

MICI3119 - Genomics I

Finlay Maguire
finlay.maguire@dal.ca
maguire-lab.github.io

Overview

- Hospital-based outbreak prevention and control
- DNA Sequencing Technologies
- Reference-based assembly
- De novo assembly
- Inferring phylogenies from genomes
- Interpreting phylogenies
- Use genomic data to respond to stop an outbreak

IPAC Investigations: UCSF NICU MRSA Outbreak



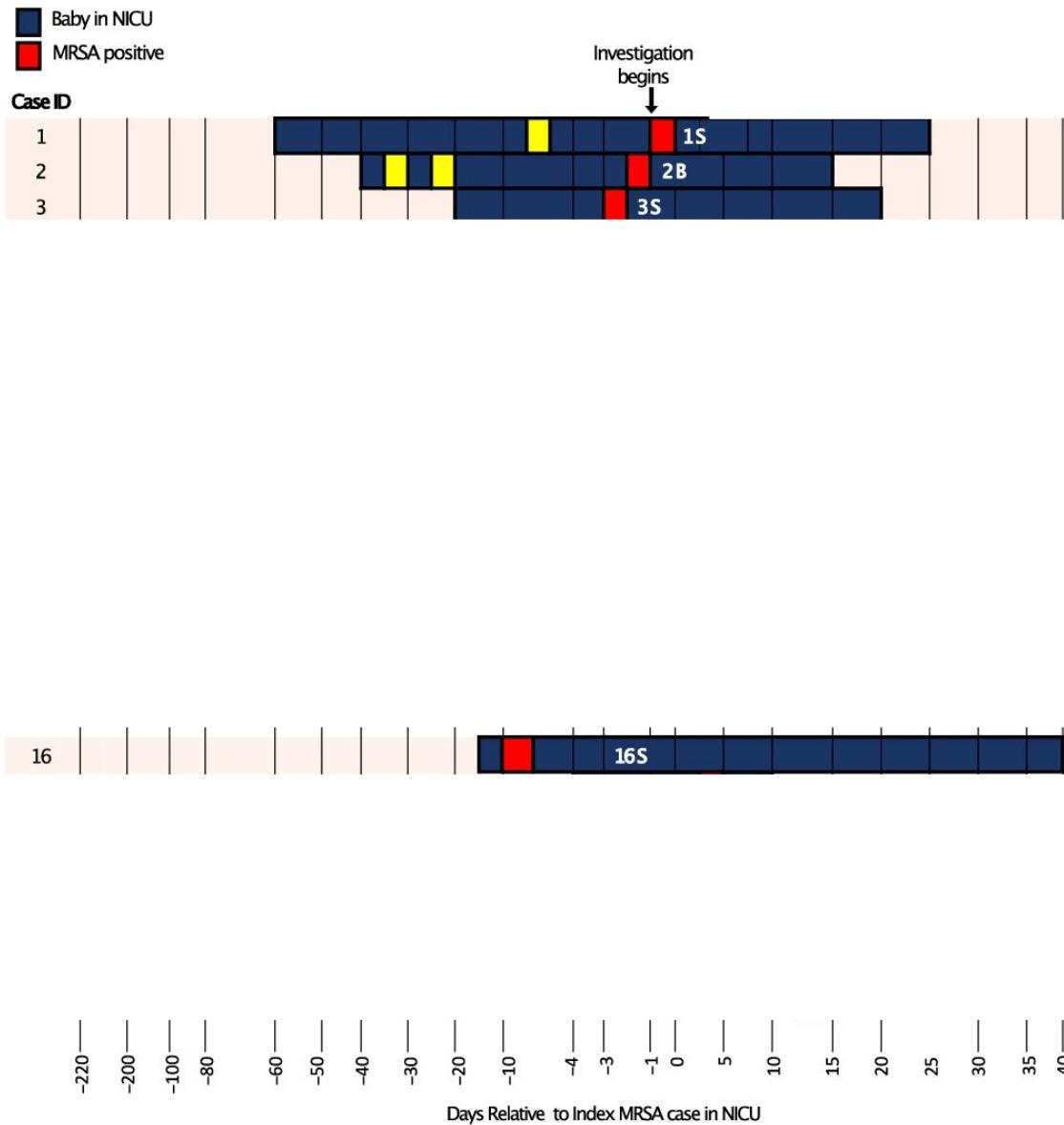
Setting: UCSF NICU (58-bed) in UCSF Benioff Children's Hospital (183-bed)

Neonatal Intensive Care Unit

<https://healthtian.com/wp-content/uploads/2020/09/neonatal-intensive-care-unit1.jpg>

Madera, Sharline, et al. "Prolonged silent carriage, genomic virulence potential and transmission between staff and patients characterize a neonatal intensive care unit (NICU) outbreak of methicillin-resistant *Staphylococcus aureus* (MRSA)." *Infection Control & Hospital Epidemiology* 44.1 (2023): 40-46.

IPAC Investigations: UCSF NICU MRSA Outbreak



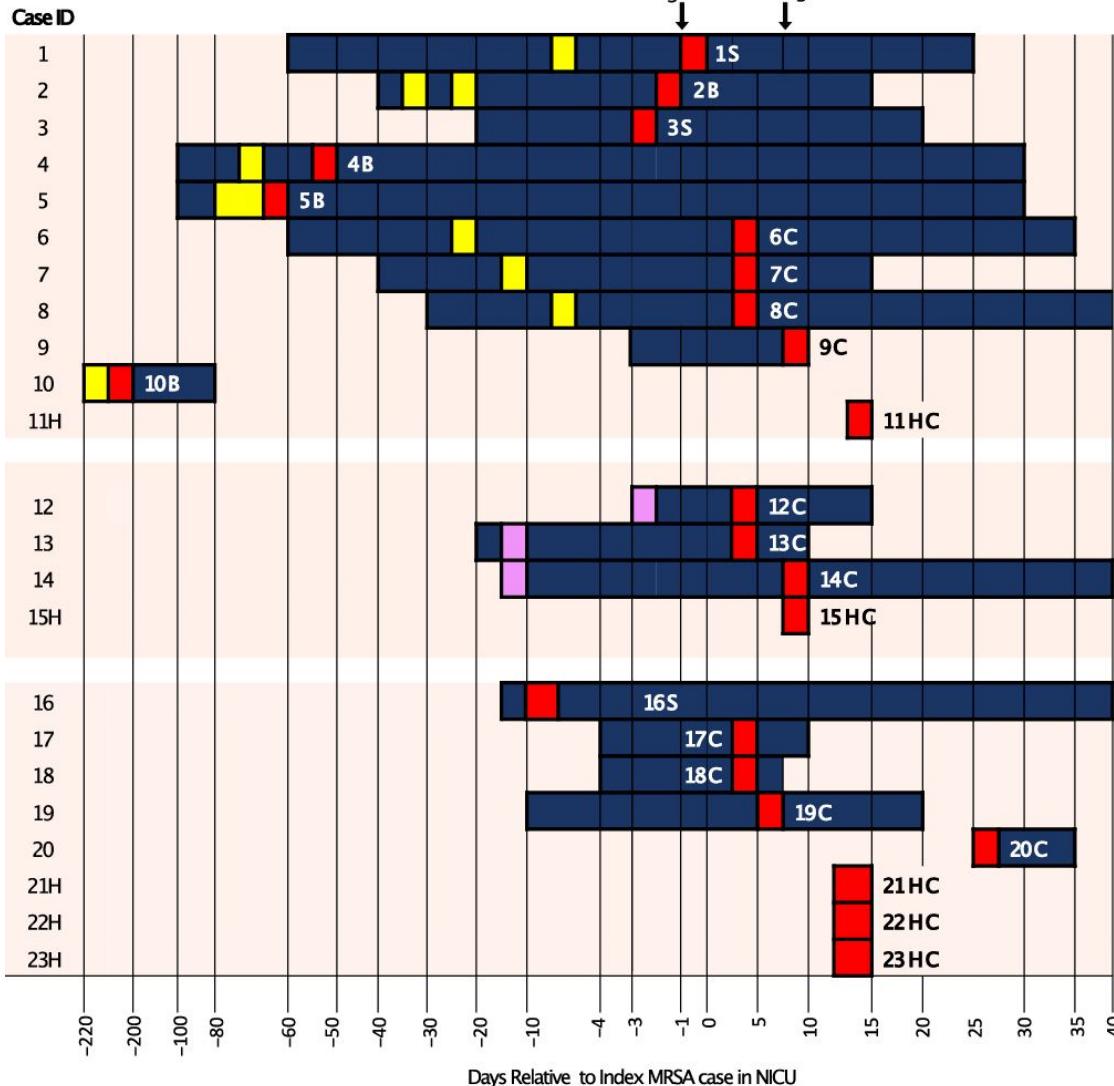
Setting: UCSF NICU (58-bed) in UCSF Benioff Children's Hospital (183-bed)

Investigation Trigger: 4 NICU infants with invasive MRSA within 8 days

Madera, Sharline, et al. "Prolonged silent carriage, genomic virulence potential and transmission between staff and patients characterize a neonatal intensive care unit (NICU) outbreak of methicillin-resistant *Staphylococcus aureus* (MRSA)." *Infection Control & Hospital Epidemiology* 44.1 (2023): 40-46.

IPAC Investigations: UCSF NICU MRSA Outbreak

 Baby in NICU
 MRSA positive



Setting: UCSF NICU (58-bed) in UCSF Benioff Children's Hospital (183-bed)

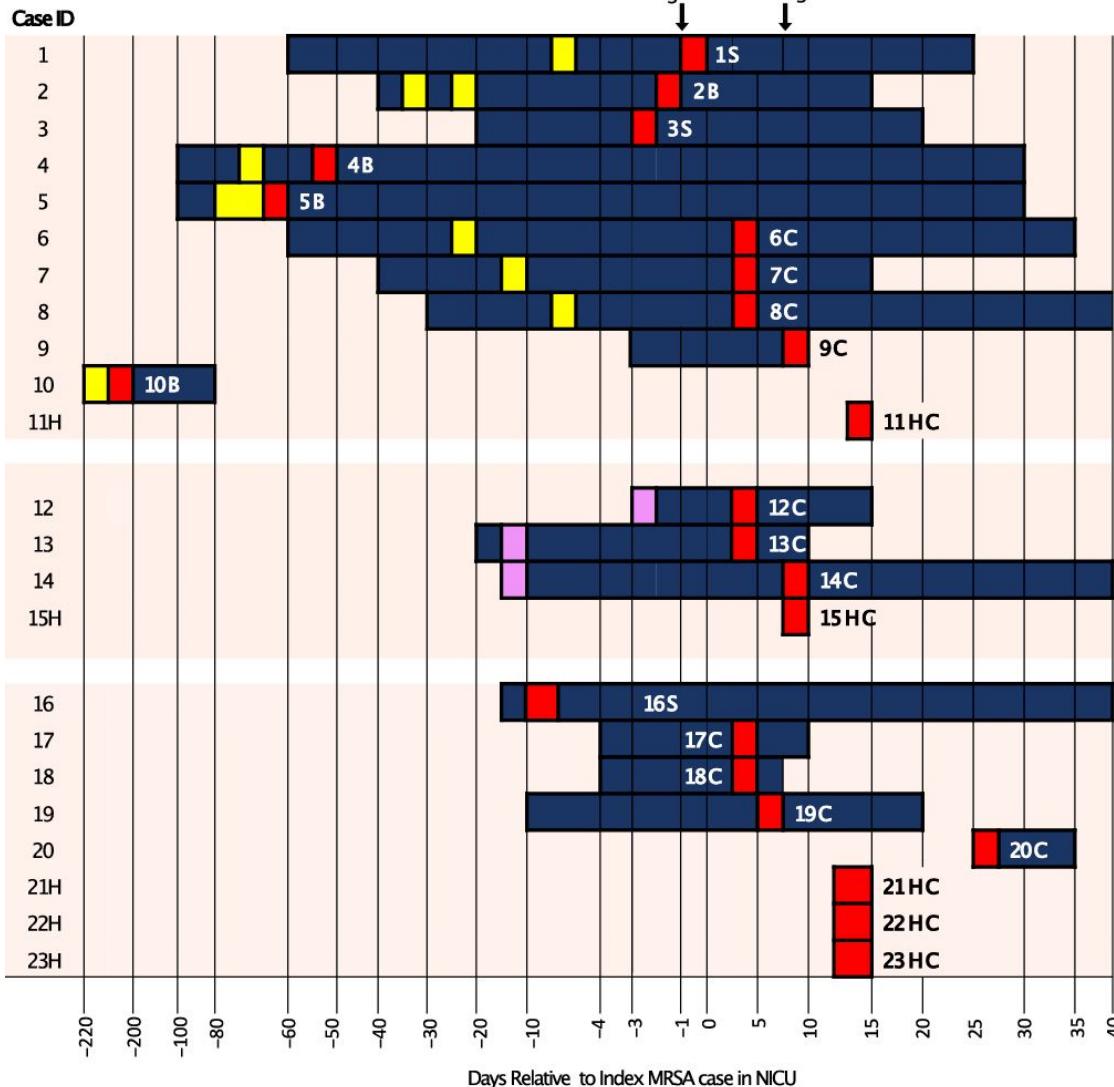
Investigation Trigger: 4 NICU infants with invasive MRSA within 8 days

Next Step: Look at recent past MRSA cases, contact trace, & expand surveillance swabbing of NICU + HCPs

Madera, Sharline, et al. "Prolonged silent carriage, genomic virulence potential and transmission between staff and patients characterize a neonatal intensive care unit (NICU) outbreak of methicillin-resistant *Staphylococcus aureus* (MRSA)." *Infection Control & Hospital Epidemiology* 44.1 (2023): 40-46.

IPAC Investigations: UCSF NICU MRSA Outbreak

Baby in NICU
MRSA positive



Setting: UCSF NICU (58-bed) in UCSF Benioff Children's Hospital (183-bed)

Investigation Trigger: 4 NICU infants with invasive MRSA within 8 days

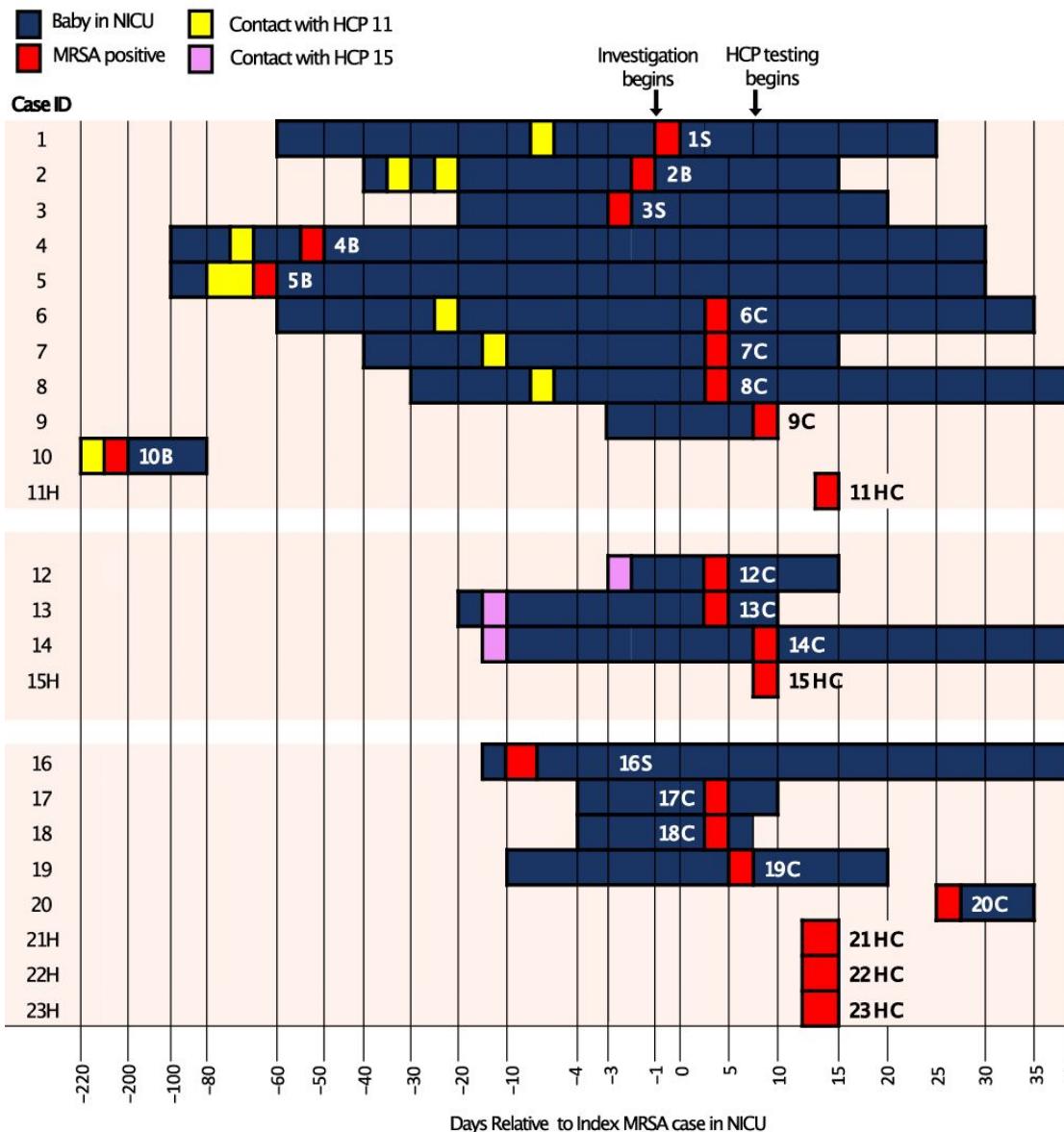
Next Step: Look at recent past MRSA cases, contact trace, & expand surveillance swabbing of NICU + HCPs

Which cases are linked?

Where did these come from?

Madera, Sharline, et al. "Prolonged silent carriage, genomic virulence potential and transmission between staff and patients characterize a neonatal intensive care unit (NICU) outbreak of methicillin-resistant *Staphylococcus aureus* (MRSA)." *Infection Control & Hospital Epidemiology* 44.1 (2023): 40-46.

IPAC Investigations: UCSF NICU MRSA Outbreak



Setting: UCSF NICU (58-bed) in UCSF Benioff Children's Hospital (183-bed)

Investigation Trigger: 4 NICU infants with invasive MRSA within 8 days

Next Step: Look at recent past MRSA cases, contact trace, & expand surveillance swabbing of NICU + HCPs

Which cases are linked?

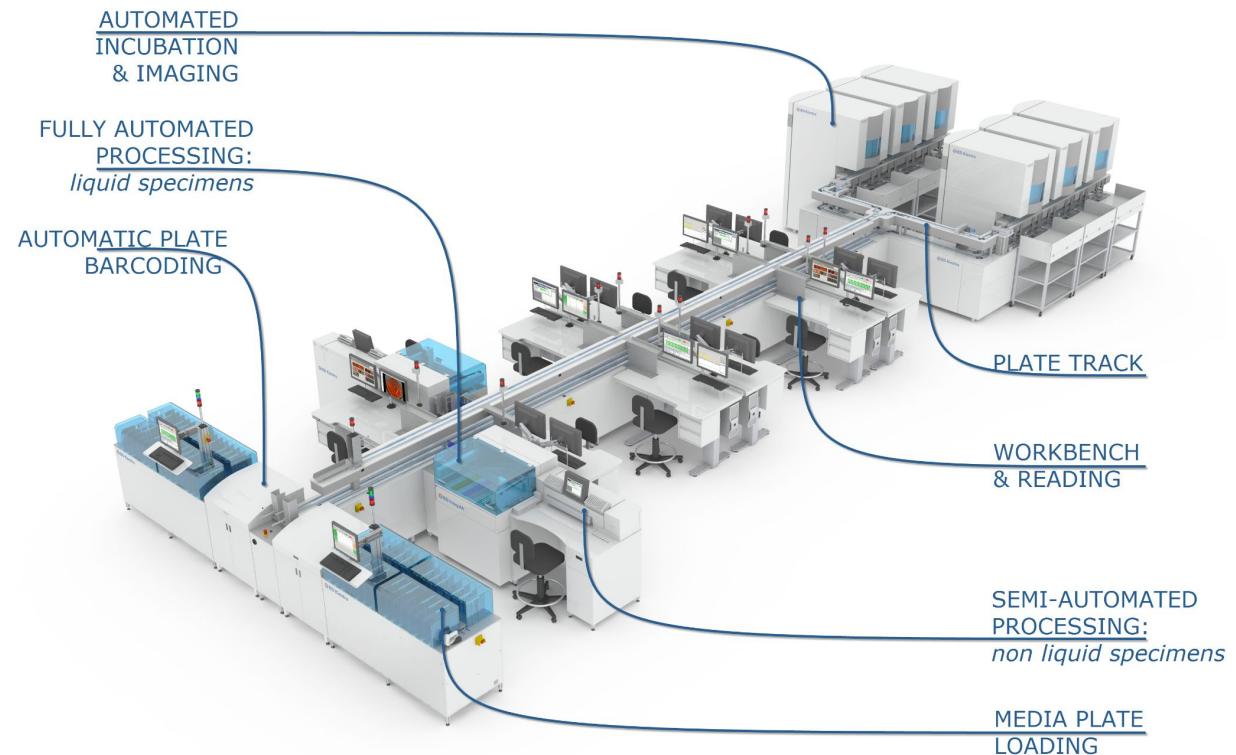
Where did these come from?

Contact tracing gives ideas but need more information => Genomes

Madera, Sharline, et al. "Prolonged silent carriage, genomic virulence potential and transmission between staff and patients characterize a neonatal intensive care unit (NICU) outbreak of methicillin-resistant *Staphylococcus aureus* (MRSA)." *Infection Control & Hospital Epidemiology* 44.1 (2023): 40-46.

So, we have lots of swabs, what now?

