



Tuberculosis Molecular Epidemiology: Deciphering Genotyping and Whole Genome Sequencing

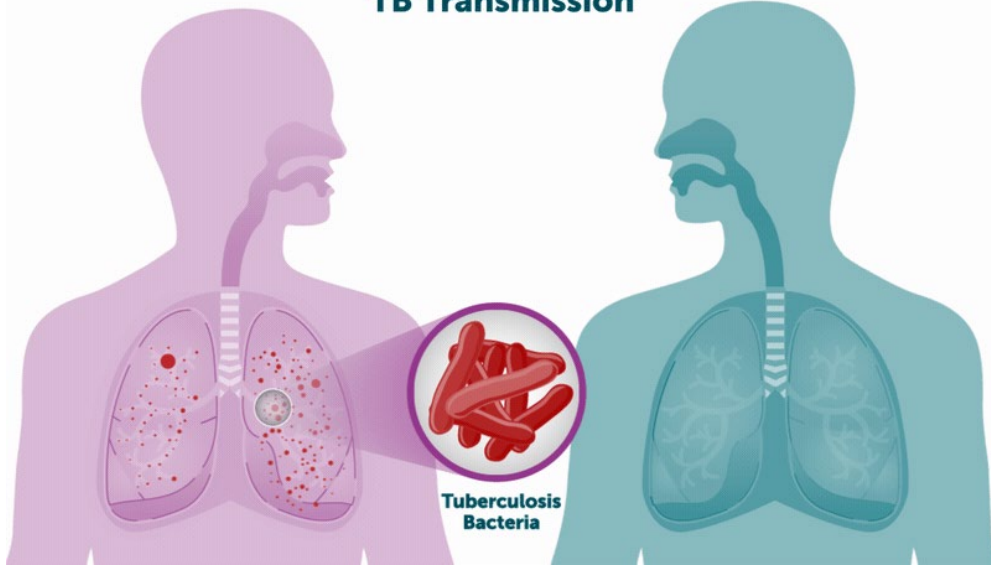
Sarah Talarico, PhD, MPH

Molecular Epidemiology and Outbreak Investigations Team
Surveillance, Epidemiology, and Outbreak Investigations Branch
Division of Tuberculosis Elimination

October 23, 2024

TB course of infection: latent TB infection and active TB disease

TB Transmission



Person with active pulmonary TB disease (infectious)




Active TB Disease

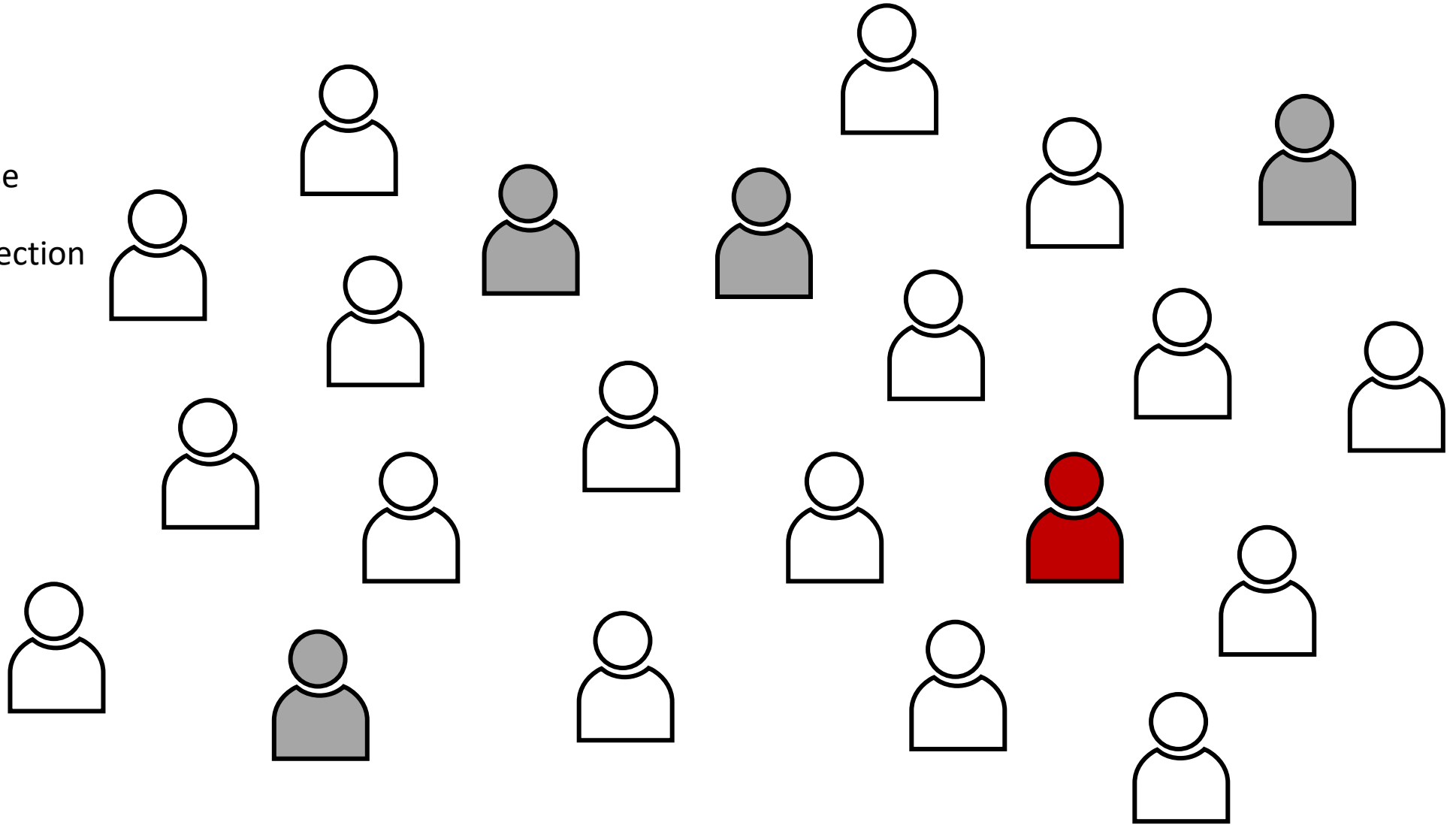
~ 5% within 2 years
(recent transmission)

~ 5% lifetime risk
(reactivation)




Latent TB Infection (not infectious)

