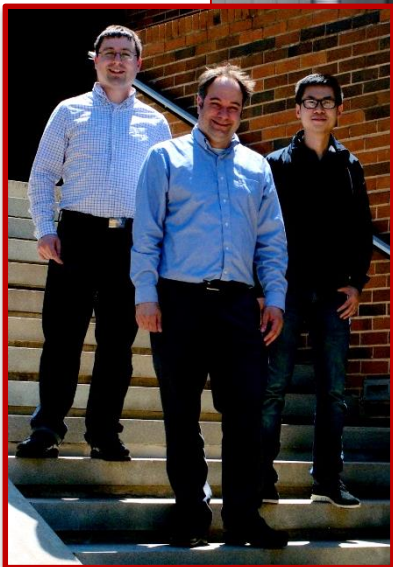# Genomic Epidemiology

**Lee Katz, Ph.D.**

Senior bioinformatician
Enteric Diseases Laboratory Branch

Computational Genomics course
Feb 11, 2020

# Acknowledgements up front

- Every single compgenomics class since 2008
- My branch at CDC
- Federal partners
- State partners

# Enteric Diseases Laboratory Branch (EDLB)



Food Safety Informatics Group,
Center for Food Safety,
University of Georgia

Enteric Diseases Bioinformatics
Team (EDBiT)

# Enteric Diseases Laboratory Branch

2011 to present

*Vibrio, Campylobacter, Escherichia, Shigella, Yersinia, Salmonella*

# PulseNet's **20**-year history of making food safer to eat

**1993** — Outbreak of a deadly form of *E. coli* infections in western states: > 700 illnesses; 4 children died. Highlights the need for a network to identify DNA fingerprints of foodborne bacteria
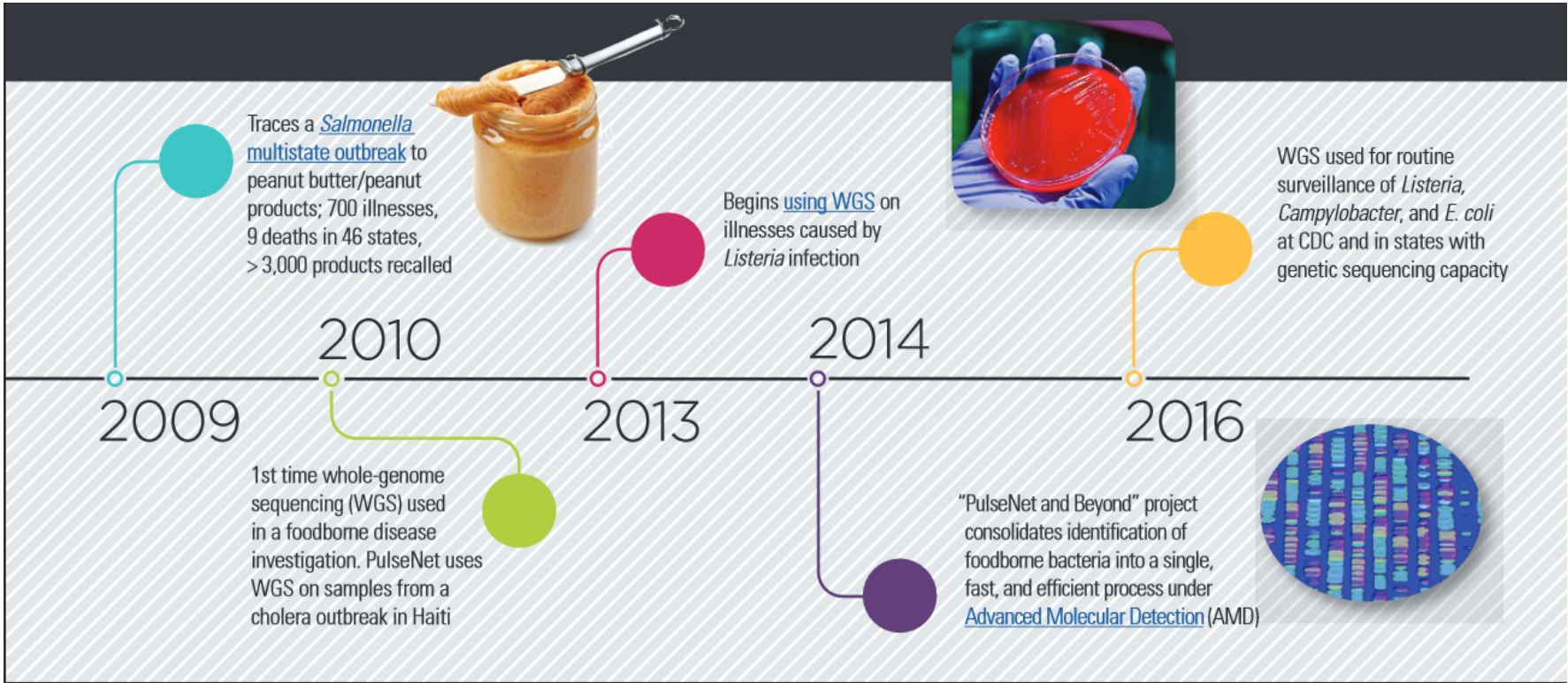
**1996** — CDC, the Association of Public Health Laboratories, federal partners, and public health labs in 4 states launch a new foodborne surveillance network called PulseNet

**2001** — PulseNet goes nationwide; all 50 state public health labs perform DNA fingerprinting of foodborne bacteria

**2002** — PulseNet wins 2nd Innovations in American Government Award (also won in 1999) for excellence and creativity in the public sector

**2006** — Identifies spinach as source of *E. coli* outbreak; 225 sickened and 5 deaths in 27 states; prompts nationwide recall

Traces a *Salmonella multistate outbreak* to peanut butter/peanut products; 700 illnesses, 9 deaths in 46 states, > 3,000 products recalled

Begins using WGS on illnesses caused by *Listeria* infection

WGS used for routine surveillance of *Listeria*, *Campylobacter*, and *E. coli* at CDC and in states with genetic sequencing capacity

**2009**

**2010**

**2013**

**2014**

**2016**

1st time whole-genome sequencing (WGS) used in a foodborne disease investigation. PulseNet uses WGS on samples from a cholera outbreak in Haiti

