

# **Implementing Whole Genome Cluster Analysis to Aid in *Salmonella* Outbreak Investigations**

William Wolfgang  
Wadsworth Center NYSDOH  
APHL  
06/01/14

# CDC statistics on Foodborne Illness in America

Each year

- 1 in 6 (48 million) get sick.
- 128,000 are hospitalized
- 3,000 die
- For about 60% the cause is unknown.

In 2011.

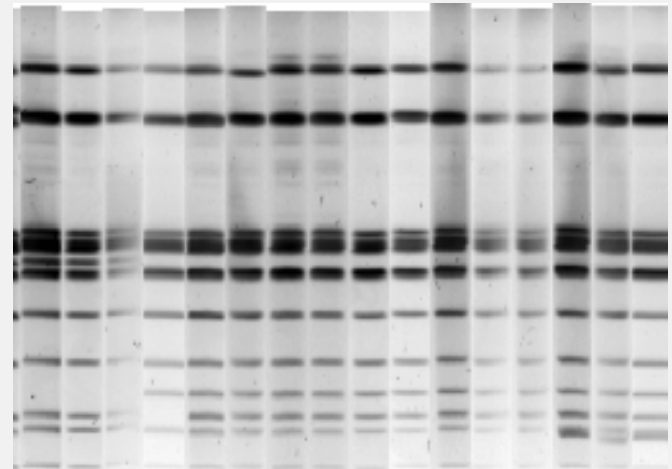
- *Salmonella* accounted for 1 million cases.
- 19,000 hospitalizations and 378 deaths.
- Rates of *Salmonella* infections are increasing.

# Surveillance of *Salmonella* in New York State

In New York all **positive patient specimens** are submitted to the Wadsworth Bacteriology Laboratory.

Bacteriology receives ~1,800 *Salmonella* patient specimens each year.

- Serotyped.
- DNA is fingerprinted by PFGE.



All data is sent to the CDC.

PFGE data bases are monitored:

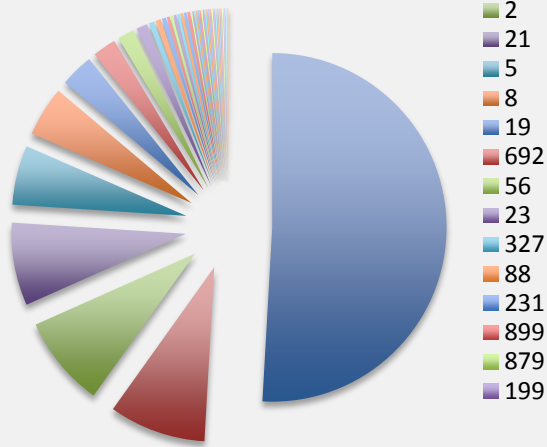
- To detect outbreaks in the patient population.
- To find source of the outbreak.



# For *Salmonella* Enteritidis

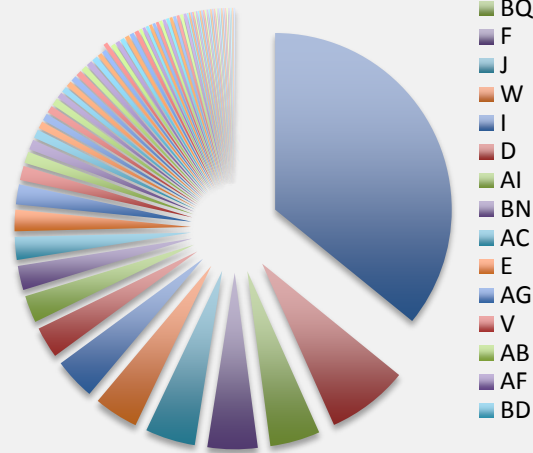
## PFGE typing methods are poor at resolving clusters

PFGE type frequency



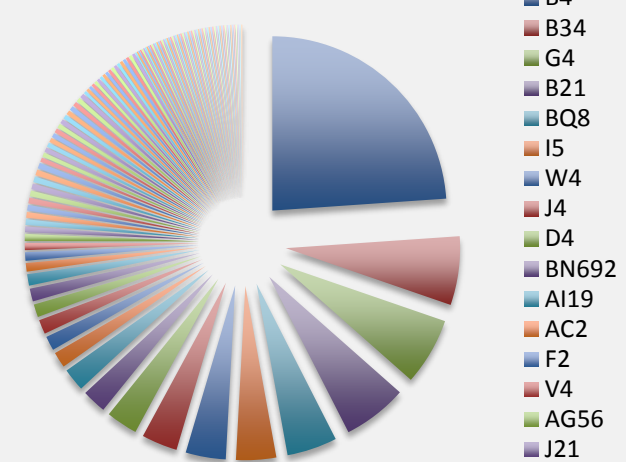
52 PFGE types

MLVA type frequency



98 MLVA types

MLVA-PFGE type frequency



163 combined  
MLVA-PFGE types

# Thus *Salmonella* Enteritidis outbreaks are rarely detected

## 2010 Outbreaks

- Montevideo - salami.
- Typhimurium – long term care facility.
- Javiana - tomatoes.
- Saintpaul - restaurant.
- **Enteritidis** - long term care facility (4).

## 2011 Outbreaks

- Typhimurium - aquatic frogs.
- Heidelberg - ground turkey.
- **Enteritidis** - food worker at a deli (692).
- Typhimurium - ground beef.
- **Enteritidis** - Turkish pine nuts (8).