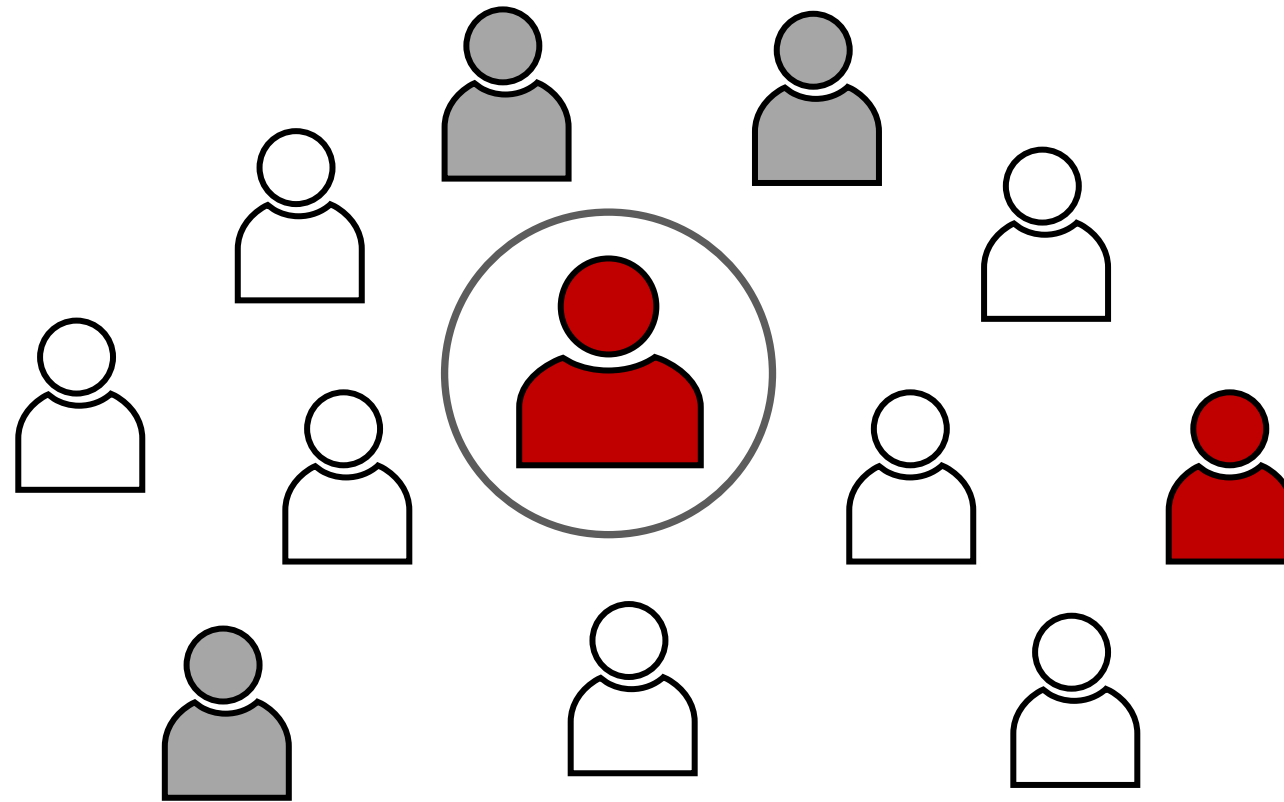
TB screening

-  Active TB case
-  Latent TB infection
-  TB negative






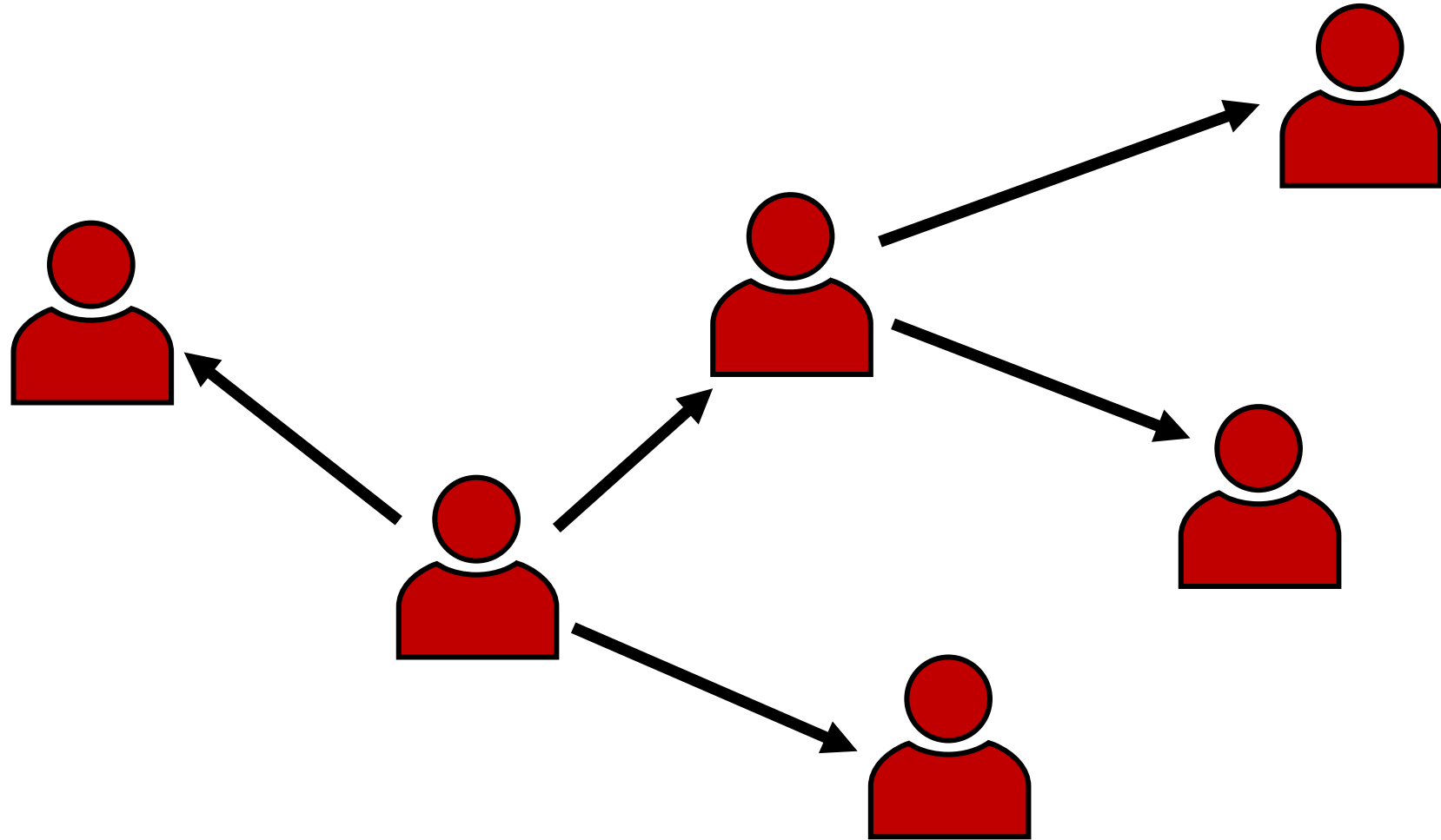
TB contact investigation

-  Active TB case
-  Latent TB infection
-  TB negative






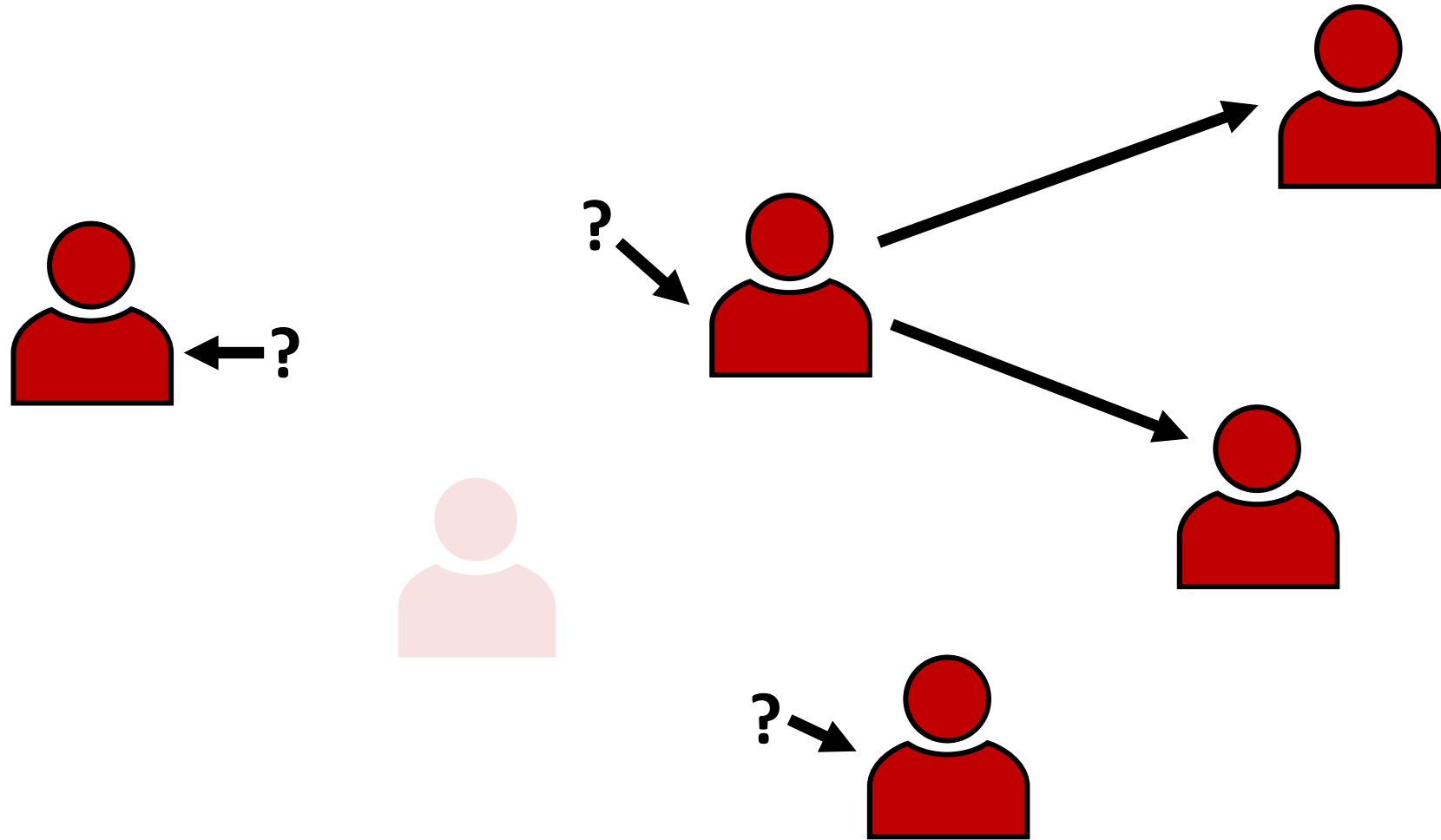
TB cluster or outbreak investigation

-  Active TB case
-  Latent TB infection
-  TB negative



TB cluster or outbreak investigation

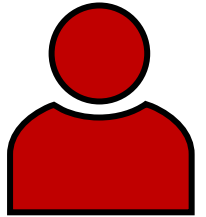
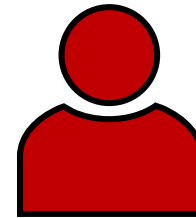
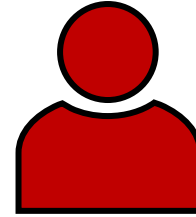
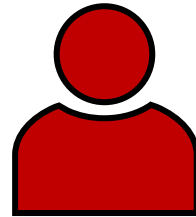
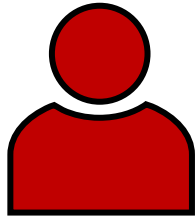
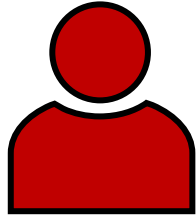
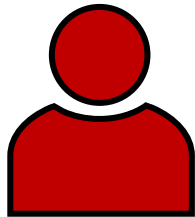
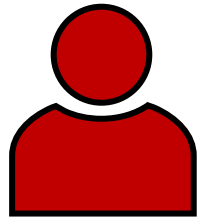
-  Active TB case
-  Latent TB infection
-  TB negative



How do we know that there is a TB cluster or outbreak occurring?

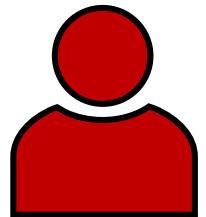
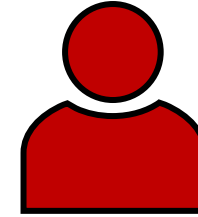
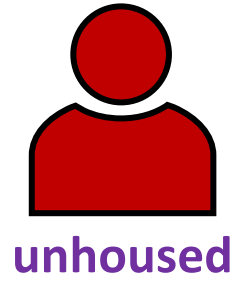
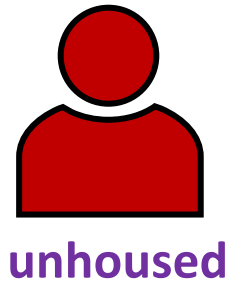
- Goal
 - Reduce the burden of TB by identifying where transmission is currently occurring and interrupting it
- Challenge
 - Distinguish recent transmission from cases infected long ago

How to identify the clusters when all we can see is cases?

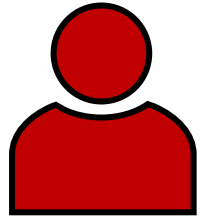
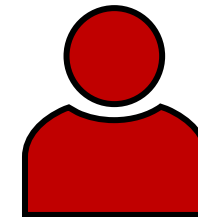
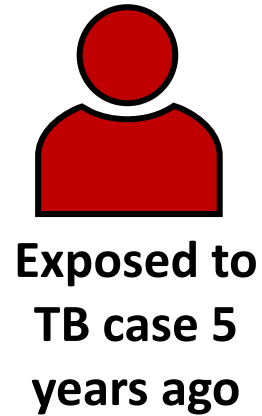
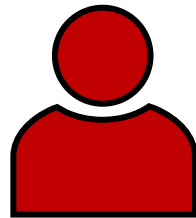


Time

Having some epi data can help...

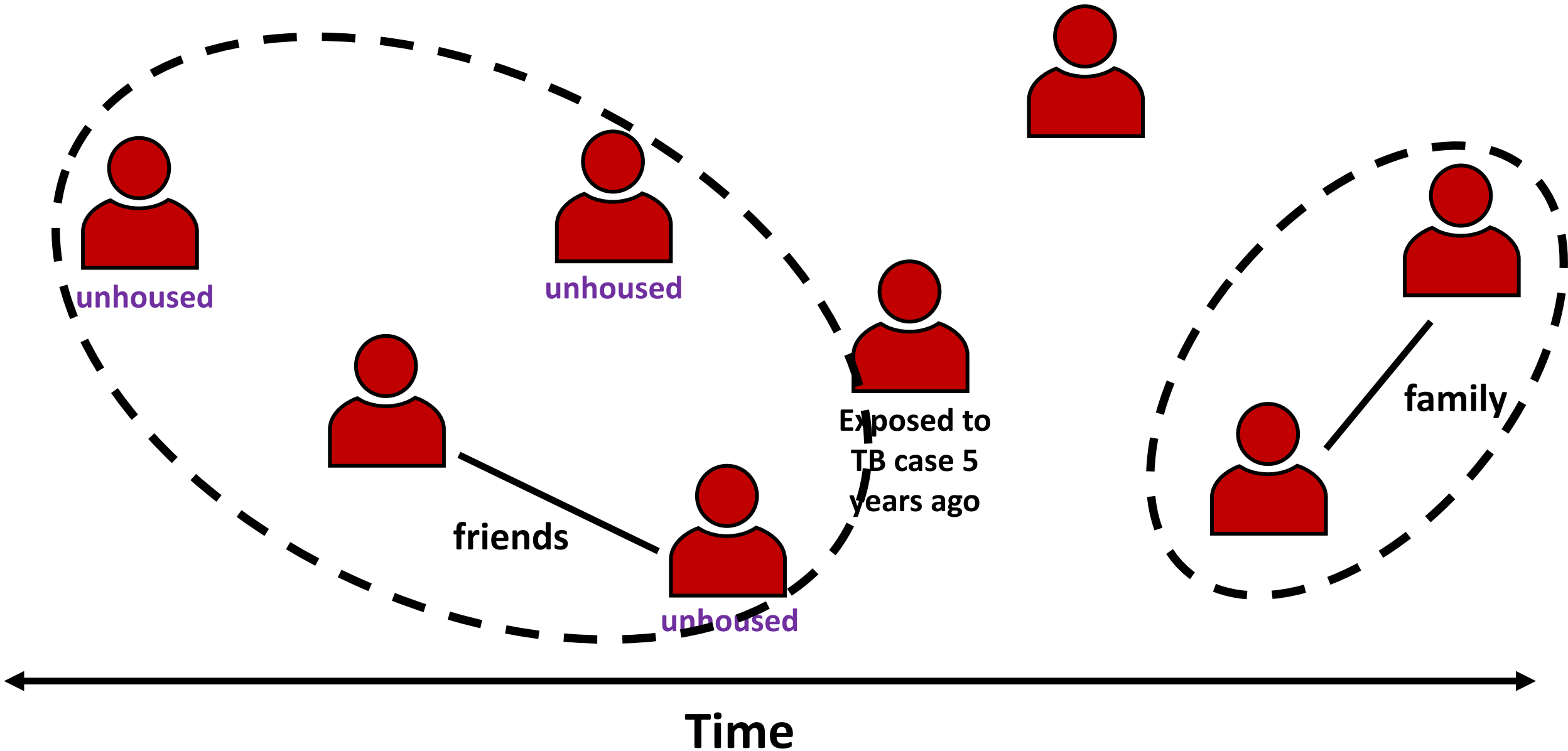


friends

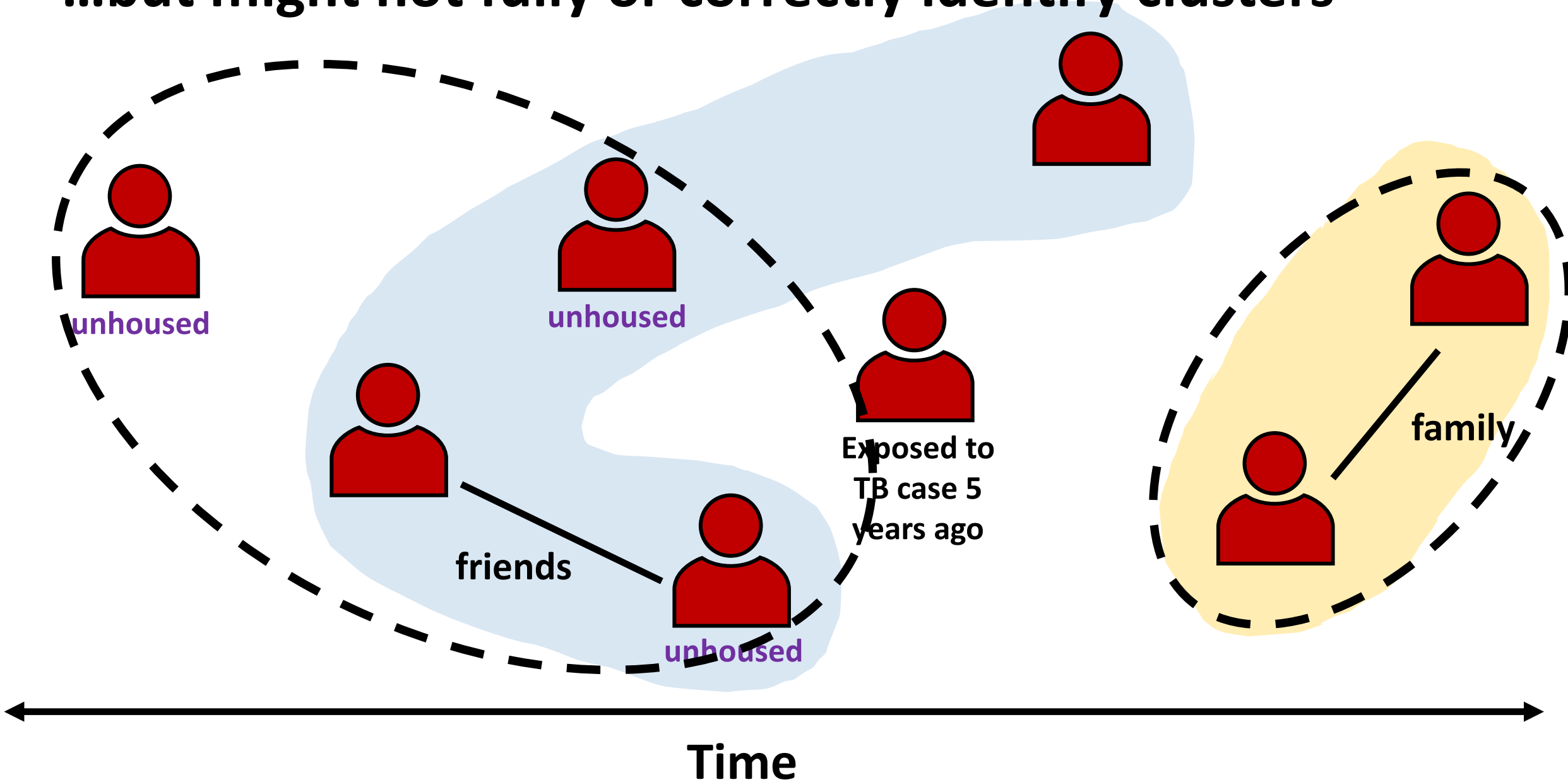


Time

Having some epi data can help...



