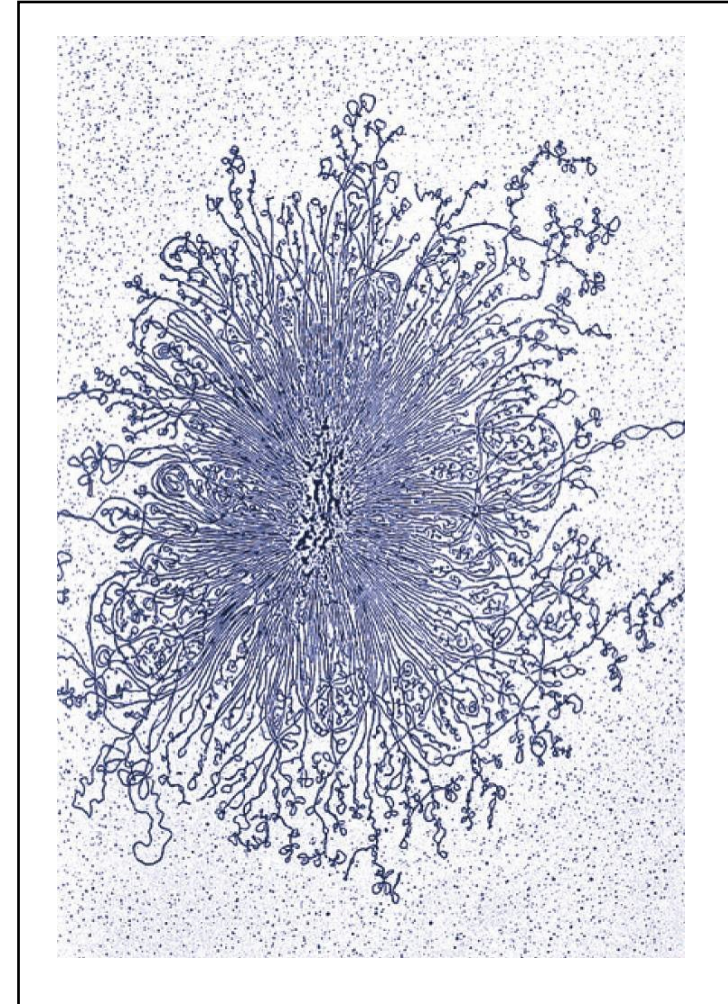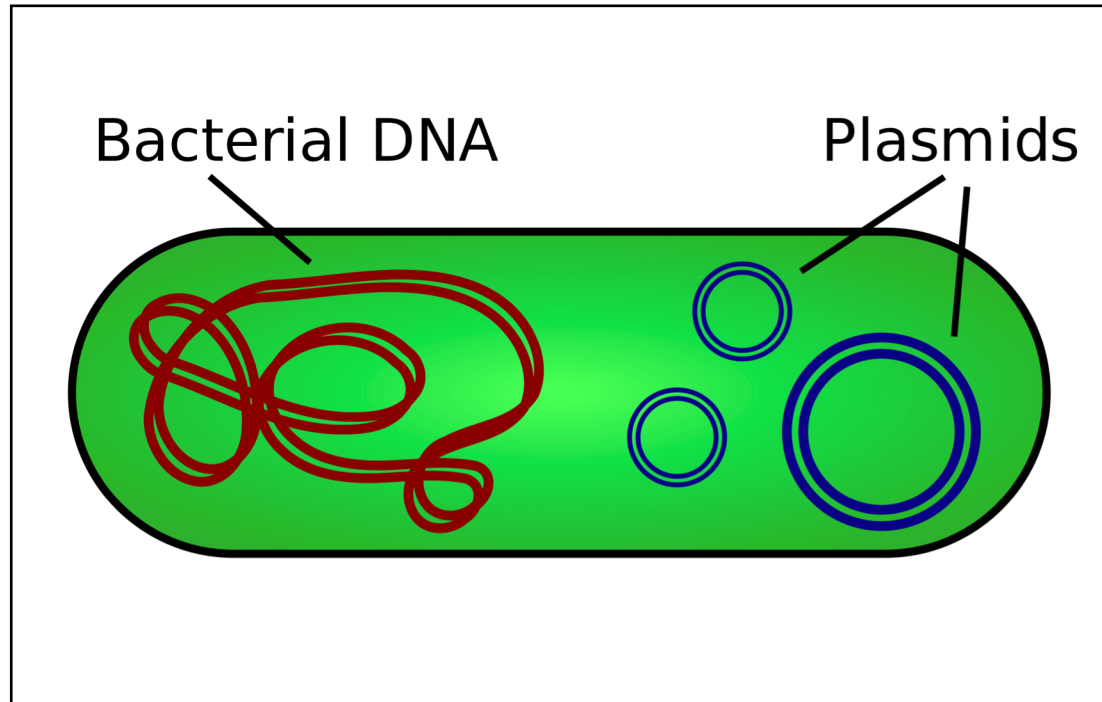# Closing
## (some of the issues with circular)
# Genomes