## 2012 Outbreaks

- Bareilly - sushi.
- Nchanga – sushi.
- Hartford - sub shop worker.
- Newport - chick and duck exposure.
- SanDiego & Poona - small turtles.
- Javiana - Mothers Day fruit baskets.
- **Enteritidis** - Cargill ground beef (9).

# Retrospective studies have shown that WGS can improve pathogen tracking and surveillance.

- MRSA
- Tb
- Drug resistant *Klebsiella pneumoniae*
- Cholera outbreak in Haiti.
- *Salmonella* Montevideo outbreak from pepper.

# A Whole-Genome Single Nucleotide Polymorphism-Based Approach To Trace and Identify Outbreaks Linked to a Common *Salmonella* *enterica* subsp. *enterica* Serovar Montevideo Pulsed-Field Gel Electrophoresis Type †

Henk C. den Bakker,<sup>1\*</sup> Andrea I. Moreno Switt,<sup>1</sup> Craig A. Cummings,<sup>2</sup> Karin Hoelzer,<sup>1</sup>  
Lovorka Degoricija,<sup>2</sup> Lorraine D. Rodriguez-Rivera,<sup>1</sup> Emily M. Wright,<sup>1</sup> Rixun Fang,<sup>2</sup>  
Margaret Davis,<sup>3</sup> Tim Root,<sup>4</sup> Dianna Schoonmaker-Bopp,<sup>4</sup> Kimberlee A. Musser,<sup>4</sup>  
Elizabeth Villamil,<sup>4</sup> HaeNa Waechter,<sup>5</sup> Laura Kornstein,<sup>5</sup>  
Manohar R. Furtado,<sup>2</sup> and Martin Wiedmann<sup>1</sup>

---

---

## Identification of a Salmonellosis Outbreak by Means of Molecular Sequencing

E. Kurt Lienau, Ph.D.

Errol Strain, Ph.D.

Charles Wang, B.S.

Jie Zheng, D.V.M., Ph.D.

Andrea R. Ottesen, Ph.D.

Christine E. Keys, M.S.

Thomas S. Hammack, M.S.

Steven M. Musser, Ph.D.

Eric W. Brown, Ph.D.

Marc W. Allard, Ph.D.

Food and Drug Administration

College Park, MD

marc.allard@fda.hhs.gov

Guojie Cao, M.S.

Jianghong Meng, D.V.M., Ph.D.

University of Maryland

College Park, MD

Robert Stones, M.S.

Food and Environment Research Agency

York, United Kingdom

# Acknowledgments

- Cornell University  
Martin Wiedmann  
Henk den Bakker
- FDA  
Eric Brown  
Peter Evans  
Marc Allard  
Errol Strain  
Ruth Timme
- Connecticut DOH  
Stacey Kinney  
John Fontana
- Minnesota DOH  
David Boxrud  
Angie Jones
- NCBI  
Bill Klimke  
Martin Shumway
- Wadsworth Bacteriology Laboratory  
Kara Mitchell
- Wadsworth Genomics Core  
Matt Shudt  
Zhen Zhang  
Charles MacGowan  
Mark Rosenthal  
Melissa Leisner  
Danielle Loranger  
Mike Palumbo  
Pascal LaPierre
- Wadsworth PulseNet Lab  
Dianna Bopp  
Deb Baker  
Lisa Thompson