...but might not fully or correctly identify clusters



Challenges to relying exclusively on epidemiologic investigation

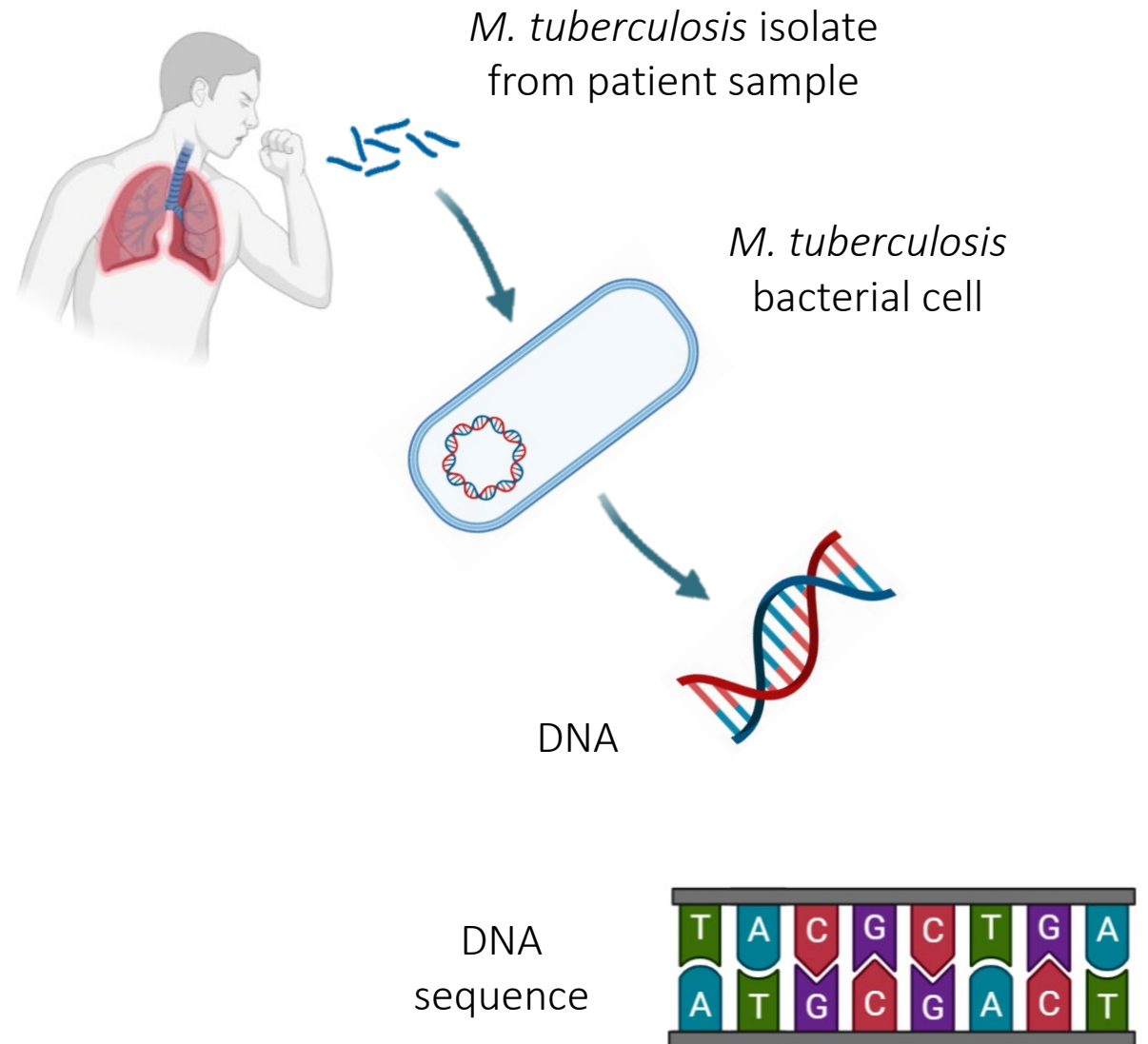
- Airborne transmission
- Exposure in congregate settings
- Long infectious periods
- Patient recall may be incomplete or unreliable
- Often in impoverished or marginalized communities

How do we know that there is a TB cluster or outbreak occurring?

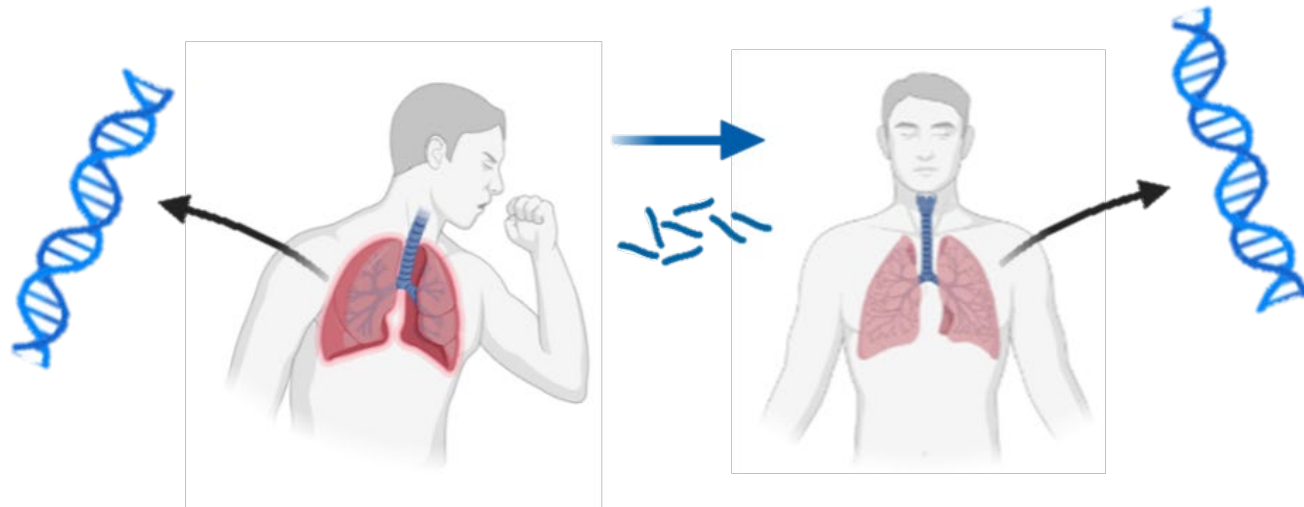
- Goal
 - Reduce the burden of TB by identifying where transmission is currently occurring and interrupting it
- Challenge
 - Distinguish recent transmission from cases infected long ago
- Approach
 - Molecular epidemiology!
 - Use molecular genotyping data, combined with clinical and epidemiologic data, to detect, investigate, and monitor recent TB transmission

TB genotyping examines the DNA of *M. tuberculosis* isolates from TB patients

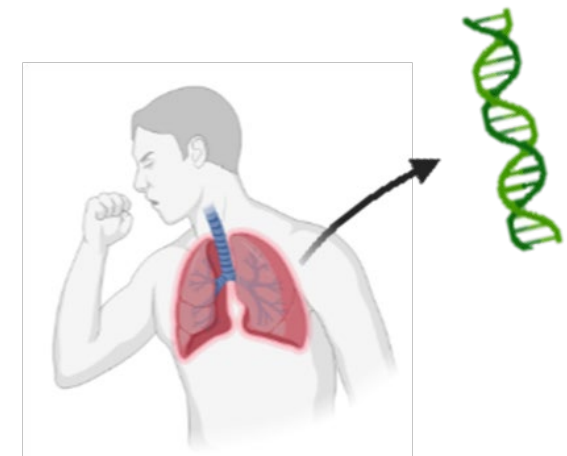
- The *M. tuberculosis* bacteria cultured from a TB patient is called the patient's isolate
- Bacteria, including *M. tuberculosis*, have DNA called a genome
- DNA is made up of four different nucleotides (abbreviated A, T, C, and G)
- The order of these nucleotides in the genome is the DNA sequence
- The genome of *M. tuberculosis* is over 4.4 million nucleotides long



Genotyping analyzes DNA to identify TB patients with similar *M. tuberculosis* genomes who are more likely to be linked by recent transmission

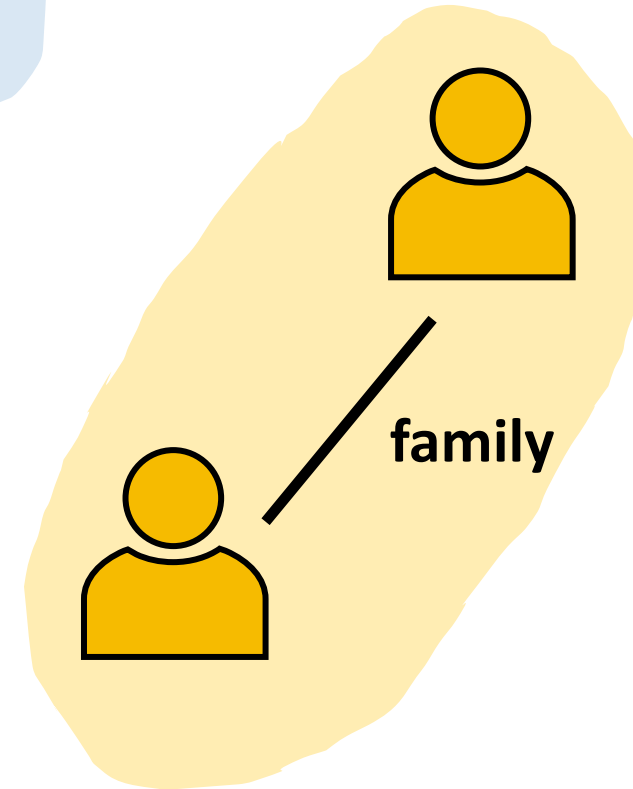
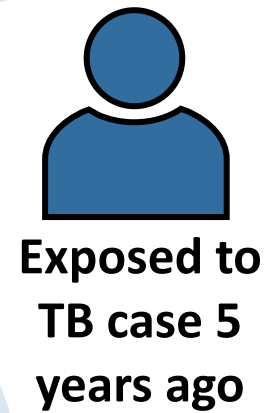
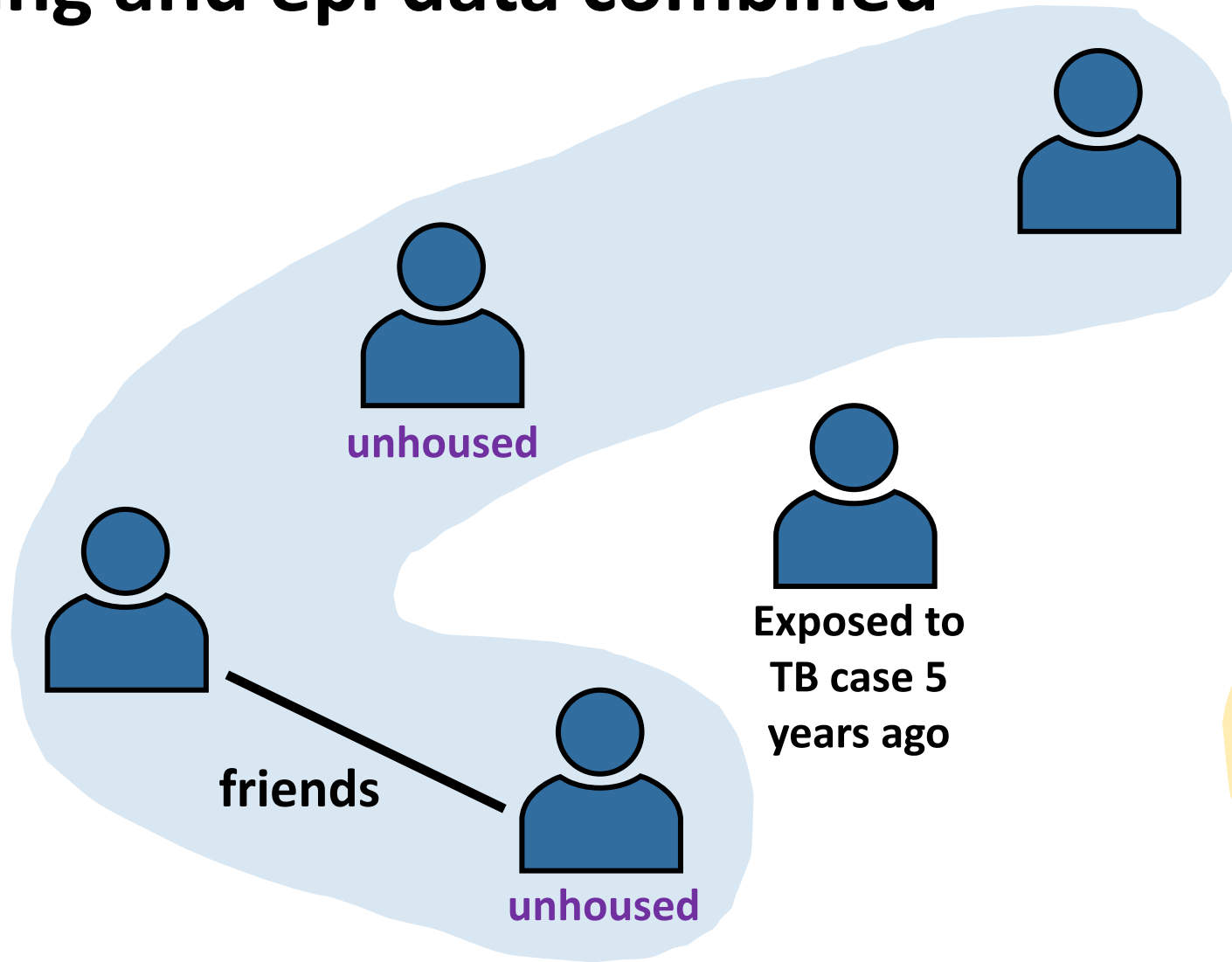
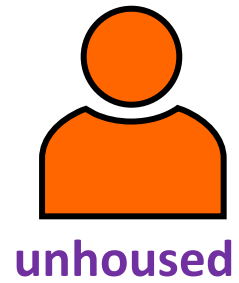


TB patients linked by recent transmission have isolates with the same genotype (blue)



TB patient not linked by recent transmission has an isolate with a different genotype (green)

Genotyping and epi data combined



TB Genotyping Methods

TB genotyping through the years

Conventional TB Genotyping

Whole-genome sequencing (WGS)

National TB
Genotyping Service
begins conventional
genotyping on all
Mtb isolates

Retrospective WGS
for select clusters
begins

National TB Molecular
Surveillance Center
begins prospective
WGS for all
Mtb isolates

Conventional
genotyping
discontinued

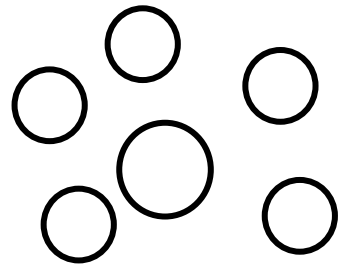
2004

2012

2018

July 2022

Whole-genome sequencing (WGS) of *M. tuberculosis*

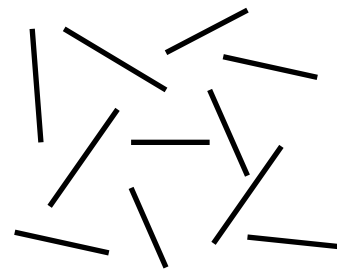


Genomic DNA

~ 4.4 million basepairs



Shear DNA



Genomic DNA fragments

~ 500 basepairs



Create library and sequence

```
@NB551186:40:H5TN5AFX:1:1110  
1:21172:1116 1:N:0:  
GGACTCCT+TATGCAGTTACGGAACC  
CAATCAGGTCCAAAGGTCTTCATCAA  
GGCGTCGGAAAGCACGTCGATAACA  
GCGTCGCTCTGTTGTTGGTTGGCTT  
+  
A/AAAEAAAAAAAAA6/EEEEAAE/A/  
/E/AEEE/EEAAE/EEEEEE/AEEAE/EEE  
//AEE6AEEA//<</<<E/EEE<<//
```

File of sequence reads

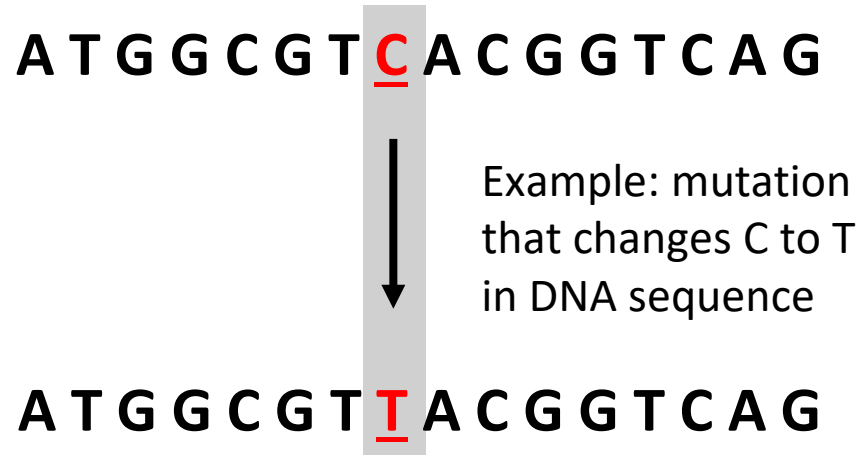


Data transfer to CDC



WGS data analysis

WGS data analysis focuses on a type of genetic variation called a single nucleotide polymorphism (SNP)