"PulseNet and Beyond" project consolidates identification of foodborne bacteria into a single, fast, and efficient process under Advanced Molecular Detection (AMD)

# Outline

- **Background**
- **Genomic Epidemiology**
  - Algorithms
  - Software
- **Example**

**The findings and conclusions in this report are those of the authors and do not necessarily represent the official position of the Centers for Disease Control and Prevention.**
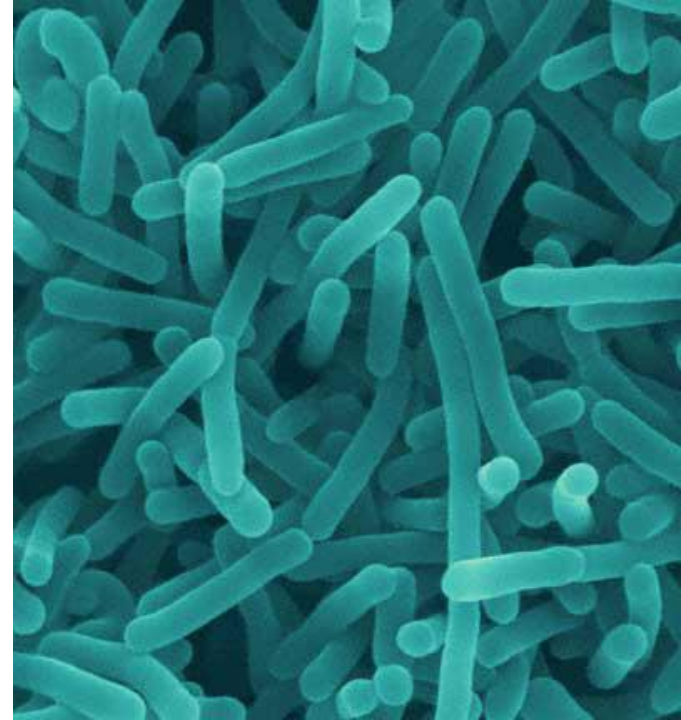
# Listeria pilot project

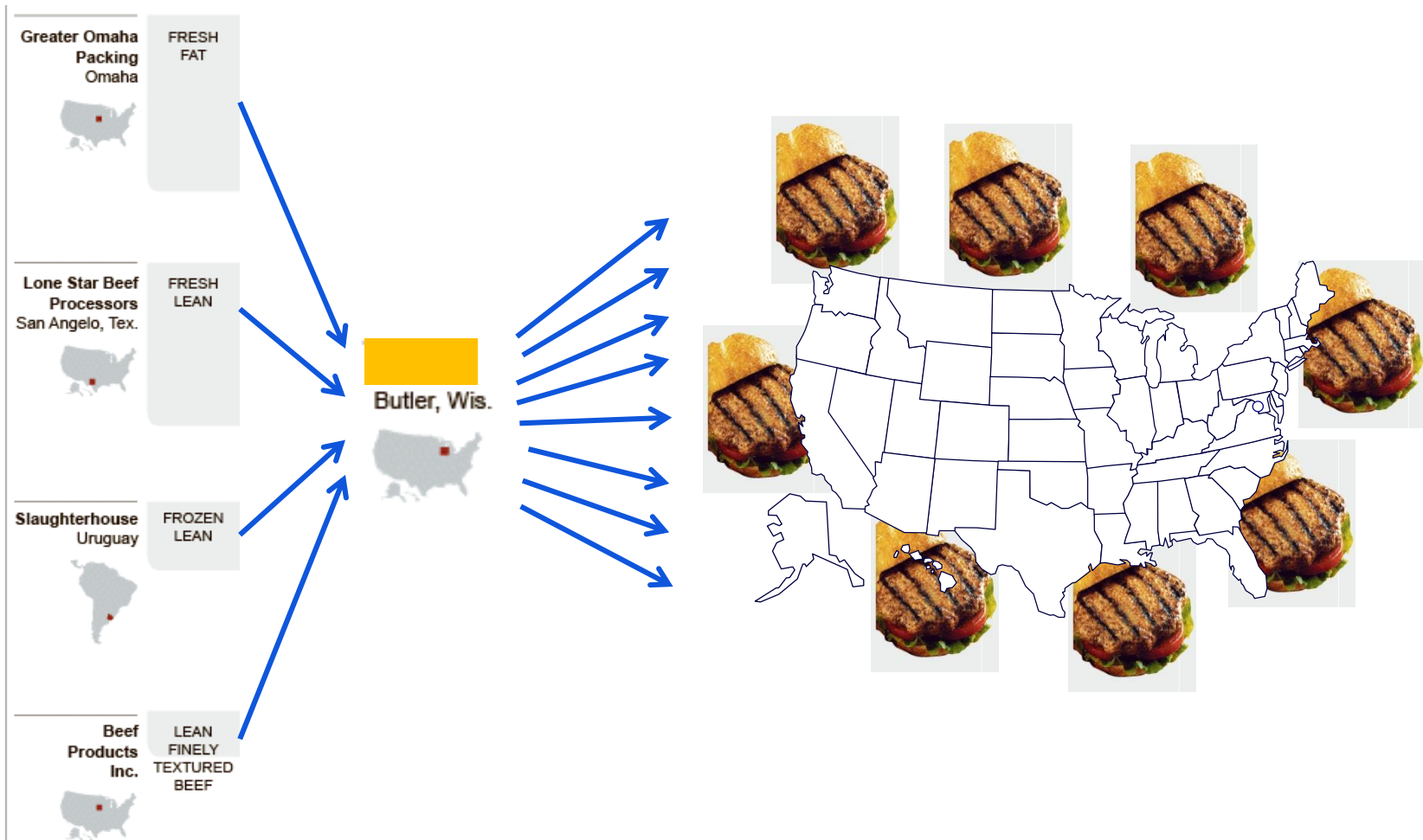**As told from a bioinformatician's perspective**



(It's an awesome perspective)

# Why *Listeria monocytogenes*?



- Illness is rare but serious, costly, and commonly outbreak associated
  - Estimated $2.8 billion in annual medical costs and lost productivity ($1.8 million/case)
- Current subtyping methods are not ideal
- Strong epidemiologic surveillance (Listeria Initiative)
- Strong regulatory component
- Listeria genome is fairly small, stable, and relatively easy to sequence and analyze. Most changes in the genome are due to point mutations and not phages.