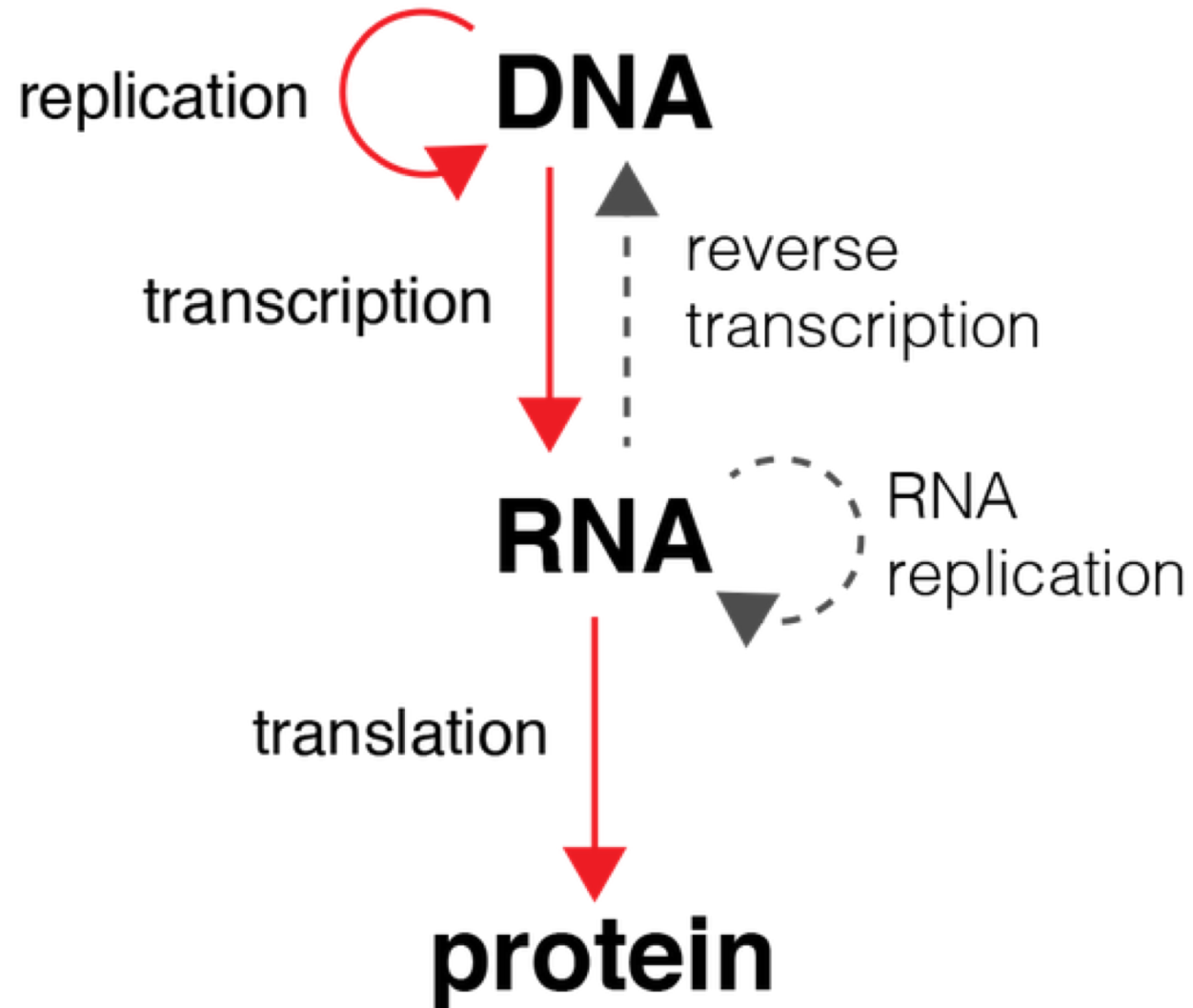Erin Young, PhD
September 17, 2021



Jurassic Park, 1993 (28 years ago)

# Bacterial genomes have a large circular chromosome made of DNA

# Central Dogma of Biology

# We sequence bacterial DNA to improve public health

## Whole Genome Sequencing

MiSeq came out in 2011

CDC is tracking and classifying illness in a new way, using advanced technology to find and stop outbreaks and combat drug-resistant germs.

### On This Page

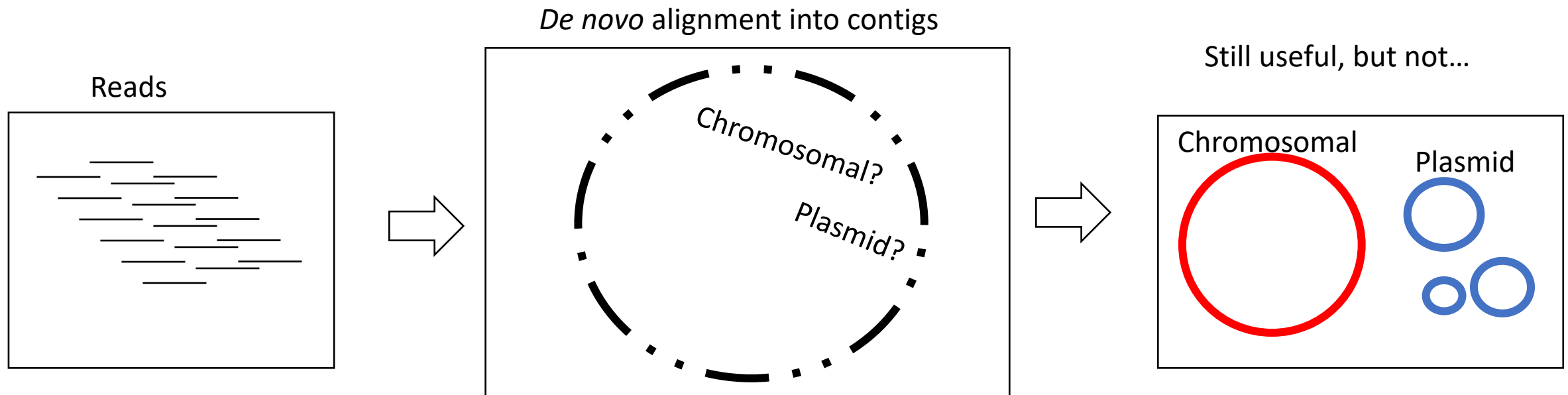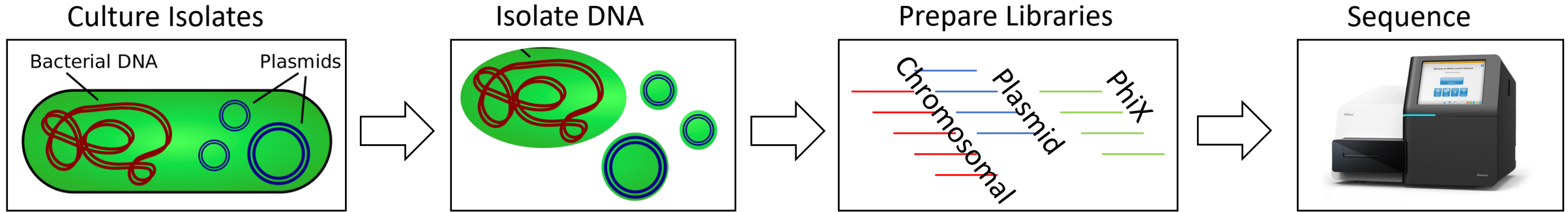Foodborne Outbreak Investigations          Beyond Foodborne Outbreaks

Success Stories          Antibiotic Resistance

Whole genome sequencing (WGS) provides detailed genetic information about germs that make people sick. CDC's Division of Foodborne, Waterborne, and Environmental Diseases uses this information to improve efforts to find, investigate, and prevent illnesses caused by bacteria, fungi, and parasites. This is especially

# Illumina WGS Sequencing

Culture Isolates

Bacterial DNA  Plasmids

Isolate DNA

Prepare Libraries

Chromosomal  Plasmid  PhiX

Sequence

Reads

*De novo* alignment into contigs

Chromosomal?

Plasmid?

Still useful, but not...