Expensively Automated sample processing and culturing

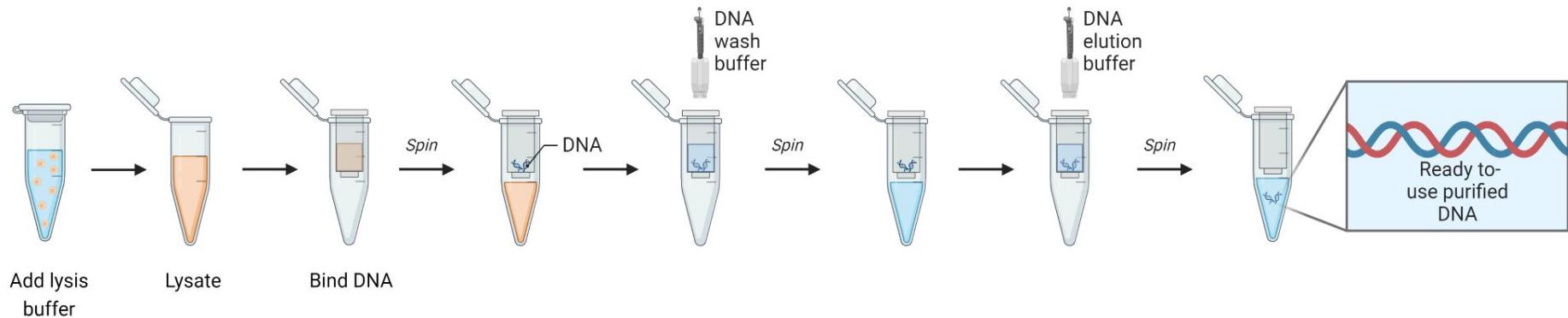


Expensively Automated sample processing and culturing

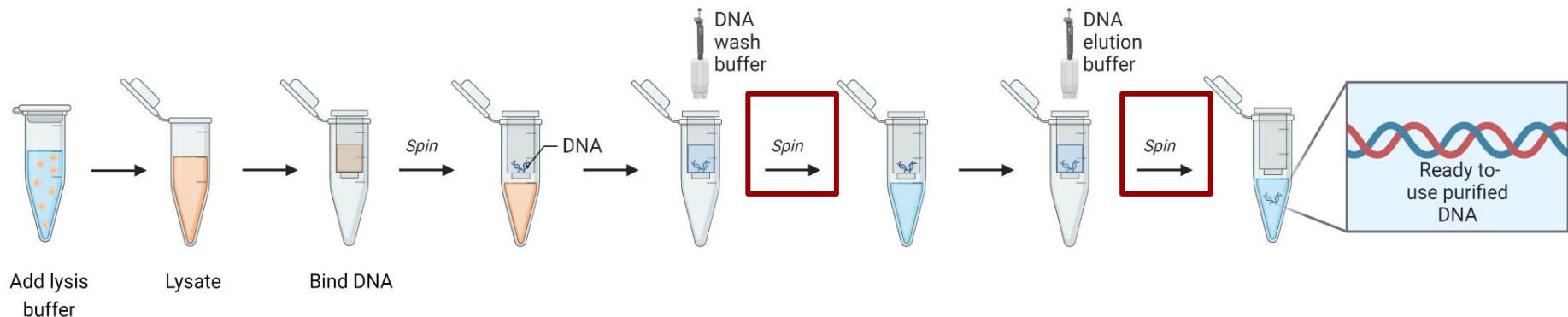


Christy Vermeiren

Automated DNA extraction with magnets + robots



Automated DNA extraction with magnets + robots



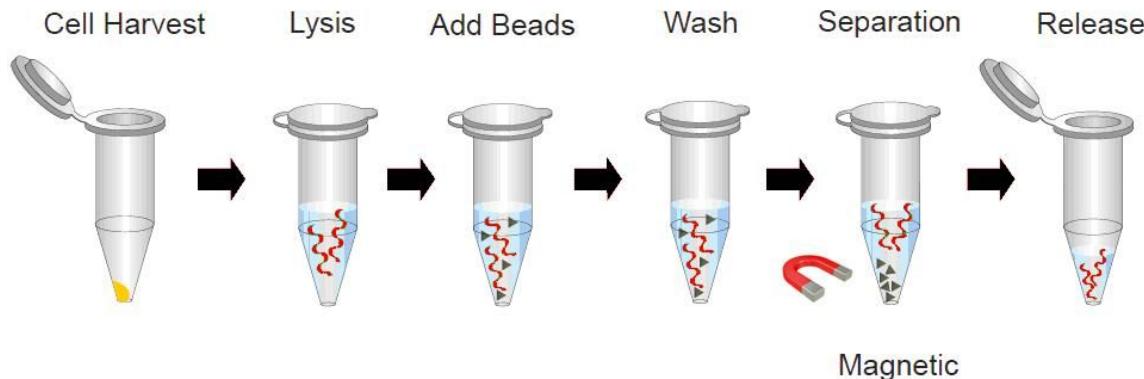
https://craft-robotics.s3.amazonaws.com/_thumbnail/Multi-Colored-Pipetting.jpg?mtime=20180919131431

<https://the-dna-universe.com/wp-content/uploads/2021/11/Hamilton-liquid-handling-robot.png>



<https://images.aatbio.com/universal/newsletter/volume-11-1/Silica%20Gen%20DNA%20Extraction%20Workflow.png>

Automated DNA extraction with magnets + robots



https://craft-robotics.s3.amazonaws.com/_thumbnail/Multi-Colored-Pipetting.jpg?mtime=20180919131431

<https://the-dna-universe.com/wp-content/uploads/2021/11/Hamilton-liquid-handling-robot.png>



<https://www.epruibiotech.com/wp-content/uploads/2021/03/Magnetic-beads-DNA-extraction-process.jpg>

Got DNA now, how do we work out what it says?

Sequencing Technology

~1972-1977

First generation

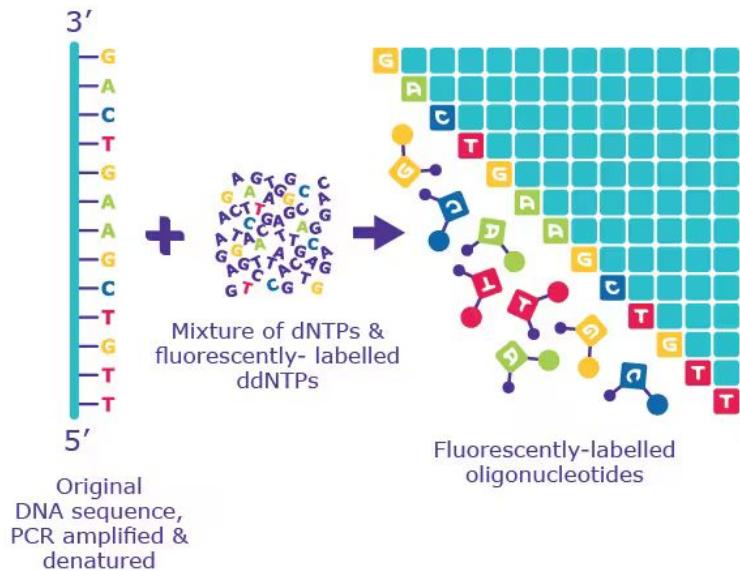


Sanger sequencing
Maxam and Gilbert
Sanger chain termination

Sanger Sequencing

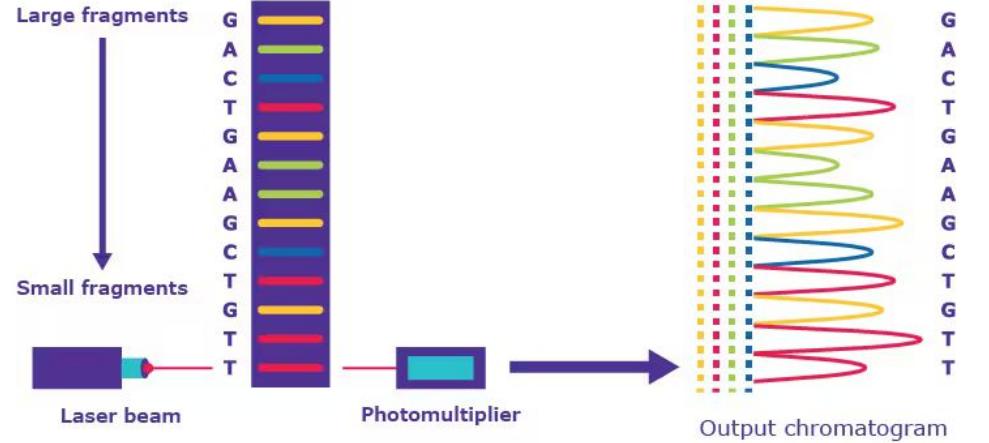
1

PCR with fluorescent, chain-terminating ddNTPs



2

Size separation by capillary gel electrophoresis



Sequencing Technology

~1972-1977

First generation



Sanger sequencing
Maxam and Gilbert
Sanger chain termination

Infer nucleotide identity using dNTPs,
then visualize with electrophoresis

500–1,000 bp fragments

Sequencing Technology

~1972-1977

~2001-2004

First generation

Second generation
(next generation sequencing)



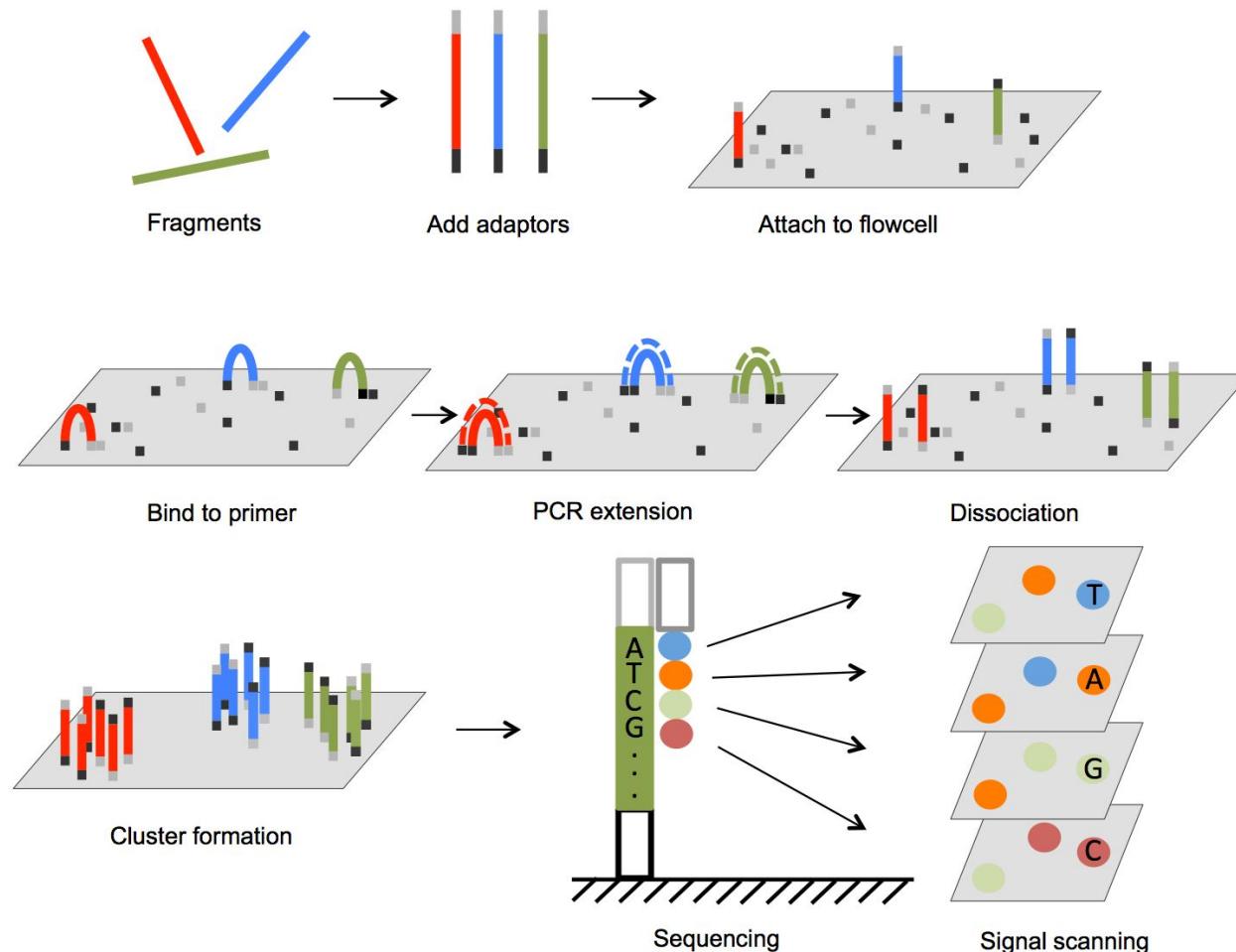
Sanger sequencing
Maxam and Gilbert
Sanger chain termination

Infer nucleotide identity using dNTPs,
then visualize with electrophoresis

500–1,000 bp fragments

454, Solexa,
Ion Torrent,
Illumina

Sequencing by Synthesis



Sequencing Technology

~1972-1977

~2001-2004

First generation

Second generation
(next generation sequencing)



Sanger sequencing
Maxam and Gilbert
Sanger chain termination

Infer nucleotide identity using dNTPs,
then visualize with electrophoresis

500–1,000 bp fragments

454, Solexa,
Ion Torrent,
Illumina

High throughput from the
parallelization of sequencing reactions

~50–500 bp fragments

Sequencing Technology

~1972-1977

~2001-2004

~2011-2015

First generation

Second generation
(next generation sequencing)

Third generation



Sanger sequencing
Maxam and Gilbert
Sanger chain termination

Infer nucleotide identity using dNTPs,
then visualize with electrophoresis

500–1,000 bp fragments

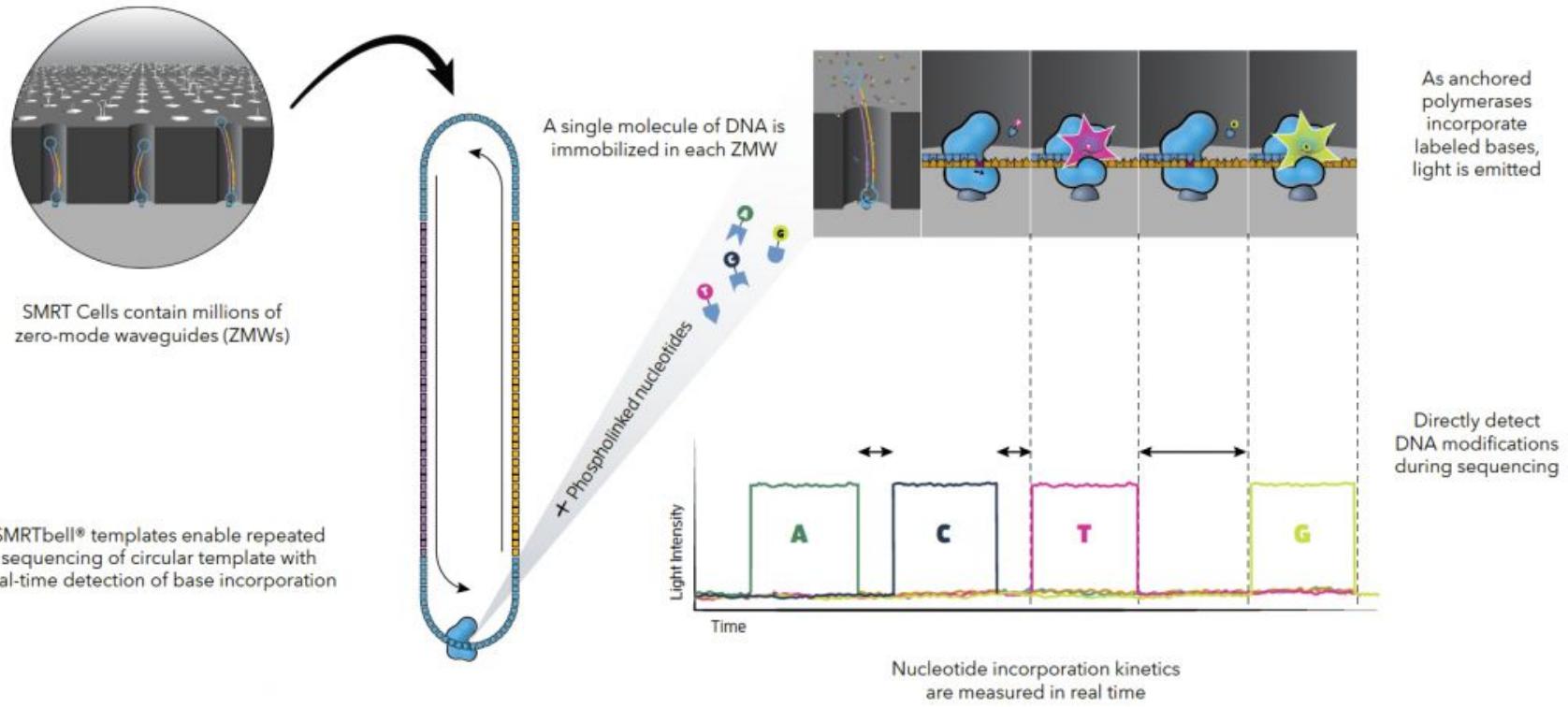
454, Solexa,
Ion Torrent,
Illumina

High throughput from the
parallelization of sequencing reactions

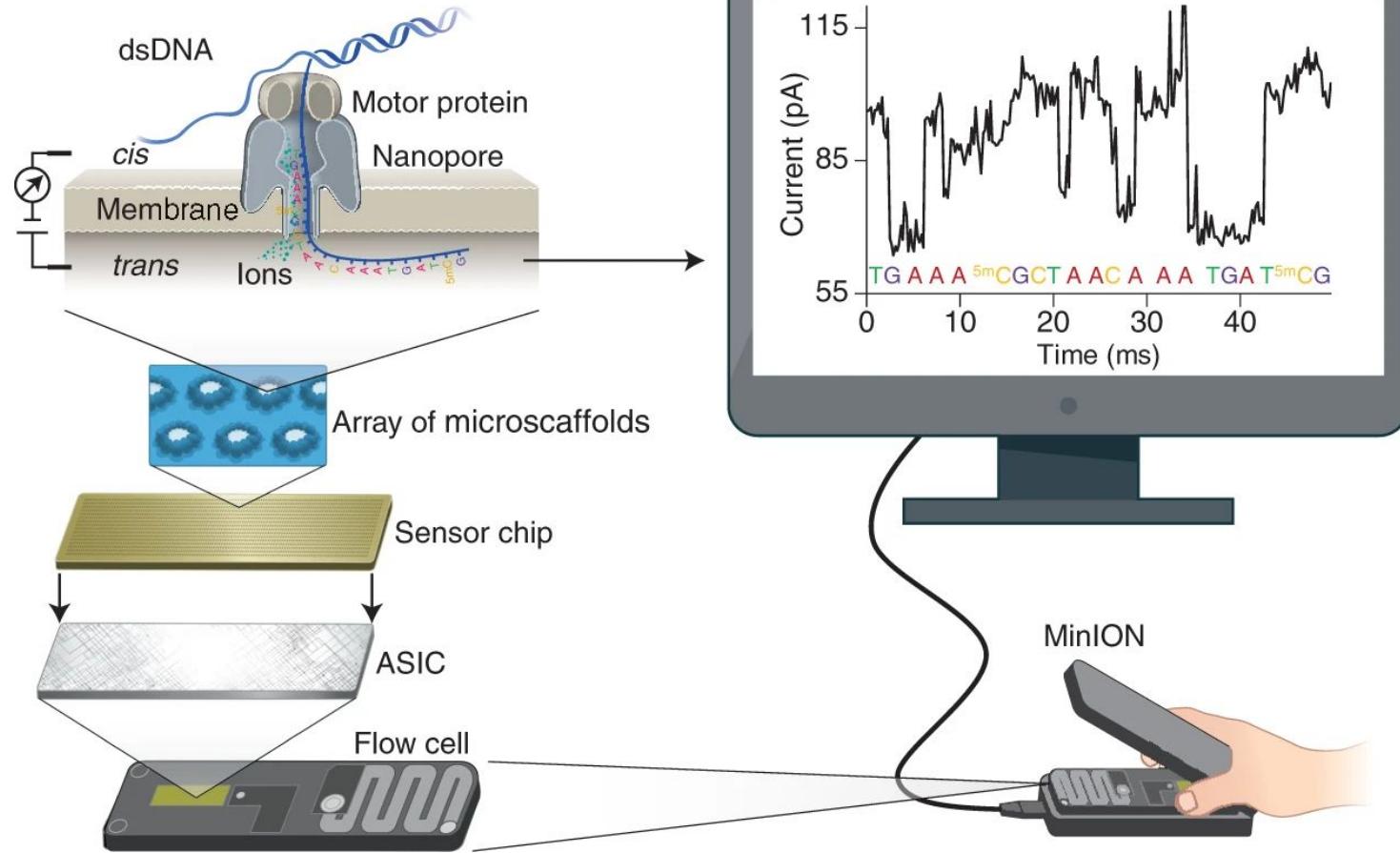
~50–500 bp fragments

PacBio
Oxford Nanopore

PacBio Sequencing



Nanopore Sequencing



Sequencing Technology

~1972-1977

~2001-2004

~2011-2015

First generation



Sanger sequencing
Maxam and Gilbert
Sanger chain termination

Infer nucleotide identity using dNTPs,
then visualize with electrophoresis

500–1,000 bp fragments

Second generation (next generation sequencing)



454, Solexa,
Ion Torrent,
Illumina

High throughput from the
parallelization of sequencing reactions

~50–500 bp fragments

Third generation



PacBio
Oxford Nanopore

Sequence native DNA in real time
with single-molecule resolution

Tens of kb fragments, on average

Sequencing Technology

~1972-1977

~2001-2004

~2011-2015

First generation



Sanger sequencing
Maxam and Gilbert
Sanger chain termination

Infer nucleotide identity using dNTPs,
then visualize with electrophoresis

500–1,000 bp fragments

Second generation (next generation sequencing)



454, Solexa,
Ion Torrent,
Illumina

High throughput from the
parallelization of sequencing reactions

~50–500 bp fragments



PacBio
Oxford Nanopore

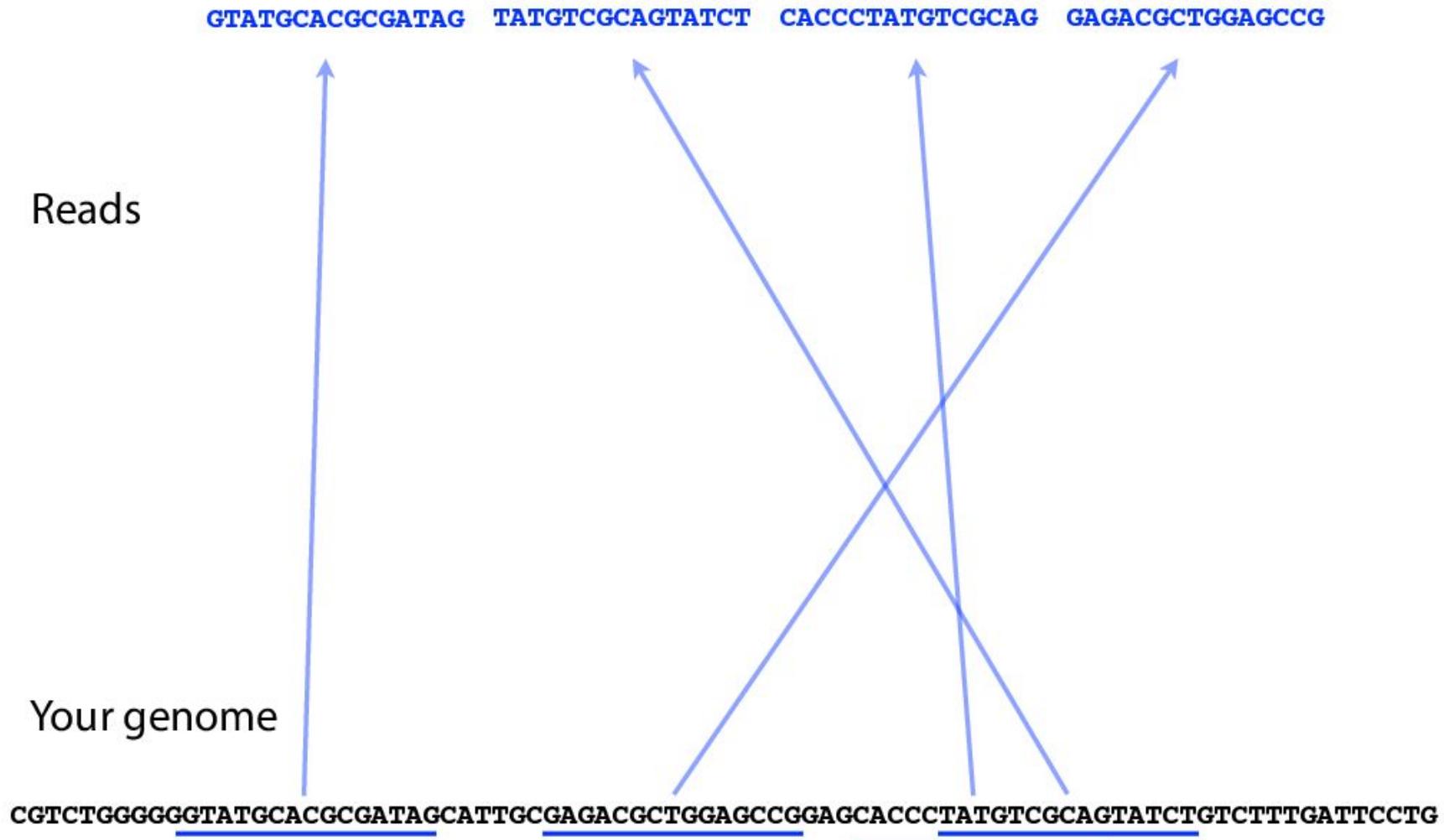
Sequence native DNA in real time
with single-molecule resolution