# The Problem: Detecting Outbreaks in an Increasingly Globalized Food System



**Greater Omaha Packing** Omaha — FRESH FAT

**Lone Star Beef Processors** San Angelo, Tex. — FRESH LEAN

**Slaughterhouse** Uruguay — FROZEN LEAN

**Beef Products Inc.** — LEAN FINELY TEXTURED BEEF

Butler, Wis.

# Limitations of
# Pulsed-Field Gel Electrophoresis (PFGE)

# Limitation: Genetically Unrelated Isolate Might Appear Same by PFGE

# Limitation: Genetically Related Isolate Might Appear Different By PFGE

# Can genomics clear up this picture?

# How do we compare genomes?

# Three major methods we use

- Kmer-based: mile-high view (shredded paper)

- MLST-based: naked eye (book pages)

- SNP-based: microscope (book letters)


- The question in this analogy:
  how similar are these two books?

# kmers



- **Kmer:** a length of DNA *k* nucleotides long

1. Shred all reads in equal sizes *k*
2. How many kmers are in common?
3. Transform into a percentage **

** Known as the jaccard distance

# Kmers, jaccard distance

**CAAAAAAAAAAAAT**          **CAAAAAAAAAAAAG**

**Here, K=12**

| | | | | |
|---|---|---|---|---|
| CAAAAAAAAAAA | 1 | 1 | CAAAAAAAAAAA | |
| AAAAAAAAAAAA | 2 | 2 | AAAAAAAAAAAA | |
| AAAAAAAAAAAG | 3 | 4 | AAAAAAAAAAAG | |

Two out of four kmers different;
Jaccard distance = 2/4 = 0.5

# Example kmer tree

- http://www.ncbi.nlm.nih.gov/pathogens/
- Software: pathogen detection pipeline at NCBI



Mile-high view
7,800 *Listeria monocytogenes* genomes in this tree

# Kmer-based software

- NCBI Pathogen Detection Pipeline
  - Not available for individual use, but the results are comprehensive and public

- Mashtree
  - Based on min-hash, implemented in Mash

- SKA
  - Split Kmer Analysis



NCBI kmer trees screen shot taken Sept 23, 2016

https://www.ncbi.nlm.nih.gov/pathogens

https://github.com/lskatz/mashtree (latest version: 1.0.4; Katz et al 2019, *JOSS*)

https://github.com/simonrharris/SKA/releases (latest version: 1.0)

# How does Mash work?

"Sketch"

# How does Mash work?

"Sketch"

## Min-hash

| | |
|---|---|
| 2<br>5<br>24<br>33<br>34 | ← "min" |

This example: just keep five hashes

| |
|---|
| 2<br>5<br>24<br>33<br>34<br>60<br>66<br>... |

← sort

May or may not keep counts

| |
|---|
| 66 – 2<br>42 – 2<br>33 – 5<br>44 – 5<br>24 – 7<br>34 – 3<br>... |

← Filter low-count

| |
|---|
| 66 – 2<br>42 – 2<br>82 – 1<br>87 – 1<br>64 – 1<br>22 – 1<br>... |

Ondov et al, "Mash", Genome Biology. http://genomebiology.biomedcentral.com/articles/10.1186/s13059-016-0997-x

# How does Mash work?

"Distance" or "dist"

**Min-hash**

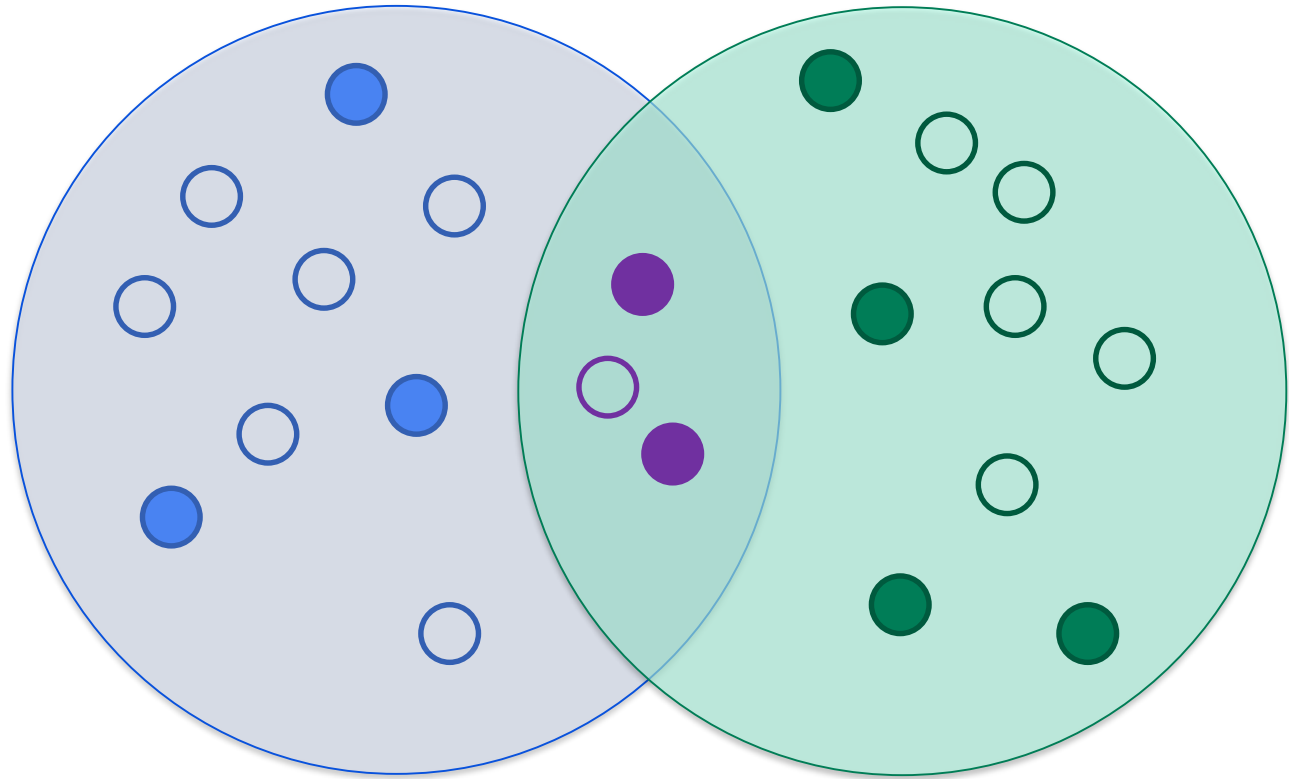| sketch1 | sketch2 |
|---------|---------|
| 2 | 2 |
| 5 | 7 |
| 24 | 24 |
| 33 | 33 |
| 34 | 34 |

Six different hashes, two differences.
Jaccard distance = 2/6 = 0.33

The resolution gets better with more hashes.