Chromosomal

Plasmid

# Long-range sequencing is less likely to have issues with troublesome regions



**Figure 1**
A schematic showing how long-read sequencing can deliver simplified, less ambiguous genome assembly. Long reads (solid arrows) have greater overlap with other reads than is provided by short reads (dashed arrows), allowing more accurate assemblies, especially in repeat regions (R). Image adapted from Schatz (2014)[4].
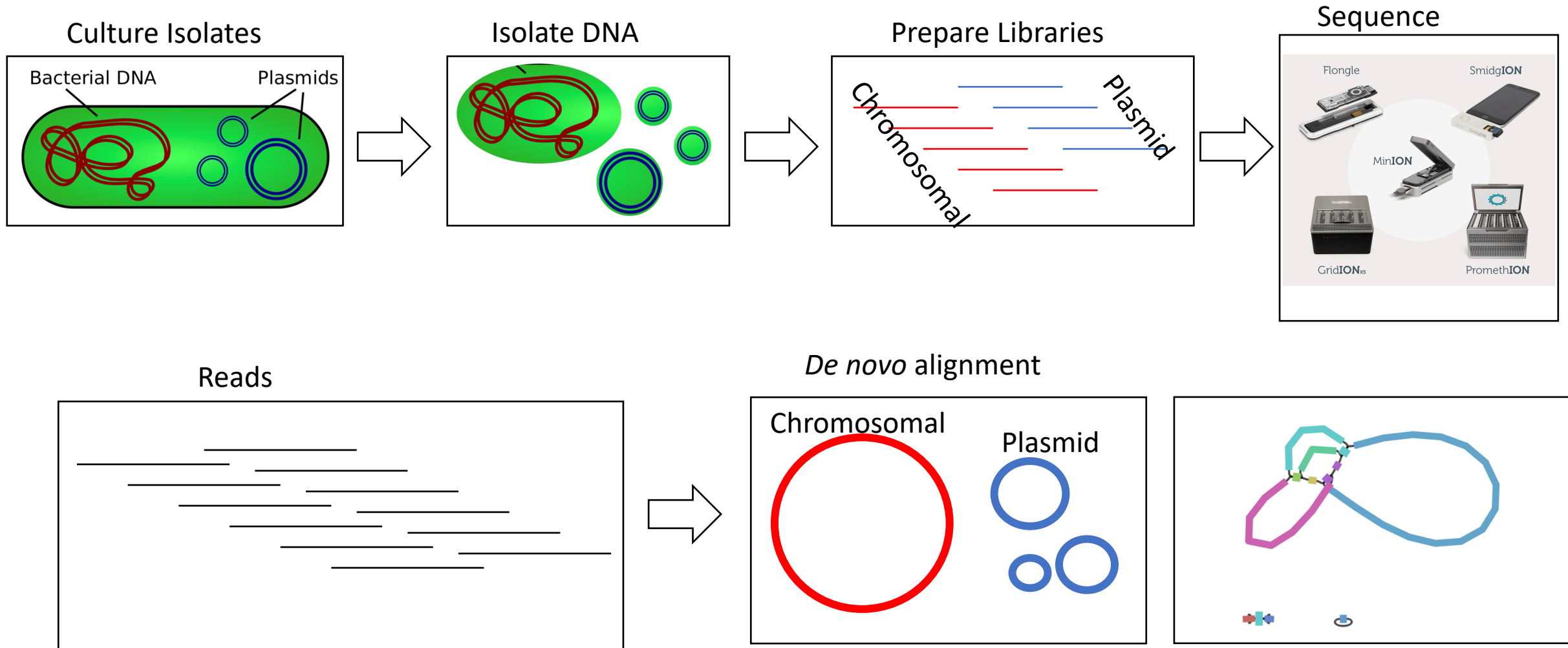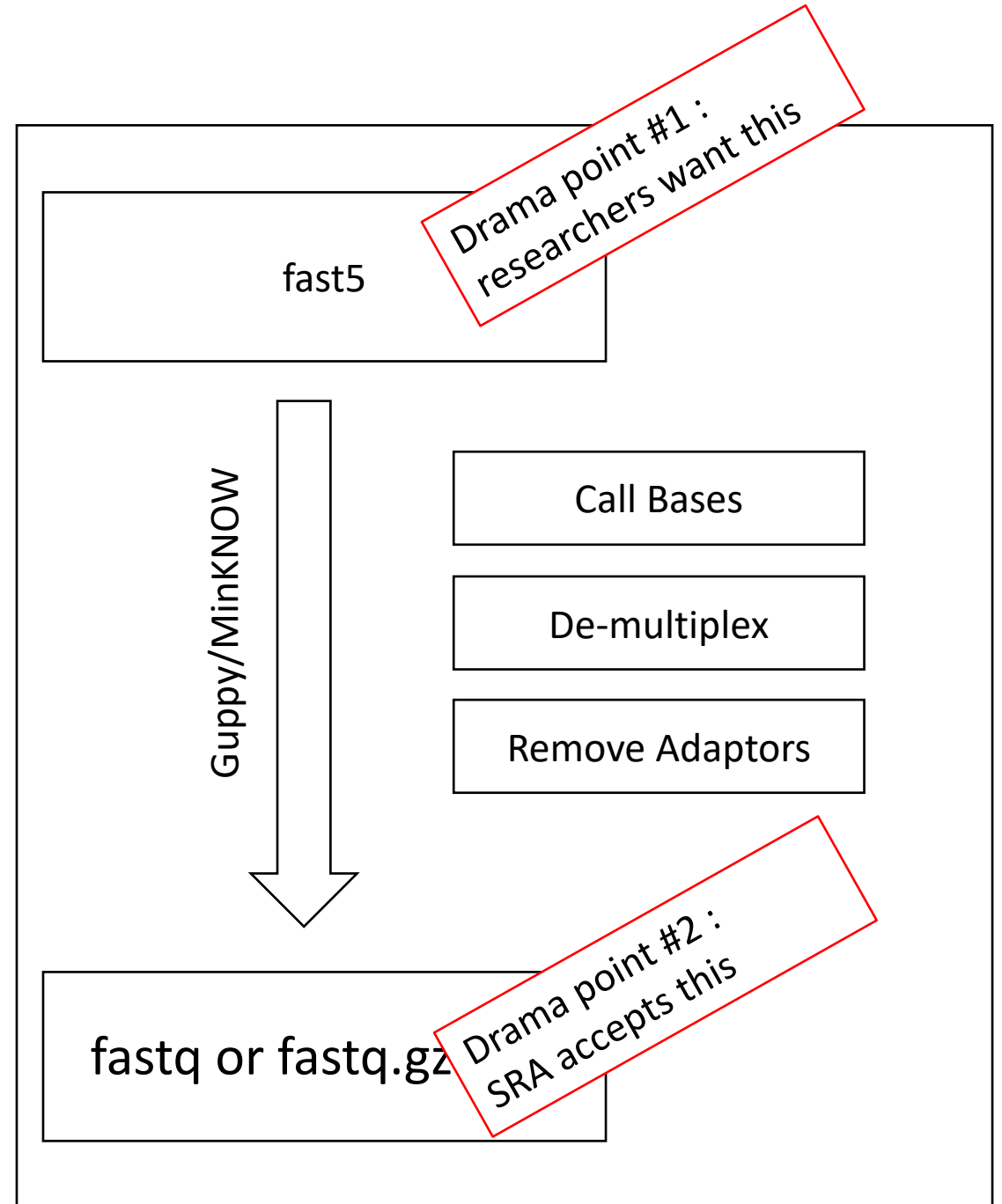
**Figure 3**
A schematic highlighting the advantages of long reads in *de novo* assembly of repetitive regions. Long read lengths are more likely to incorporate the whole repetitive region (shown in red) allowing more accurate assembly with fewer gaps. Image adapted from Sam Demharter[7].
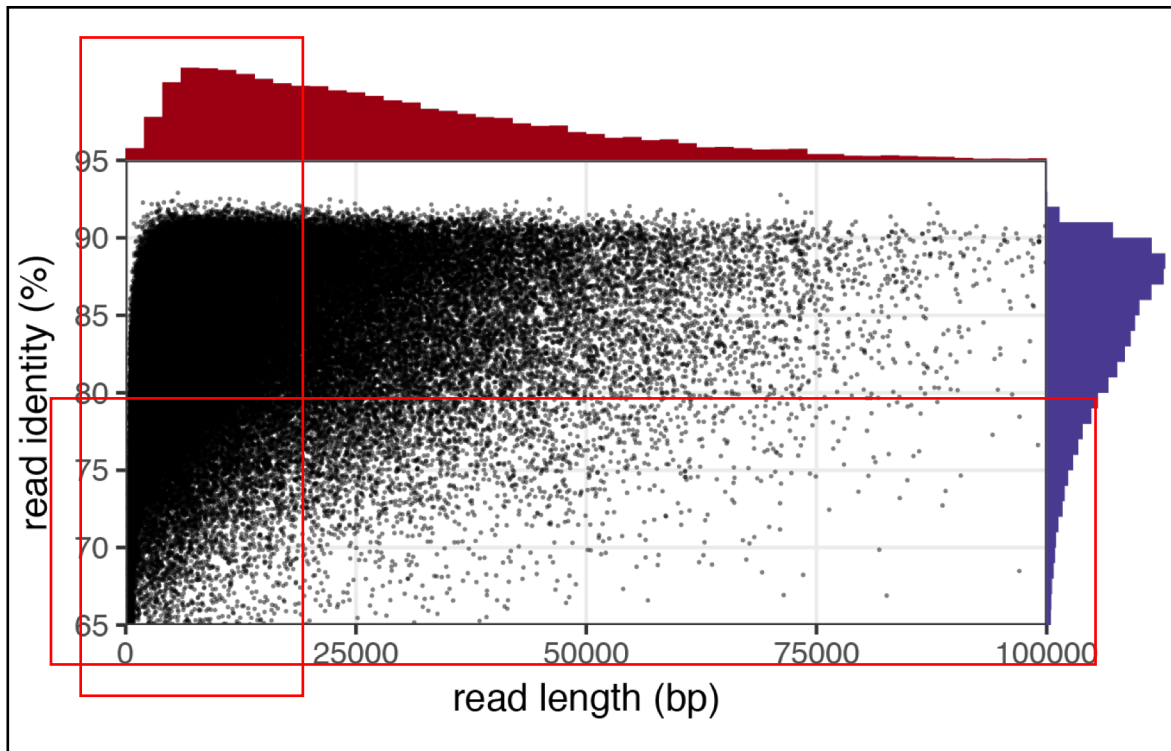
# Oxford Nanopore WGS Sequencing



