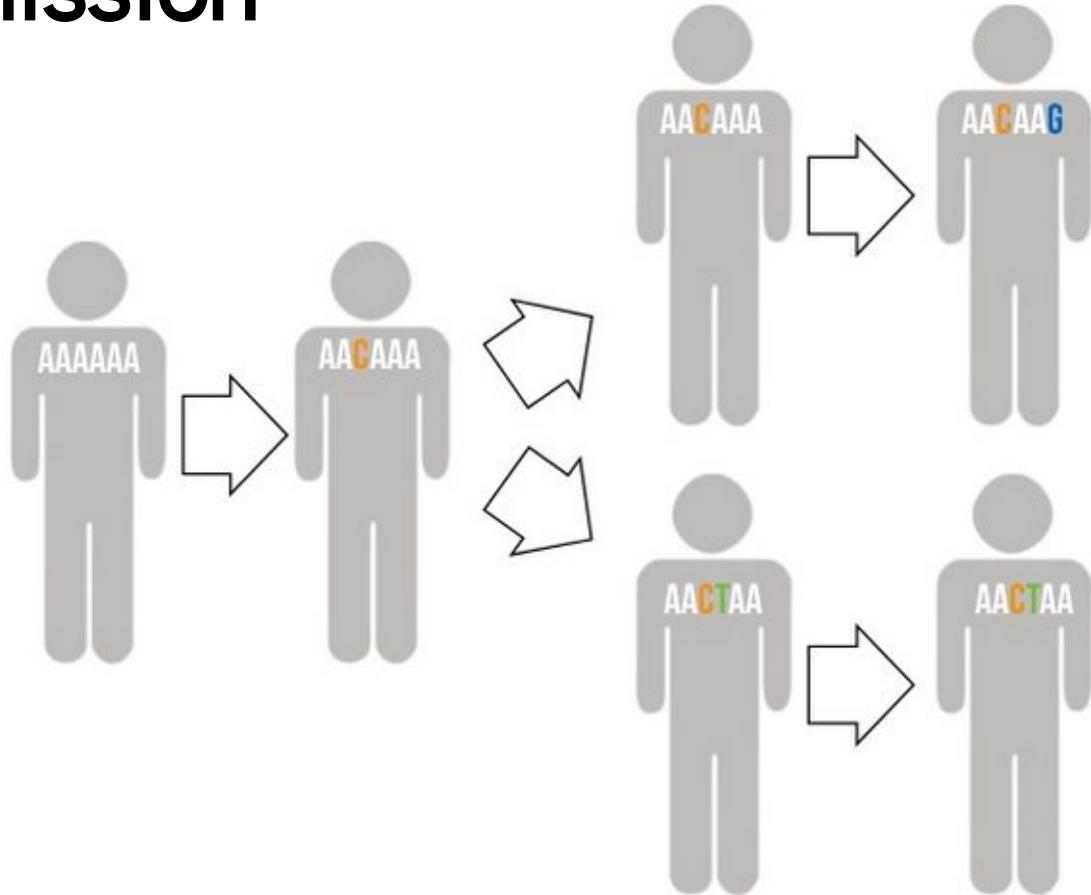


# Phylogenomics



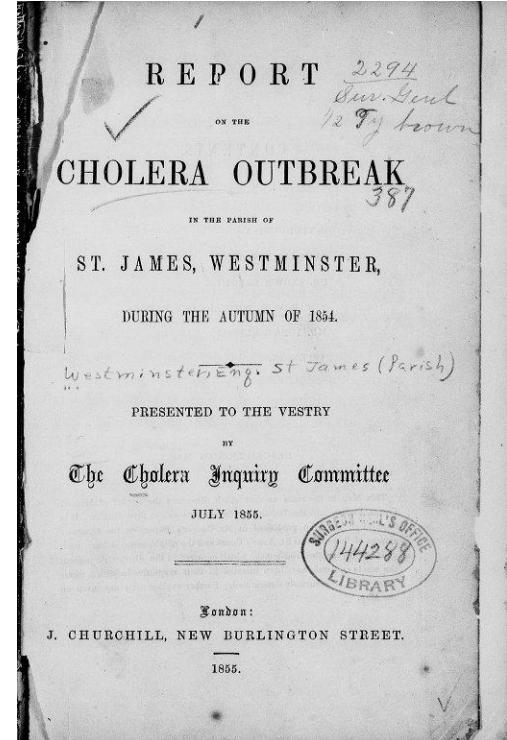
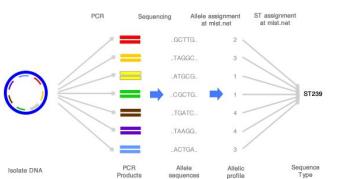
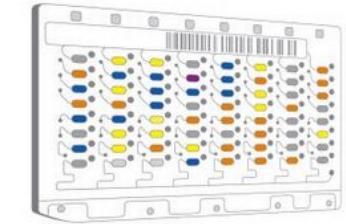
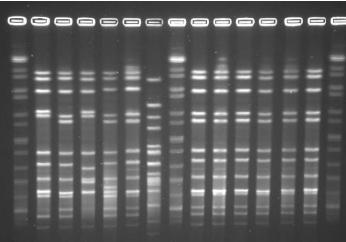
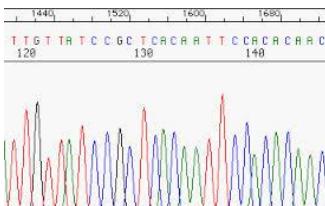
# Inferring transmission

- :: Identical sequence does not imply transmission
- :: Easier to rule out than in

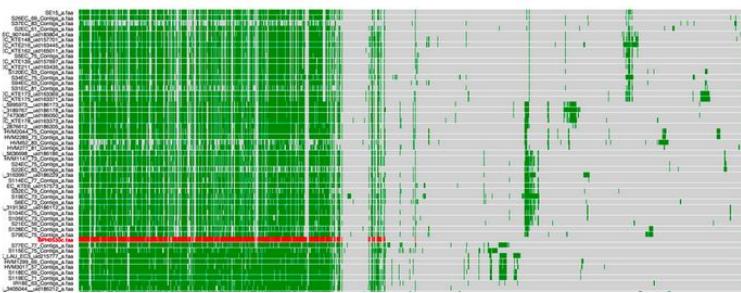
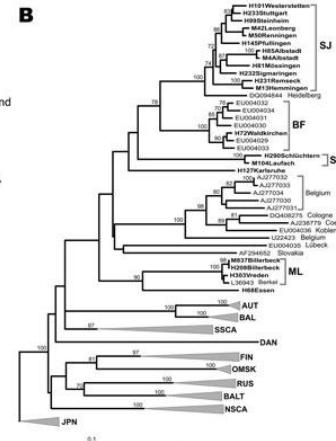


# Conclusions

# Traditional workflow



# Modern workflow



A chimpanzee is shown from the chest up, wearing black-rimmed glasses and a light blue button-down shirt. It is seated at a white desk, looking directly at the camera with a neutral expression. Its hands are resting on the desk, which has some papers and a pen on it. In the background, there's a dark object, possibly a bookshelf or a wall.

Thank you for listening