Ondov et al, "Mash", Genome Biology. http://genomebiology.biomedcentral.com/articles/10.1186/s13059-016-0997-x

# Min-hash visualization

# Mashtree

## What it is and what it isn't

| Is | Isn't |
| --- | --- |
| Builds trees | Infers phylogeny |
| Fast | Slow |

## When to use it

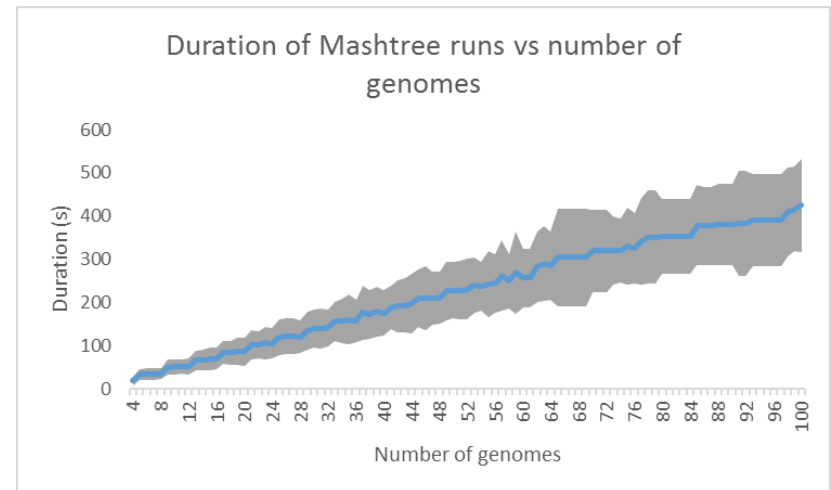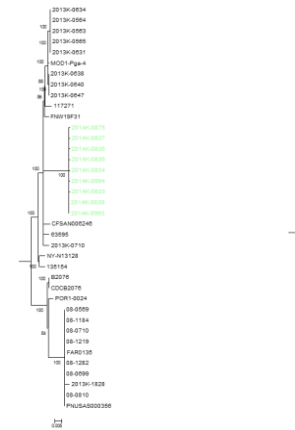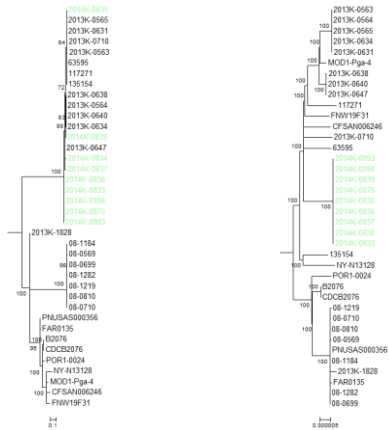| Use it when | Don't use it when |
| --- | --- |
| Need fast estimate | Need solid results |
| Need to know a good reference genome | Inferring phylogenetic relatedness |
| Large, diverse dataset | Not diverse or not large dataset |

# Mashtree is fast

- I had a tree of > 1500 genomes and ran Mashtree on the genomes of every clade with fewer than 101 taxa.

- The forward Illumina read of every genome was analyzed.

- Grey shading indicates the range of durations. (next slide)



0.050

# Mashtree is fast

- I had a tree of > 1500 genomes and ran Mashtree on the genomes of every clade with fewer than 101 taxa.

- The forward Illumina read of every genome was analyzed.

- Grey shading indicates the range of durations.

Duration of Mashtree runs vs number of genomes
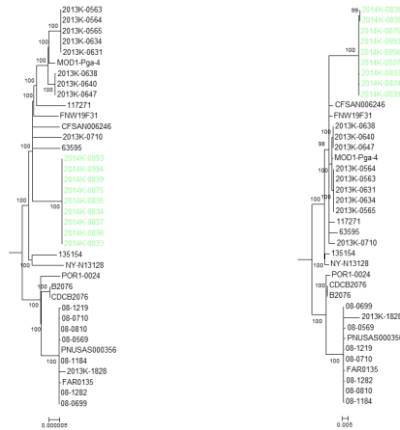
# The Mashtree v0.06 accuracy



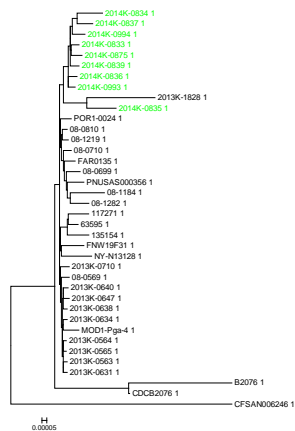Lyve-SET
Sn = 100%
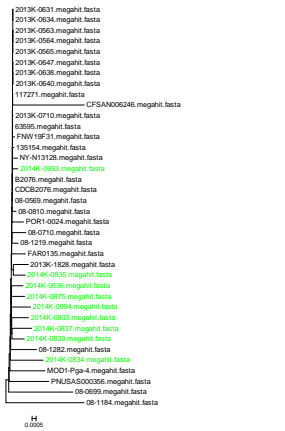Sp = 100%

kSNP3
Sn = 100%
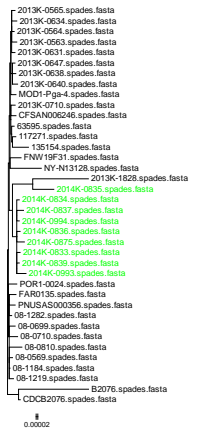Sp = 58%

RealPhy
Sn = 100%
Sp = 100%

Snp-Pipeline
Sn = 100%
Sp = 100%

Mash v0.06
Raw reads
min_depth: 5x
Sn = 100%
Sp = 97%

Mash v0.06
Megahit asm
59-3141 contigs
Sn = 78%
Sp = 100%

Mash v0.06
SPAdes asm
23-46 contigs
Sn = 100%
Sp = 97%

1409MLJN6-1
$n_{pos}$ = 9
$n_{neg}$ = 29
Dataset from *Katz et al*, "Lyve-SET",
2017, MGEN