Culture Isolates
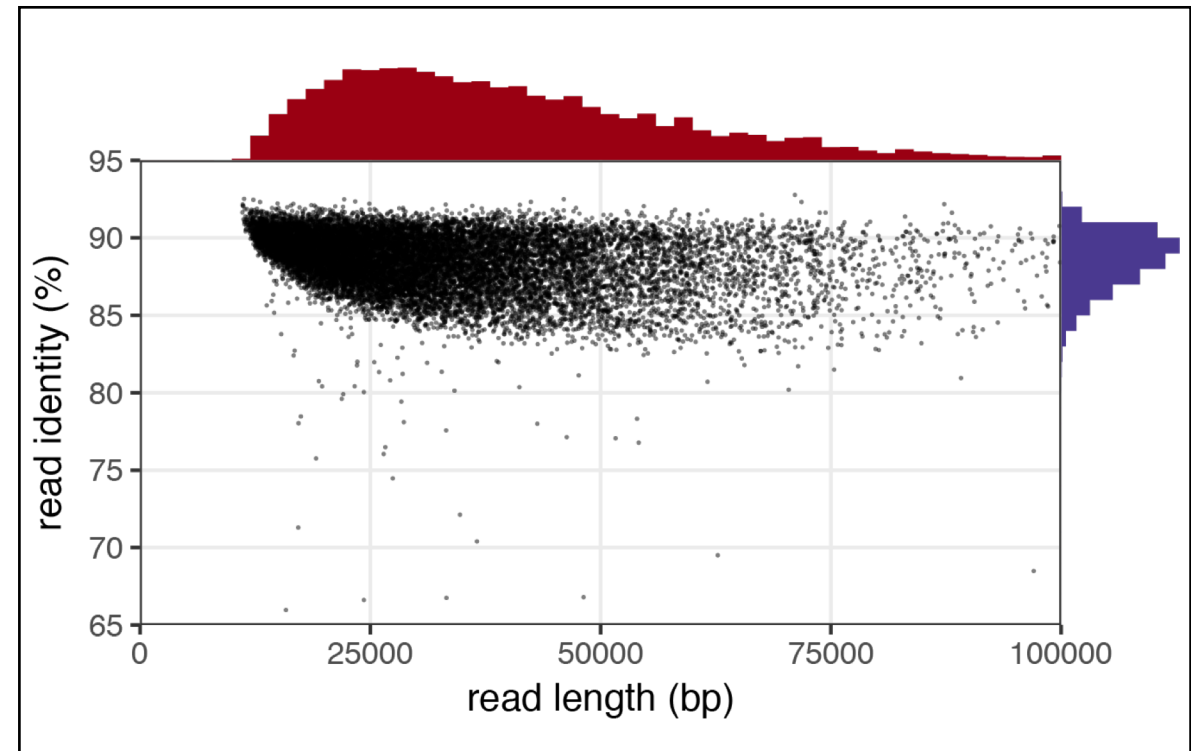
Isolate DNA

Prepare Libraries

Sequence

Reads

*De novo* alignment

# Nanopore Output



fast5

Drama point #1 : researchers want this

Guppy/MinKNOW

Call Bases

De-multiplex

Remove Adaptors

fastq or fastq.gz

Drama point #2 : SRA accepts this

https://en.wikipedia.org/wiki/Nanopore_sequencing

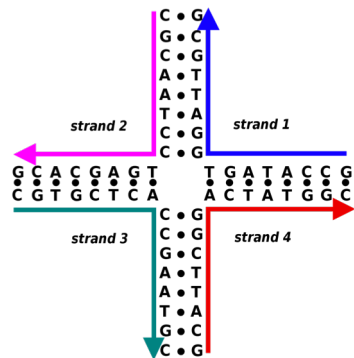# Filtlong

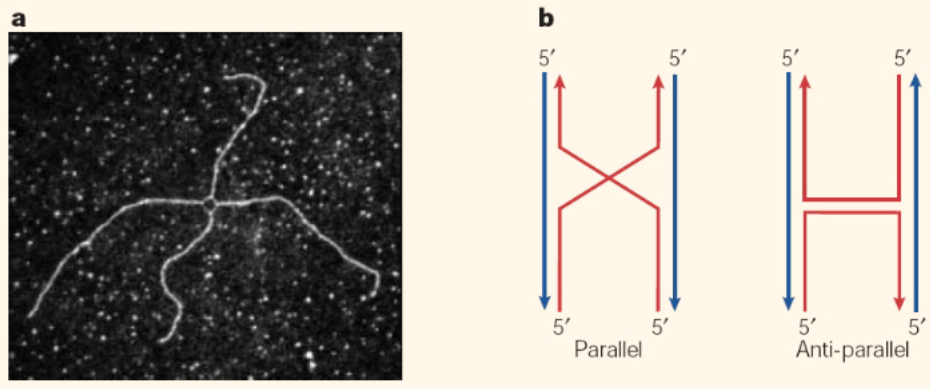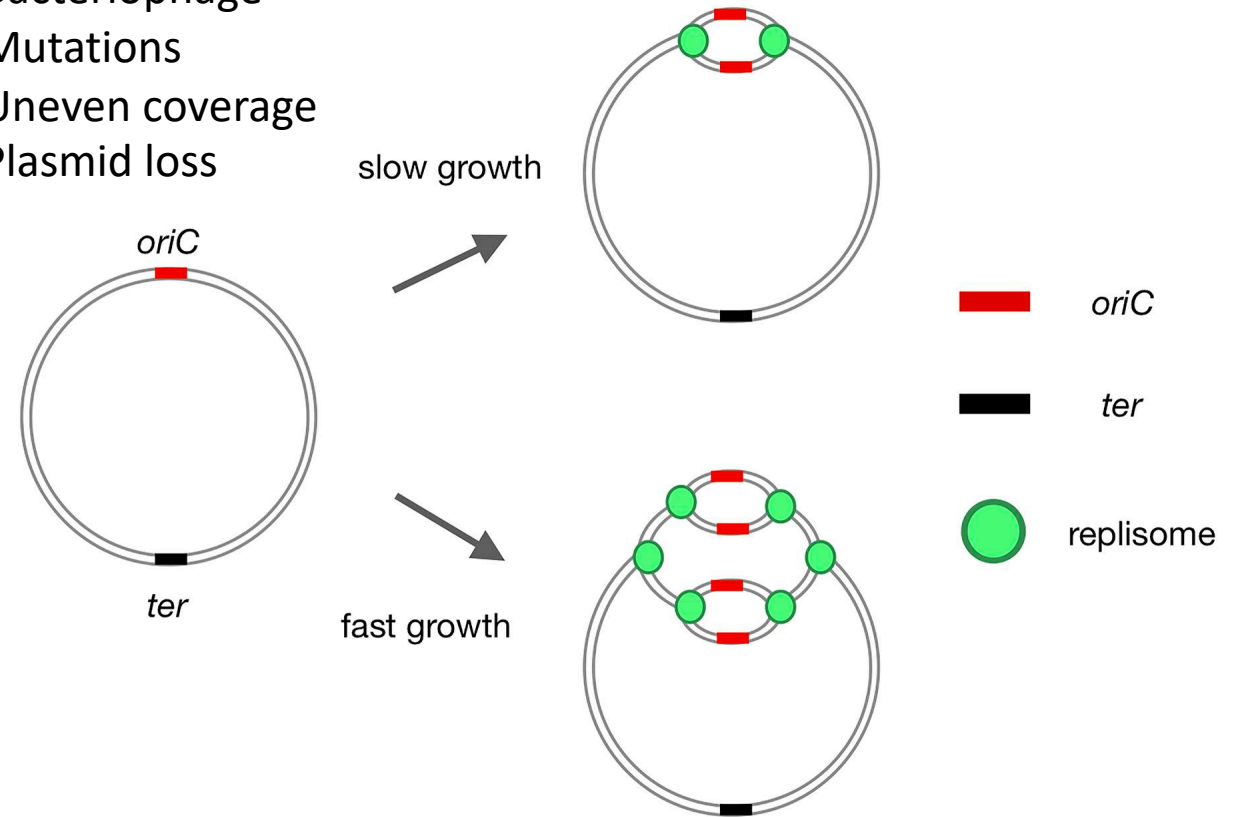You don't want all the reads (trust me)