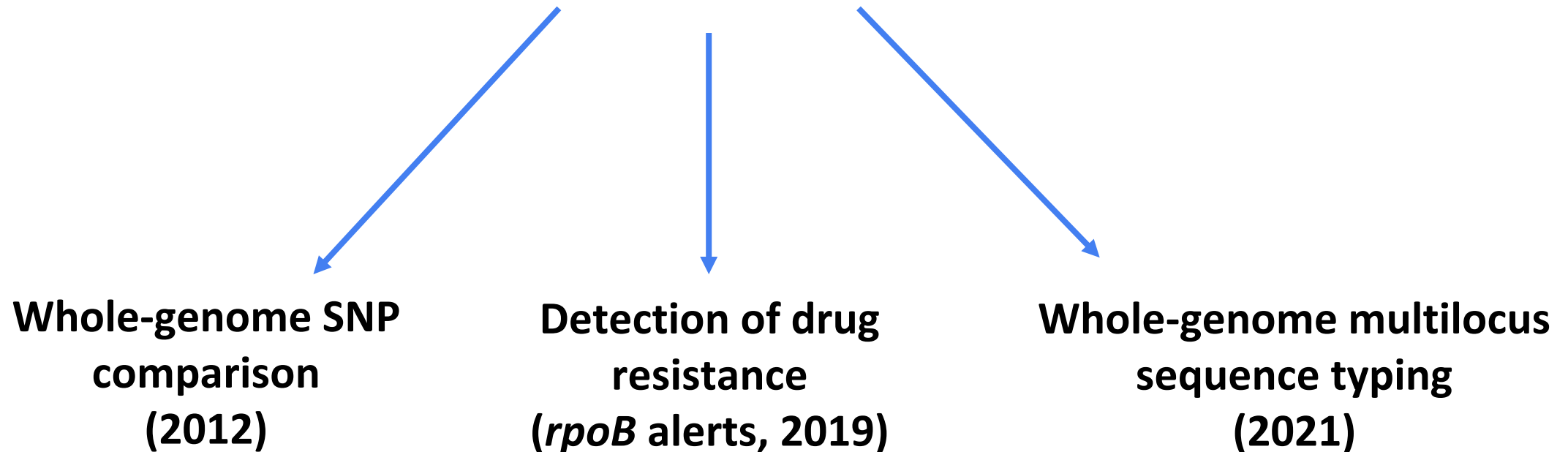


Single nucleotide
polymorphisms (SNPs)
throughout the genome

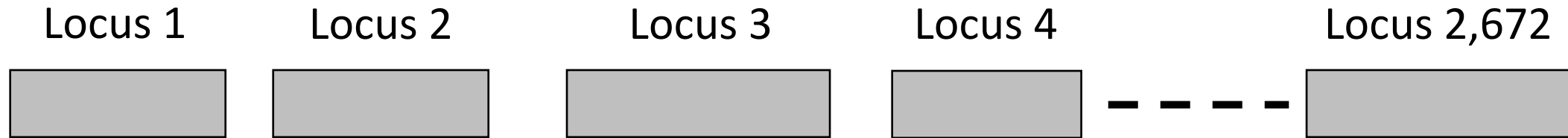
WGS data can be used for many different types of analyses

```
@NB551186:40:H5TN5AFX:1:1110  
1:21172:1116 1:N:0:  
GGACTCCT+TATGCAGTTACGGAACC  
CAATCAGGTCCAAGGCTTCATCAA  
GGCGTCGGAAAGCACGTCGATAACA  
GCGTCGCTCTGTTGTTGGTTGGCTT  
+  
A/AAAEAAAAAAAAA6/EEEEAAE/A/  
/E/AEEE/EEAAE/EEEEEE/AEEAE/EEE  
//AEE6EAEEA//<</<<E/EEE<<//
```

WGS data

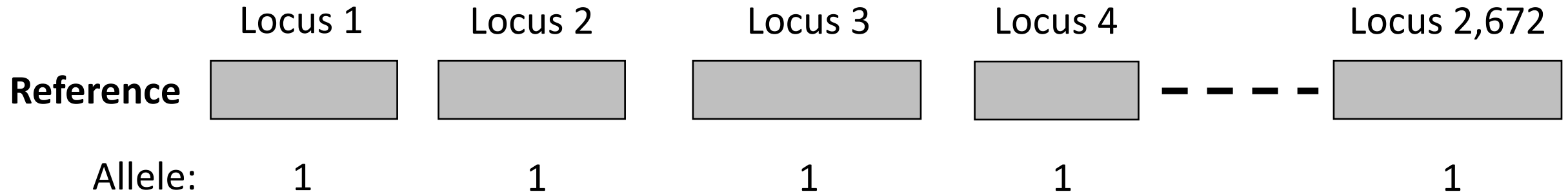


Whole-genome multilocus sequence typing (wgMLST)



- **Compares sequence at 2,672 loci throughout the genome**
 - Covers ~70% of the genome
- **Locus: location in the genome**
 - In this case, each locus is an individual gene

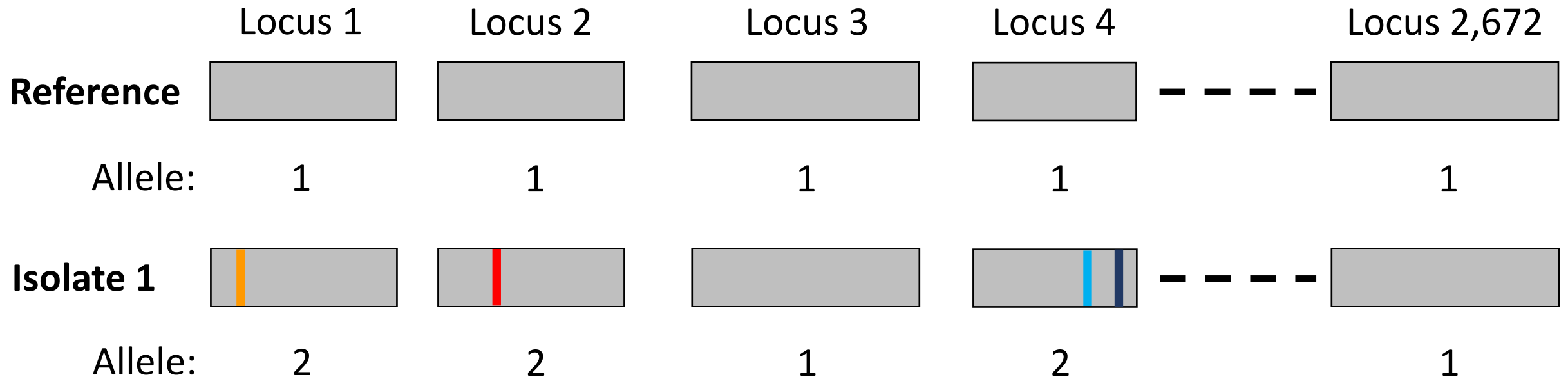
Whole-genome multilocus sequence typing (wgMLST)



Locus: location in the genome; in this case, each locus is an individual gene

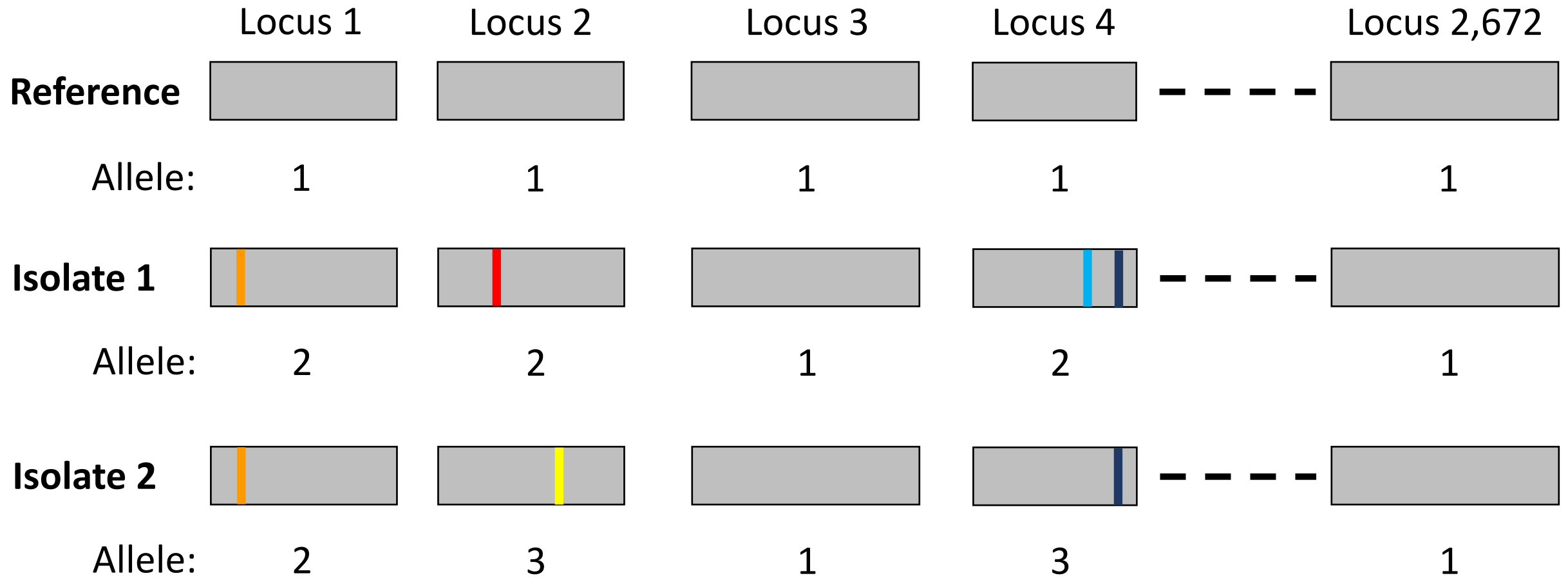
Allele: variant form of a gene

Whole-genome multilocus sequence typing (wgMLST)



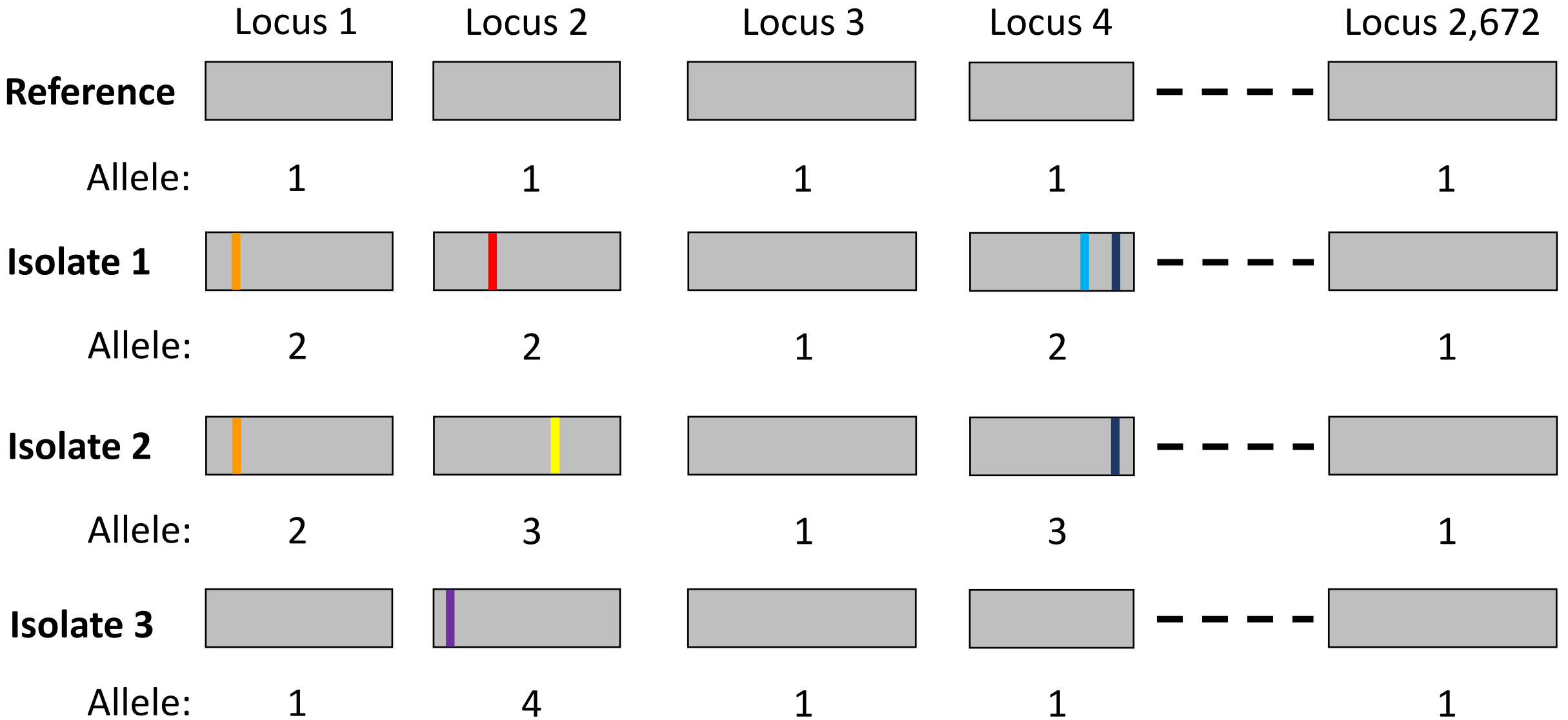
Colored lines indicate a single nucleotide difference compared to the reference

Whole-genome multilocus sequence typing (wgMLST)