◼ part of outbreak

# Mashtree is command line

```
# Installation
$ cpanm -L ~ Mashtree
$ export PERL5LIB=$PERL5LIB:$HOME/lib/perl5

# Usage
$ mashtree.pl --help

# Execution
$ mashtree.pl --numcpus 12 --genomesize 4700000 \
  *.fastq.gz \
  [*.fasta] [*.gbk] [*.fasta.gz] [*.gbk.gz] \
  > mashtree.dnd
```

https://github.com/lskatz/mashtree

# MLST



- **MLST:** multilocus sequence typing
- **Locus:** a place in a genome. Plural: **loci**

- Identify a set of loci (genes) in the genome
- Compare each locus in a genome against the set of loci
- Count differences and the number of loci compared
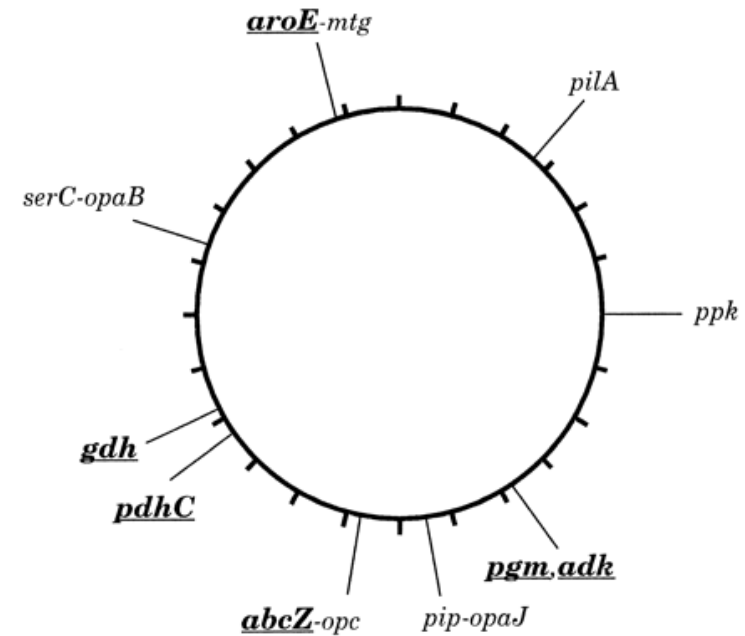
Different kinds
- 7-gene MLST
- wgMLST (whole genome MLST)
- cgMLST (core genome MLST)
- … and more

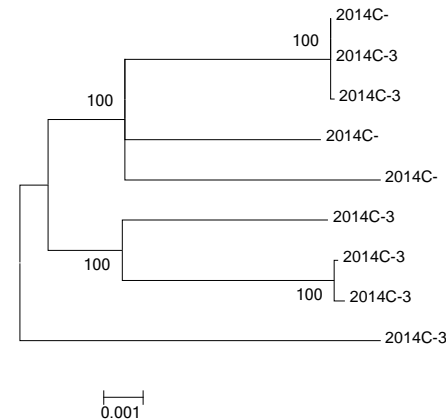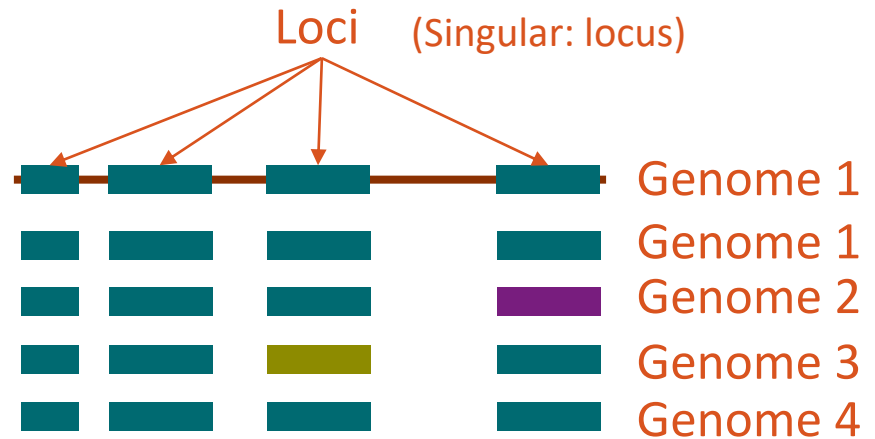Image credit: Wikipedia.org
Software: BioNumerics

# 7-gene MLST

- Choose about seven loci in the genome
- Compare all genomes based on these seven loci
- This profile of alleles is called a **sequence type (ST)**



Maiden et al 1998 *PNAS*

# Animation of MLST

0. Assemble the genome
1. Identify the loci
2. Call alleles
3. Compare with other genomes and their alleles
4. Create a phylogeny

- Note: many methods do not require an assembly and these are called **assembly-free methods**.



Loci    (Singular: locus)

Genome 1

Genome 1
Genome 2
Genome 3
Genome 4

Phylogeny

2014C-
2014C-3
2014C-3
2014C-
2014C-
2014C-3
2014C-3
2014C-3
2014C-3

100
100
100
100

0.001

# Whole-genome MLST

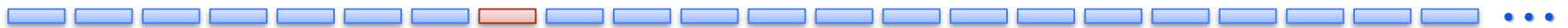~one locus per 1,000 nucleotides (nt) in the genome.