Tens of kb fragments, on average

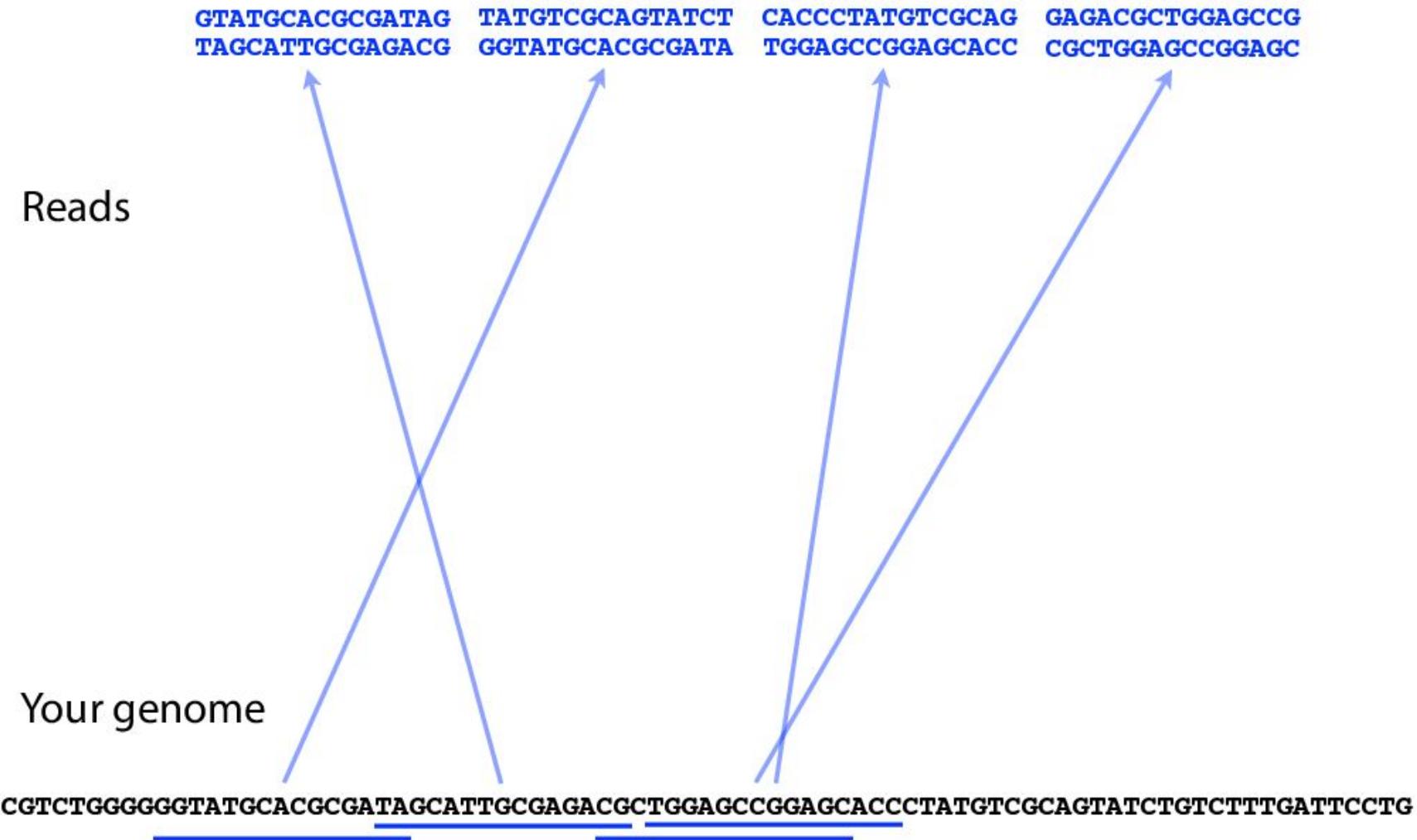
Short-read sequencing

Long-read sequencing

Reads are randomly(-ish) sampled from the DNA



Reads are randomly(-ish) sampled from the DNA



Reads are randomly(-ish) sampled from the DNA



Reads are randomly(-ish) sampled from the DNA

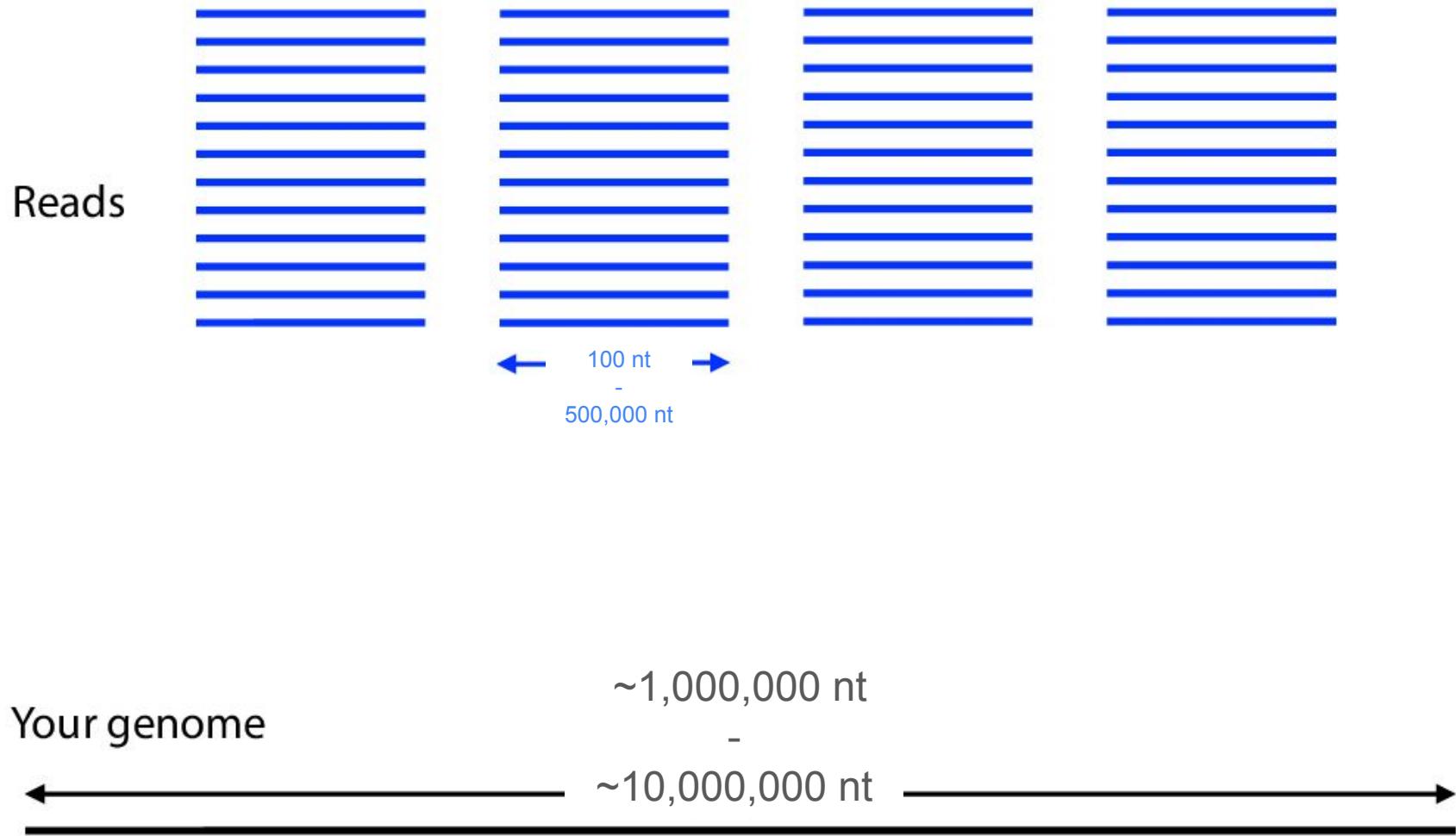
Reads

GTATGCACGCGATAG	TATGTCGCA GTATCT	CACCCTATGTCGCAG	GAGACGCTGGAGCCG
TAGCATTGCGAGACG	GGTATGCA CGCGATA	TGGAGCCGGAGCACC	CGCTGGAGCCGGAGC
TGTCTTGATT CCTG	CGCGATAGCATTGCG	GCATTGCGAGACGCT	CCTATGTCGAGTAT
GACGCTGGAGCCGGA	GCACCCCTATGTCGA	GTATCTGTCTTGAT	CCTCATCCTATTATT
TATCGCACCTACGTT	CAATATTGATCATG	GATCACAGGTCTATC	ACCCTATTAACCACT
CACGGGAGCTCTCCA	TGCATTGGTATTTT	CGTCTGGGGGTATG	CACGCGATAGCATTG
GTATGCACGCGATAG	ACCTACGTTCAATAT	TATTATCGCACCTA	CCACTCACGGGAGCT
GCGAGACGCTGGAGC	CTATCACCCCTATTAA	CTGTCCTTGATTCC	ACTCACGGGAGCTCT
CCTACGTTCAATATT	GCACCTACGTTCAAT	GTCTGGGGGTATGC	AGCCGGAGCACCCCTA
GACGCTGGAGCCGGA	GCACCCCTATGTCGA	GTATCTGTCTTGAT	CCTCATCCTATTATT
TATCGCACCTACGTT	CAATATTGATCATG	GATCACAGGTCTATC	ACCCTATTAACCACT
CACGGGAGCTCTCCA	TGCATTGGTATTTT	CGTCTGGGGGTATG	CACGCGATAGCATTG

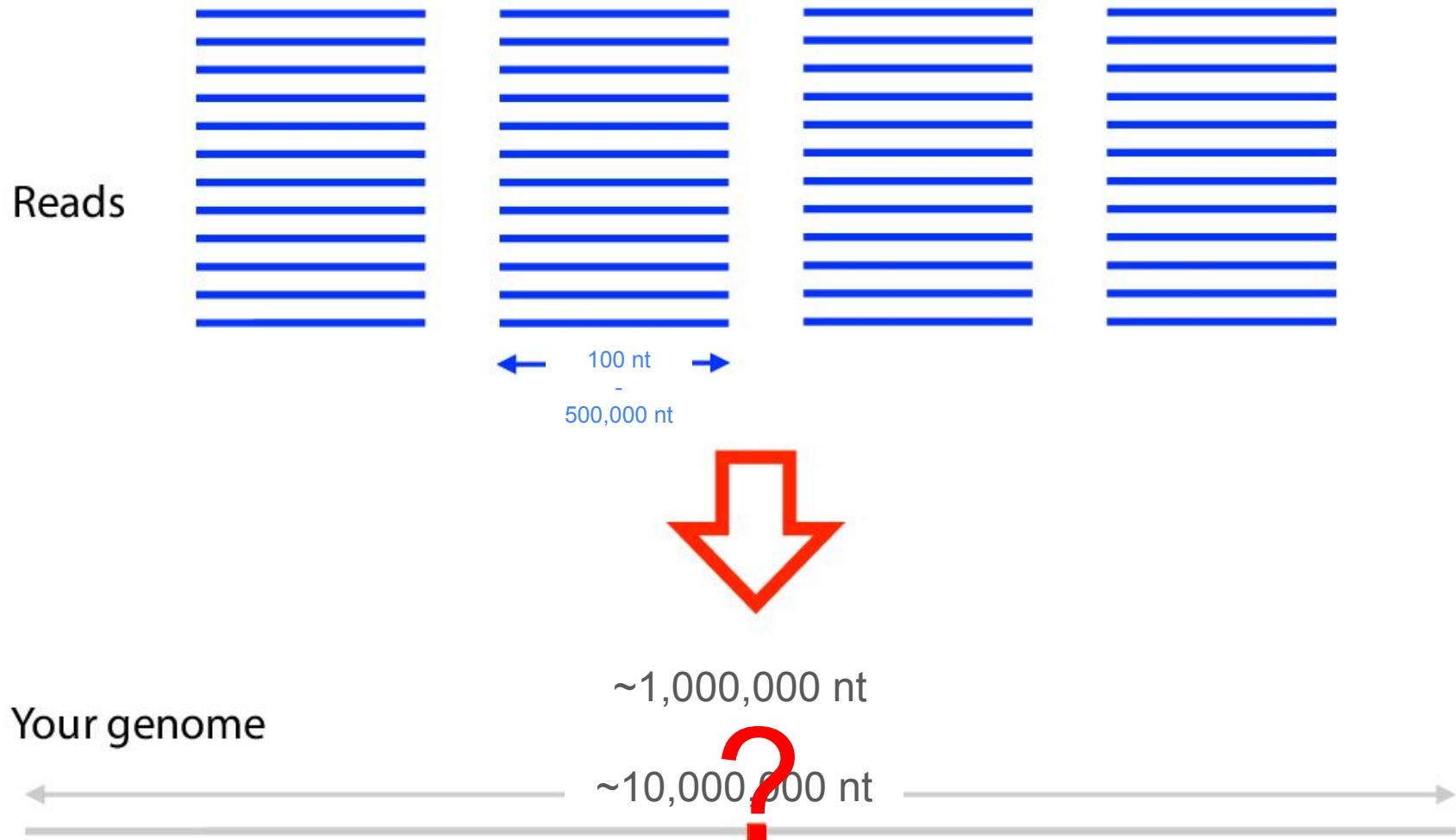
Your genome

CGTCTGGGGGTATGCACGCGATAGCATTGCGAGACGCTGGAGCCGGAGCACCCTATGTCGAGTATCTGTCTTGATT CCTG

Reads are randomly(-ish) sampled from the DNA



Reads are randomly(-ish) sampled from the DNA



Reads also have measurement error: FASTQ

The diagram illustrates a FASTQ sequence record. It consists of two main parts: a header (Label) and a sequence (Sequence). The header is '@FORJUSP02AJWD1' and the sequence is 'CCGTCAATTCAATTAAAGTTTAACCTTGCAGGCCGTACTCCCCAGGGCGGT'. A blue bracket labeled 'Label' points to the first part of the line, and another blue bracket labeled 'Sequence' points to the rest of the line.

```
@FORJUSP02AJWD1  
CCGTCAATTCAATTAAAGTTTAACCTTGCAGGCCGTACTCCCCAGGGCGGT
```

Reads also have measurement error: FASTQ

The diagram illustrates a FASTQ sequence record. It starts with a label '@FORJUSP02AJWD1' followed by a sequence of bases: CCGTCAATTCAATTAAAGTTTAACCTTGCGGCCGTACTCCCCAGGCGGT. Two blue boxes with arrows point to specific parts: one labeled 'Label' points to the start of the sequence, and another labeled 'Sequence' points to the sequence itself.

@FORJUSP02AJWD1
CCGTCAATTCAATTAAAGTTTAACCTTGCGGCCGTACTCCCCAGGCGGT

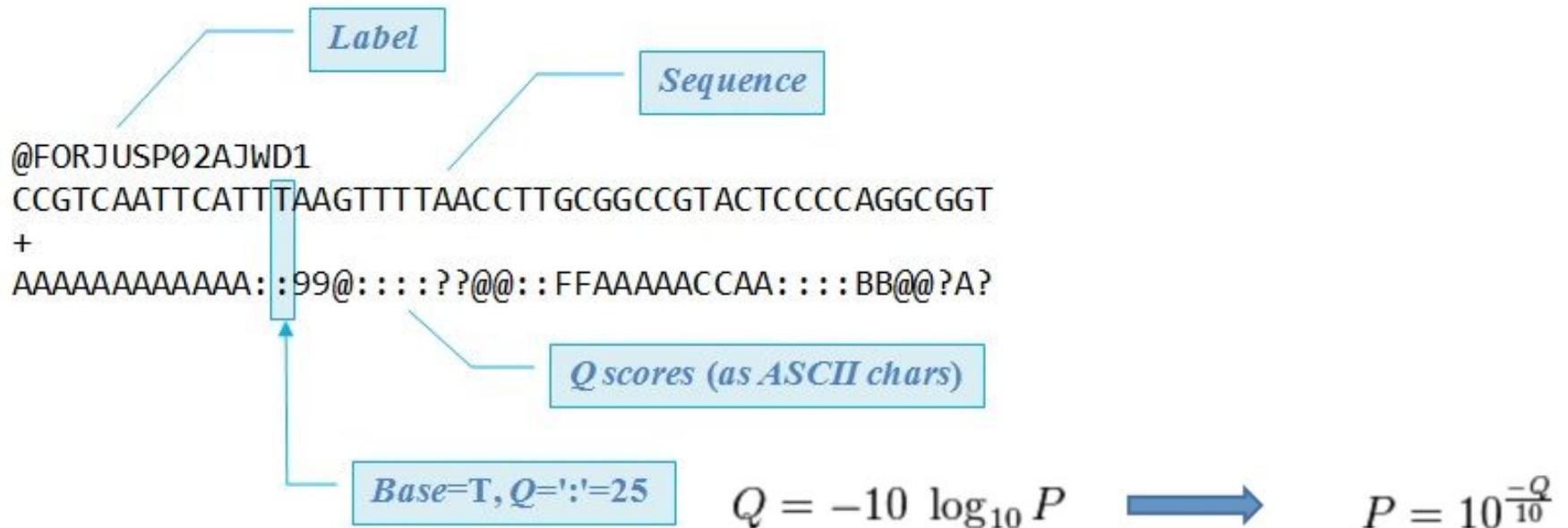
$$Q = -10 \log_{10} P \longrightarrow P = 10^{-\frac{Q}{10}}$$

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 10000	99.99%
50	1 in 100000	99.999%

https://www.drive5.com/usearch/manual/fastq_files.html

<https://learn.genome.edu/ngs-file-formats/quality-scores/>

Reads also have measurement error: FASTQ



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 10000	99.99%
50	1 in 100000	99.999%

https://www.drive5.com/usearch/manual/fastq_files.html

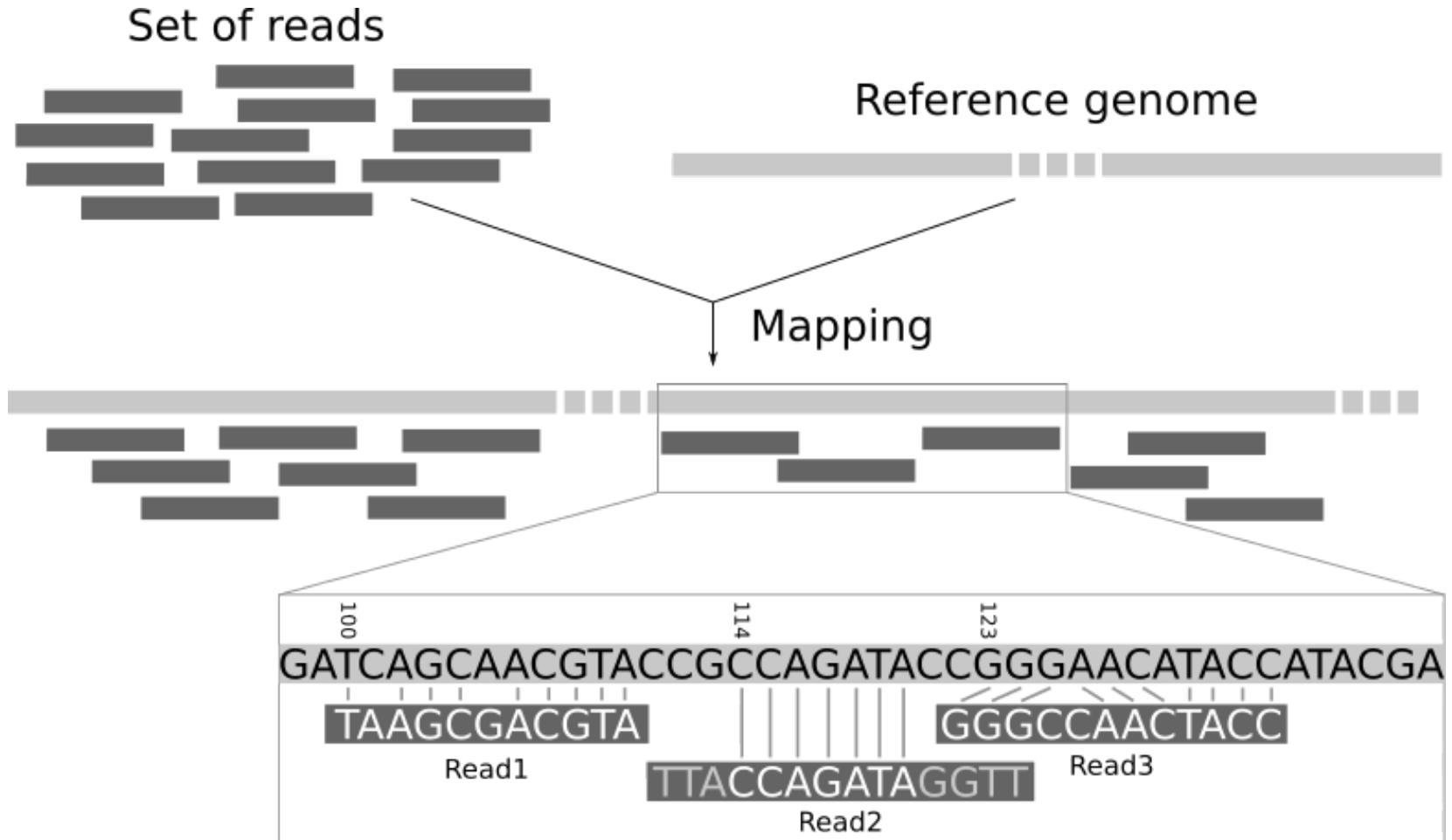
<https://learn.genome.edu/ngs-file-formats/quality-scores/>

What do we do with these reads?