Before

After

# Filtlong cannot filter out all issues

## Holliday Junction artefacts



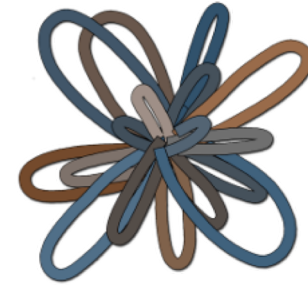## Evolution does not stop in the lab
- Bacteriophage
- Mutations
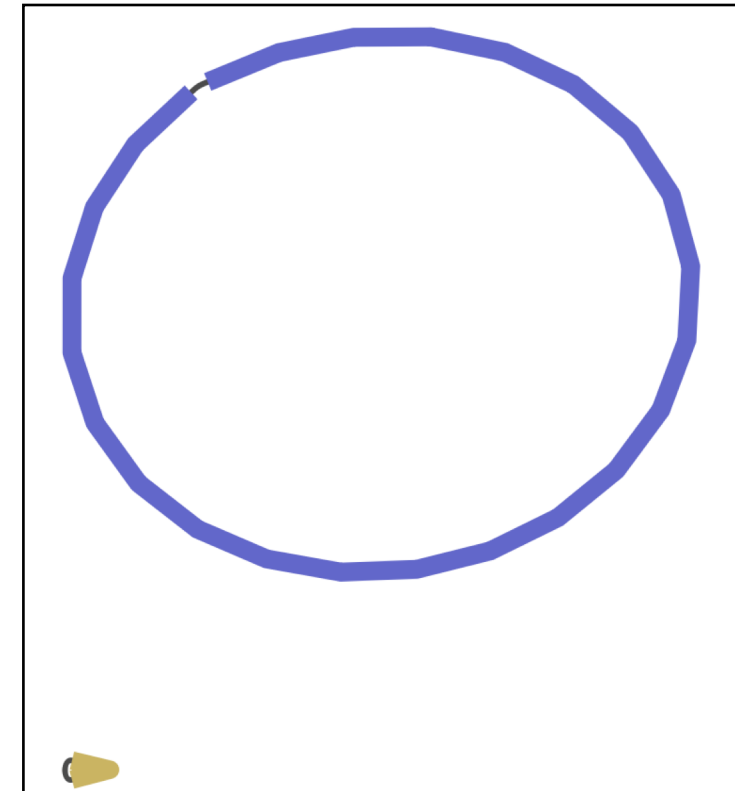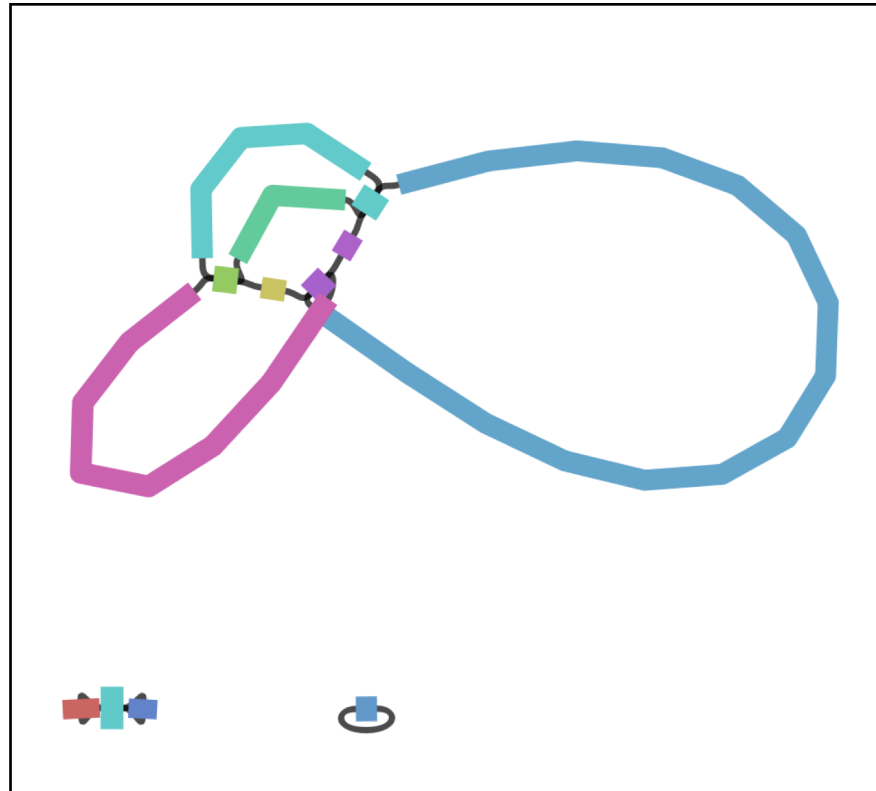- Uneven coverage
- Plasmid loss