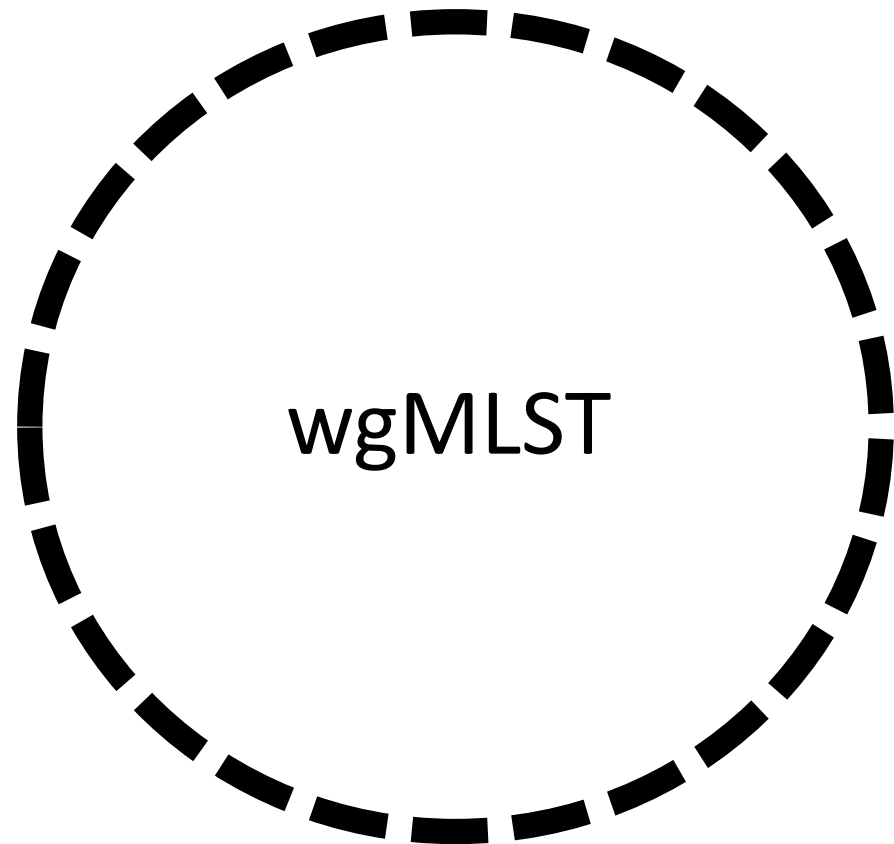
Colored lines indicate a single nucleotide difference compared to the reference

Whole-genome multilocus sequence typing (wgMLST)



Colored lines indicate a single nucleotide difference compared to the reference

Whole-genome multilocus sequence typing (wgMLST)



Compares sequence at 2,672 loci
(covers about 70% of genome)

$\geq 99.7\%$ of loci have
matching sequence



wgMLSType cluster
(example: MTBC123456)

$< 99.7\%$ of loci have
matching sequence



MTBCunique

wgMLSType results

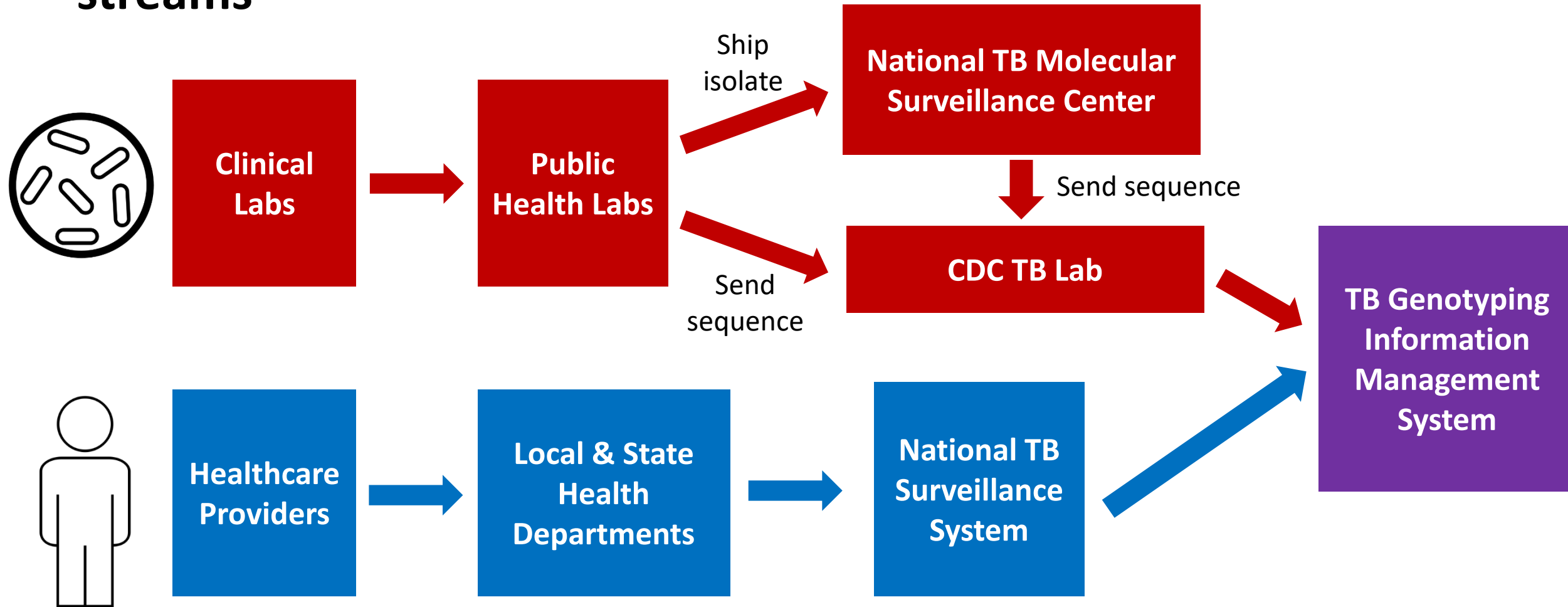
<i>Mtb</i> Isolate	wgMLSType
Isolate 1	MTBC099909
Isolate 2	MTBC099909
Isolate 3	MTBC099909
Isolate 4	MTBC099909
Isolate 5	MTBC099909
Isolate 6	MTBCunique
Isolate 7	MTBC000015
Isolate 8	MTBC000015

wgMLSType results

<i>Mtb</i> Isolate	wgMLSType	
Isolate 1	MTBC099909	Cluster 1: MTBC099909
Isolate 2	MTBC099909	
Isolate 3	MTBC099909	
Isolate 4	MTBC099909	
Isolate 5	MTBC099909	
Isolate 6	MTBCunique	> 7 SNPs from any other isolate nationally
Isolate 7	MTBC000015	Cluster 2: MTBC000015
Isolate 8	MTBC000015	

Combining TB genotyping data with TB surveillance data

TB Genotyping Information Management System (TB GIMS): combines genotyping and patient data flowing in from two streams



- Search ▼
- Genotype Results
- Patient Results
- Records** ▼
- Submitted Isolates
- Submitted SRS
- New Isolates
- Edit Isolates
- Submit Isolates
- Find Duplicates
- Import Data
- wgMLSType Change
- Reports and Tools** ▼
- Watch List
- Cluster Snapshot
- Generate Reports
- Templates
- Export Data
- Recent Transmission
- Custom Cluster List
- Alert Tracking List
- LITT Analysis Tools** ▼
- Build LITT Case Data File
- Run LITT

Tuberculosis Genotyping Information Management System

The last TB GIMS Surveillance Upload includes data transmitted to CDC through: **02/21/2023**

Searches and reports will only include data reported to CDC by the state and included in the latest TB GIMS surveillance upload.

Announcements:

No New Announcements.


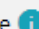
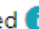


State: ALL ▼ Submit

Genotyping Surveillance Coverage

Year	2020	2021	2022	2023 *
National (%)	98.0	97.8	88.9	19.7

*Year to date. NA-Not Available. Source: NTIP Data Processed Weekly

Timeliness of Genotyping - by Isolate

Time From → To	Median number of days		Goal (days)
	National		
	2022	2023 *	
Specimen collection → Isolate shipped to genotyping lab 	54	65	NA
Receipt at genotyping lab → Genotype create date 	9	12	14
Genotype create date → State Case No. entered 	0	1	56
Genotype create date → Isolate Linked 	2	5	90
Specimen collection → Isolate Linked 	79	97	NA

Search

Genotype Results

Patient Results
 Records

Submitted Isolates

Submitted SRS

New Isolates

Edit Isolates

Submit Isolates

Find Duplicates

Import Data

wgMLSType Change

 Reports and Tools

Watch List

Cluster Snapshot

Generate Reports

Templates

Export Data

Recent Transmission

Custom Cluster List

Alert Tracking List


 LITT Analysis Tools

Patient Results







To view surveillance data on TB patients, enter the search criteria then click **Find**

For more information, refer to **Patient Results** in the online help.

 Information * Required

Basic Options

State *	County	Region	
<input type="text"/>	<input type="text"/>	<input type="text"/>	
wgMLSType 	GENType 	PCRTYPE	Cluster Name
<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>
State Case #	Submitter #	Accession #	Cluster Name2
<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>
Spoligotype	MIRU	MIRU2	
<input type="text"/>	<input type="text"/>	<input type="text"/>	
Date Type	Start Date	End Date	
<input type="text"/>	MM/DD/YYYY 	MM/DD/YYYY 	

Advanced Options

 Find

 Clear

 Create Watch List Item


Search Results

Displaying 1 to 7 of 7 Records

<input type="checkbox"/>	CCID	STATE	STATE CASE #	Accession #	RPTDATE	CNTDATE	SEX	RACEHISP	AGE3	CITY	COUNTY	ZIPCODE	ORIGIN	CNTRYLN	YRSIN_US2
<input type="checkbox"/>		ST	2023ST000000001	23RF0001	01/26/2023	01/14/2023	F	WHITE	65+	ANYTOWN	COUNTY A		USBORN		
<input type="checkbox"/>		ST	2023ST000000002	23RF0002	01/11/2023	01/14/2023	M	HISP	65+	ANYTOWN	COUNTY A		USBORN		
<input type="checkbox"/>		ST	2023ST000000003	23RF0003	01/10/2023	01/14/2023	M	WHITE	45-64	TOWNSVILLE	COUNTY A		USBORN		
<input type="checkbox"/>		ST	2023ST000000004	23RF0004	01/10/2023	02/04/2023	F	BLACK	65+	METRO HEIGHTS	COUNTY A		USBORN		
<input type="checkbox"/>		ST	2023ST000000005	23RF0005	01/09/2023	01/14/2023	F	BLACK	25-44	ANYTOWN	COUNTY A		USBORN		
<input type="checkbox"/>		ST	2023ST000000006	23RF0006	01/03/2023	01/07/2023	M	ASIAN	45-64	TOWNSVILLE	COUNTY A		USBORN		
<input type="checkbox"/>		ST	2023ST000000007	23RF0007	12/20/2022	01/14/2023	F	MULT	45-64	ANYTOWN	COUNTY A		USBORN		

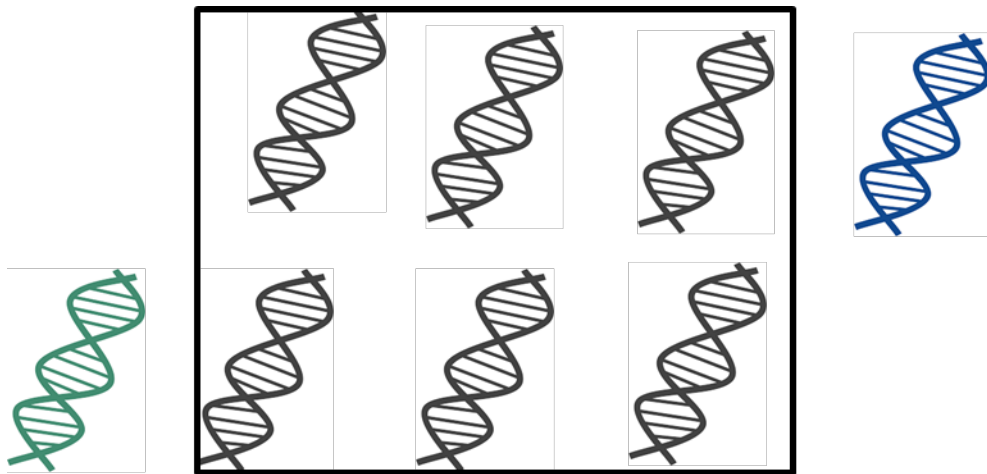
Some of the functions in TB GIMS

- Create isolate records before shipping to genotyping lab
- Search and view genotyping and patient results
- Link isolate to patient surveillance record
- Run reports and summaries of clusters
- Create watch lists to be notified of new cases in a cluster
- Run and track molecular surveillance algorithms
 - Weekly cluster alerts
 - Large outbreak surveillance
 - Yearly estimates of recent transmission
- Request additional WGS analysis (whole-genome SNP comparison)

Whole-genome SNP comparison

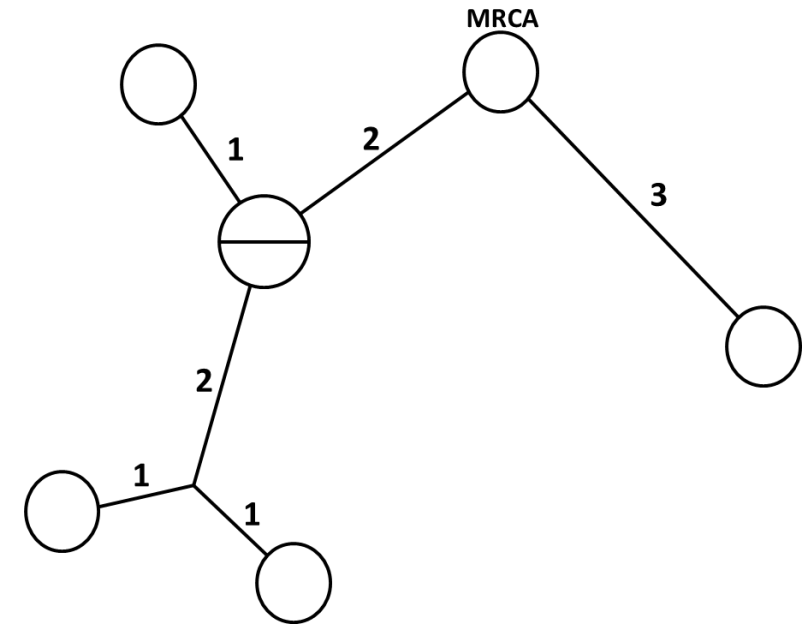
Isolates in a genotype cluster can be further differentiated by whole-genome SNP (wgSNP) comparison