Reference-based assembly



Reference-based assembly



Variant Calling / Consensus generation

reference	GATCCATGTAGTACCAT T AGTACAGTACCATATAT GATCCATGTAGTACCAT T AGTACAGTACCATATA ATCCATGTAGTACCAT T AGTACAGTACCATATAT CATGTAGTACCAT T AGTACAGTACCATATAT GTAGTACCAT T AGTACAGTACCATATAT GTAGTACCAT T AGTACAGTACCATATAT
reads	TAGTACCAT T AGTACAGTACCATATAT TAGTACCAT T AGTACAGTACCATATAT AGTACCAT T AGTACAGTACCATATAT AGTACCAT T AGTACAGTACCATATAT GTACCAT T AGTACAGTACCATATAT
consensus	GATCCATGTAGTACCAT T AGTACAGTACCATATAT

Variant Calling / Consensus generation

- *variants or mutations* are detected when the aligned reads differ from the reference

reference	GATCCATGTAGTACCAT T AGTACAGTACCATATAT GATCCATGTAGTACCAT C AGTACAGTACCATATA ATCCATGTAGTACCAT C AGTACAGTACCATATAT CATGTAGTACCAT C AGTACAGTACCATATAT GTAGTACCAT C AGTACAGTACCATATAT GTAGTACCAT C AGTACAGTACCATATAT
reads	TAGTACCAT C AGTACAGTACCATATAT TAGTACCAT C AGTACAGTACCATATAT AGTACCAT C AGTACAGTACCATATAT AGTACCAT C AGTACAGTACCATATAT GTACCAT C AGTACAGTACCATATAT
consensus	GATCCATGTAGTACCATCAGTACAGTACCATATAT

Variant Calling / Consensus generation

- *ambiguous or mixed bases* are detected when the aligned reads have evidence for more than one type of base

reference

GATCCATGTAGTACCAT**T**AGTACAGTACCATATAT

GATCCATGTAGTACCAT**T**AGTACAGTACCATATA

ATCCATGTAGTACCAT**C**AGTACAGTACCATATAT

CATGTAGTACCAT**C**AGTACAGTACCATATAT

GTAGTACCAT**C**AGTACAGTACCATATAT

GTAGTACCAT**T**AGTACAGTACCATATAT

TAGTACCAT**C**AGTACAGTACCATATAT

TAGTACCAT**C**AGTACAGTACCATATAT

AGTACCAT**T**AGTACAGTACCATATAT

AGTACCAT**T**AGTACAGTACCATATAT

GTACCAT**C**AGTACAGTACCATATAT

reads

GATCCATGTAGTACCAT**Y**AGTACAGTACCATATAT

consensus



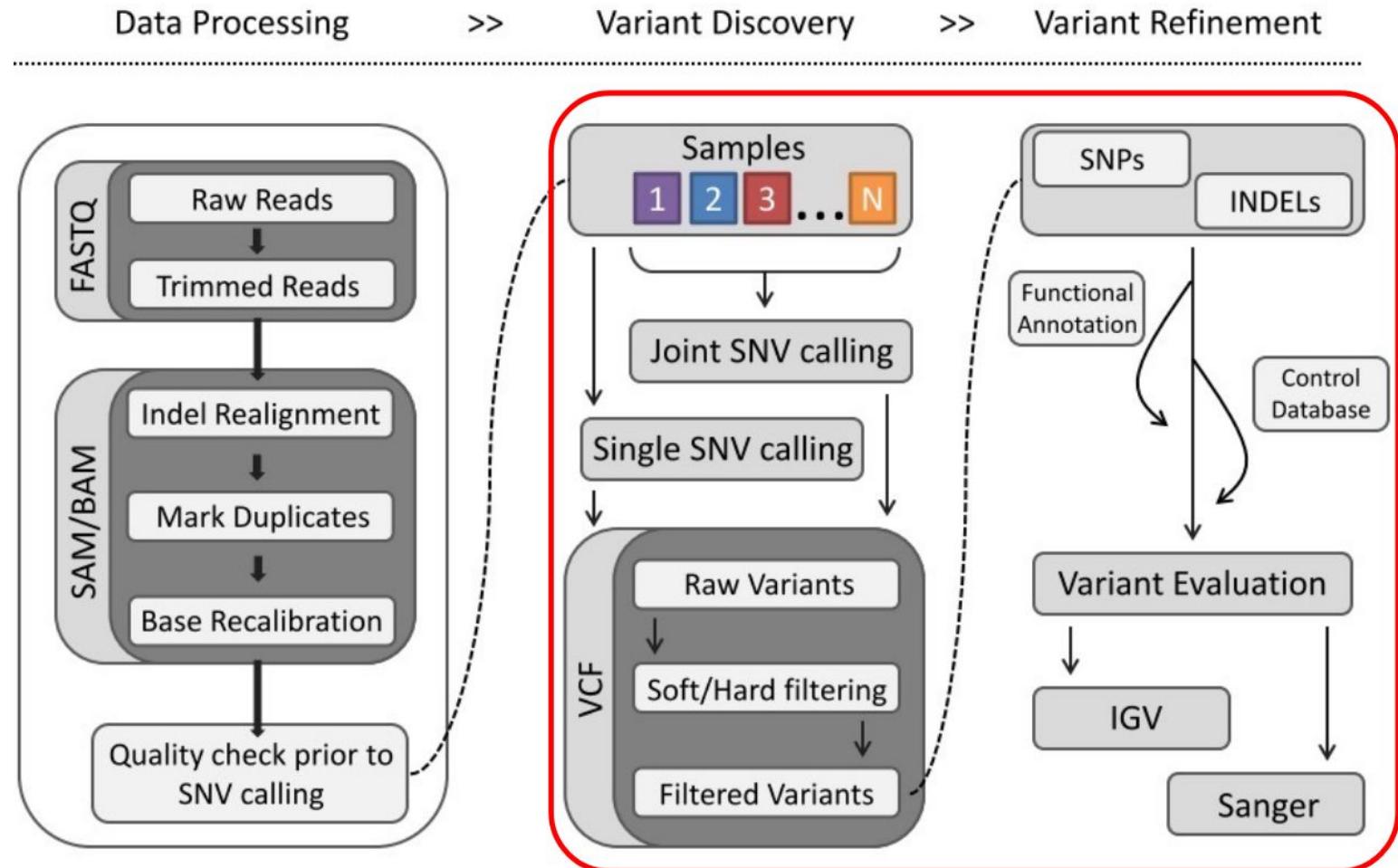
IUPAC ambiguity

Distinguishing real variants from error is non-trivial

GT_TA_CTGTCGTTGTAATAC_TCCAC_GATGTC
GT_TA_CTGTCGTTGTAATAC_TCCAC_GATGTC
GT_TA_CTGTCGTTGTAATAC_TCCAC_GATGTC
GT_TA_CTGTCGTTGTAATAC_TCCAC_GATGTC
GT_TA_CTGTCGTTGTAATgCTCCACGATGTC
GT_TA_CTGTCGTTGTAATAC_TCCAC_GAATGTC
GT_TA_CTGTCGTTGTAATAC_TCCAC_GATGTC
GT_TA_CTGTCGTGGTAATAC_TCCAC_GATGTC
GT_TA_CTGTCGTTGTAATAC_TCCAC_GATGTC
GT_TA_aTGTCGTTGTAATAC_TCCAC_GATGTC
GT_TA_CTGTCGTTGTA_cTAC_TCCAC_GATGTC
GT_TA_CTGTCGTTGTAATAC_TCCAC_GATGTC

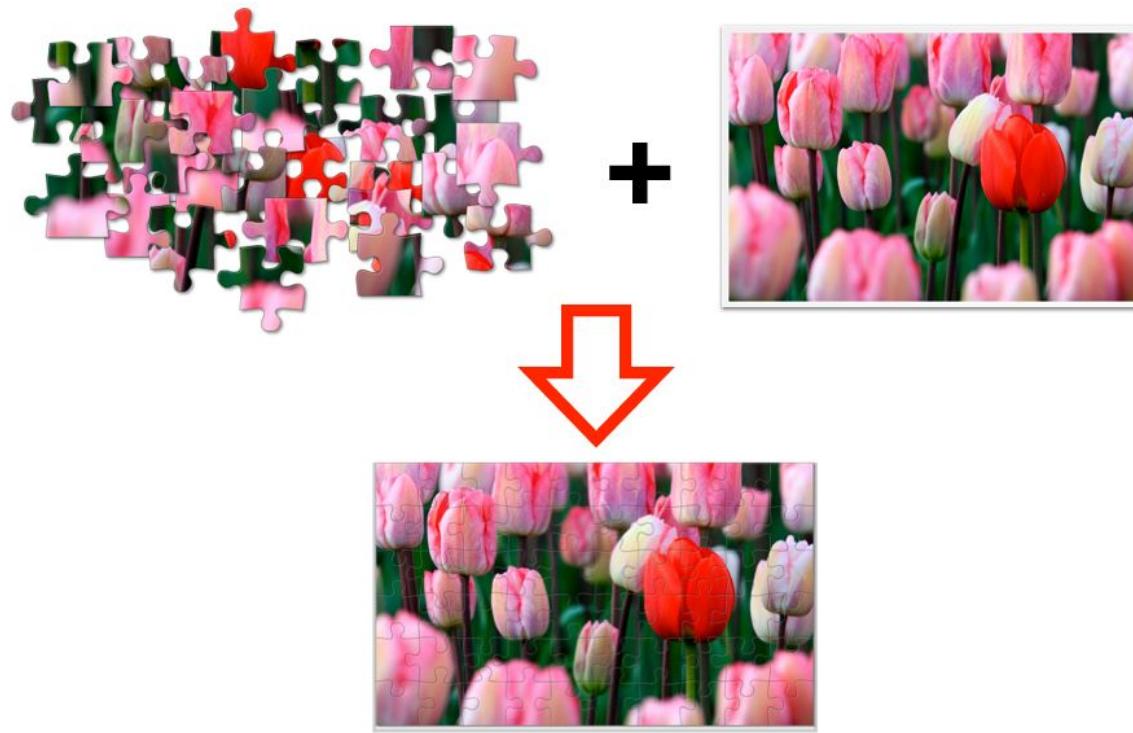
sequencing errors SNP

Lots of steps to get accurate reference-based data

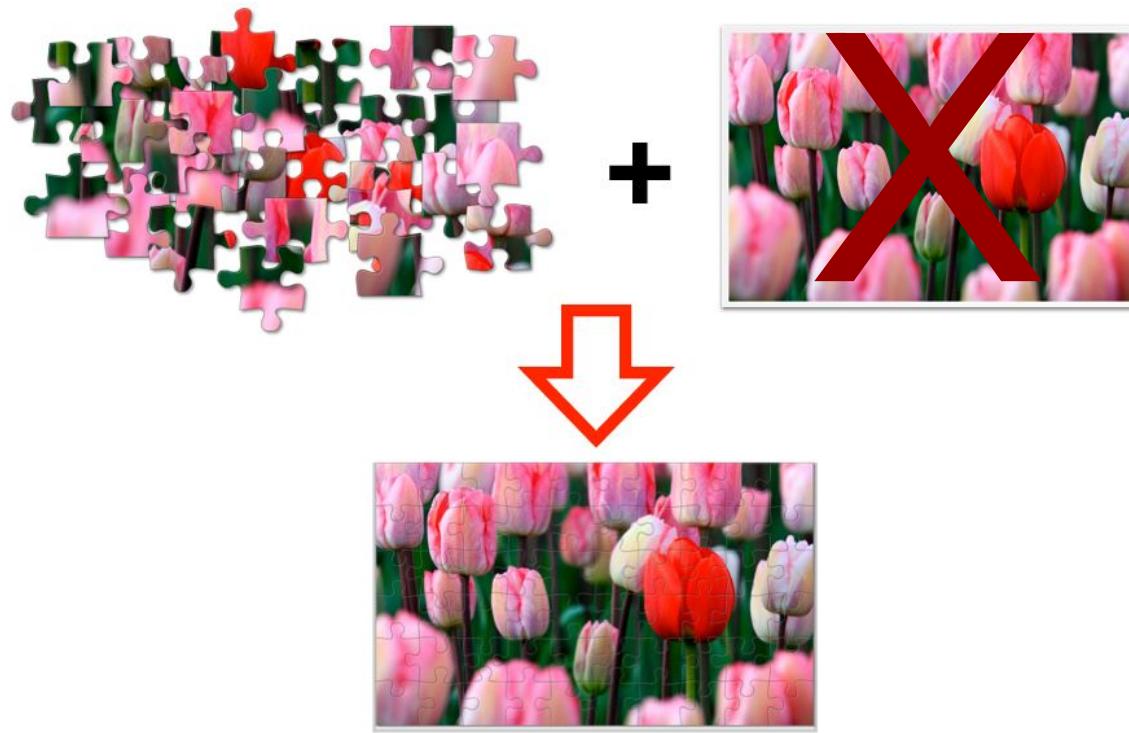


What if we don't have a reference?

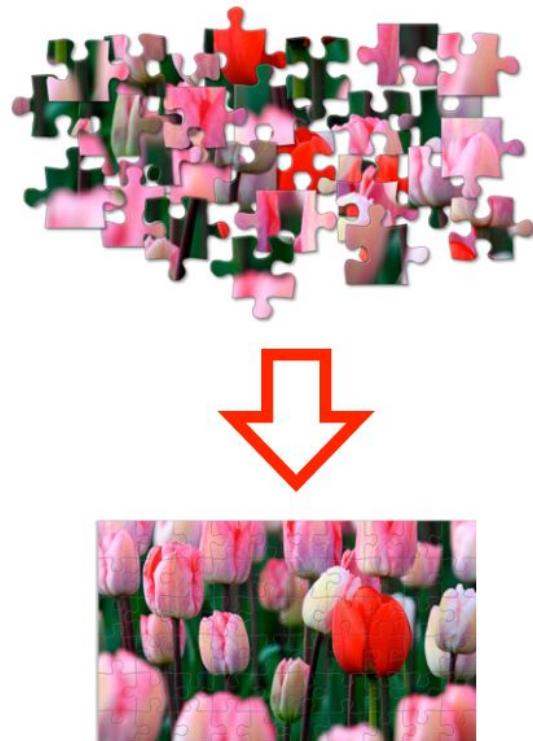
de novo Assembly



de novo Assembly



de novo Assembly



de novo Assembly

Reconstruct this

CTAGGCCCTCAATTTTT
CTCTAGGCCCTCAATTTTT
GGCTCTAGGCCCTCATTTTT
CTCGGCTCTAGCCCCTCATT
TATCTCGACTCTAGGCCCTCA
TATCTCGACTCTAGGCC
TCTATATCTCGGCTCTAGG
GGCGTCTATATCTCG
GGCGTCGATATCT
GGCGTCTATATCT
GGCGTCTATATCTCGGCTCTAGGCCCTCATT

From
these

→ GGCGTCTATATCTCGGCTCTAGGCCCTCATT

de novo Assembly

Reconstruct this

???

CTAGGCCCTCAATTTT
GGCGTCTATATCT
CTCTAGGCCCTCAATTTT
TCTATATCTCGGCTCTAGG
GGCTCTAGGCCCTCATTTTT
CTCGGCTCTAGCCCCTCATTTT
TATCTCGACTCTAGGCCCTCA
GGCGTCGATATCT
TATCTCGACTCTAGGCC
GGCGTCTATATCTCG

From
these

Enough sampling (depth) and reads will overlap

If the end of 1 read matches the start of another:

The diagram illustrates two DNA sequences. The top sequence is **TCTATATCTCGGCTCTAGG**, and the bottom sequence is **TATCTCGACTCTAGGCC**. Vertical lines connect the bases at positions 1 through 11, indicating a perfect match between the two sequences across their entire lengths.

Enough sampling (depth) and reads will overlap

If the end of 1 read matches the start of another:

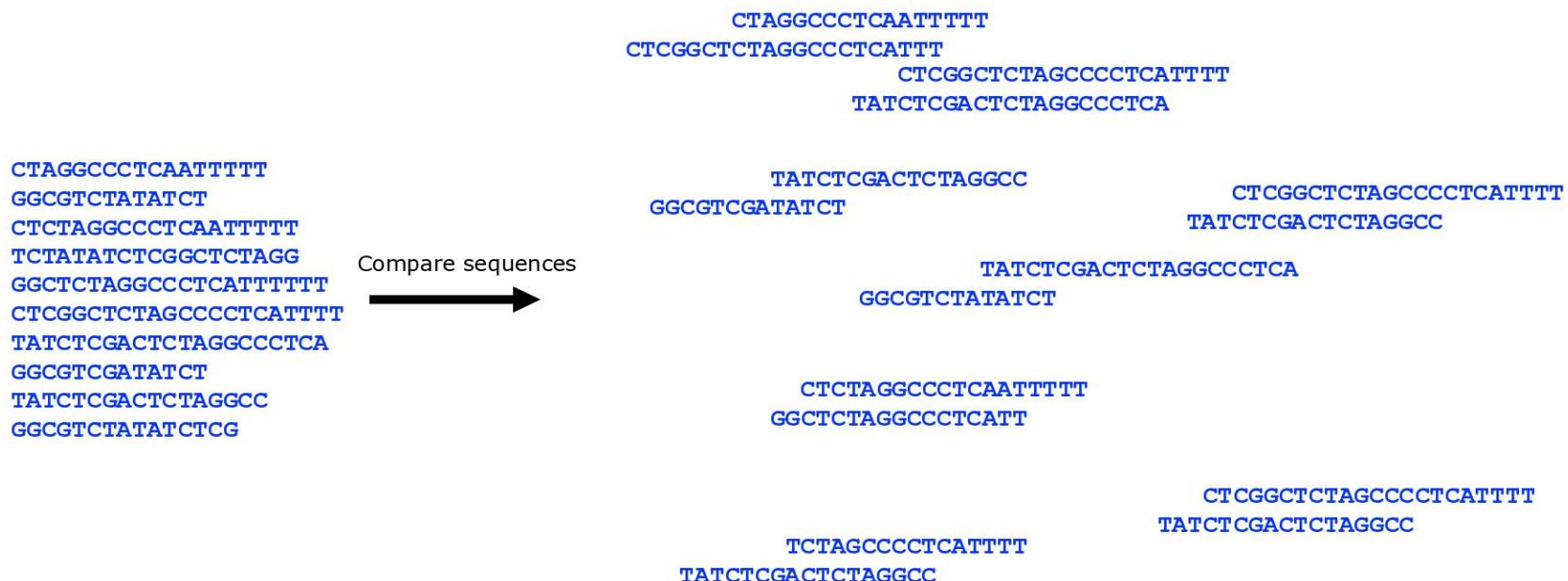
TCTATATCTCGGCTCTAGG
| | | | | | | | | |
TATCTCGACTCTAGGCC

Then they **MIGHT** be from overlapping bits of the genome:

TCTATATCTCGGCTCTAGG
GGCGTCTATATCTCGGCTCTAGGCCCTCATTTTT
TATCTCGACTCTAGGCC

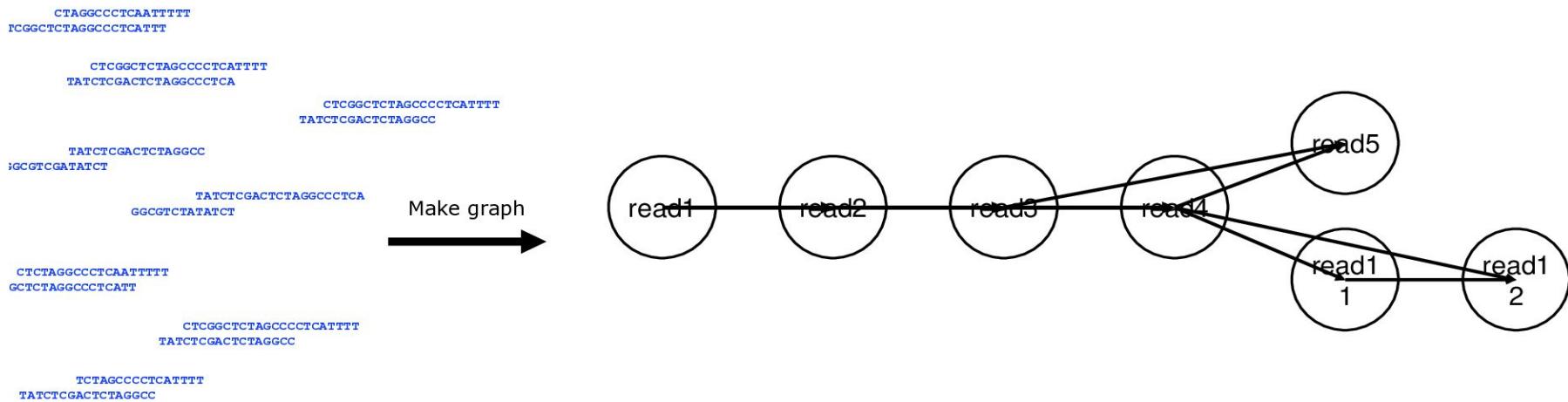
Use shared sequences to get longer fragments (contigs)

- Step 1: Find all overlapping pairs of reads



Use shared sequences to get longer fragments (contigs)

- Step 2: Construct a graph representing read-read overlaps

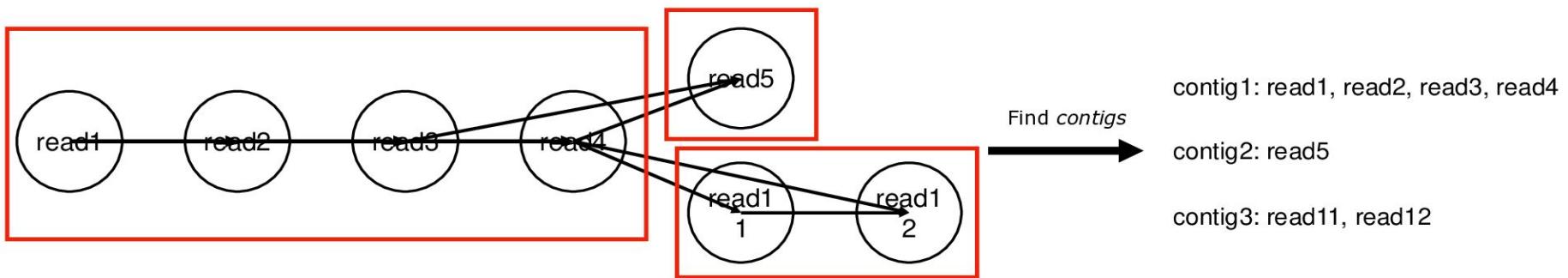


Every read is a vertex

Every overlapping pair of reads is connected with an edge

Use shared sequences to get longer fragments (contigs)

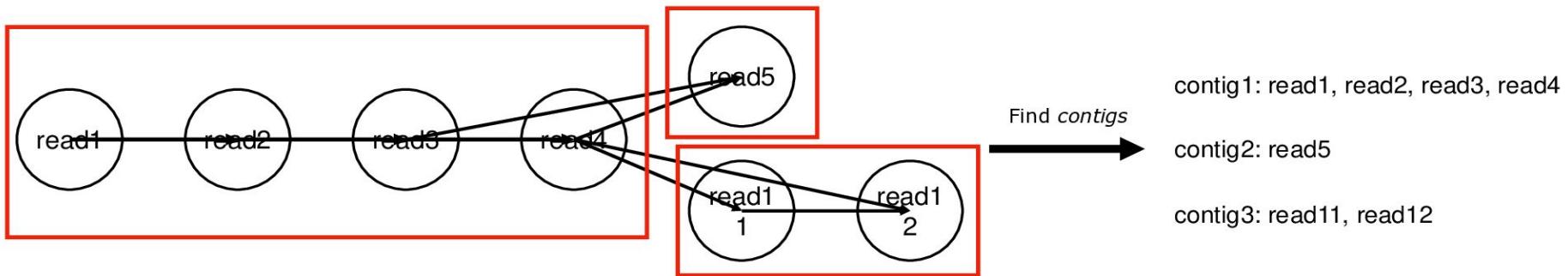
- Step 3: Analyze graph structure to find possible repeat boundaries



Key term: *contig*. Short for “contiguous sequence” it is the result of assembling reads together

Use shared sequences to get longer fragments (contigs)

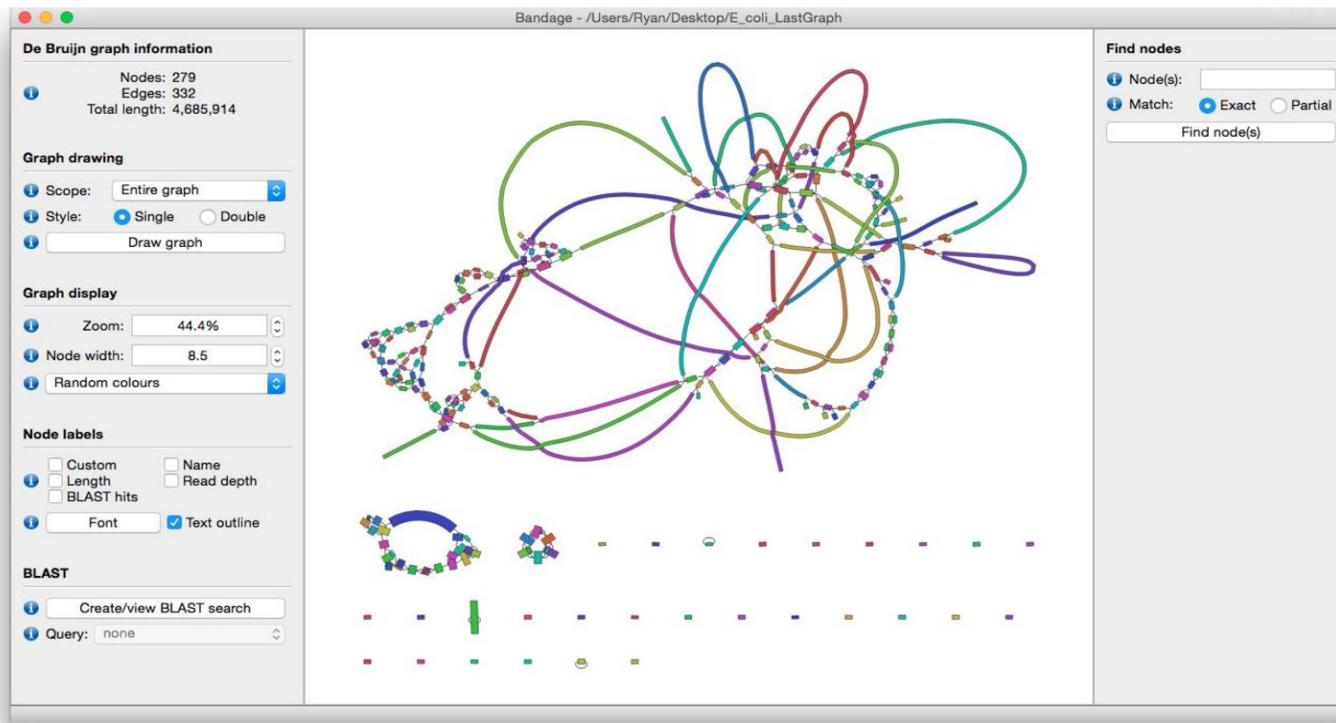
- Step 3: Analyze graph structure to find possible repeat boundaries



Key term: *contig*. Short for “contiguous sequence” it is the result of assembling reads together

Note: assembly is often done with a variant of this which uses exact matches of even shorter bits of reads (called k-mers)

Real assembly graphs are messy and complicated!