# Cannoli outbreak

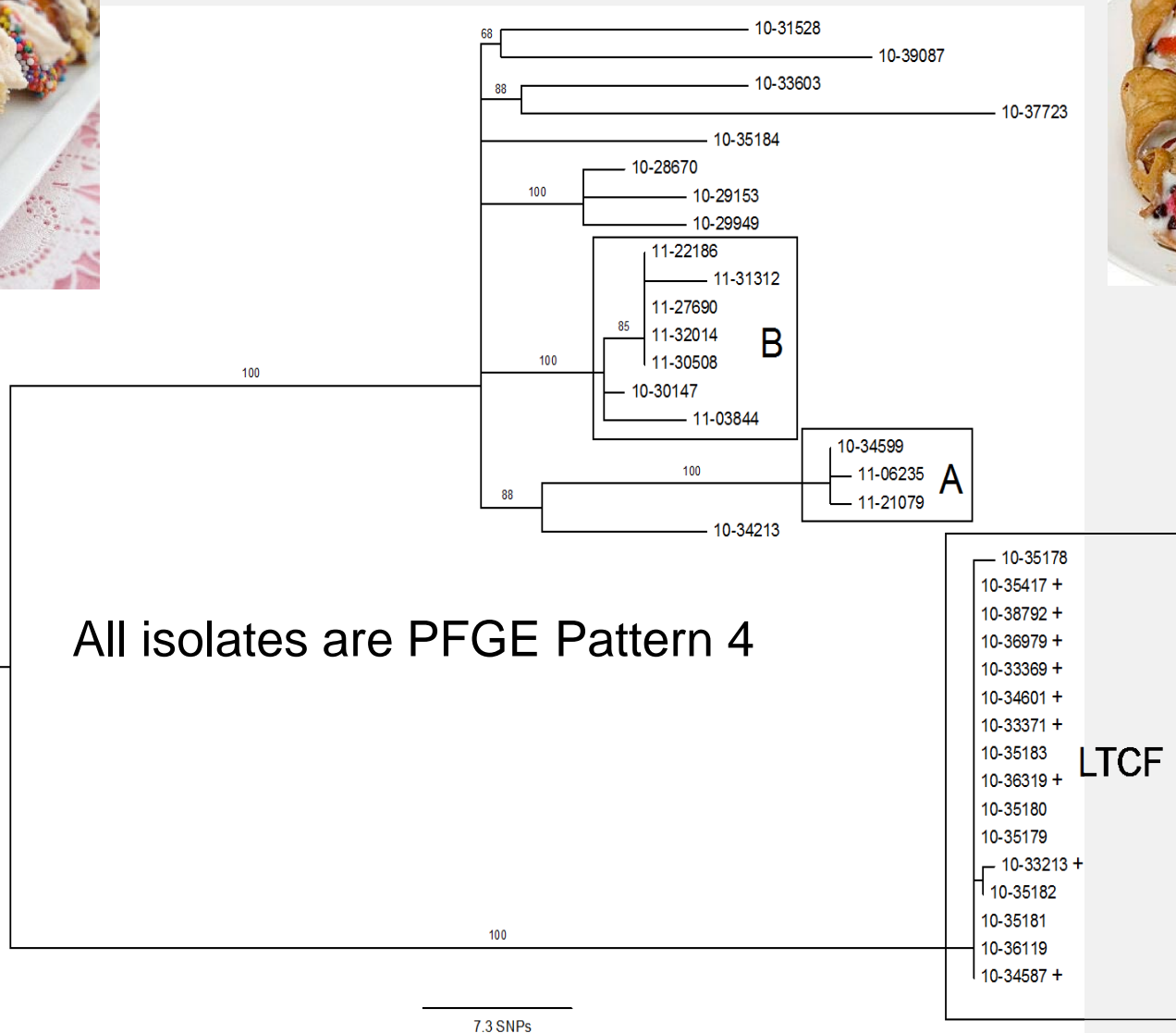
- Sept. 2010 Connecticut Dept. of Health identifies a *Salmonella* outbreak in a long term care facility.
- Outbreak was linked to cannoli from a Westchester bakery.
- Both NY and CT cases consumed cannoli's.
- Isolates had the most common PFGE pattern, JEGX01.0004.



# Retrospective cohort

Key	County	Date	PFGE-MLVA Combined
IDR1000029153	Cattaraugus	8/10/10	JEGX01.0004W
IDR1000031528	Rockland	8/26/10	JEGX01.0004W
IDR1000033213	Putnam	9/10/10	JEGX01.0004W
IDR1000033369	Putnam	9/10/10	JEGX01.0004W
IDR1000033371	Putnam	9/11/10	JEGX01.0004W
IDR1000034601	Washington	9/13/10	JEGX01.0004W
IDR1000034587	Westchester	9/20/10	JEGX01.0004W
IDR1000035417	Putnam	9/22/10	JEGX01.0004W
IDR1000035178	Westchester	9/13/10	JEGX01.0004W
IDR1000035179	Greenwich CT	9/12/10	JEGX01.0004W
IDR1000035180	Westchester	9/12/10	JEGX01.0004W
IDR1000035181	Westchester	9/13/10	JEGX01.0004W
IDR1000035182	Westchester	9/12/10	JEGX01.0004W
IDR1000035183	Greenwich CT	9/16/10	JEGX01.0004W
IDR1000036119		9/17/10	JEGX01.0004W
<b>IDR1100035184</b>	<b>Westchester</b>	<b>9/16/10</b>	<b>JEGX01.0004AE</b>
IDR1000036319	Putnam	9/28/10	JEGX01.0004W
IDR1000036979	Putnam	10/8/10	JEGX01.0004W
IDR1000038792	Nassau	10/29/10	JEGX01.0004W
IDR1000034599	Orange	9/15/10	JEGX01.0004W
IDR1100006235	Westchester	2/21/11	JEGX01.0004W
IDR1100021079	Rockland	7/13/11	JEGX01.0004W
IDR1000030147	Out-Of-State	8/22/10	JEGX01.0004W
IDR1100003844	Onondaga	2/1/11	JEGX01.0004W
IDR1100022186	Yates	7/22/11	JEGX01.0004W
IDR1100027690	Erie	9/6/11	JEGX01.0004W
IDR1100030508	Madison	10/9/11	JEGX01.0004W
IDR1100031312	Suffolk	10/5/11	JEGX01.0004W
IDR1100032014	Onondaga	10/22/11	JEGX01.0004W
IDR1000028670	Nassau	8/8/10	JEGX01.0004B
IDR1000029949	Suffolk	8/16/10	JEGX01.0004B
IDR1000033603	Erie	9/14/10	JEGX01.0004B
IDR1000034213	Erie	9/13/10	JEGX01.0004B
IDR1000037723	Westchester	10/4/10	JEGX01.0004B
IDR1000039087	Westchester	10/27/10	JEGX01.0004B

# WGCA can identify an outbreak cluster not detected by PFGE



# Implementing WGCA in real-time.

- Evaluate WGCA efficacy compared to PFGE.
  - Speed
  - Actionable Clusters
  - Cost
- Develop in house bioinformatic pipeline.
- Develop communication pipeline to epidemiologists.
- Determine cluster parameters that represent an outbreak from a single source.
- Acquire a real data set to evaluate evolving informatic methods.