Step 1. Detect a genotype cluster



Whole-genome
multilocus sequence
typing

Step 2. Examine genetic relationship among isolates in the cluster



Whole-genome single nucleotide
polymorphism comparison

Example: wgMLSType versus wgSNP comparison

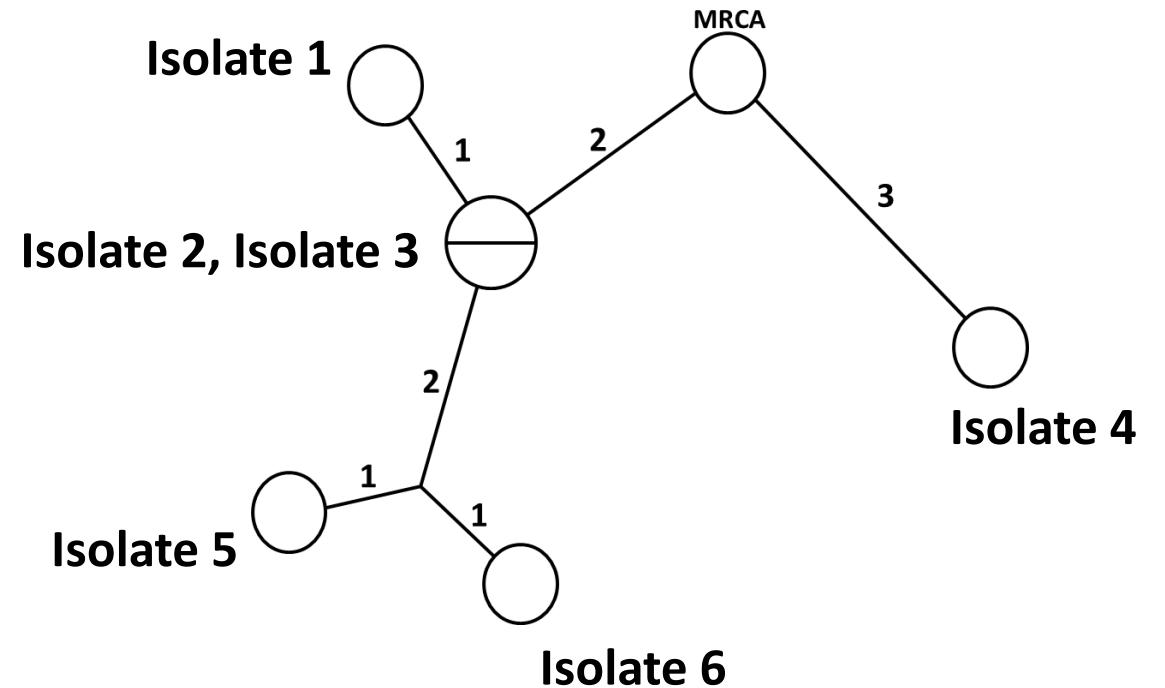
Step 1. wgMLSType

Cluster of six isolates with the same genotype MTBC001089 is detected

Isolate 1	MTBC001089
Isolate 2	MTBC001089
Isolate 3	MTBC001089
Isolate 4	MTBC001089
Isolate 5	MTBC001089
Isolate 6	MTBC001089

Step 2. wgSNP

Six MTBC001089 isolates are further differentiated and genetic relationships are examined

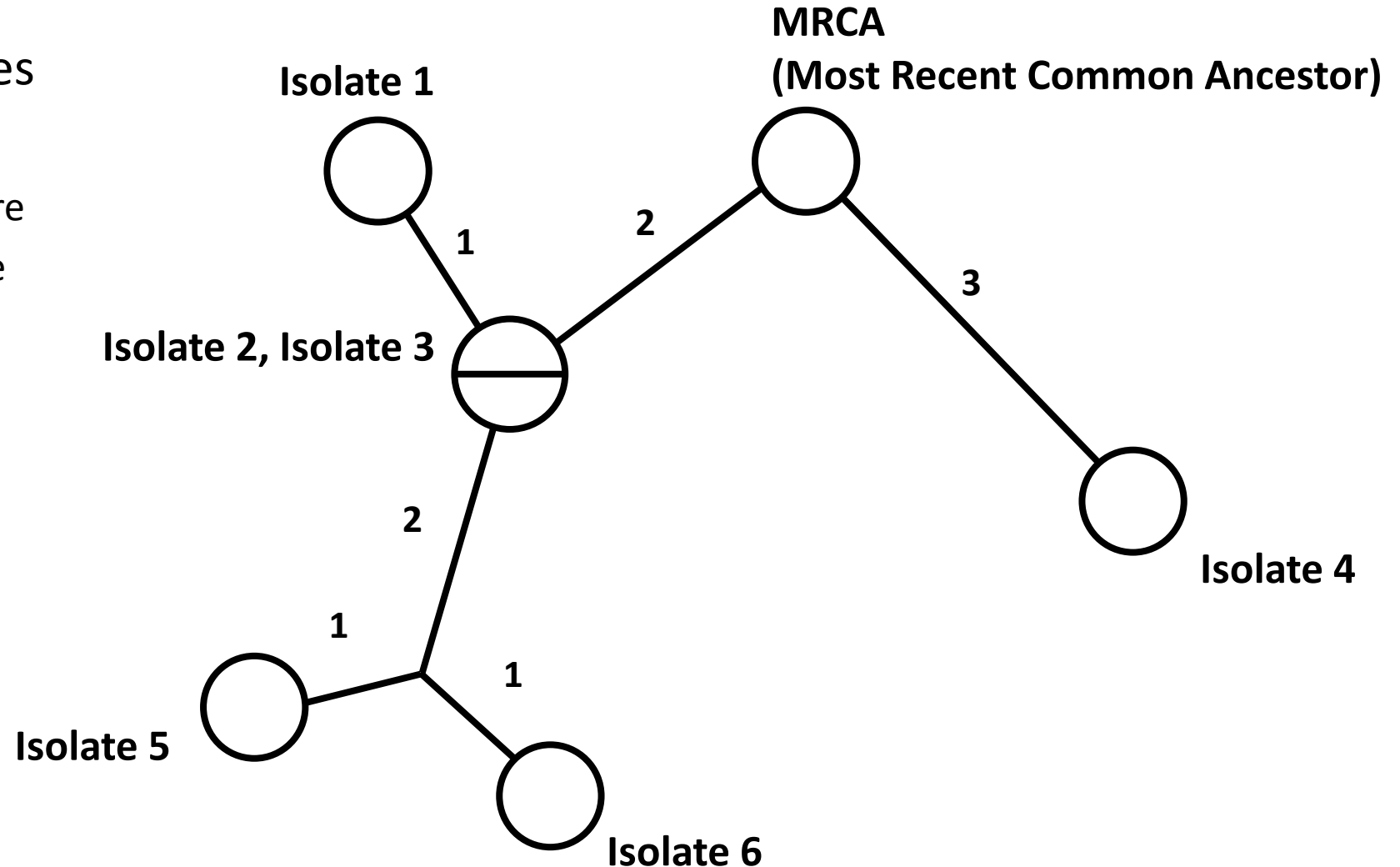


Whole-genome SNP comparison

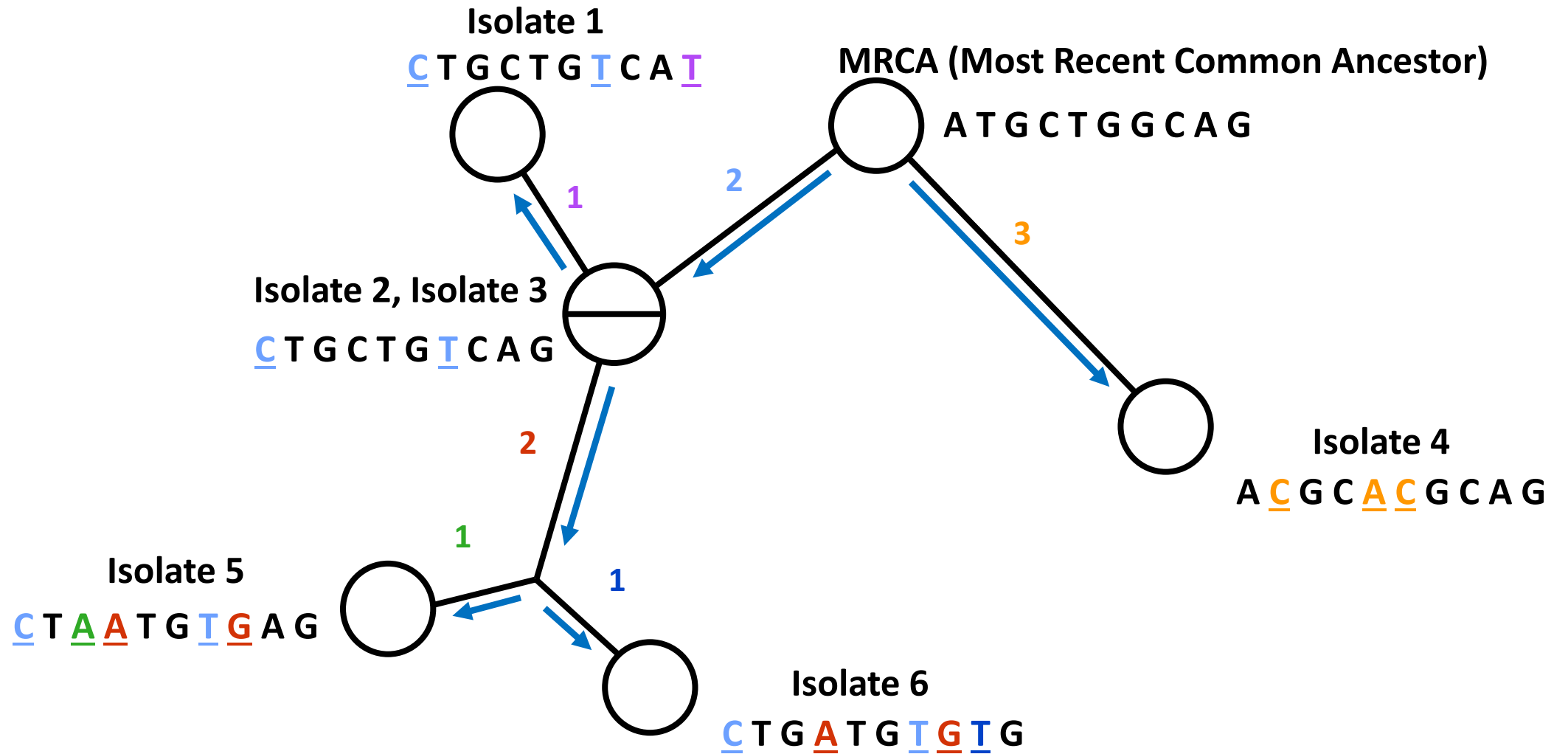
Reference	A	T	G	C	T	G	G	C	A	G
Isolate 1	<u>C</u>	T	G	C	T	G	<u>T</u>	C	A	<u>T</u>
Isolate 2	<u>C</u>	T	G	C	T	G	<u>T</u>	C	A	G
Isolate 3	<u>C</u>	T	G	C	T	G	<u>T</u>	C	A	G
Isolate 4	A	<u>C</u>	G	C	<u>A</u>	<u>C</u>	G	C	A	G
Isolate 5	<u>C</u>	T	<u>A</u>	<u>A</u>	T	G	<u>T</u>	<u>G</u>	A	G
Isolate 6	<u>C</u>	T	G	<u>A</u>	T	G	<u>T</u>	<u>G</u>	<u>T</u>	G

Whole-genome SNP comparison

- Isolates are shown as circles (called nodes)
 - Isolates that differ by 0 SNPs are displayed together in one node
- Lines are proportional in length to the number of SNPs that differ between the isolates
- Lines are labeled with the number of SNPs

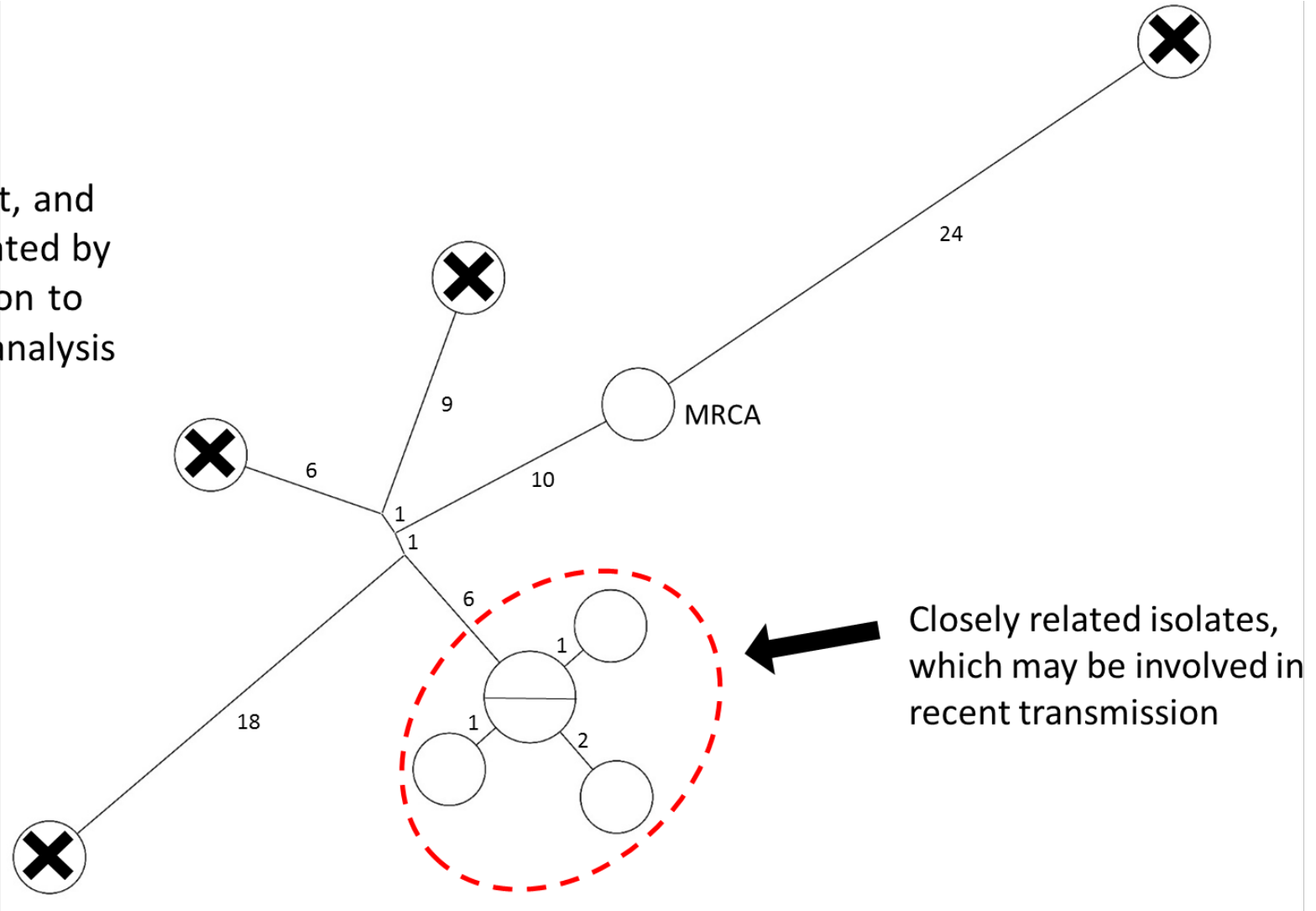


Whole-genome SNP comparison



Phylogenetic trees can be used to inform epidemiologic investigations

X = genetically distant, and unlikely to be related by recent transmission to other isolates in analysis