Simplified graph of *E. Coli* (4.6Mbp) genome, viewed with Bandage: https://rrwick.github.io/E_coli_LG.html

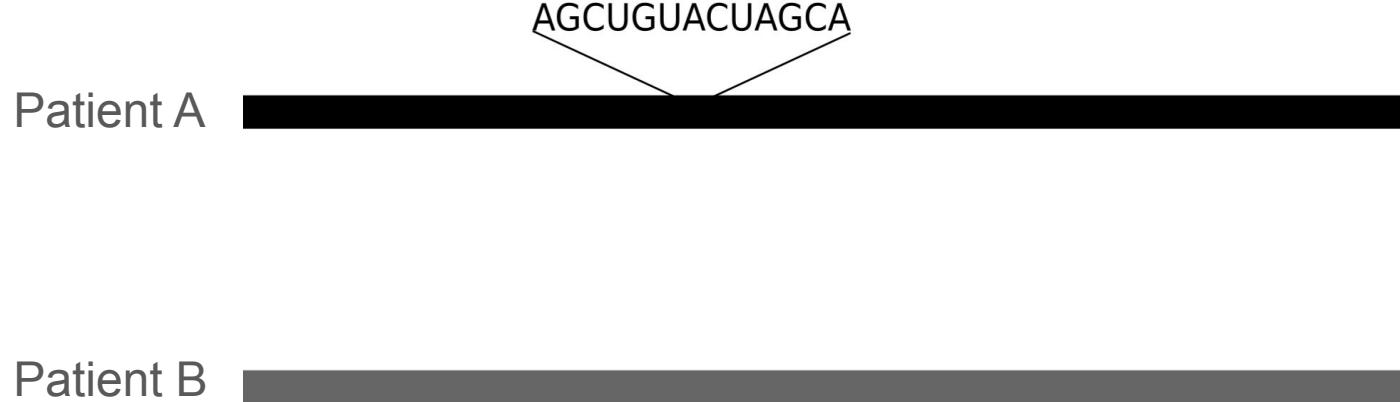
Got genomes, but how do we use them to
stop our outbreak?

Compare genomes from each patient

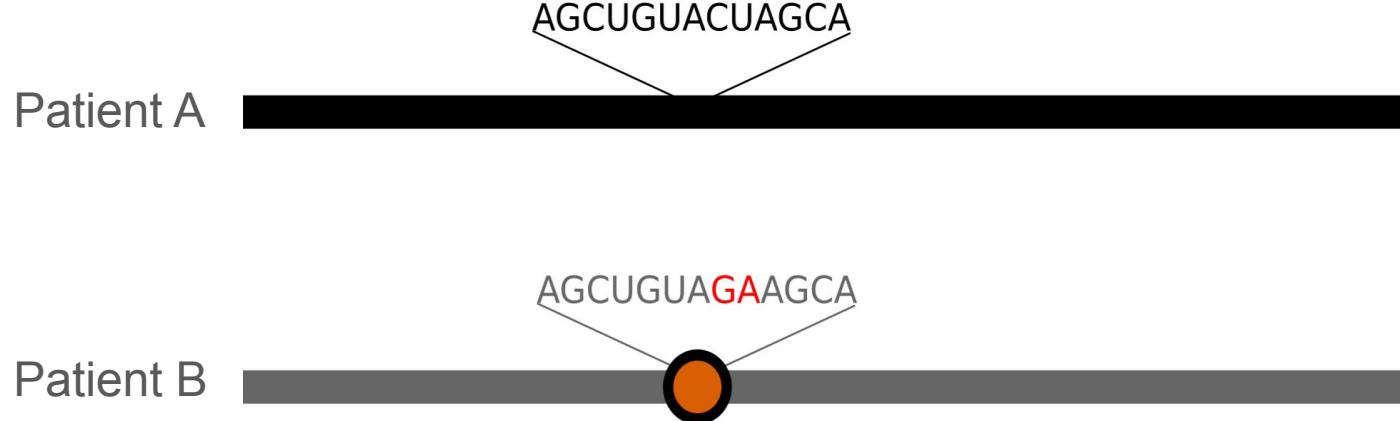
Patient A

Patient B

Compare genomes from each patient



Compare genomes from each patient



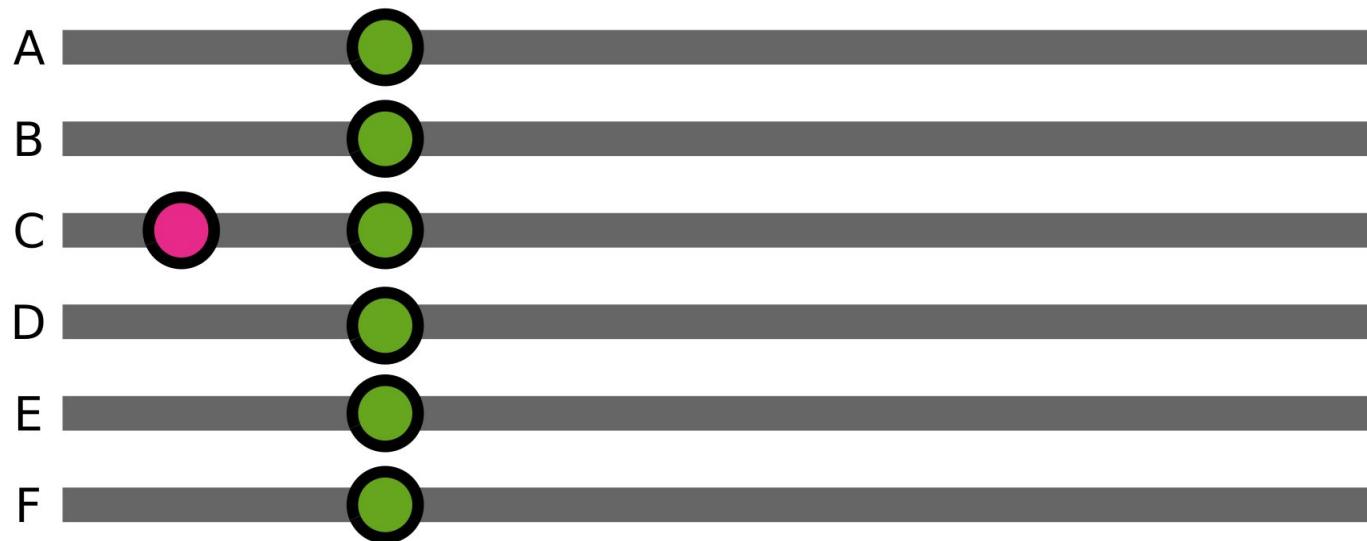
Compare (align) genomes from each patient



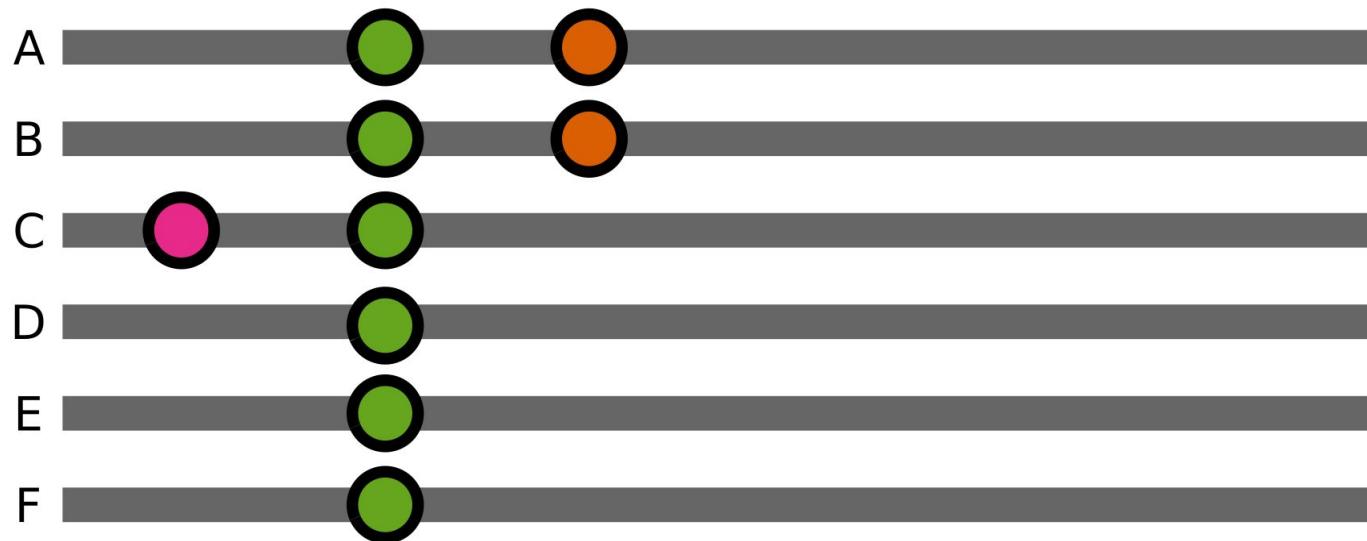
Compare (align) genomes from each patient



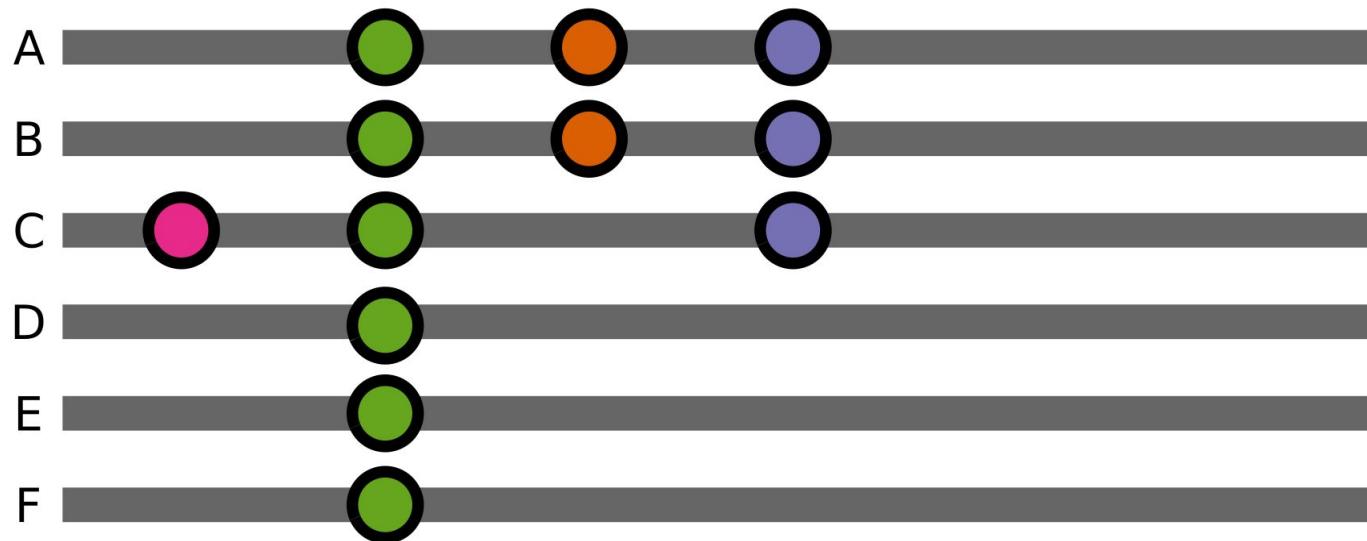
Compare (align) genomes from each patient



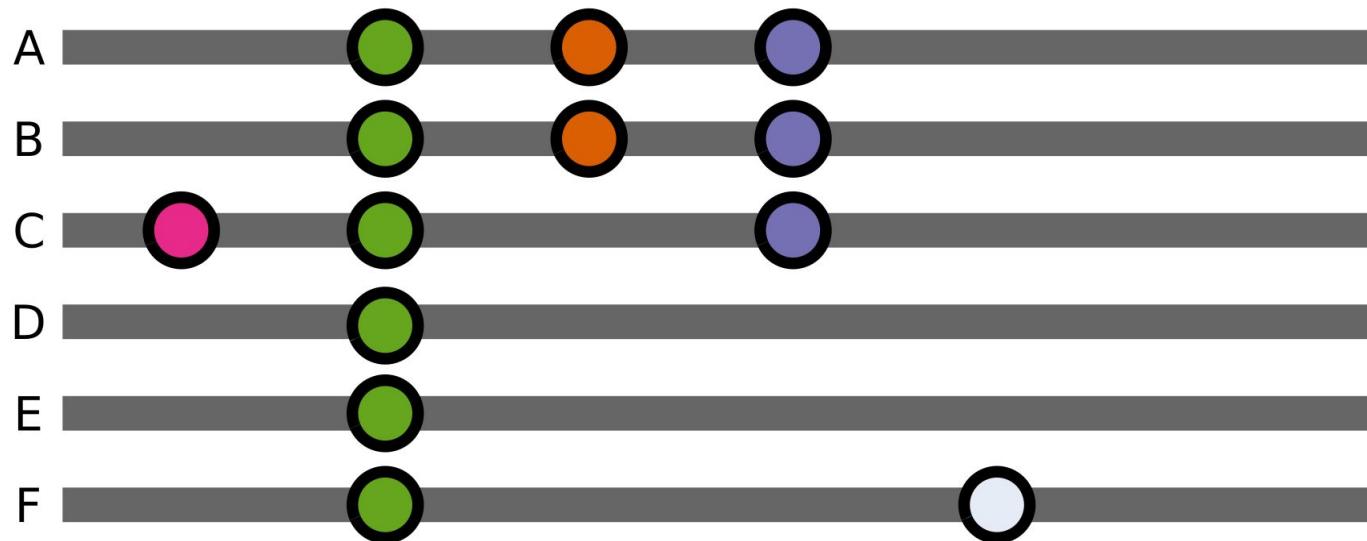
Compare (align) genomes from each patient



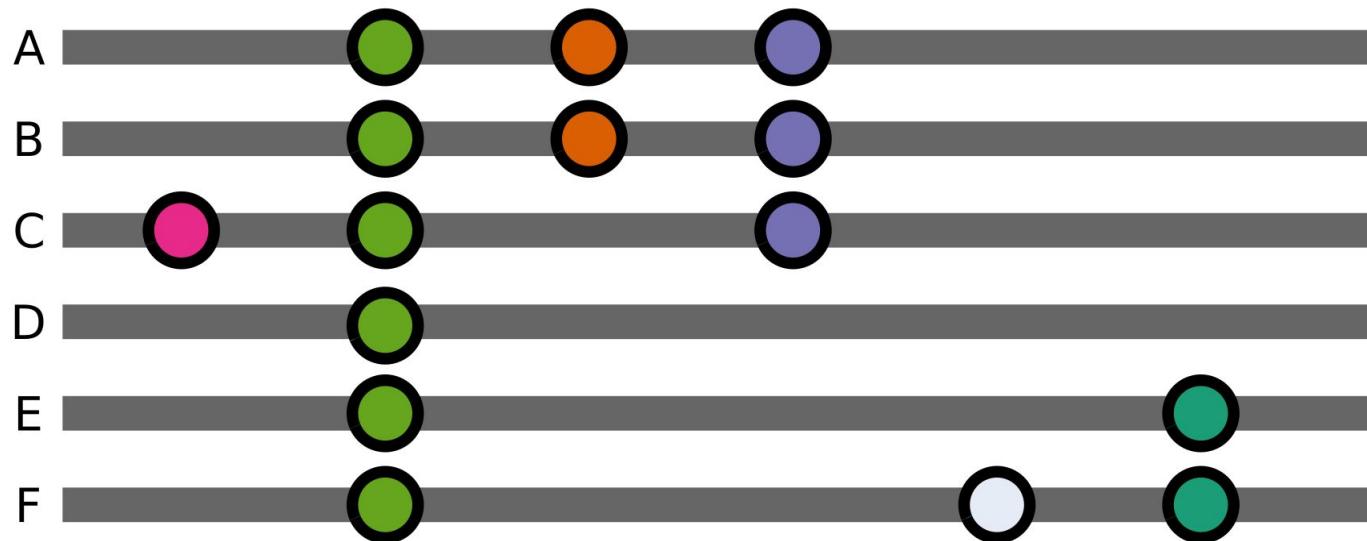
Compare (align) genomes from each patient



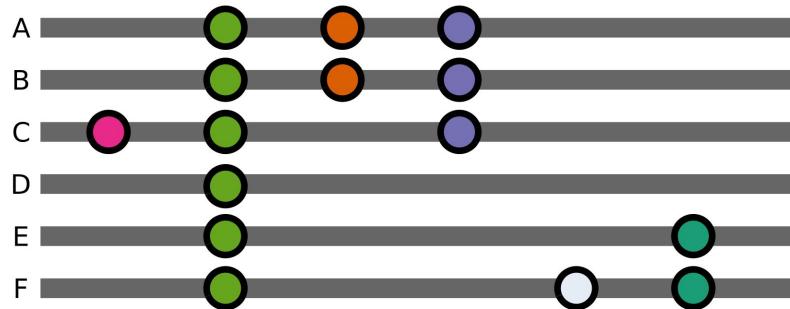
Compare (align) genomes from each patient



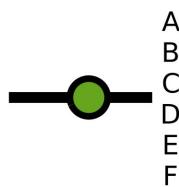
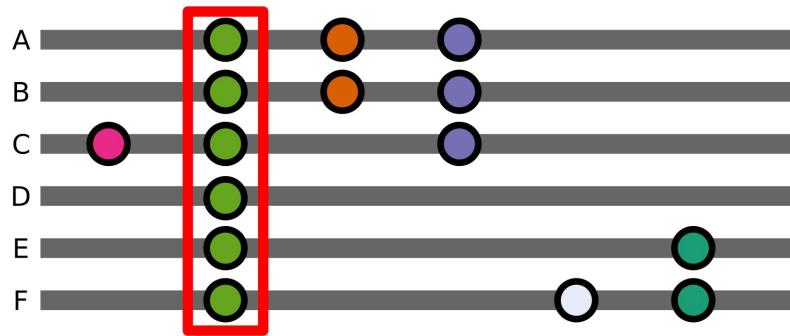
Compare (align) genomes from each patient



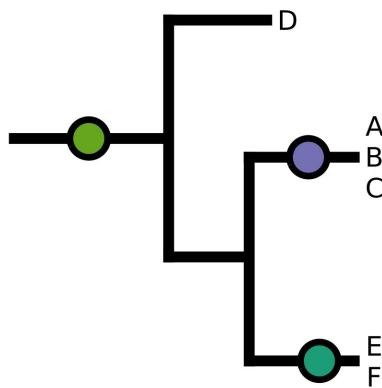
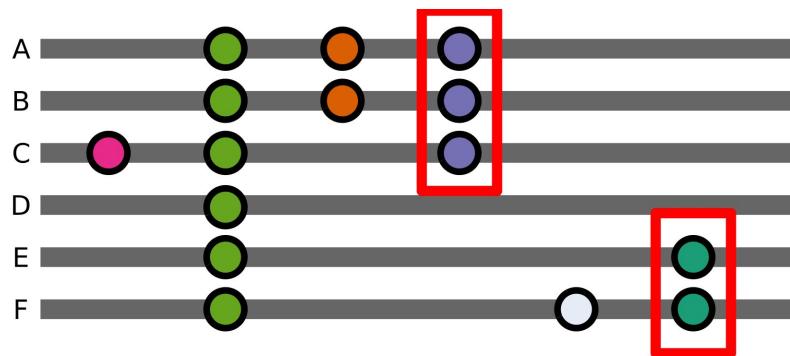
Use pattern of similarities/differences to infer relationships