Data is analyzed using a portal developed by Informatics core

The screenshot displays the Galaxy web interface at <https://galaxy.wadsworth.org/>. The top navigation bar includes 'Galaxy', 'Analyze Data', 'Workflow', 'Shared Data', 'Visualization', 'Help', and 'User'. A search bar on the left shows results for 'Bill's Salmonella Pipelines', with a red box highlighting the search results. Below the search bar is a list of tool categories such as 'Text Manipulation', 'Filter and Sort', 'Join, Subtract and Group', 'Convert Formats', 'Extract Features', 'Fetch Sequences', 'Fetch Alignments', 'Get Genomic Scores', 'Operate on Genomic Intervals', 'Statistics', 'Wavelet Analysis', 'Graph/Display Data', 'Regional Variation', 'Multiple regression', 'Multivariate Analysis', 'Evolution', 'Motif Tools', 'Multiple Alignments', 'Metagenomic analyses', 'FASTA manipulation', and 'NGS: QC and manipulation'. The main workspace shows a heatmap visualization with a dendrogram on the left and a color-coded matrix. The right-hand 'History' panel lists recent jobs: '20: heatmap.pdf' (133.5 KB), '19: Log file.txt' (28 lines), '18: PhyML\_tree.txt' (1 line), and '17: Statistics.txt' (7 lines). The 'Statistics.txt' job output is visible at the bottom of the history panel, showing a table of metrics for various samples.

Sample	% reads mapped	% correctly	
swgs1194	98.70%	98.42%	99.1
swgs1195	97.93%	97.67%	99.1
swgs1196	98.69%	98.37%	99.1
swgs1197	97.00%	96.71%	99.1
swgs1198	95.28%	95.03%	99.1

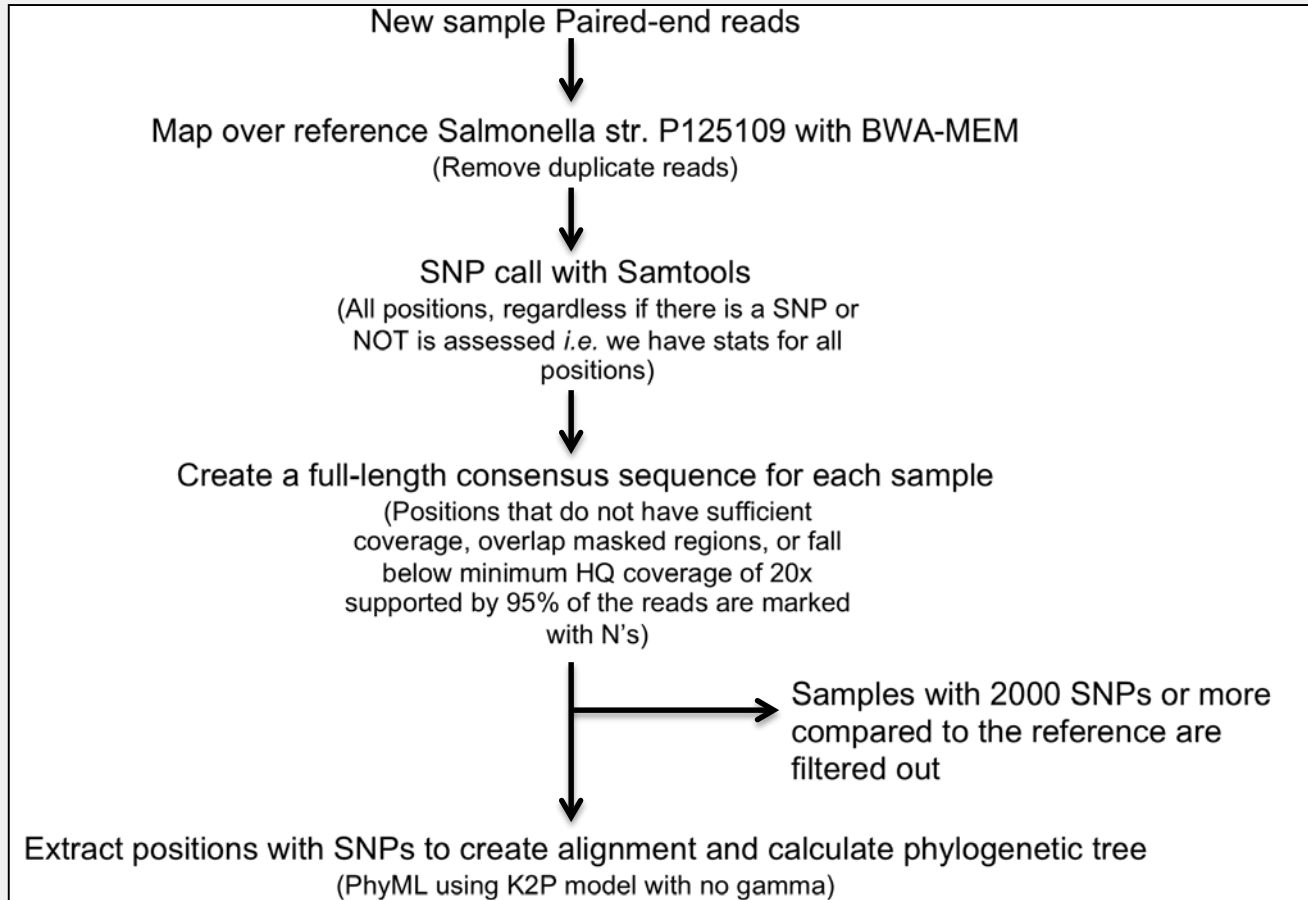
### Bill's Salmonella Pipelines

Minnesota Data Pipeline SNP finder and Tree builder for Minnesota Salmonella samples