Different species have different sizes

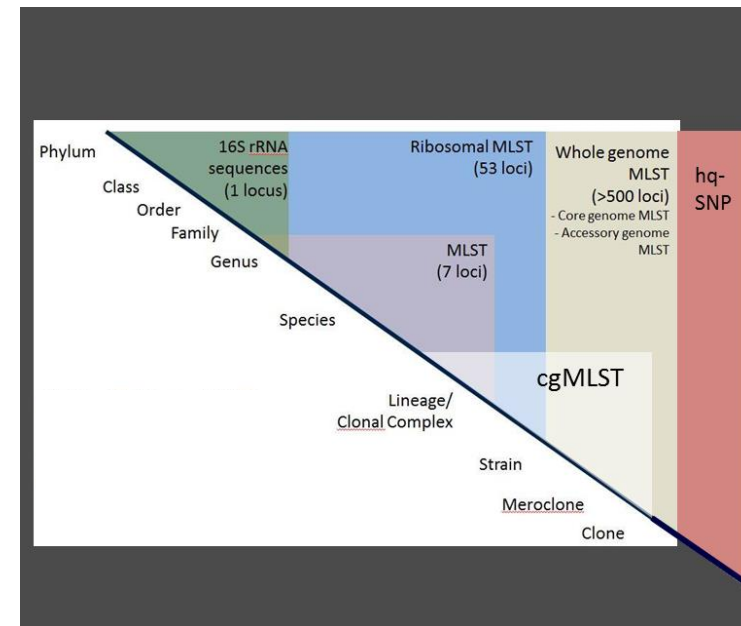e.g., *L. monocytogenes* has ~3,000,000 nt and ~3,000 loci

Strain A

Strain B

Strain C

# Flavors of multilocus sequence type analysis



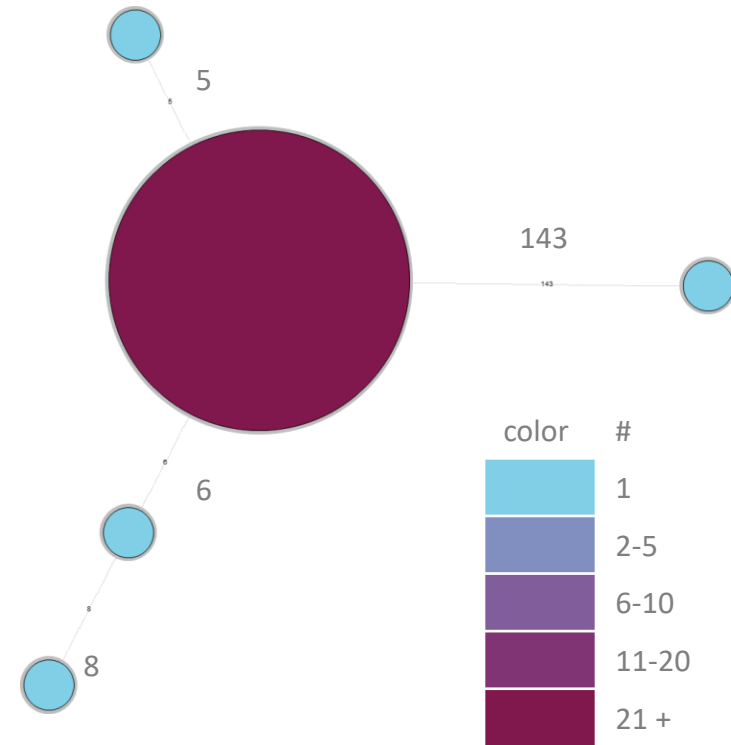Maiden et al *Nat Rev Microbiol.* 2013 **11:**728-36

- Subsets of genes can be used to identify genus/species and lineage (rMLST/ MLST)

- Core genome MLST are the genes that are in common in vast majority of genomes belonging to a genus species (for Listeria – 1748 genes belong to core and are present in ~98% of isolates tested)

# Example wgMLST tree

- Larger circles represent more with the same sequence type (ST)
- 4800 loci represented
- Distances shown on the connecting lines

- The style of tree shown is called a **minimum spanning tree**
- wgMLST can also be displayed in a conventional tree



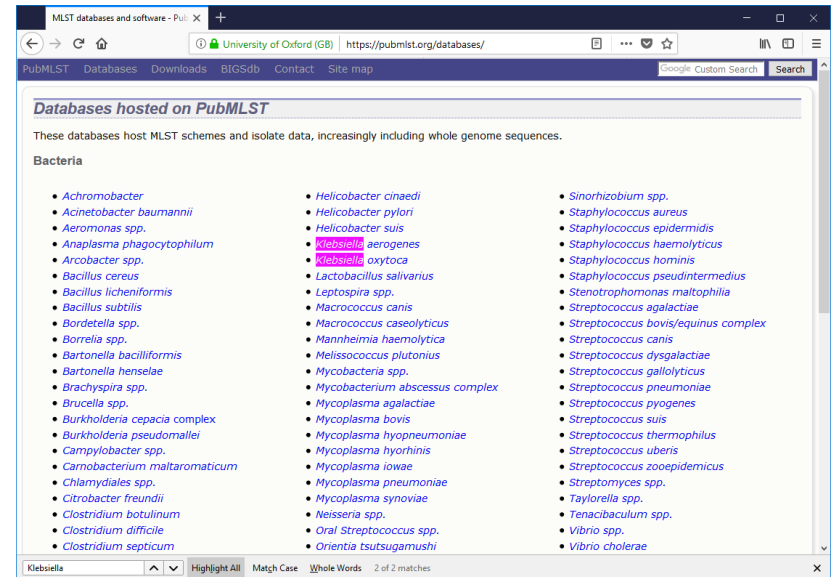| color | # |
|---|---|
|  | 1 |
|  | 2-5 |
|  | 6-10 |
|  | 11-20 |
|  | 21 + |

# MLST software

- StringMLST
  - Compare kmers of raw reads against a database

- BioNumerics
  - Graphical user interface

- SRST2, Ariba
  - Map raw reads onto database

- *mlst*
  - BLAST genome assembly against database

- Mentalist
  - Command line, meant for wgMLST schemes



Image taken from http://www.applied-maths.com/applications/wgmlst

For more information: Page et al 2017, "Comparison of Multi-locus Sequence Typing software for next generation sequencing data."