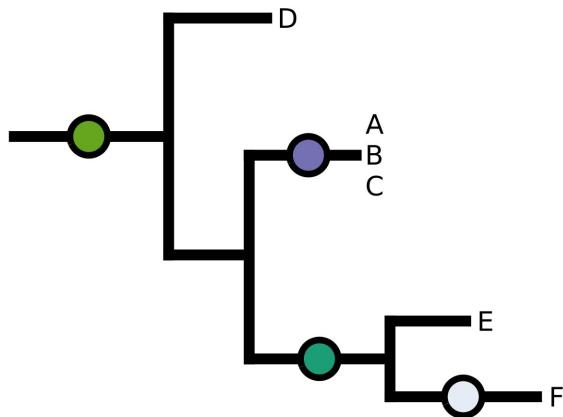
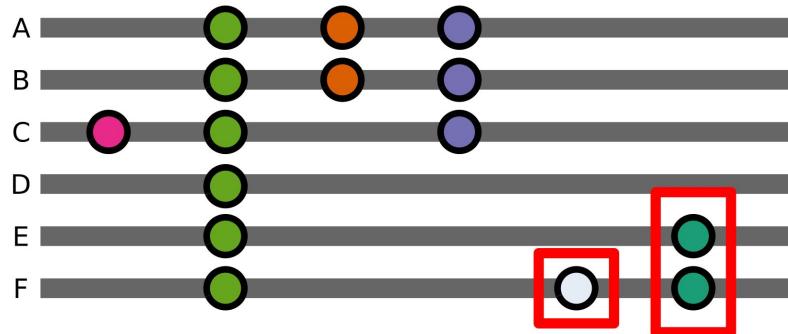
Use pattern of similarities/differences to infer relationships



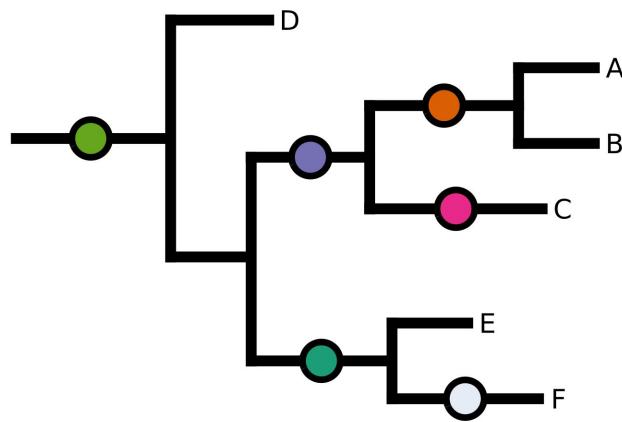
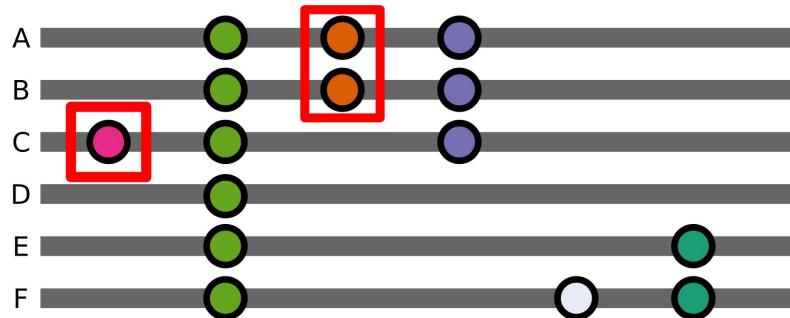
Use pattern of similarities/differences to infer relationships



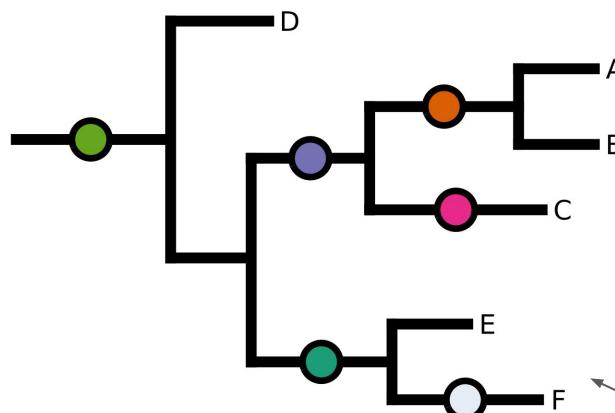
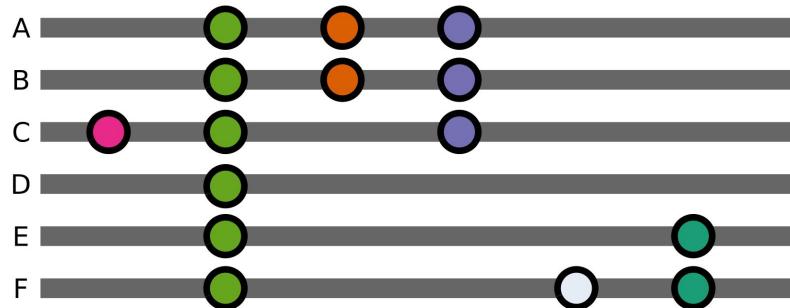
Use pattern of similarities/differences to infer relationships



Use pattern of similarities/differences to infer relationships



Use pattern of similarities/differences to infer relationships

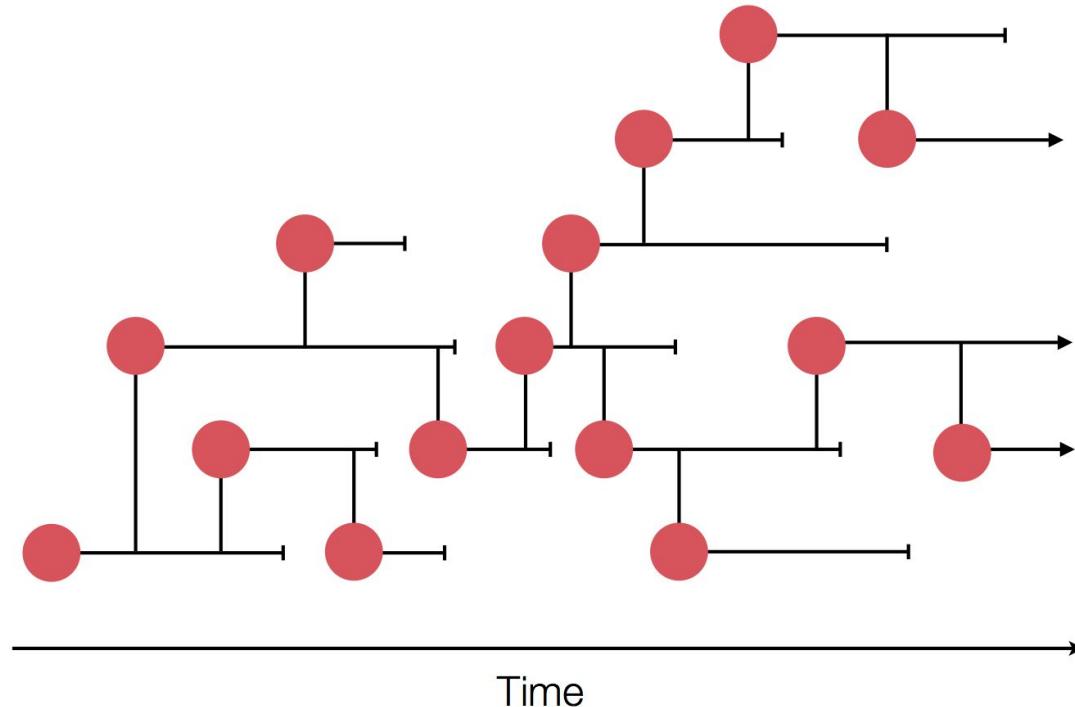


Phylogenetic Tree

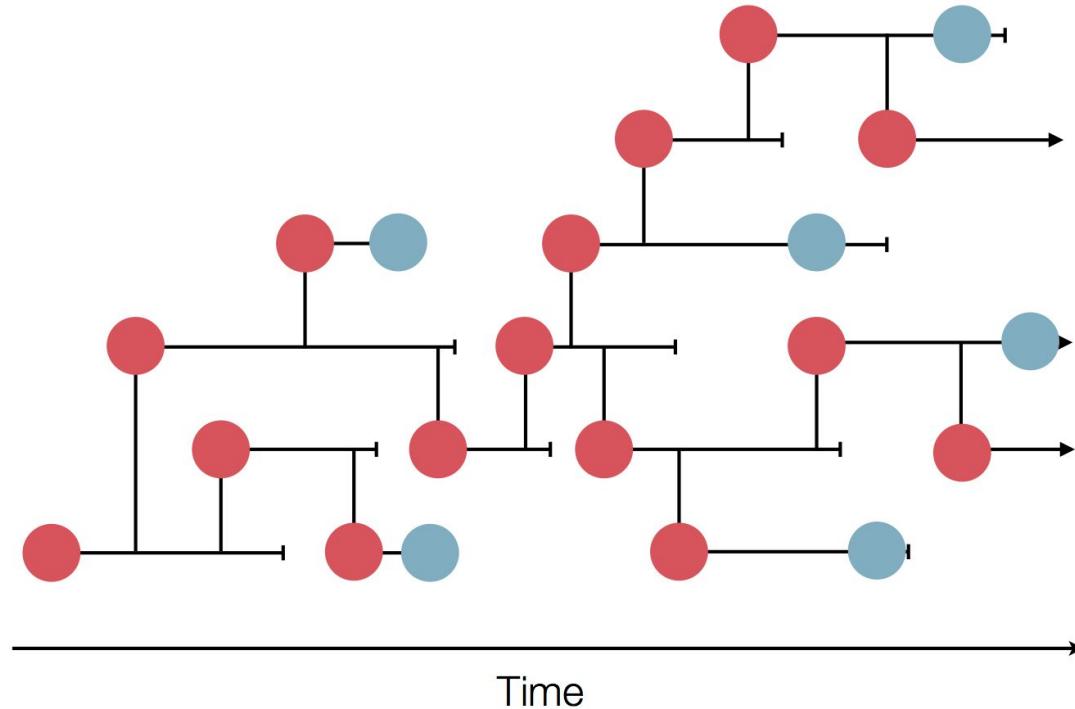
Note: this generally uses statistical models which incorporate differences in substitution rates and mutation rates

How does this tree help us?

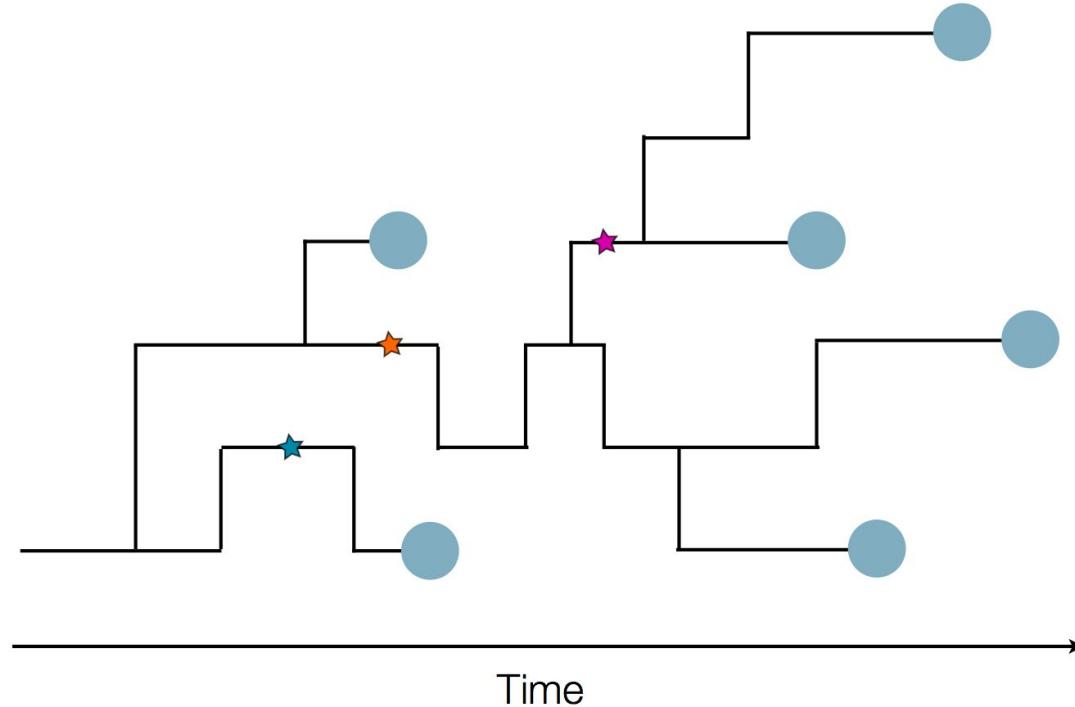
Trees sample the underlying transmission network



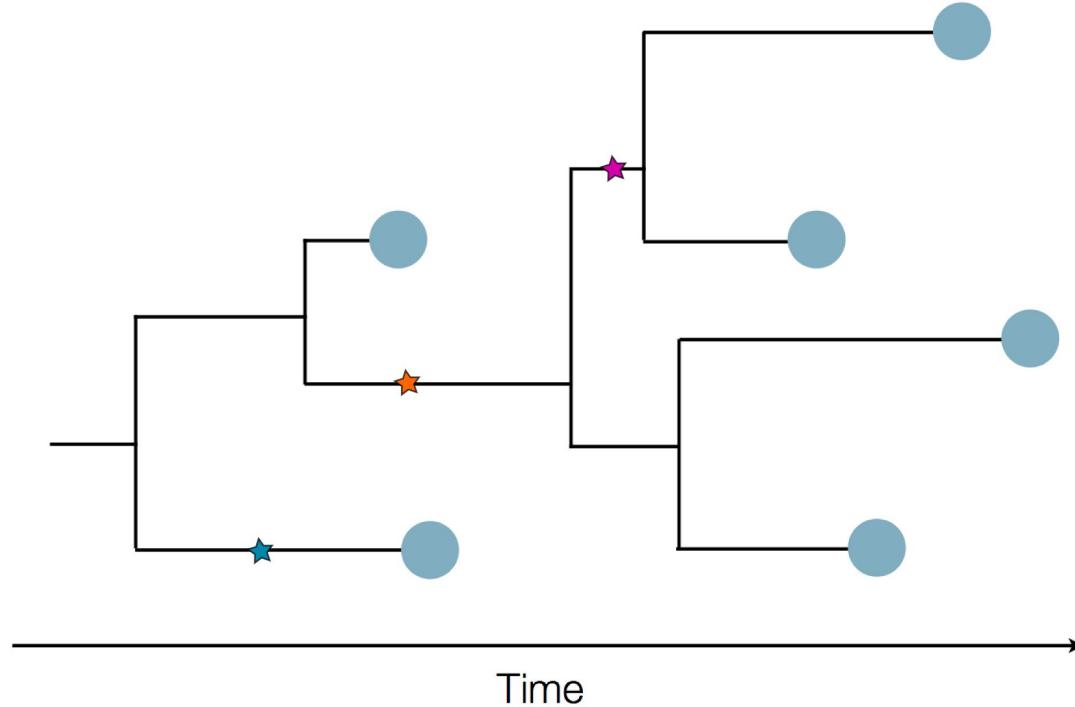
Trees sample the underlying transmission network



Trees sample the underlying transmission network

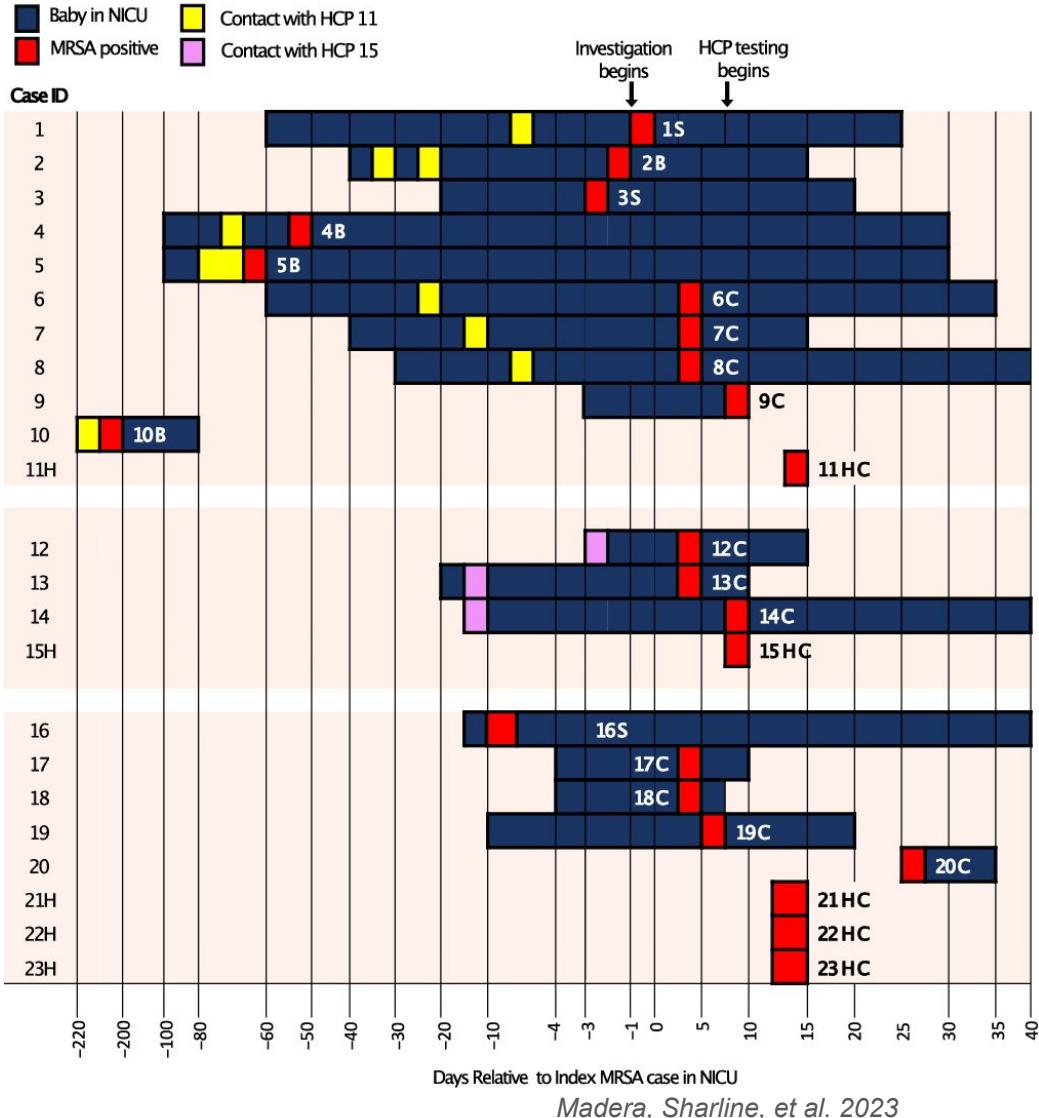


Trees sample the underlying transmission network

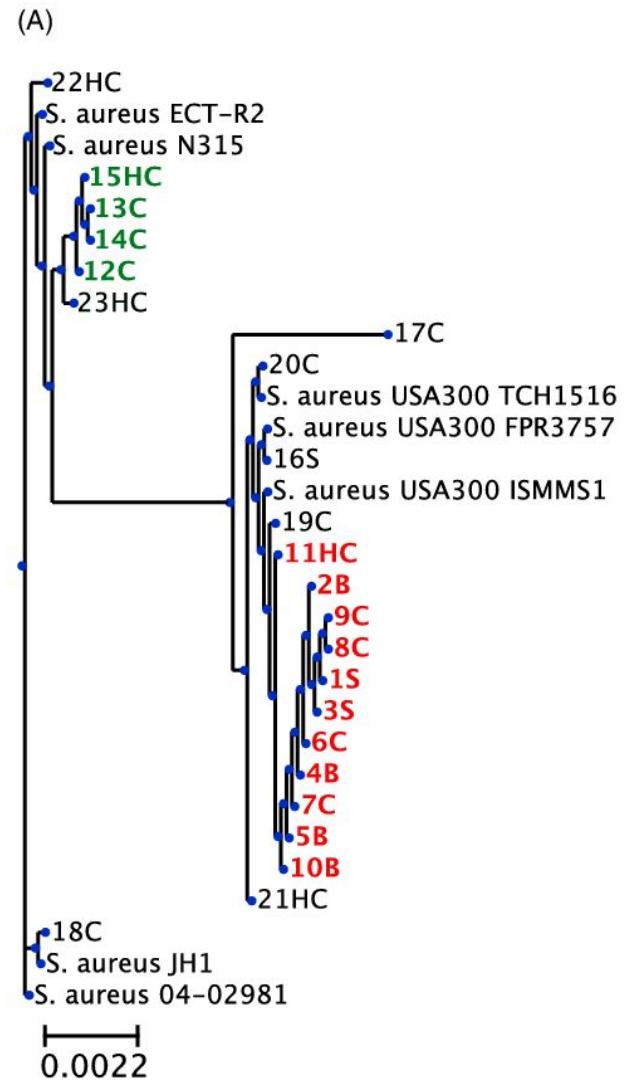
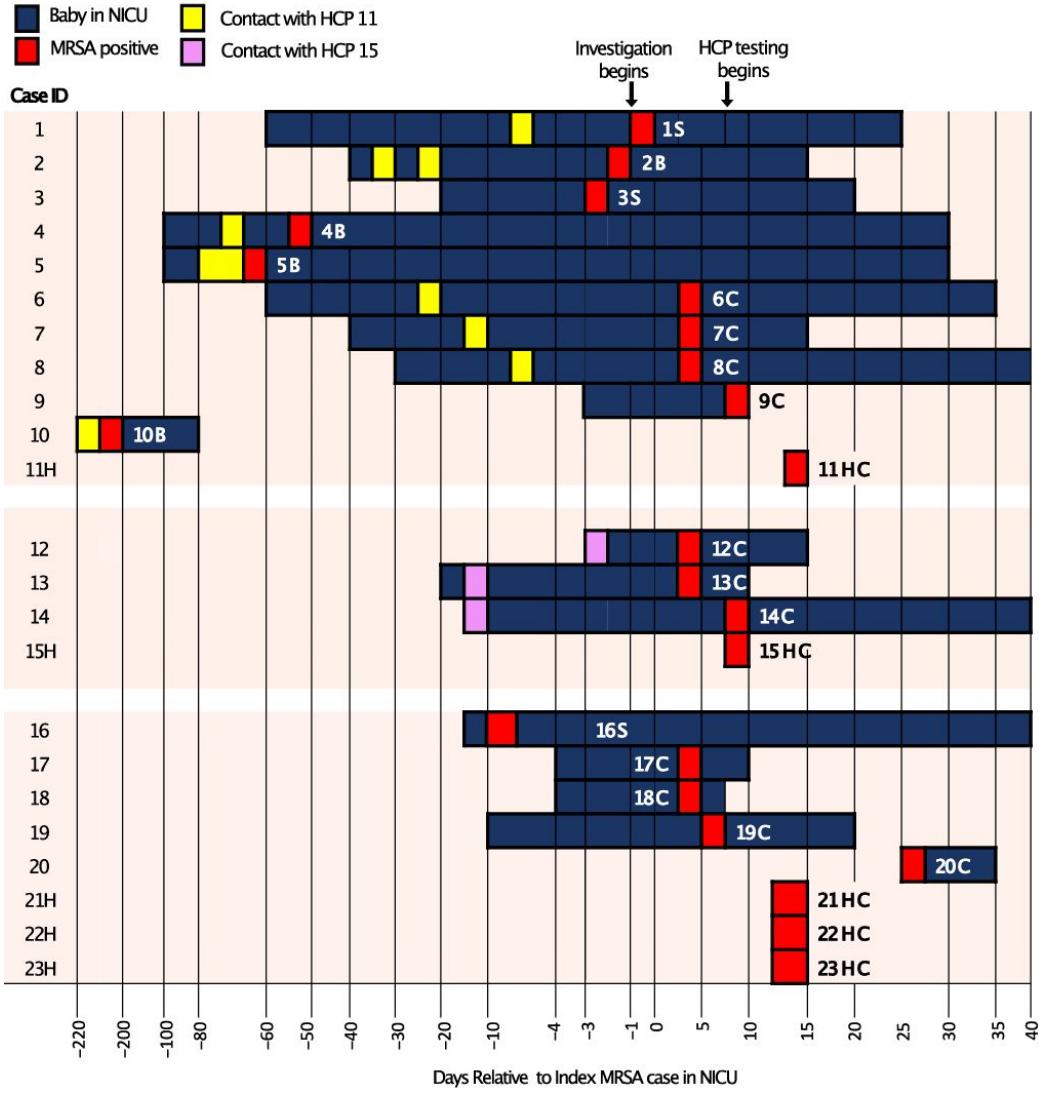


What do our genomes & trees tell us about
our NICU-MRSA outbreak?