Pipeline Merger Create Tree and Heatmap from Salmonella merged data

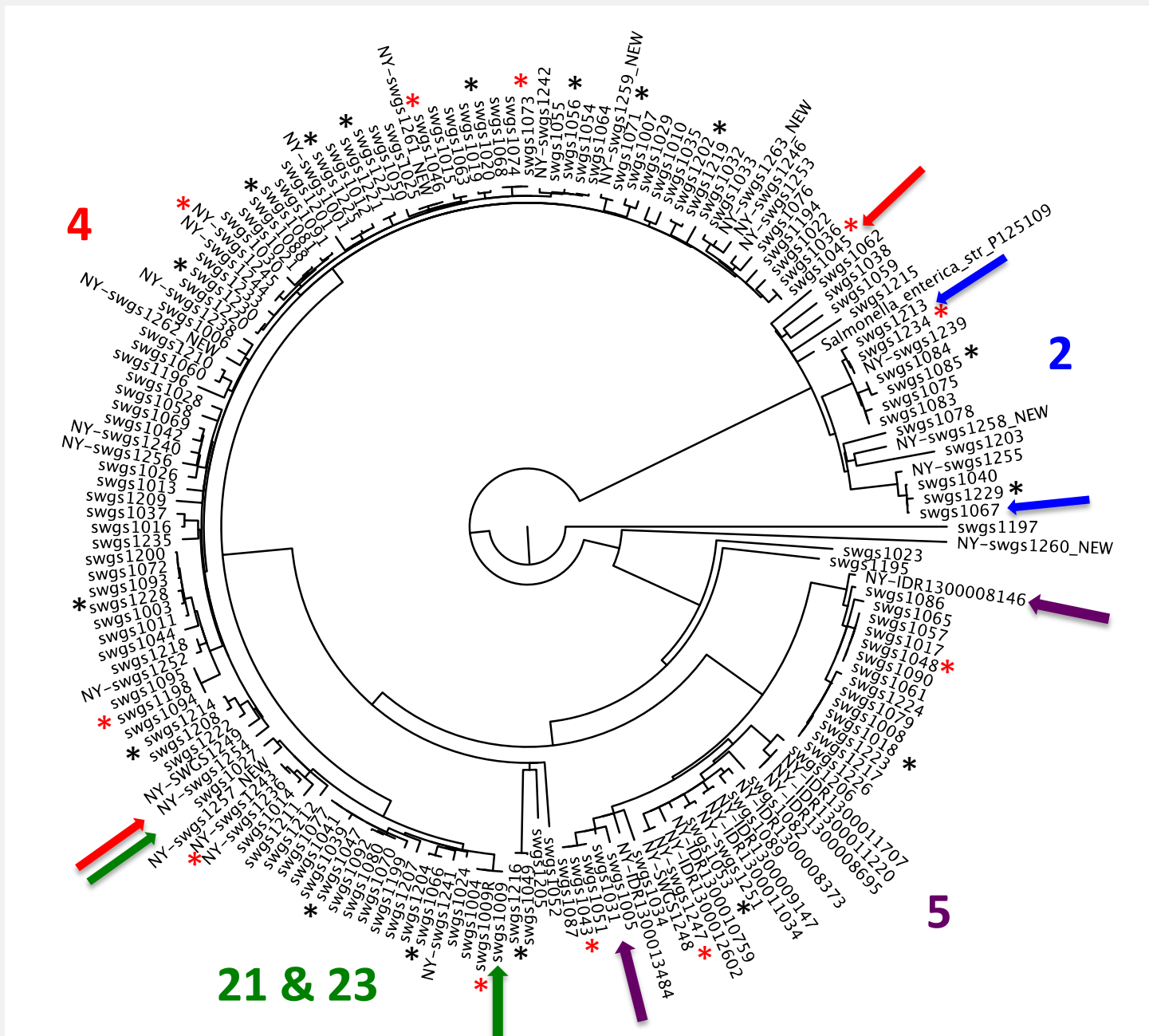
Salmonella Pipeline SNP finder and Tree builder for Salmonella

# Pipeline





# Clusters vs. PFGE







# In 6 months:

156 isolates have been sequenced and analyzed.

- .ca 6 isolates / week

28 clusters were reported to epidemiologists.

- 1 cluster every week

12 clusters have **zero** snp differences.

- Collected up to 6 months apart.

10 clusters acquired one or more new isolates.