# MLST Resources

- **Main MLST site:**
  **https://pubmlst.org/**

- **BigsDB manual:**
  **http://bigsdb.readthedocs.io/en/latest/concepts.html**

- **API:**
  **https://pubmlst.org/rest/**

- **Also see:**
  - https://enterobase.warwick.ac.uk/
  - http://bigsdb.web.pasteur.fr/listeria/

- Jolley & Maiden 2010, *BMC Bioinformatics* **11:**595
- Jolley *et al.* (2017) *Database* **2017:** bax060

# SNPs



- Compare individual letters in a query genome against the reference genome
- hqSNP: high-quality SNP (ie, high confidence)
- hqSNP indicates some high threshold, e.g.,
- 10x coverage
- 75% consensus

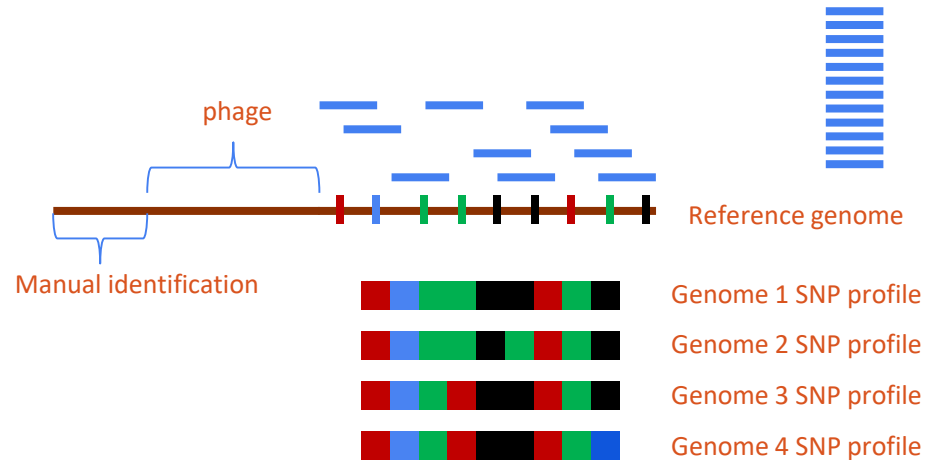Image taken from geneious.com

# SNP analysis

0. Pre-processing
   a) Identification of troublesome regions
   b) Read cleaning

1. Mapping

2. SNP calling
   a) % consensus
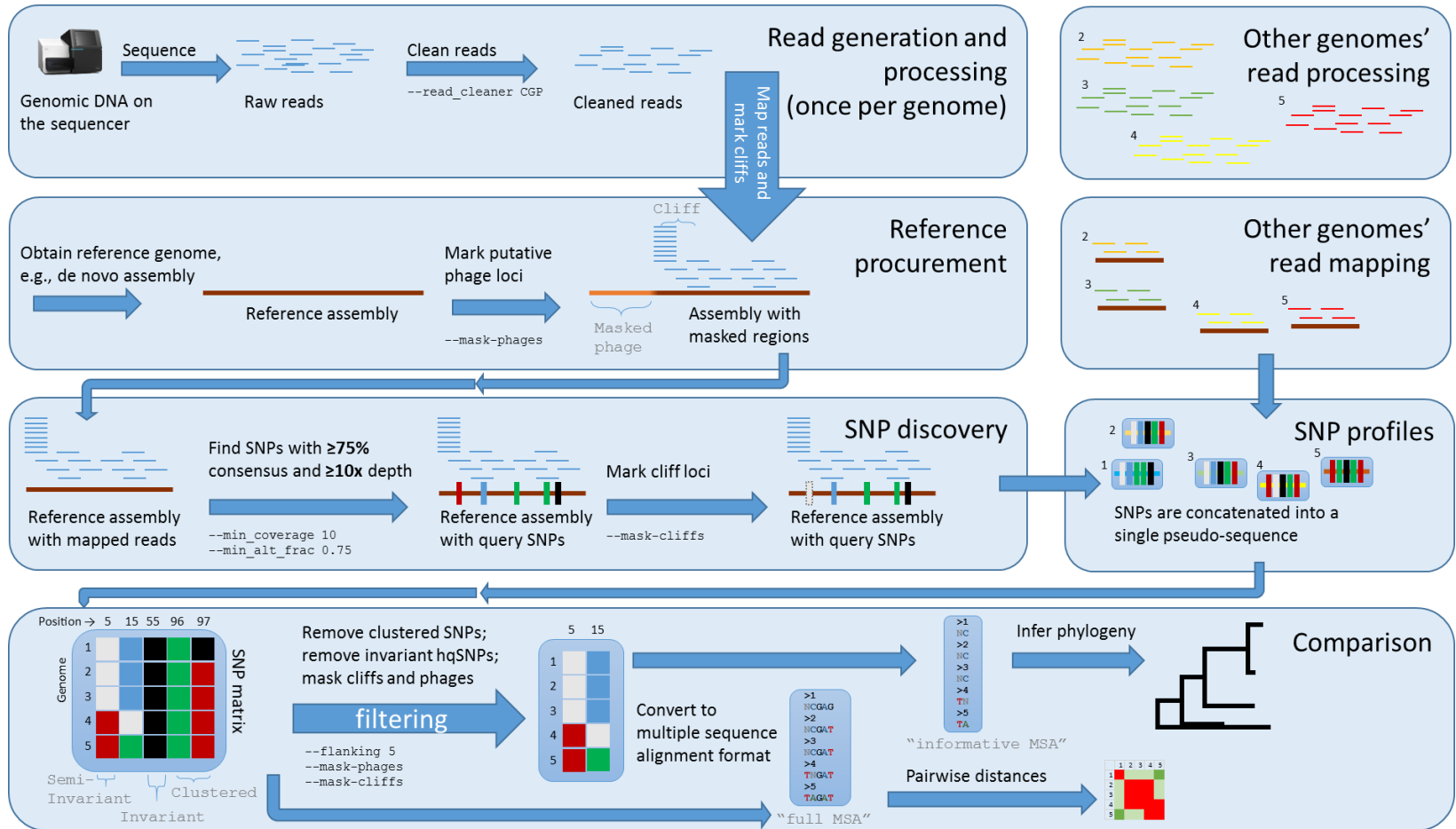   b) x depth
   c) Other filters

3. Phylogeny inference

phage

Manual identification

Reference genome

Genome 1 SNP profile

Genome 2 SNP profile

Genome 3 SNP profile

Genome 4 SNP profile

2014C-
100
2014C-3
2014C-3
100
2014C-
2014C-
Phylogeny
2014C-3
2014C-3
100
2014C-3
100
2014C-3

0.001

# More details



https://github.com/lskatz/lyve-SET
Katz et al. (2017) A comparative analysis of the Lyve-SET phylogenomics pipeline for genomic epidemiology for foodborne pathogens. Frontiers in Microbiology 8: 375.