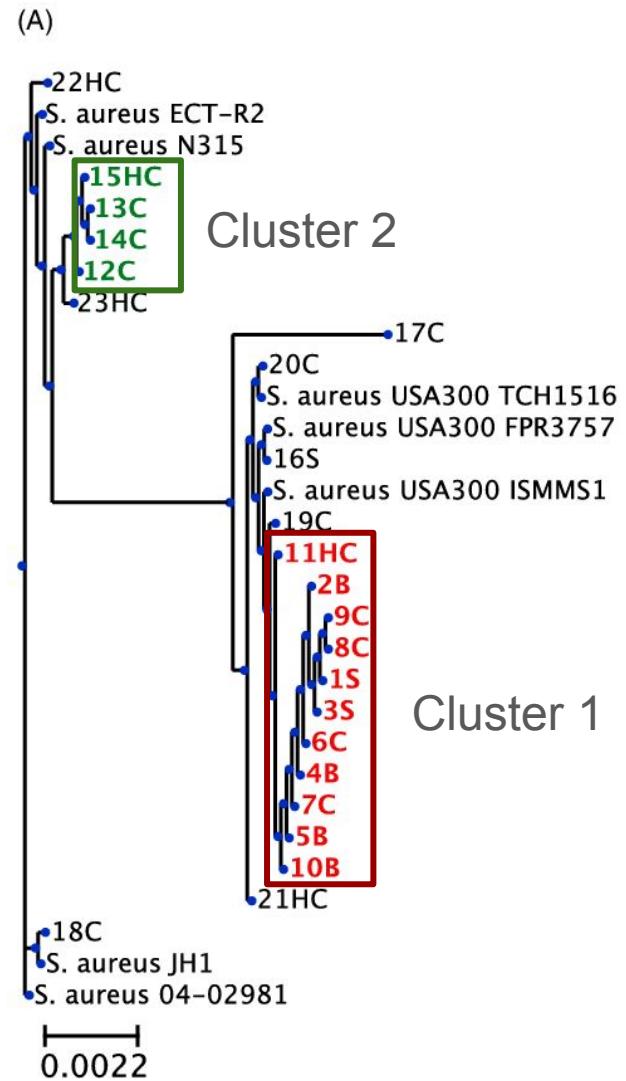
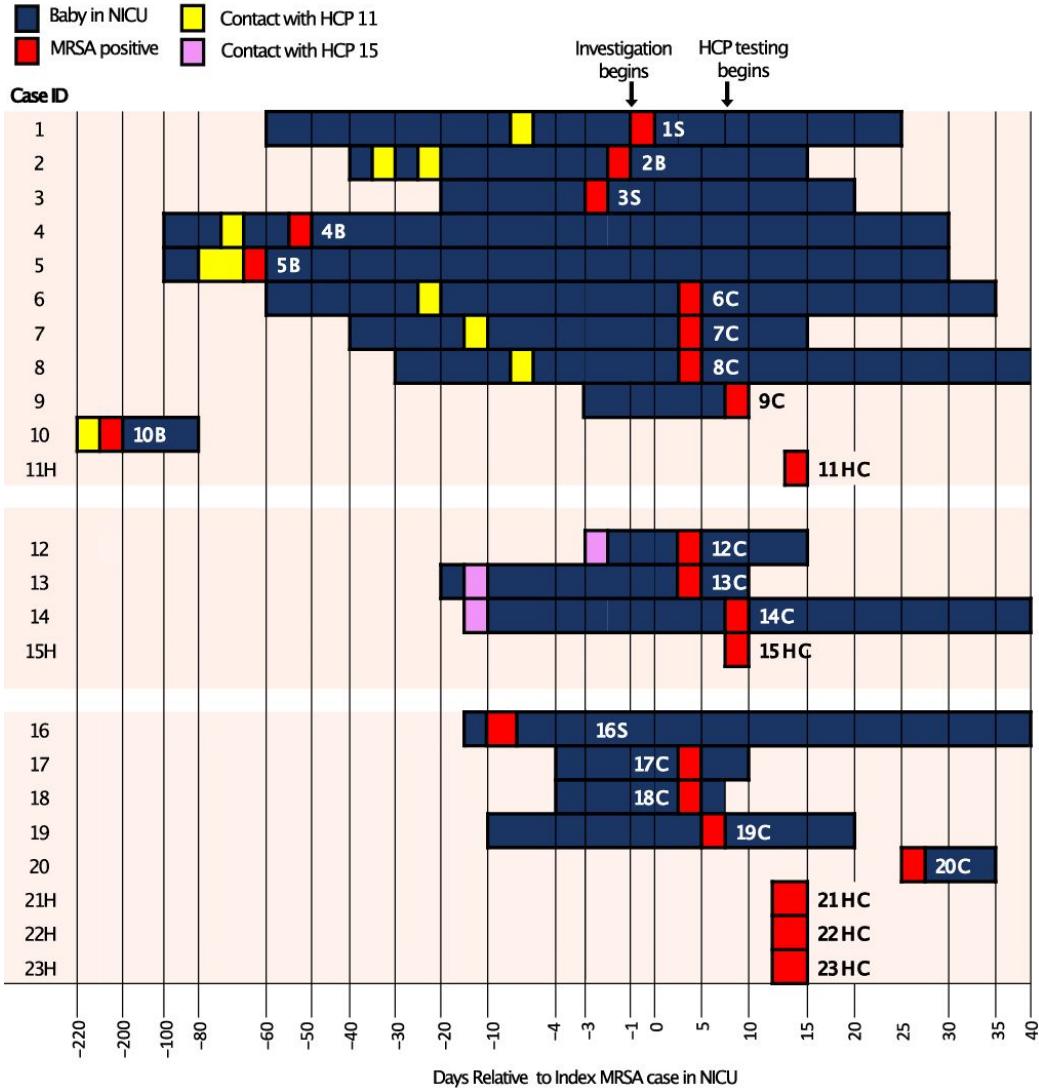
Genomics tells us there are 2 outbreaks!



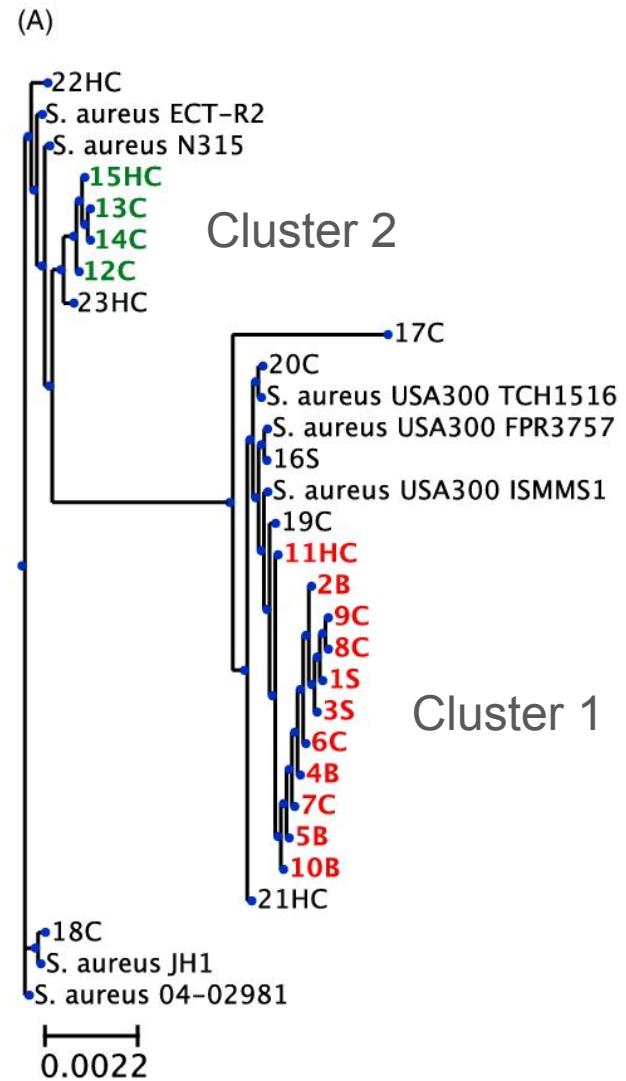
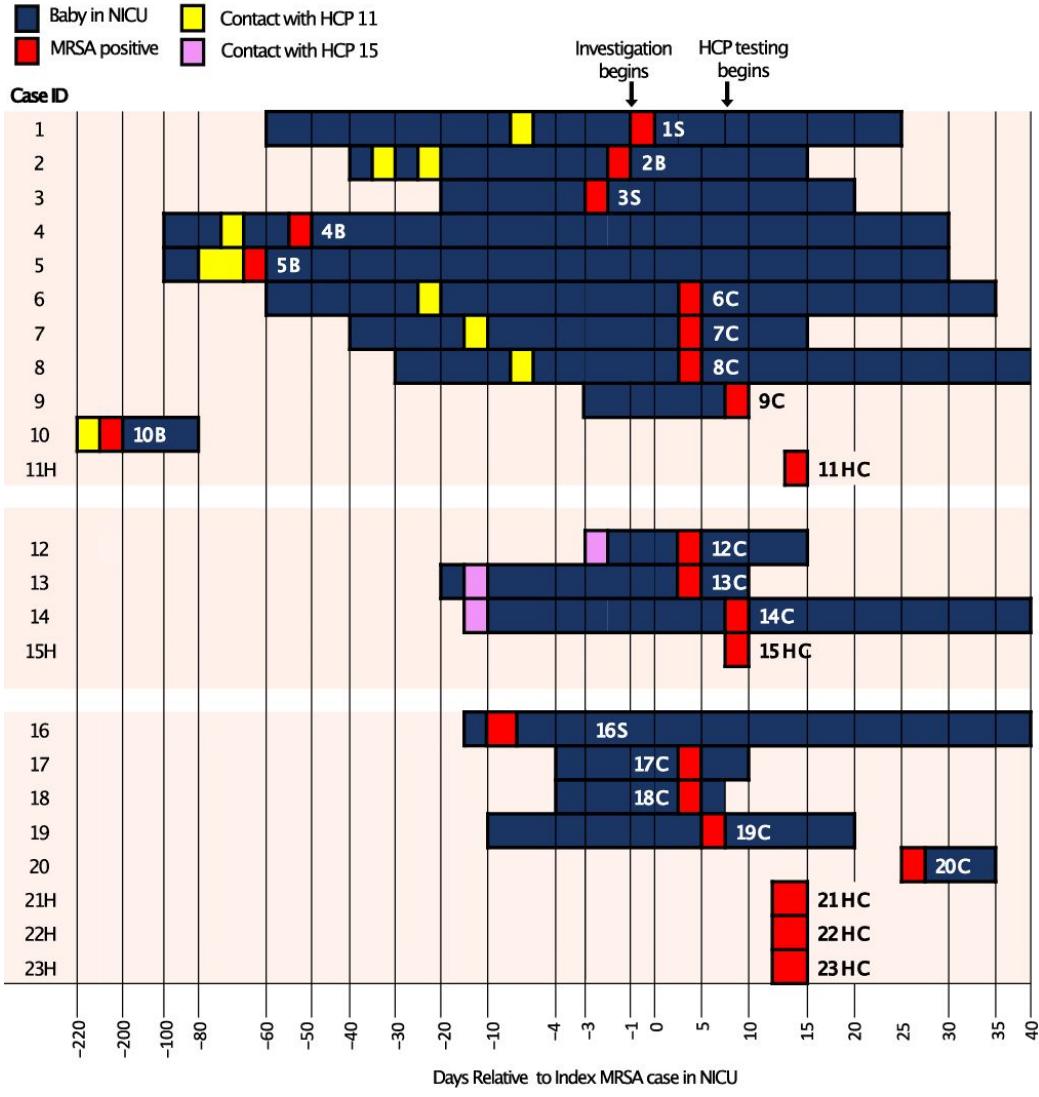
Genomics tells us there are 2 outbreaks!



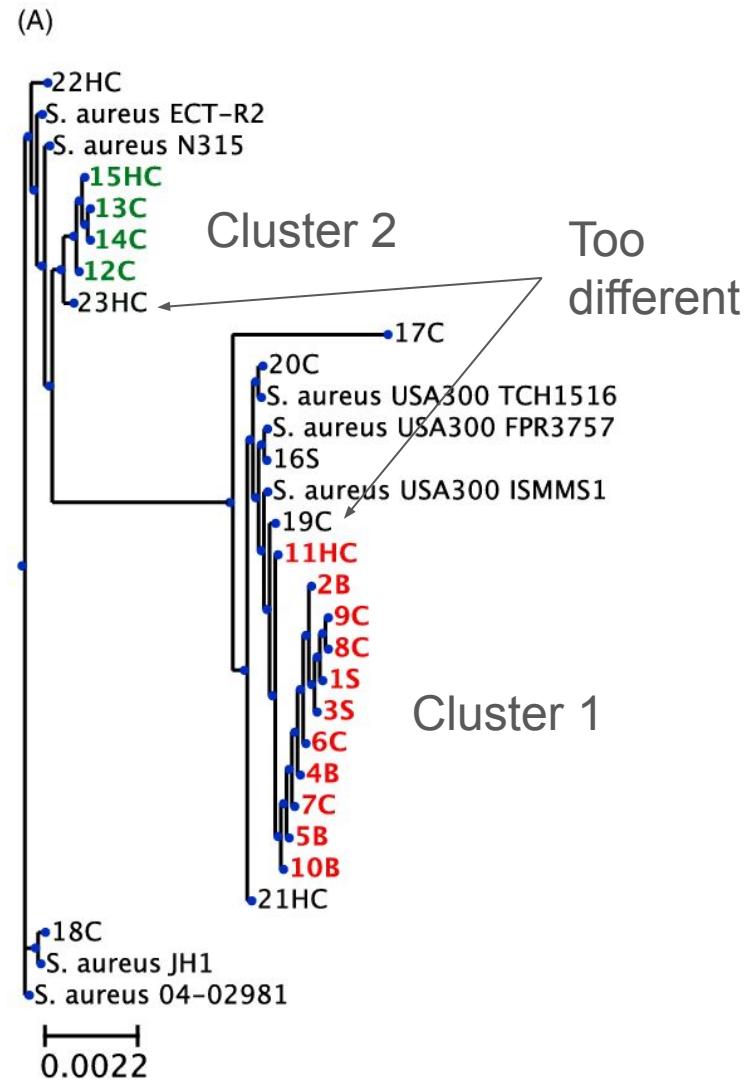
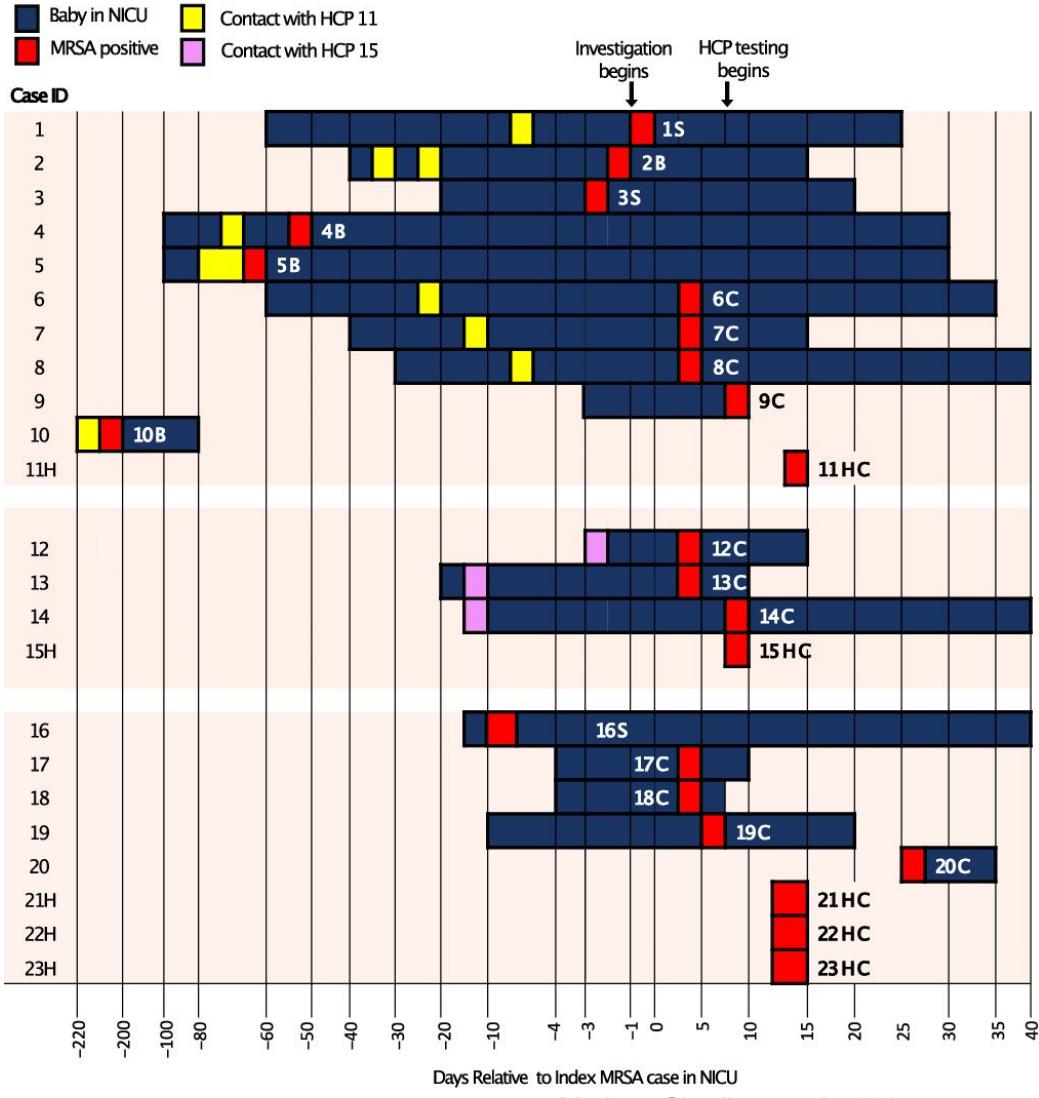
Genomics tells us there are 2 outbreaks!



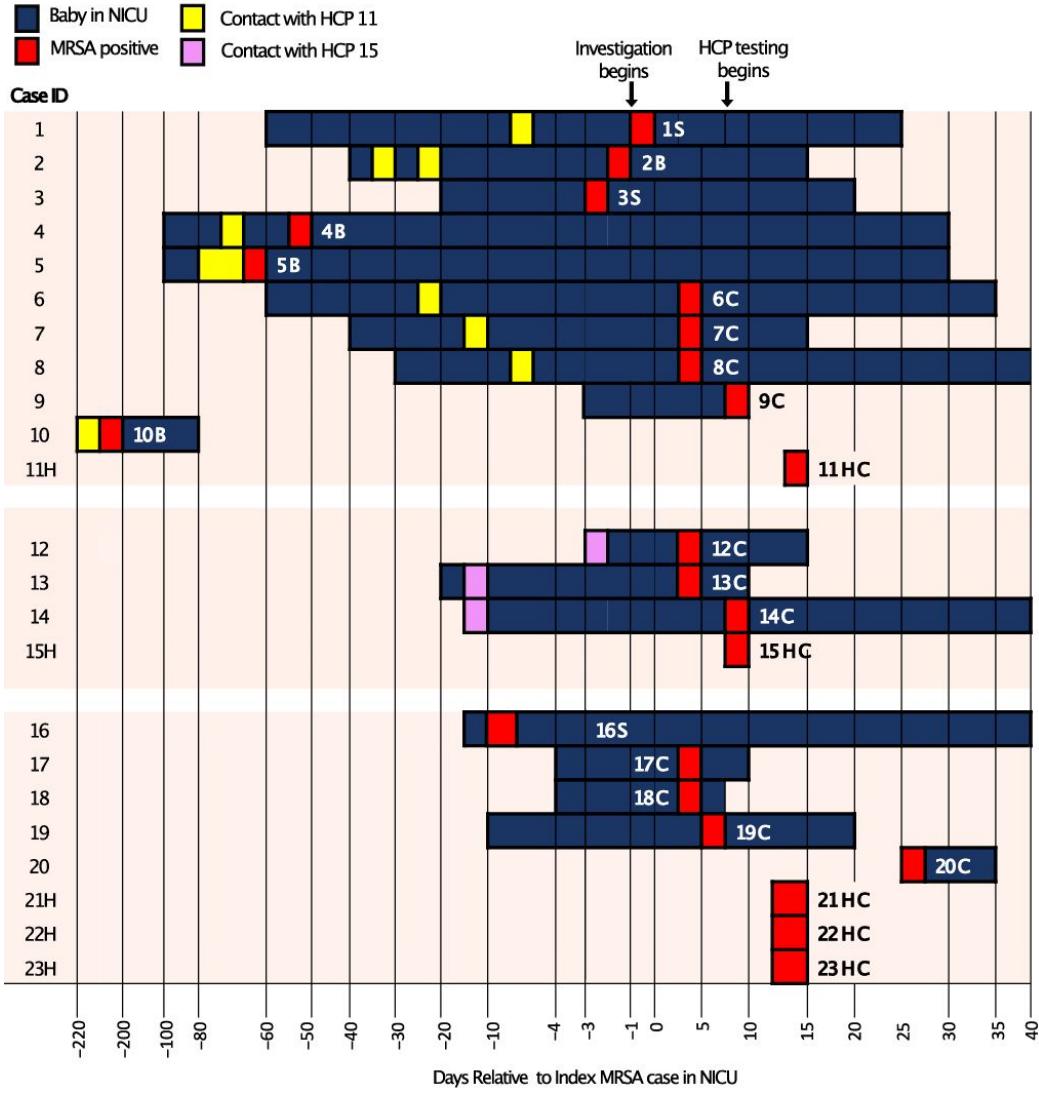
Genomics tells us there are 2 outbreaks!