- 7 isolates in the largest cluster

# Clusters with identical SNPs

cluster	# of isolates	temporal distance	spatial distance	epidemiological report
GC-13	3	2d	metro	family thanksgiving meal
GC-02	2	4d	same county	shopped at same grocery store
GC-27	2	6d	same county	sibs
GC-09	4	42d	2 in one county; 2 distant	husband wife pair; other two not related
GC-07	3	6mo	metro	no common source identified

# PFGE vs. WGCA for surveillance

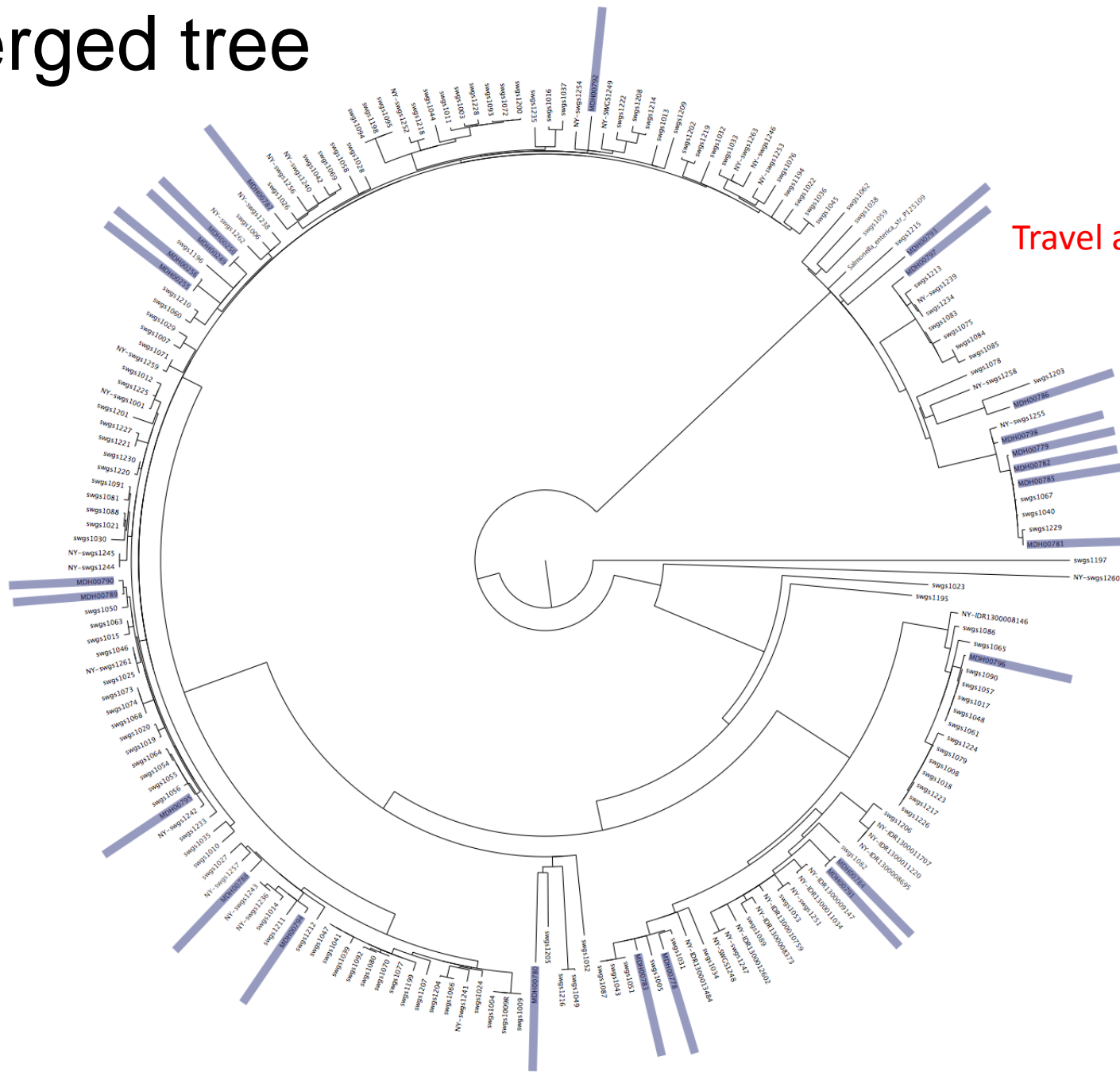
metric	PFGE	WGCA
TAT from isolation	2 days	7 days
Cost	\$69	\$294
Technician time	8h	10h
Actionable clusters	2	28

# Two State Network

- Collaborating with Minnesota.
  - Currently no informatics in house.
- We pull their sequences of Basespace.
- Run through our pipeline.