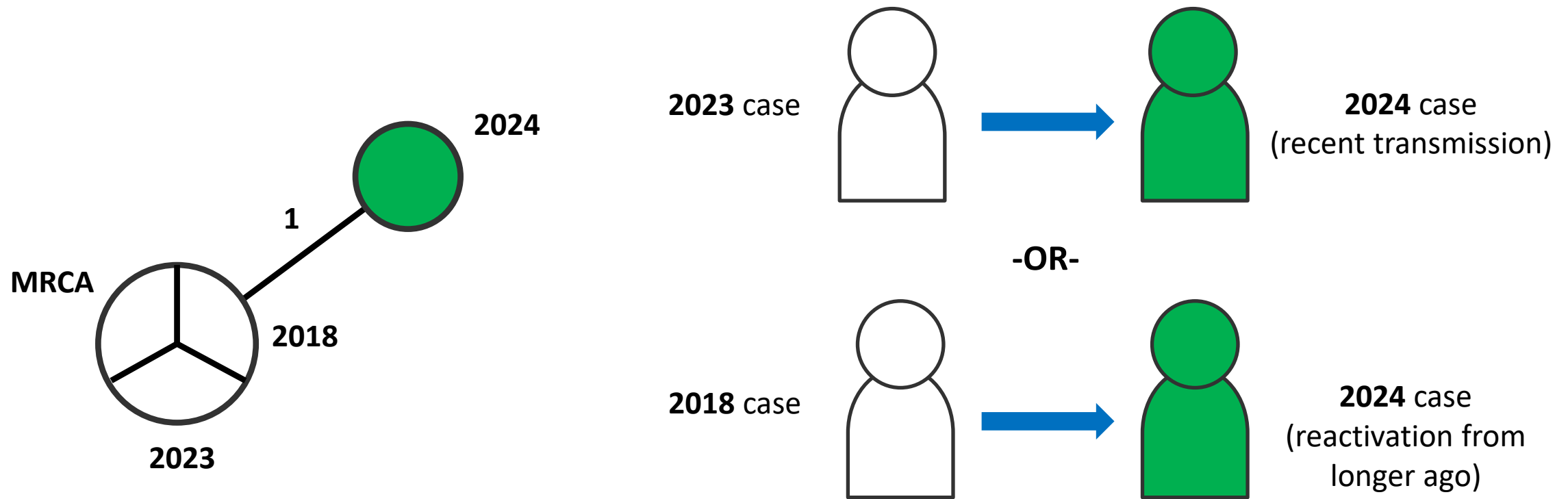


Limitations of whole-genome SNP comparison for understanding TB transmission

- Recent transmission is easier to rule out than to confirm with WGS
 - Even isolates that are closely related or identical by wgSNP can be due to reactivation
 - This is because mutations may not occur as frequently during latent infection and therefore SNPs may not accumulate

Limitations of whole-genome SNP comparison for understanding TB transmission

- Recent transmission is easier to rule out than to confirm with WGS



Limitations of whole-genome SNP comparison for understanding TB transmission

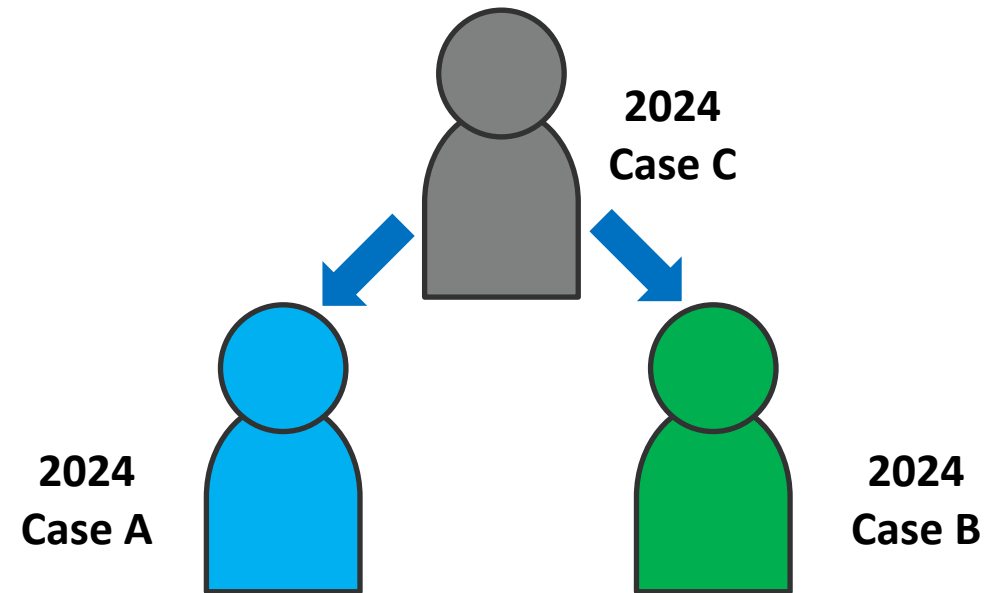
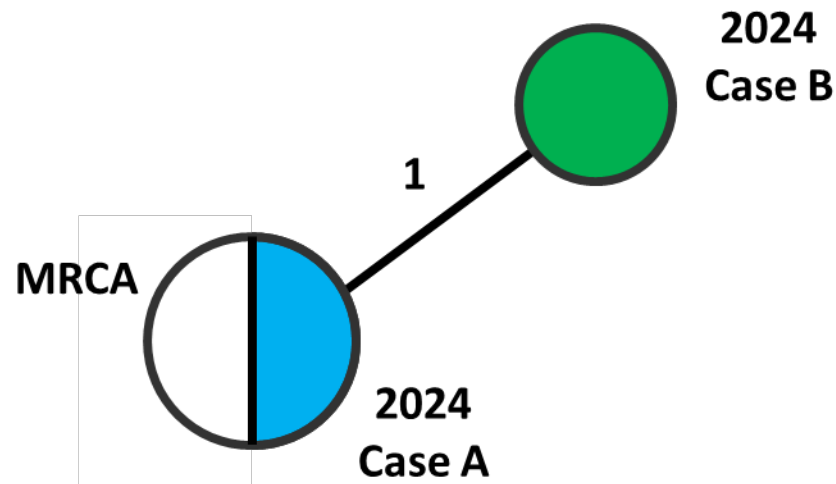
- A phylogenetic tree shows how isolates are genetically related to each other, but is not the same as a transmission diagram

Limitations of whole-genome SNP comparison for understanding TB transmission

- A phylogenetic tree shows how isolates are genetically related to each other, but is not the same as a transmission diagram
 - Isolates might be missing from the analysis
 - Examples: Cases that are not yet diagnosed, not culture-confirmed, contaminated isolates, out-of-country cases

Limitations of whole-genome SNP comparison for understanding TB transmission

- A phylogenetic tree shows how isolates are genetically related to each other, but is not the same as a transmission diagram

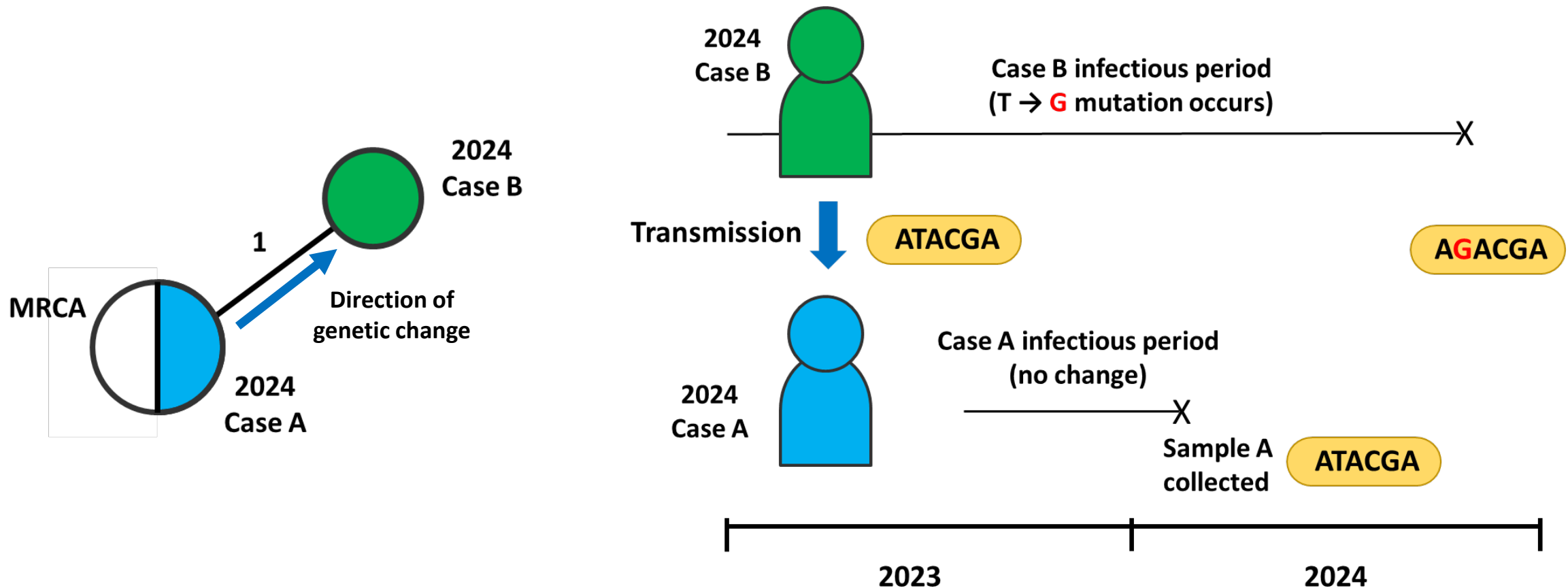


Limitations of whole-genome SNP comparison for understanding TB transmission

- A phylogenetic tree shows how isolates are genetically related to each other, but is not the same as a transmission diagram
 - Might be differences in sequence between bacterial population transmitted and bacterial population in sample that gets sequenced
 - Direction of genetic change depicted on tree might not be the same as direction of transmission

Limitations of whole-genome SNP comparison for understanding TB transmission

- A phylogenetic tree shows how isolates are genetically related to each other, but is not the same as a transmission diagram



Limitations of whole-genome SNP comparison for understanding TB transmission

- **WGS alone should not be used to infer direction of TB transmission**
- **The phylogenetic tree should be used in conjunction with clinical and epidemiologic information to assess recent transmission and infer direction of TB transmission**

WHOLE-GENOME SNP COMPARISON CASE STUDIES

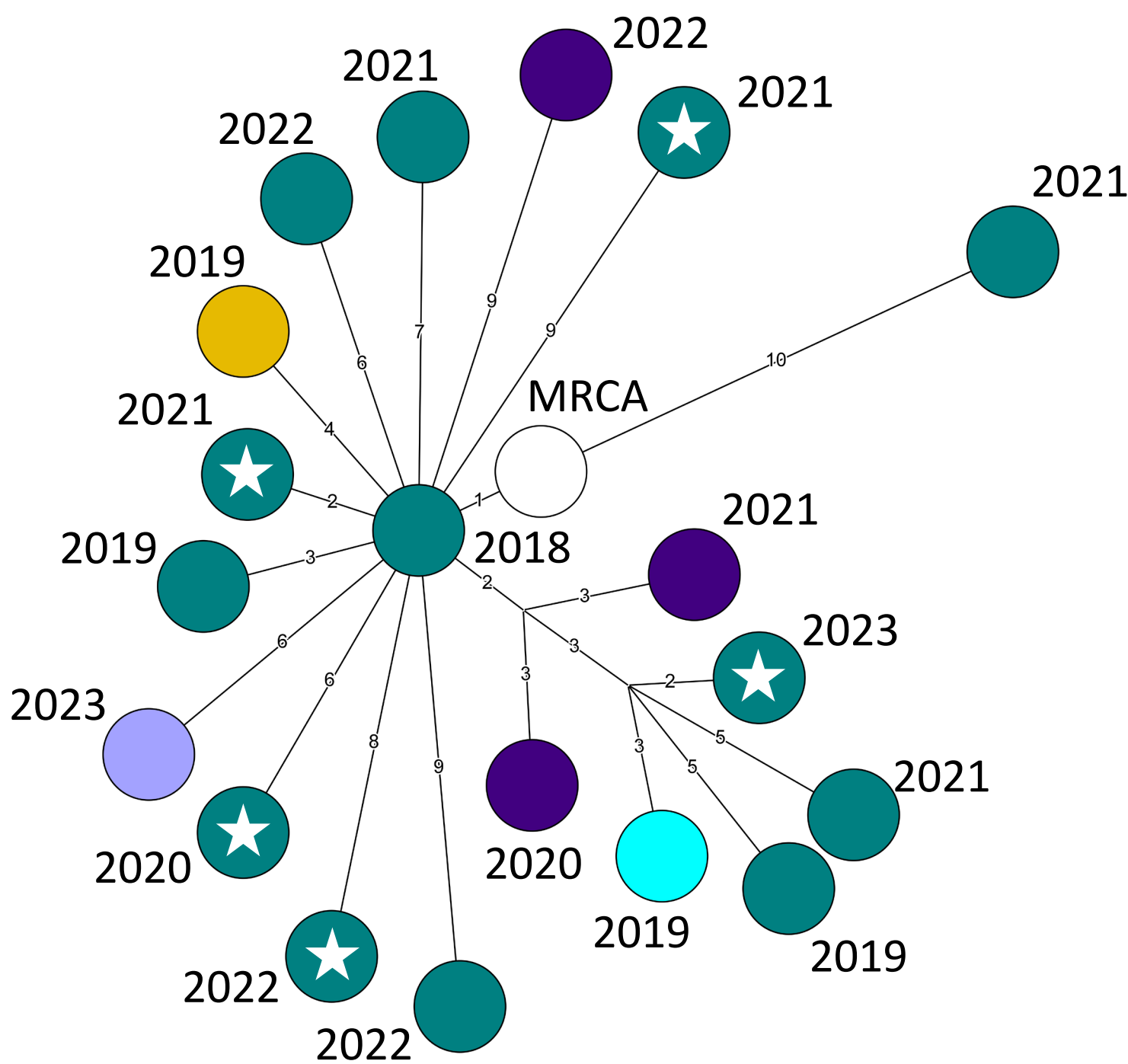
A Couple Case Studies

- 1. Refuting recent transmission among cases in a genotype-matched cluster**
- 2. Separating cases into subclusters of possible recent transmission**

Case Study 1: Refuting recent transmission among cases in a genotype matched cluster

Background

- **County cluster alert with 5 genotype matched cases**
- **Only 19 cases nationally with this wgMLSType since 2018**
 - 13 of these in the state with the cluster alert
- **All 19 cases in non-U.S.–born patients originating from the same country**