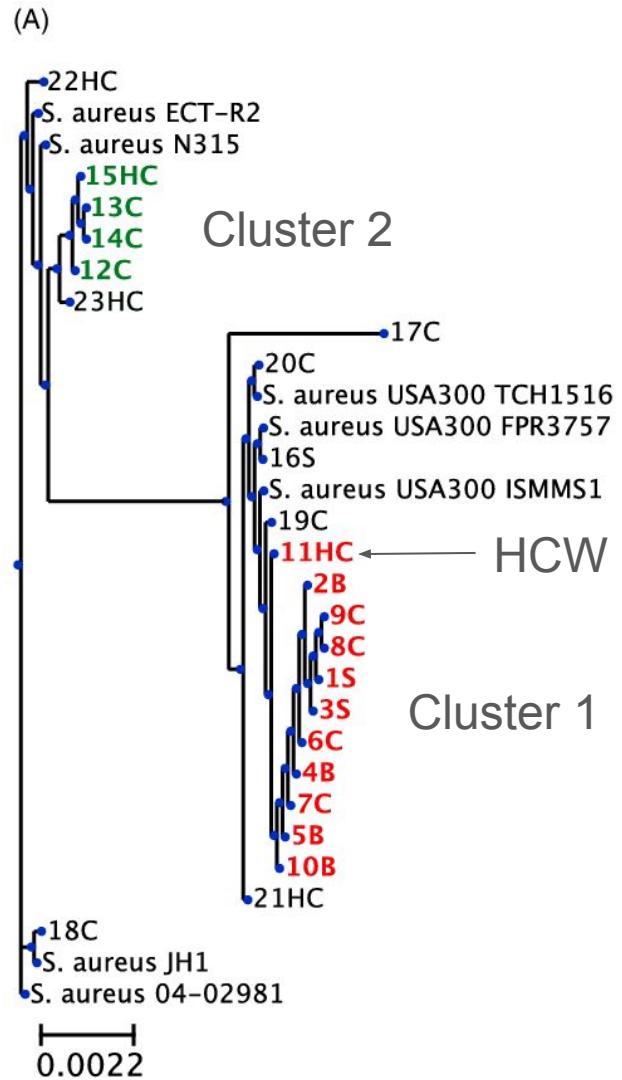
Genomics tells us there are 2 outbreaks!



Genomics tells us HCW11 likely source of Cluster 1



Madera, Sharline, et al. 2023

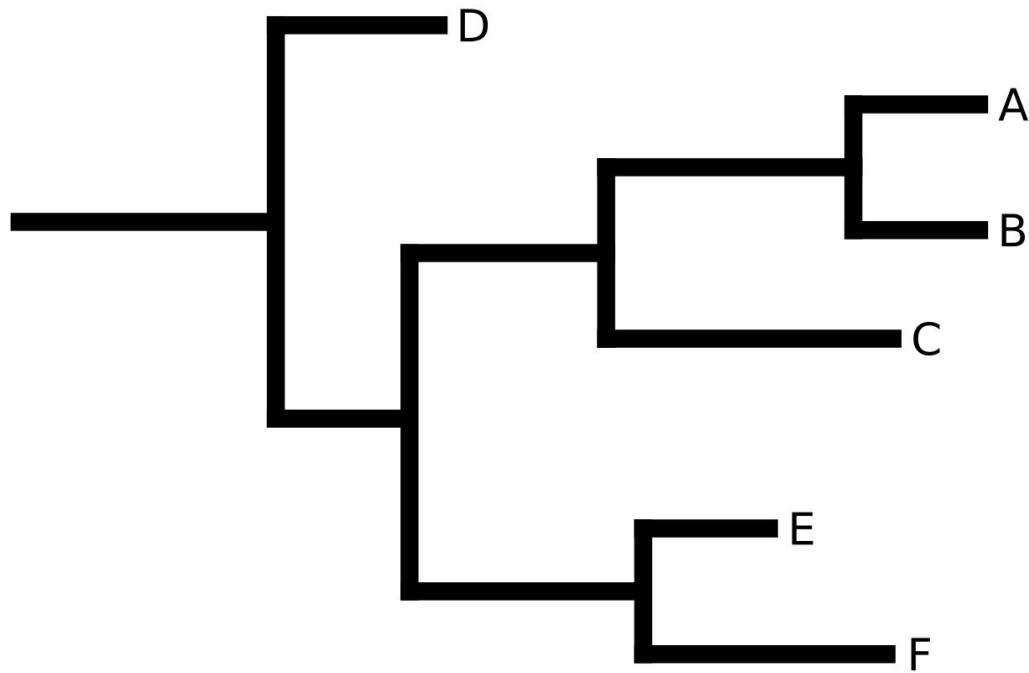


Conclusions

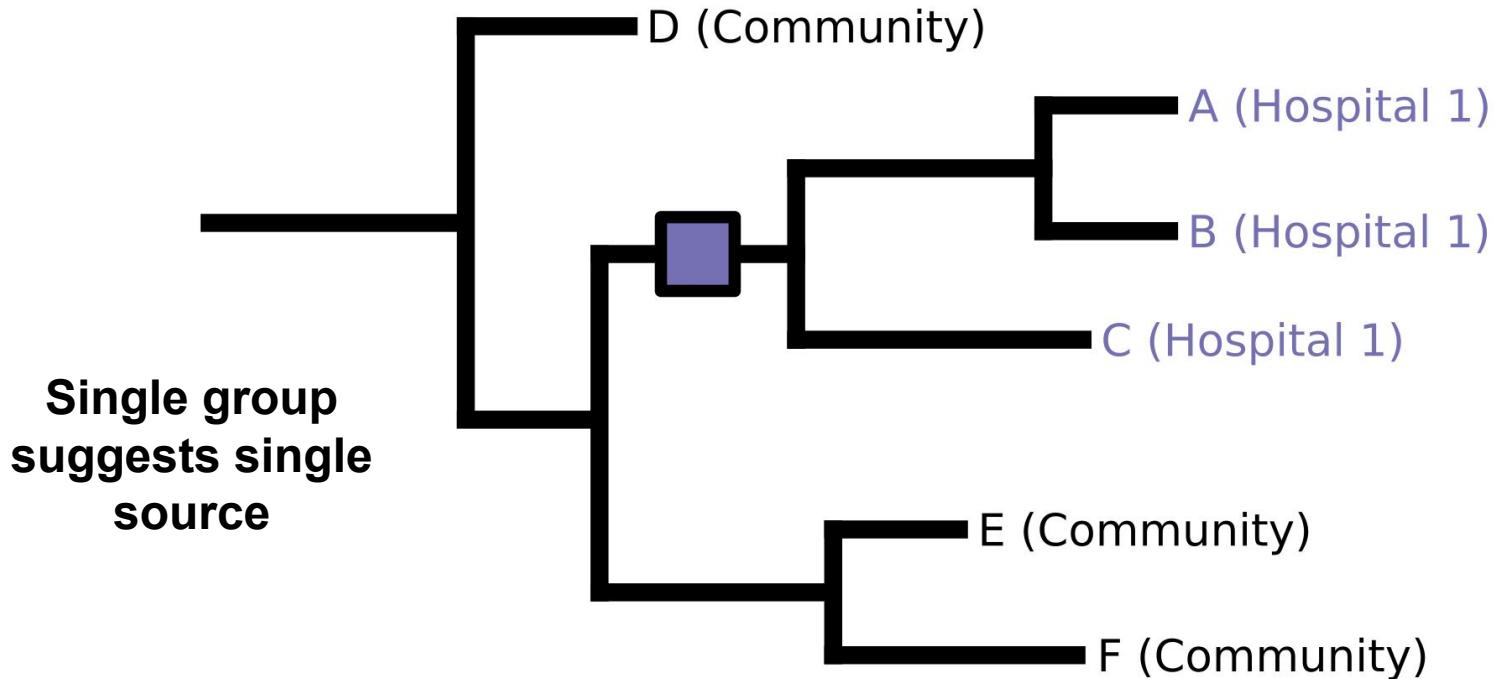
- Infection Prevention and Control (IPAC) increasingly incorporates genomics
- Clinical microbiology can involve differences in scale to normal lab work
- Sequencing involves randomly* sampling many short (2nd generation) or fewer long noisy (3rd generation) “reads” from DNA
- Sequencing is a physical process so involves measurement error
- Reference-based assembly: comparing reads to reference and trying to distinguish real changes from errors.
- *de novo* assembly: stitching together reads* which share sequences.
- Comparing genomes by alignment lets you find shared/different base pairs.
- Phylogenetic trees can be inferred from these patterns.
- Trees represent a sampling from the underlying evolving population.
- Genomes and trees can be used to track and stop outbreaks (among many other things!).

Extra Slides

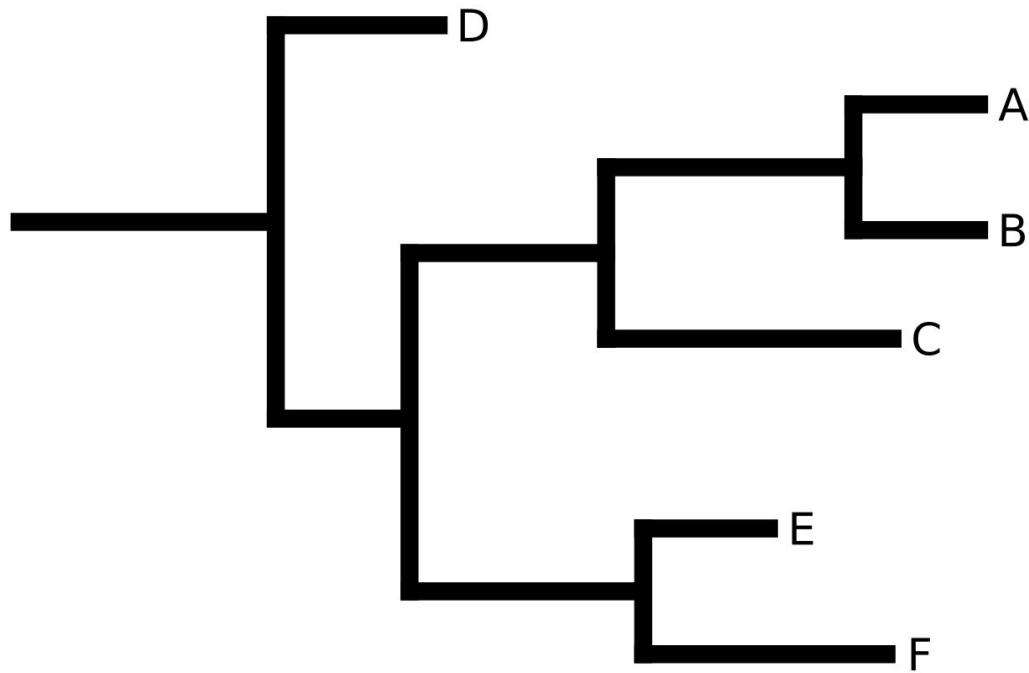
So, tree tells us about possible outbreak source



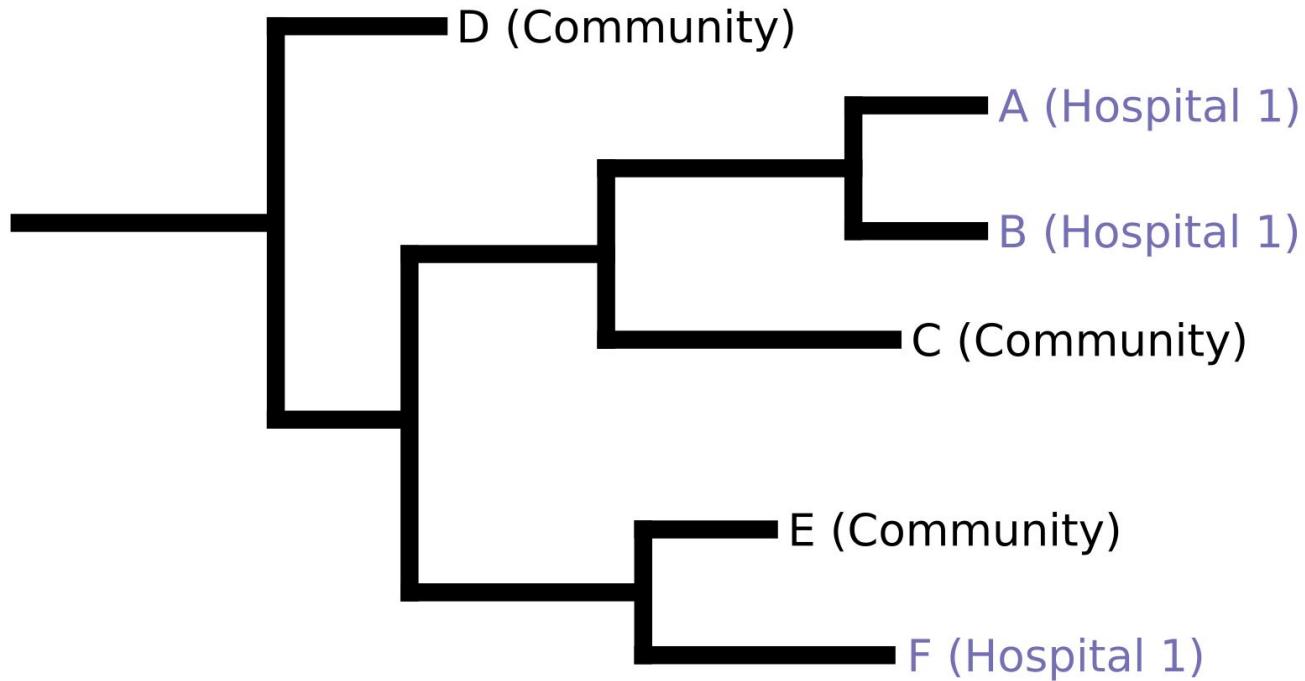
So, tree tells us about possible outbreak source



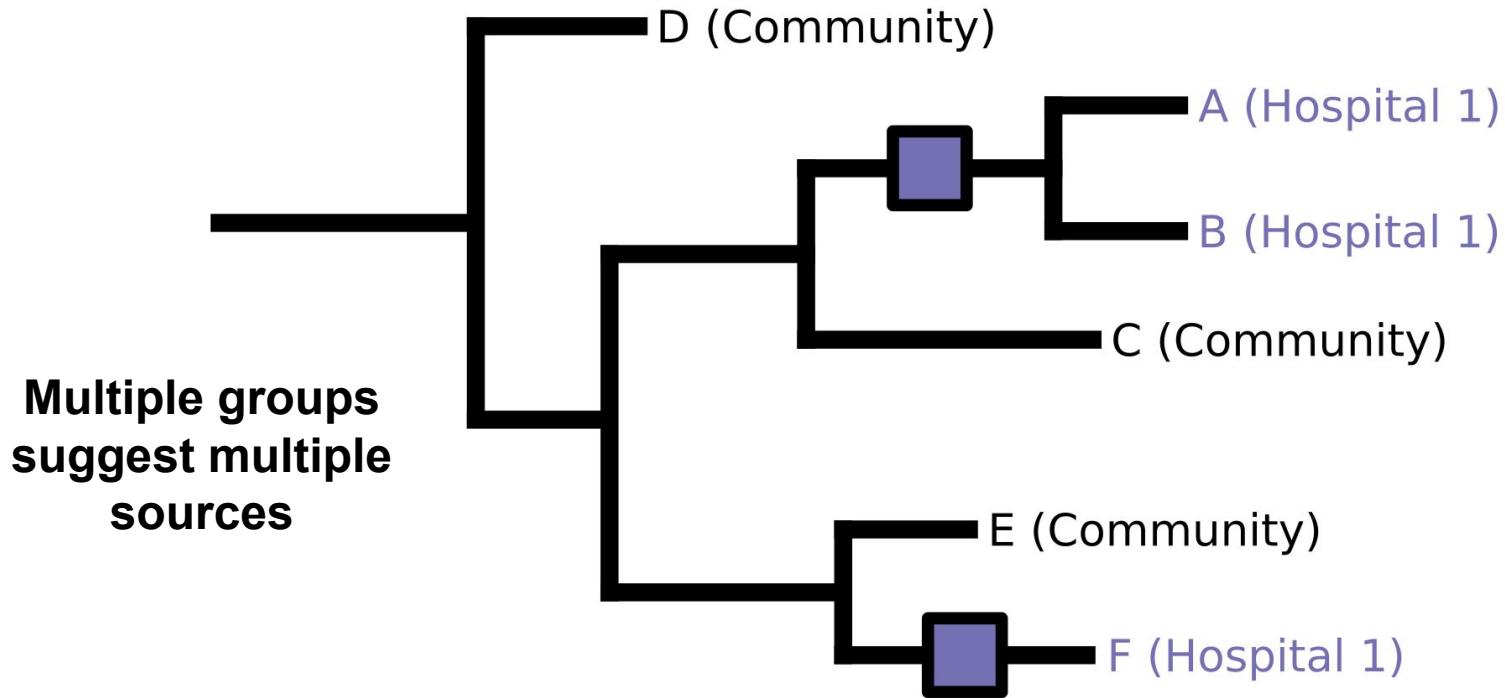
So, tree tells us about possible outbreak source



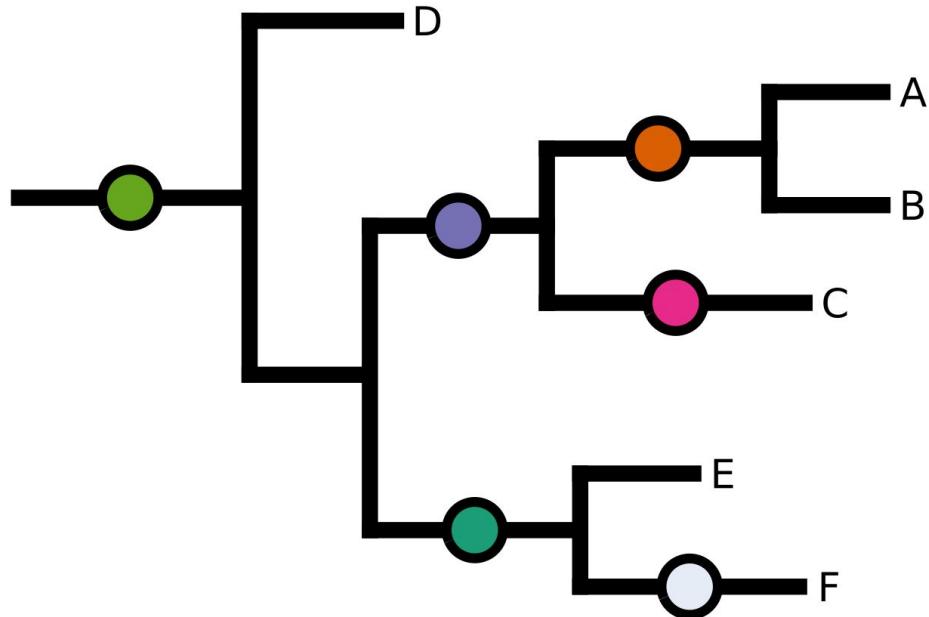
So, tree tells us about possible outbreak source



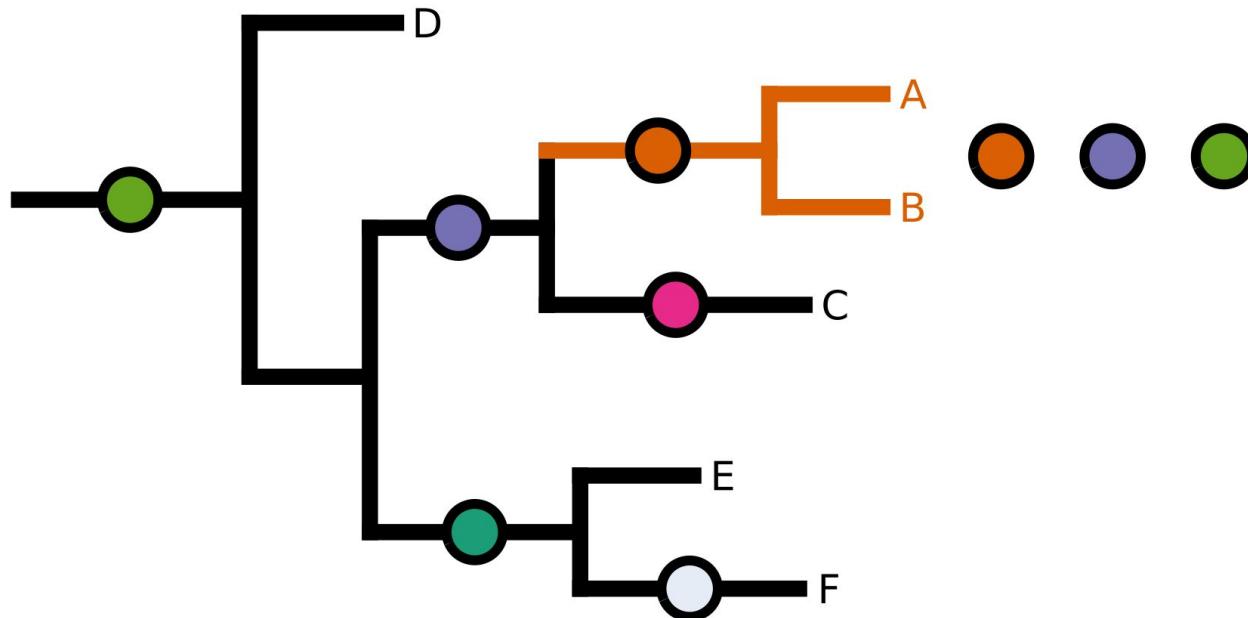
So, tree tells us about possible outbreak source



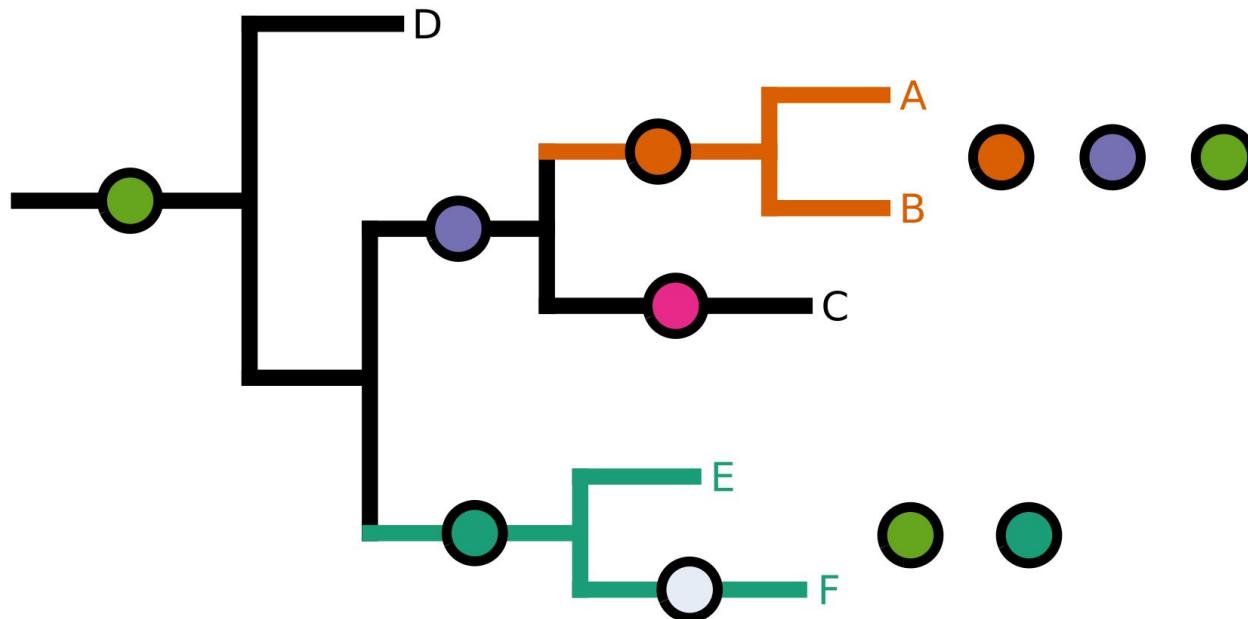
Define lineages (groups) of pathogens



Can also define lineages (groups) of pathogens



Can also define lineages (groups) of pathogens



Lineages, typing, and phenotyping