# Merged tree



# FDA Genome Trackr network

## State Health labs

- New York
- Florida
- Arizona
- Washington
- Minnesota
- Virginia
- Maryland

## FDA labs

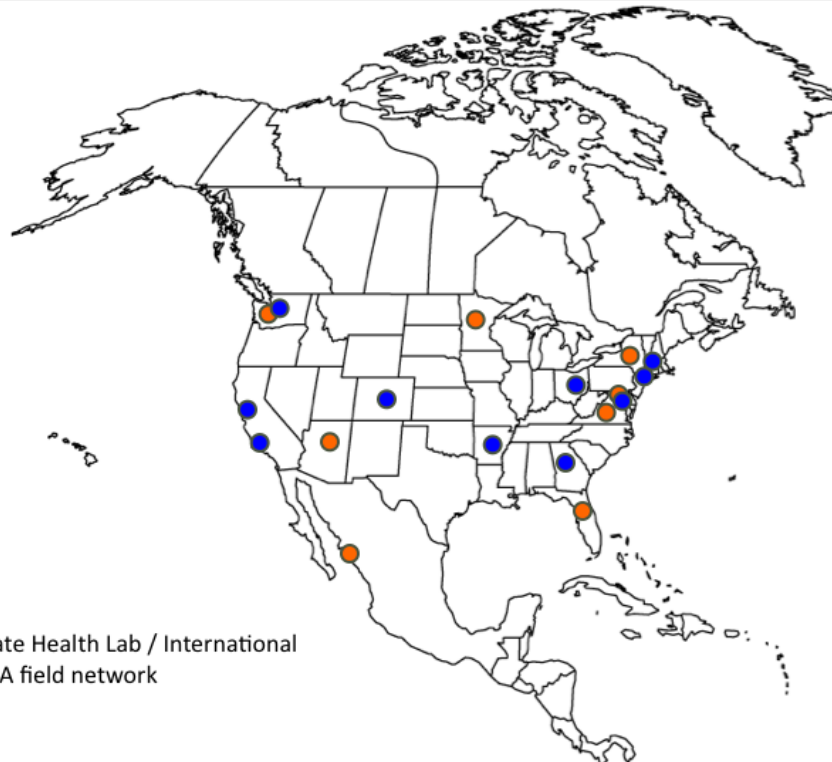
- 9 FDA field labs
- CFSAN - MOD1
- CFSAN - Wiley
- IEH (contracting lab)

## International labs

- Mexico
- Ireland
- UK (FERA)
- Columbia

## Contributors

- Turkey
- Brazil
- Italy



# Genomic Surveillance Machine

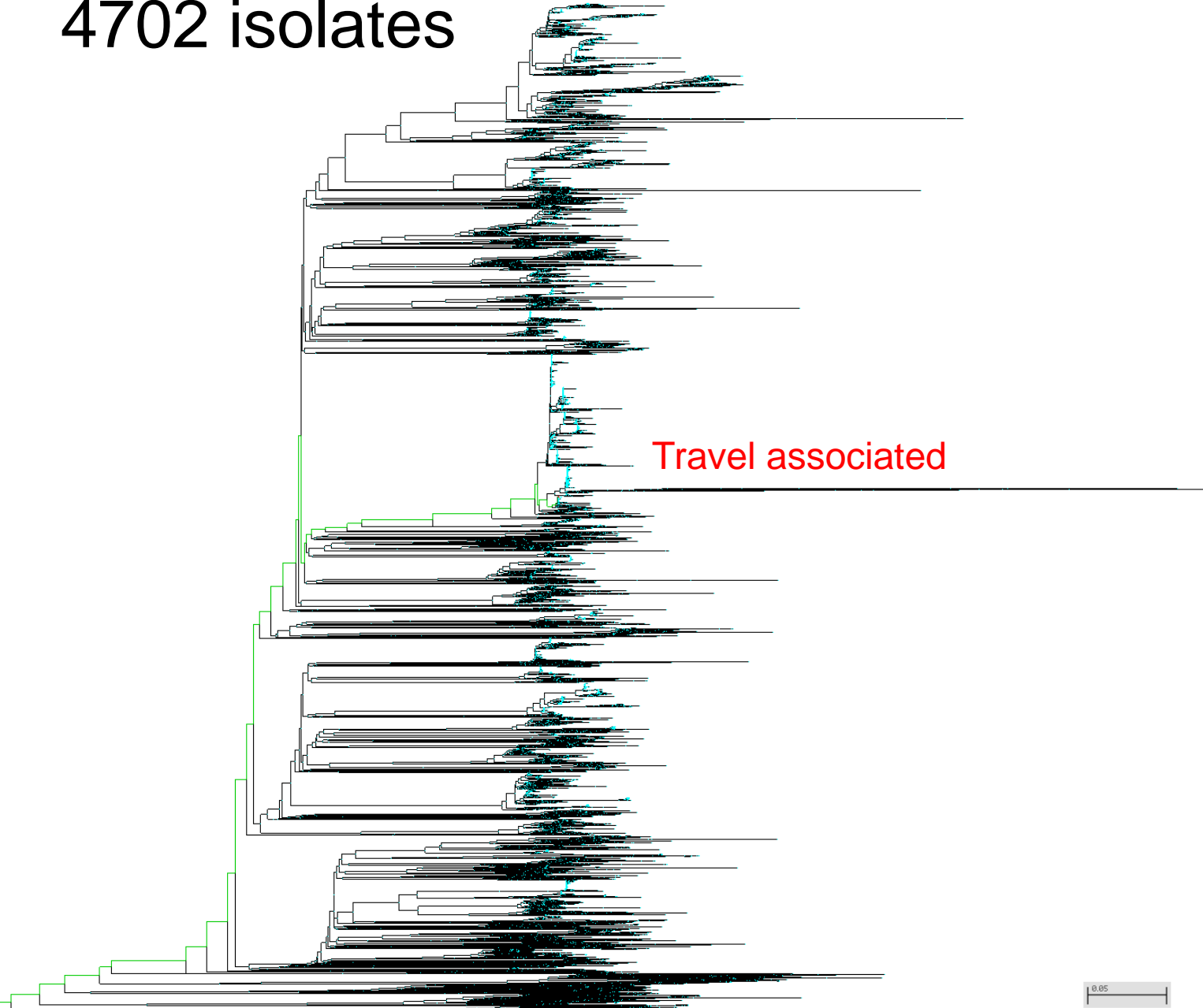
- State labs feed the machine by uploading sequences from isolates received through surveillance.
- Federal and other support for reagents and equipment.
- NCBI to analyze the products of this machine and reports results to state and federal agencies.