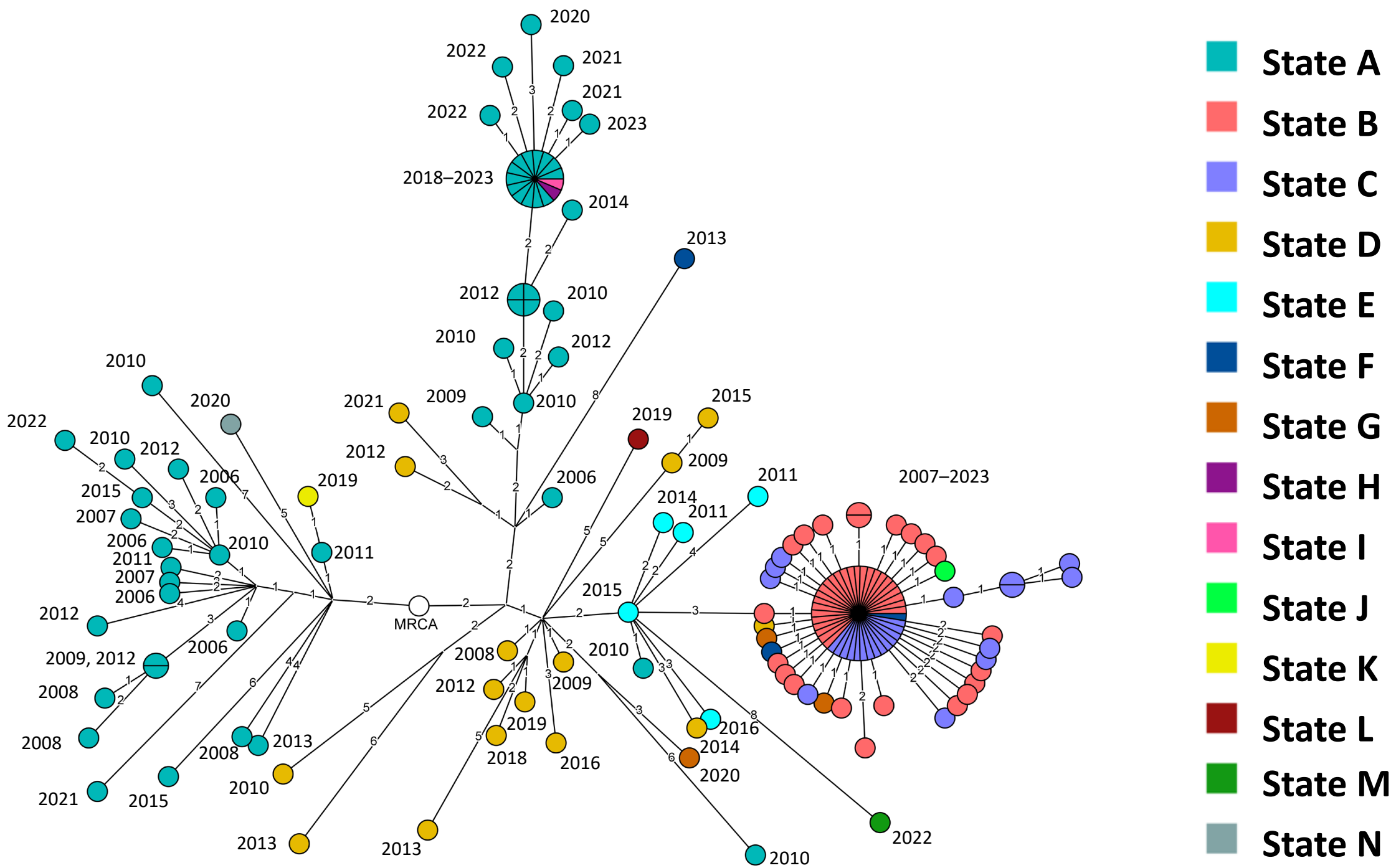


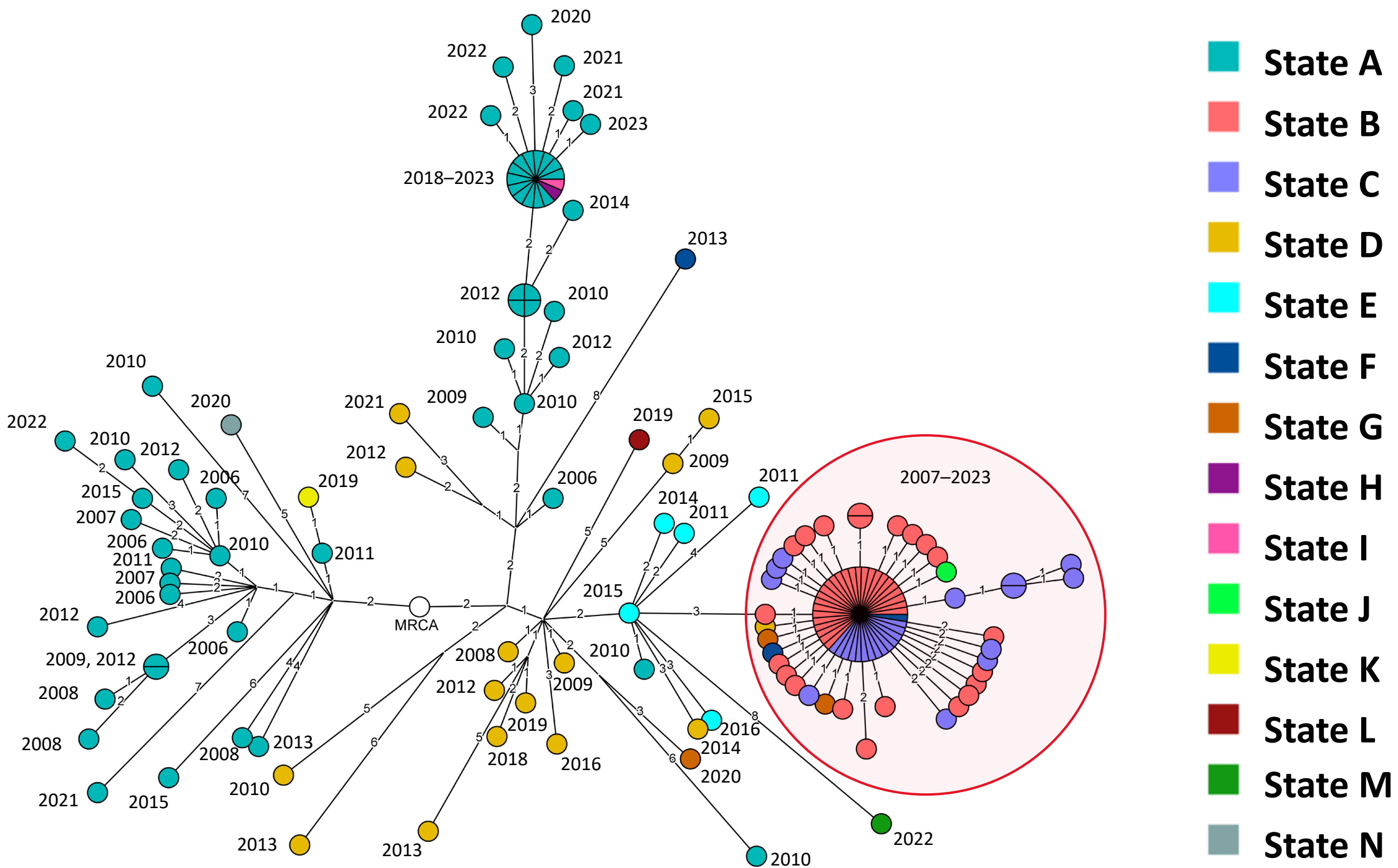
**Isolate from
case in the
cluster alert**

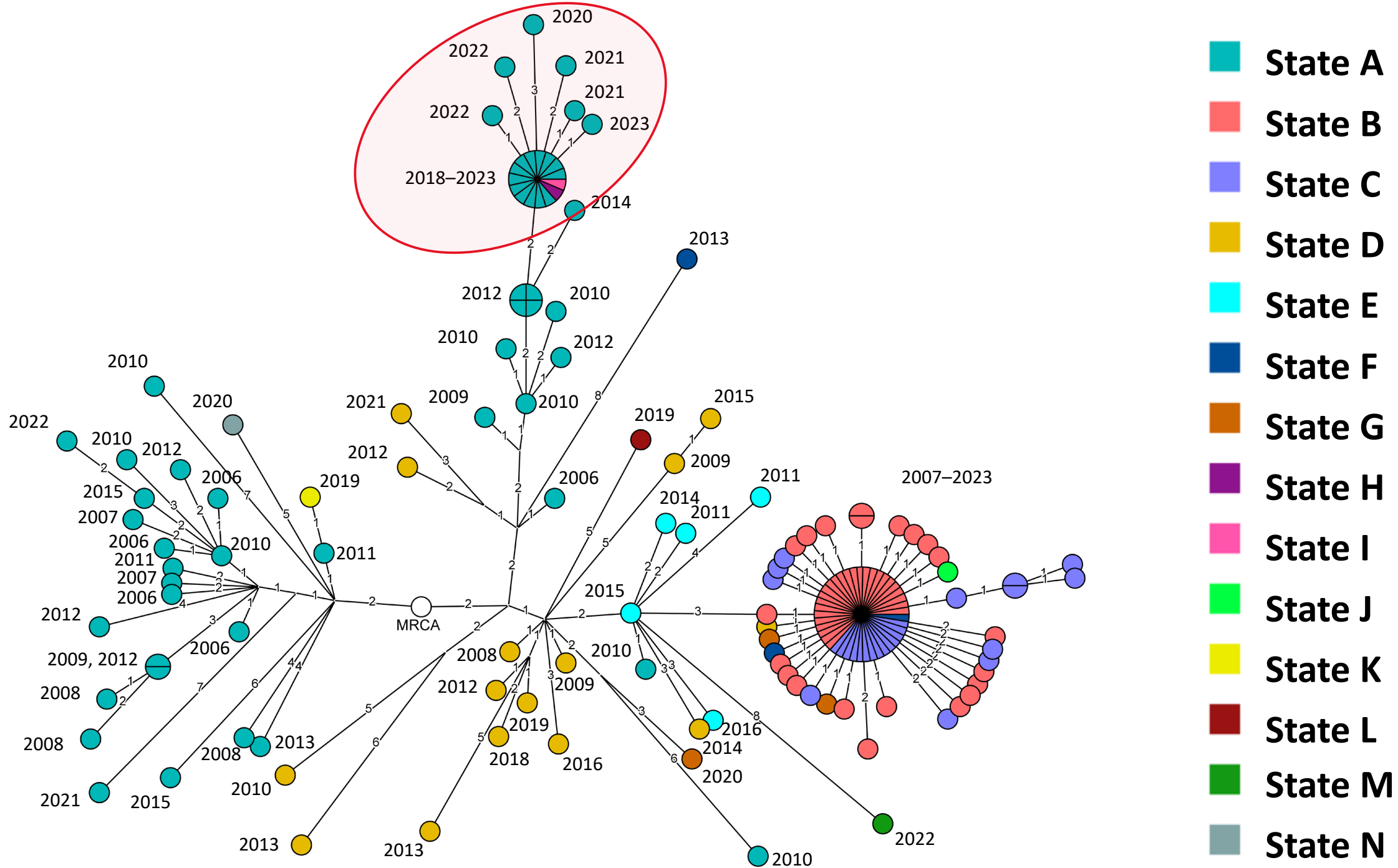
Case Study 2: Separating cases into subclusters of possible recent transmission

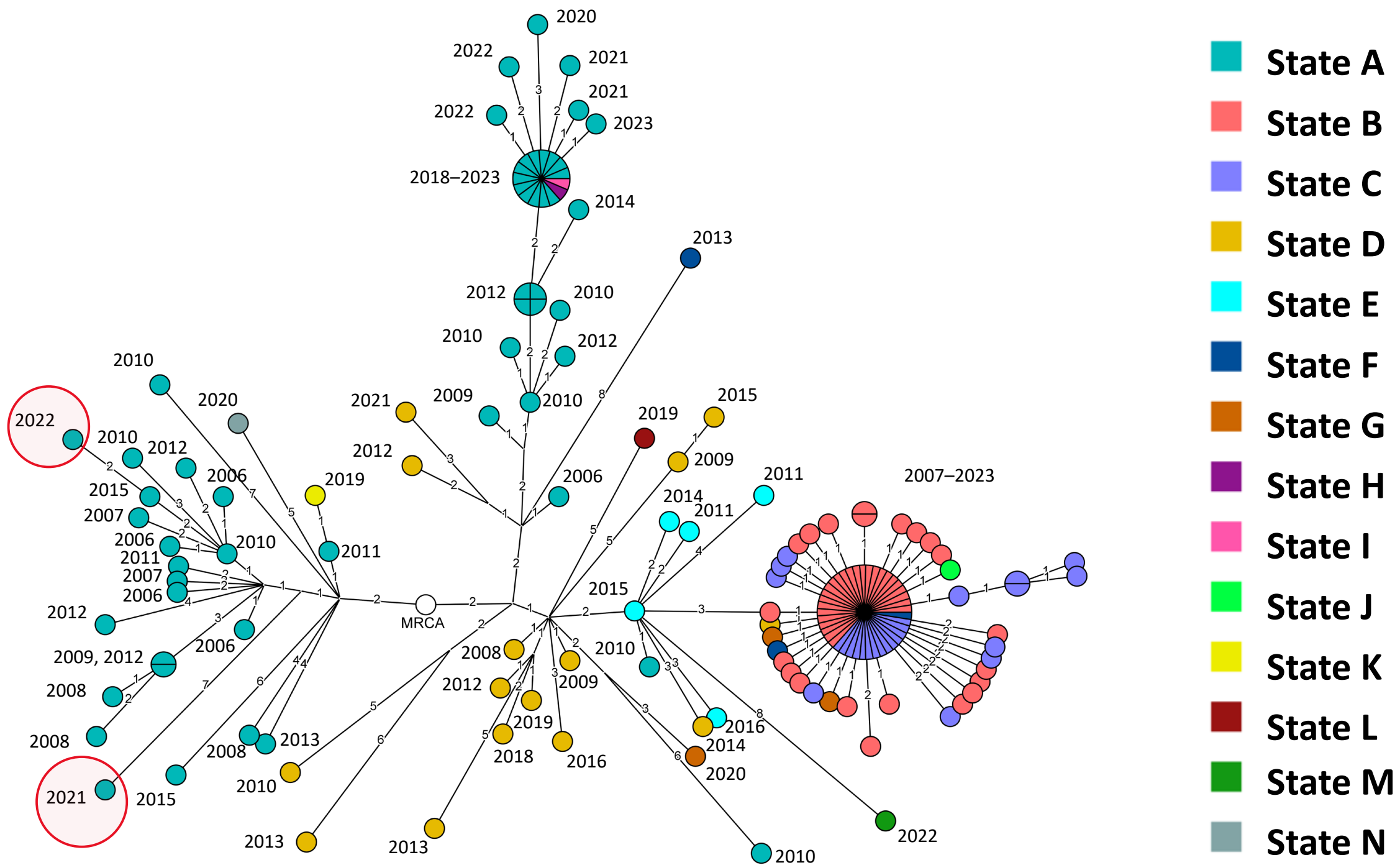
Background

- **Common genotype that has generated numerous county cluster and large outbreak alerts going back to 2012**
 - 15 county cluster alerts (4 counties in 3 states)
 - 2 separate large outbreak alerts
- **46 cases with this wgMLSType since 2018**
 - Additional 123 cases in this cluster prior to 2018
 - 18 states









Summary

- We use the genotyping data to identify clusters of TB cases that might represent recent transmission or an outbreak
 - Opportunity for public health action to interrupt further transmission
- *M. tuberculosis* isolates are first assigned a wgMLSType that is used for cluster detection and in molecular surveillance algorithms
- Whole-genome SNP comparison is then used to further assess which cases in a genotype-matched cluster might be related by recent transmission
 - Can be used to inform cluster and outbreak investigation

QUESTIONS?

For more information, contact CDC

1-800-CDC-INFO (232-4636)

TTY: 1-888-232-6348 [cdc.gov](https://www.cdc.gov)

Follow us on X (Twitter) [@CDCgov](https://twitter.com/CDCgov) & [@CDCEnvironment](https://twitter.com/CDCEnvironment)

The findings and conclusions in this report are those of the authors and do not necessarily represent the official position of the U. S. Centers for Disease Control and Prevention.