# *De novo* long read assembly

- Flye
- Miniasm/minipolish
- Raven
- Canu/Canu2
- RedBean
- Unicycler (hybrid)
- And more!
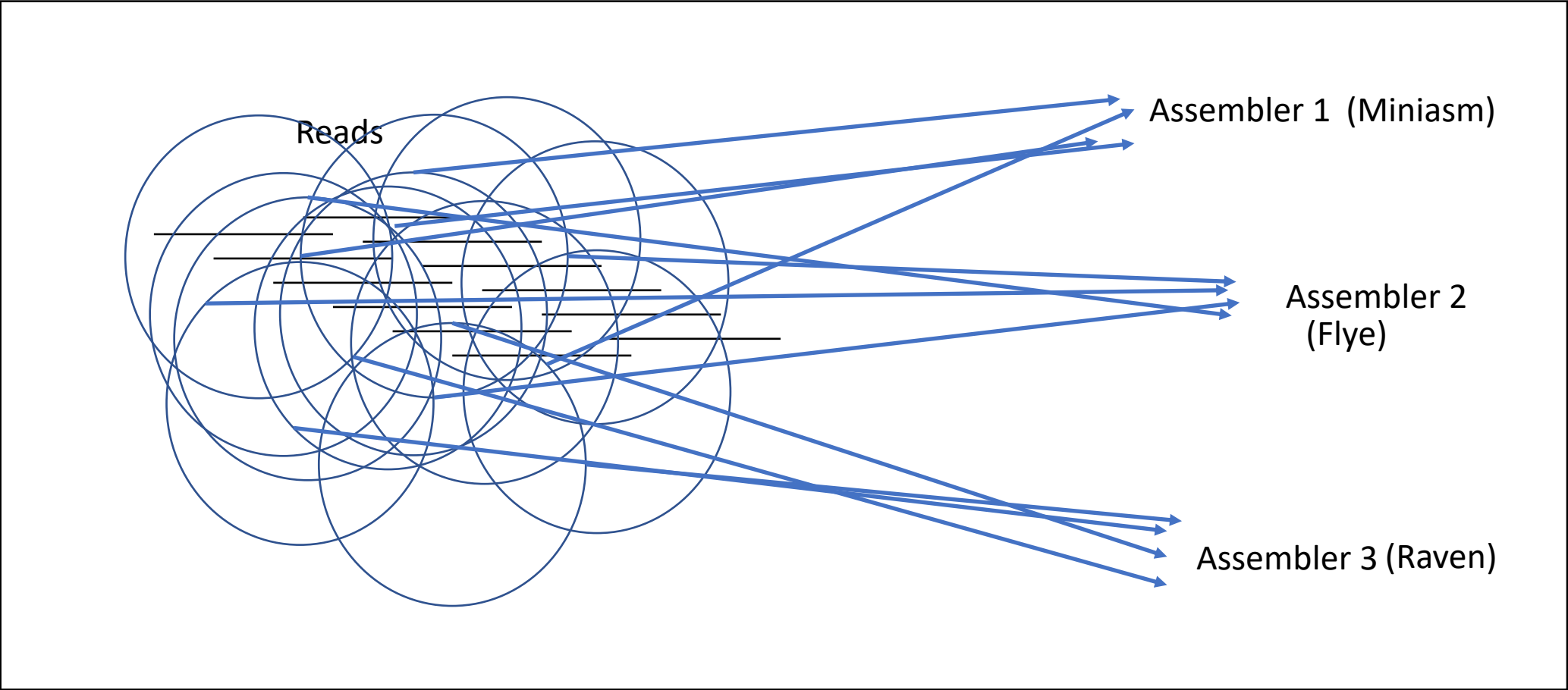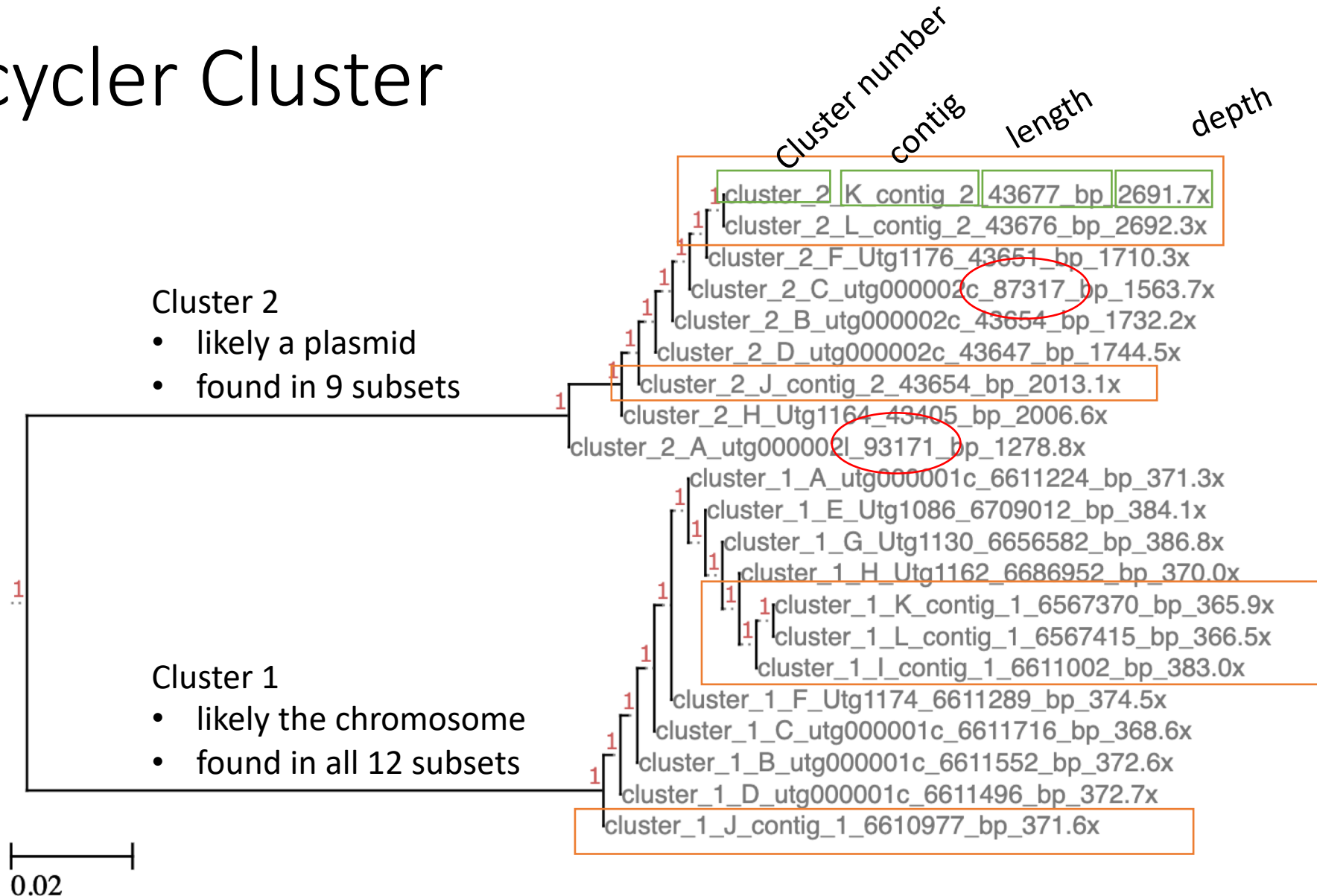
But ... which one is BEST?

# Trycycler: consensus by agreement

- Create multiple assemblies

- Resolve the differences between them

# Trycycler : SubSample & Assemble

# Trycycler Cluster



Cluster 2
- likely a plasmid
- found in 9 subsets

Cluster 1
- likely the chromosome
- found in all 12 subsets

In theory, all contigs in a cluster will have similar lengths and depth

# Trycycler reconcile

- All contigs in a cluster should have
  - Similar depth
  - Similar length
  - Similar sequence
- The end user must remove contigs that are not similar "enough"
- Reconcile
  - Ensure sequences on the same strand
  - Fix circularization
  - Rotate to common start