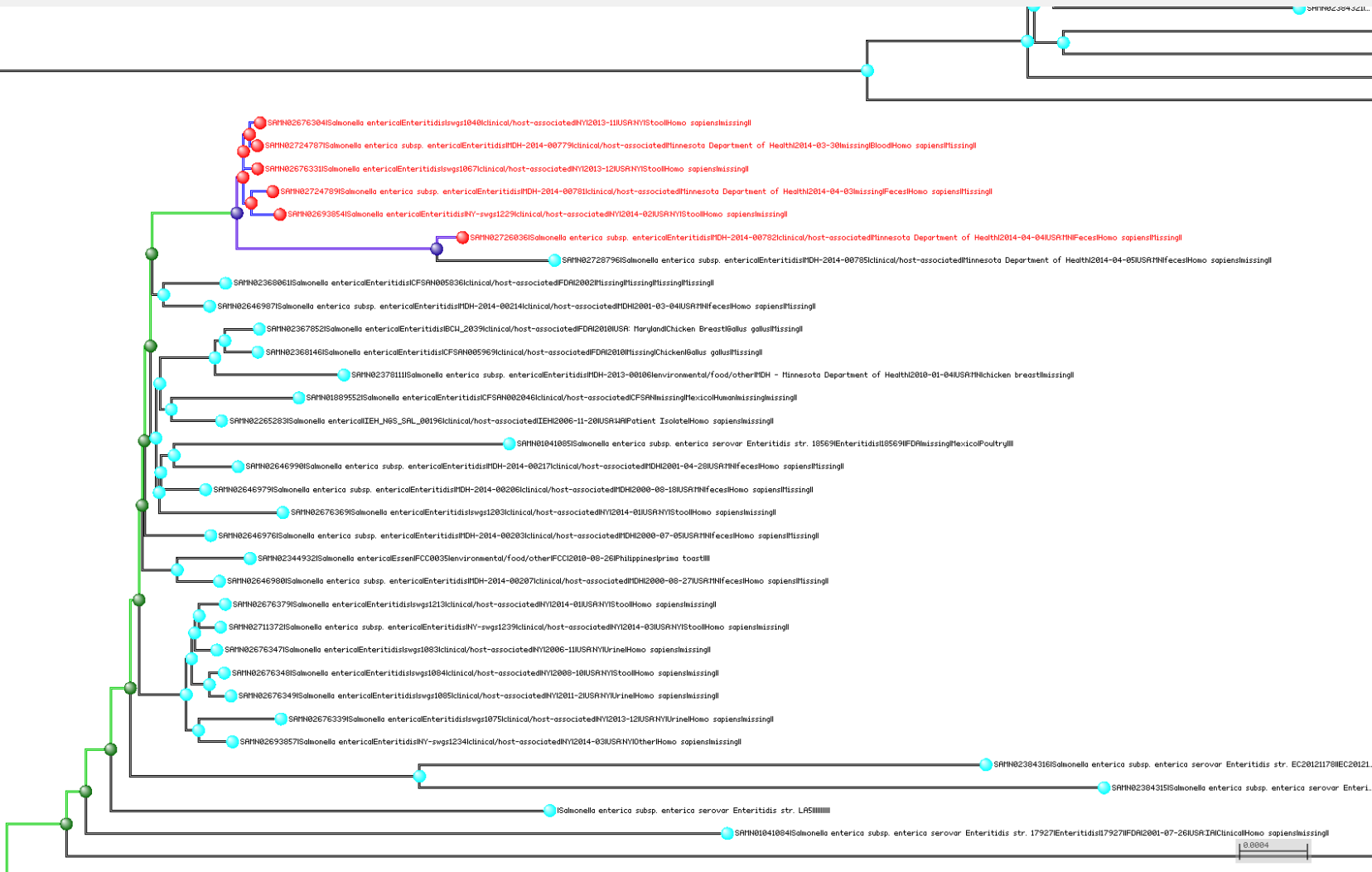




# NCBI *Salmonella* tree

4702 isolates





0.0004

# Expected Outcomes for WGS surveillance

## Laboratory

- Improve outbreak cluster detection.
- Clusters will be detected more rapidly and from fewer isolates.

## Epidemiologists

- Allow identification of clusters within **endemic** patterns.
- Allow more efficient use of resources by focusing on highly genetically related clusters.
- Solve more clusters.

## Public Health

- More efficient identification and removal of pathogen sources.

# Challenges exist in Creating a Network

- Increasing amounts of data.
- Metadata: how much should be public?
  - In real time?
  - What elements?
- Transitioning:
  - Integration with PFGE typing.
  - Integrating surveillance at a national level.
  - Paying.
- As sequencing technology and bioinformatics evolve:
  - Need to maintain backward compatibility

# THE FUTURE

## Near term

- Universal *Salmonella* tree
- Hands off data submission to NCBI
- Hands off data analysis

## Further out

- Entire process automated
  - Sample preparation
  - Sequencing
  - Identification of clusters
  - Reporting

# Summary

- WGS can improve surveillance activities and outbreak traceback.
- It is practical to develop network.
- It is likely the transition will be gradual – first to go live.
  - Pathogens with more stable genomes.
  - Pathogens with greatest Public Health Impact.