# Trycycler MSA aligns sequences in a cluster

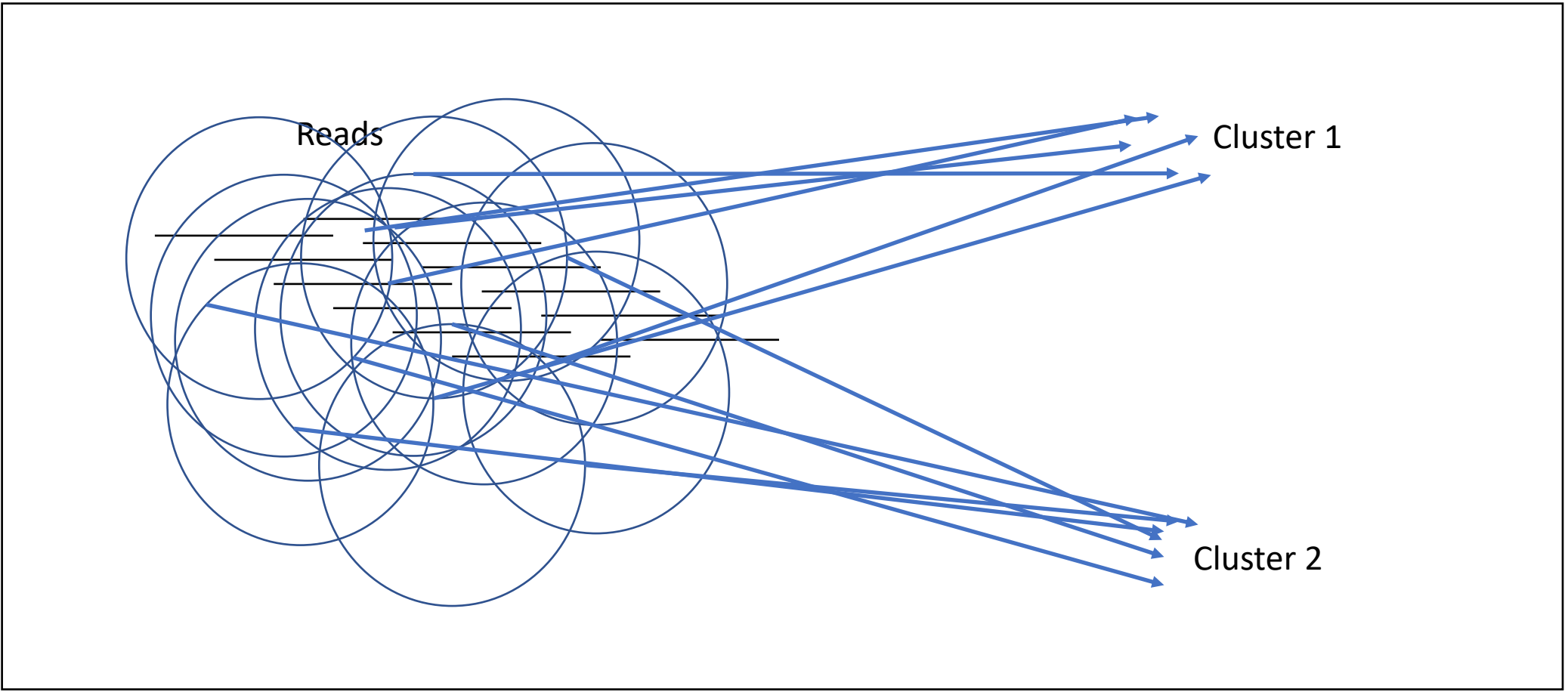For example, it would take sequences like this:

```
GGCAGAGCGACGTAAATTACGAGTAAAGGAGGGGAGAGCATTAAGCATGCCTAAACTG
GGCAGAGCGCGACGTAAATTACGAGTAAAAGGAGGGAGGAGCATTAAGCCATGCCTACTG
GGCAGAGCGCGACTAAATTTACGAGTAAAGGAGGGAGGAGCATAGCCATGCCTAAACTG
```

And produce an alignment like this:

```
GGCAGAG––CGACGTAAA–TTACGAGT–AAAGGAGGGGA–GAGCATTAAG–CATGCCTAAACTG
GGCAGAGCGCGACGTAAA–TTACGAGTAAAAGGA–GGGAGGAGCATTAAGCCATGCCT––ACTG
GGCAGAGCGCGAC–TAAATTTACGAGT–AAAGGA–GGGAGGAGCAT––AGCCATGCCTAAACTG
```

# Trycycler partition assigns reads to cluster

# Trycycler consensus



The Simpsons (1989)

For example, it would take sequences like this:

```
GGCAGAGCGACGTAAATTACGAGTAAAGGAGGGGAGAGCATTAAGCATGCCTAAACTG
GGCAGAGCGCGACGTAAATTACGAGTAAAAGGAGGGAGGAGCATTAAGCCATGCCTACTG
GGCAGAGCGCGACTAAATTTACGAGTAAAGGAGGGAGGAGCATAGCCATGCCTAAACTG
```

And produce an alignment like this:

```
GGCAGAG––CGACGTAAA–TTACGAGT–AAAGGAGGGGA–GAGCATTAAG–CATGCCTAAACTG
GGCAGAGCGCGACGTAAA–TTACGAGTAAAAGGA–GGGAGGAGCATTAAGCCATGCCT––ACTG
GGCAGAGCGCGAC–TAAATTTACGAGT–AAAGGA–GGGAGGAGCAT––AGCCATGCCTAAACTG
```

# Polishing : Because we are not done, yet

- Polishing is using prior reads to "correct" errors in the final assembly
  - Nanopolish : polishes raw ONT reads
  - Medaka : polishes assembly with ONT reads
  - Racon : polishes assembly with Illumina or ONT reads
  - Pilon : polishes assembly with Illumina reads
- Many assemblers include a polishing step
- Over-polishing is a thing

# Donut falls : A Trycycler Nextflow Workflow

https://github.com/UPHL-BioNGS/Donut_Falls



Once guppy has called bases, removed adapters, and demultiplexed
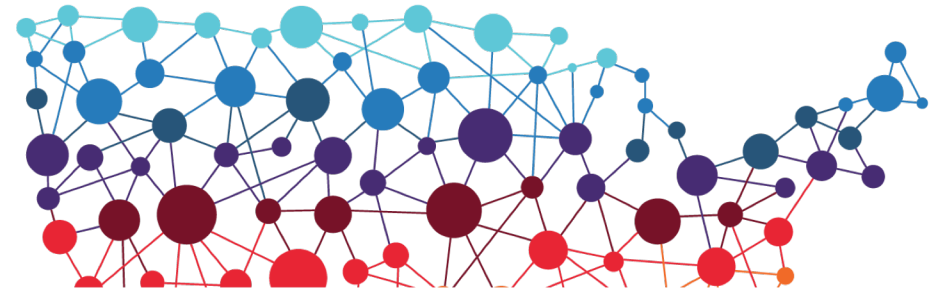- Create a sample key that links barcode and sample_id and Illumina fastq files
- Run phase 1 :

`nextflow run Donut_Falls.nf -c configs/singularity.config`

- Examine tree (at http://etetoolkit.org/treeview/)
- Remove problematic clusters
- Run Phase 2:

`nextflow run Donut_Falls.nf -c configs/phase2_singularity.config`

- Examine clusters with Bandage
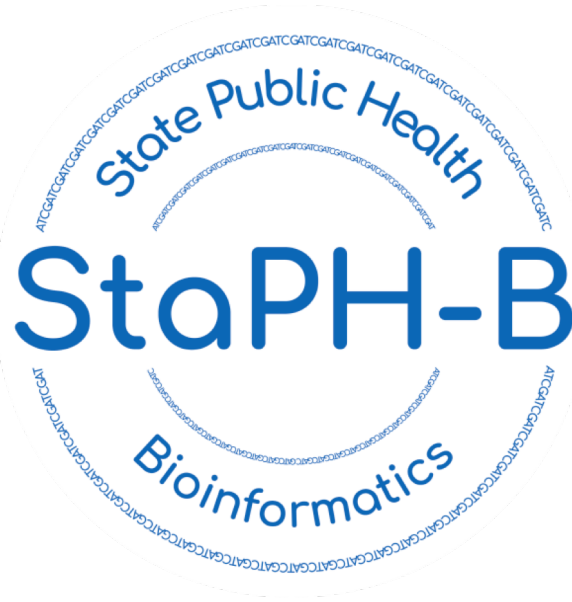- Find AMR genes, submit to repositories, etc.

Questions?

Erin Young, PhD

🐦 : https://twitter.com/ErinYoun
🐙 : https://github.com/erinyoung