# SNP software

- Lyve-SET
  - Optimized for outbreak surveillance.

- SNP-Pipeline
  - FDA SNP pipeline.  Optimized for regulatory workflow.  Optimized for speed and accuracy of SNPs.

- SNVPhyl
  - Public Health Agency of Canada.  Graphical User Interface in Galaxy.

**Each bioinformatician to have their own personal short-read aligner by 2016**

Posted on March 23, 2015 by jovialscientist

OXFORD, UK.  The Bioinformatics Society ("BS" for short) have declared that they will reach their aim of every bioinformatician having their own personal short-read aligner by the end of 2016, *The ScienceWeb* have learned.

There are approximately 28,362 scientists globally who identify themselves as being "bioinformaticians" or "computational biologists" (those who identify themselves as "bioinformagicians" have been excluded – not just from this analysis, but from life in general).  A recent survey of short-read aligners identified 23,872 different software tools, all of which basically do the same thing.

"We're almost there!" exclaimed base-pair hyper-bot Hang Li from the Broad Institute.  "As soon as I published that paper on the Ferris Bueller transform, I knew the field would take off!  And it has – we have one valuable publication and 23,871 incremental improvements" finished the Hang Li AI, a 7-dimensional intelligence that exists only in the minimal amount of memory need to represent a human.

The field of bioinformatics sequence analysis has been criticised by other areas of science for basically solving the same 3 problems over and over again, sometimes with only a marginal improvement and often with a marked deterioration in quality.

https://thescienceweb.wordpress.com/2015/03/23/each-bioinformatician-to-have-their-own-personal-short-read-aligner-by-2016/
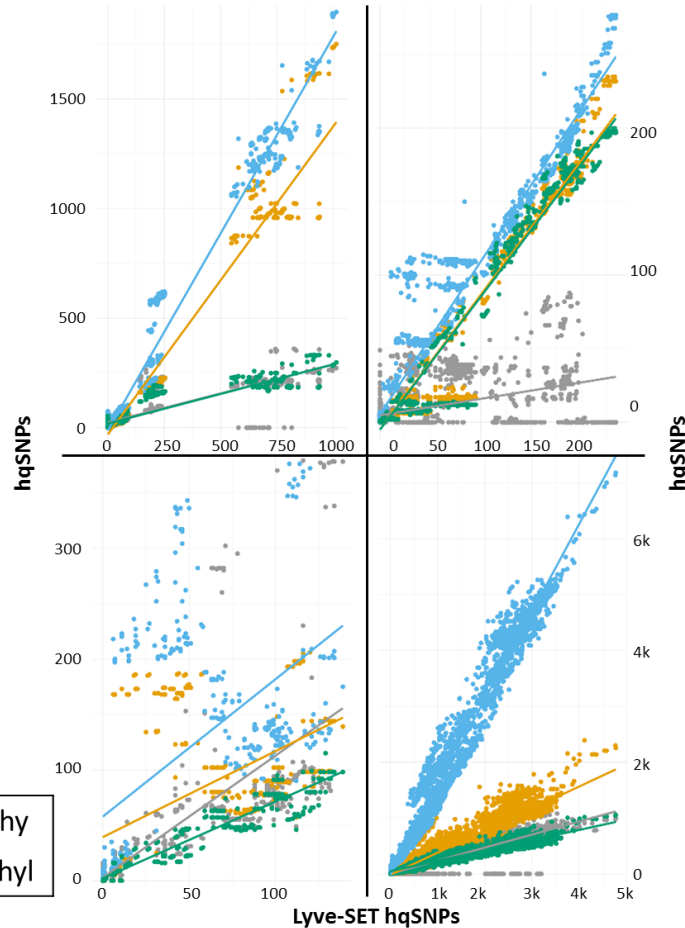
# Installation and sample run

```
$   cd ~/bin/
$   git clone https://github.com/lskatz/lyve-SET
$   cd Lyve-SET
$   git checkout v1.1.4f
$   make install
$   export PATH=$PATH:~/bin/lyve-SET/scripts
#  You may also add this to your bash profile
$   echo >> ~/.bash_profile "export PATH=$PATH:~/bin/lyve-SET/scripts"
$   which launch_set.pl
$   set_test.pl lambda lambda --numcpus 4
# Takes about two minutes
$   ls lambda/msa/tree.dnd
```

# Comparison of Lyve-SET with other SNP pipelines

**L. monocytogenes**

| Pipeline | y=mx+b | R² |
|---|---|---|
| kSNP | y=0.26+24 | 0.69 |
| RealPhy | y=1.14+31 | 0.96 |
| SNP-Pipeline | y=1.8x-13 | 0.97 |
| SNVPhyl | y=0.27x+19 | 0.58 |

**S. enterica**

| Pipeline | y=mx+b | R² |
|---|---|---|
| kSNP | y=0.11x+4.7 | 0.23 |
| RealPhy | y=0.92x-5 | 0.95 |
| SNP-Pipeline | y=1.0x+5.4 | 0.96 |
| SNVPhyl | y=0.91-5.1 | 0.94 |

**E. coli**

| Pipeline | y=mx+b | R² |
|---|---|---|
| kSNP | y=1.1x+2.9 | 0.43 |
| RealPhy | y=0.78+39 | 0.27 |
| SNP-Pipeline | y=1.2x+58 | 0.3 |
| SNVPhyl | y=0.69x+2.1 | 0.92 |

**C. jejuni**

| Pipeline | y=mx+b | R² |
|---|---|---|
| kSNP | y=0.23x+4 | 0.89 |
| RealPhy | y=0.4x-15 | 0.88 |
| SNP-Pipeline | y=1.6x-17 | 0.97 |
| SNVPhyl | y=0.18+49 | 0.92 |

Legend:
- kSNP
- RealPhy
- SNP-Pipeline
- SNVPhyl

Each data point is a SNP distance as determined by Lyve-SET (x-axis) and the distance of an alternative SNP pipeline (y-axis). The slope indicates the number of SNPs per Lyve-SET SNP.

# SNPs overlayed on MLST loci

# Comparison with whole-genome MLST (Listeria monocytogenes only)



| Max Lyve-SET hqSNPs | y=mx+b | R² |
|---|---|---|
| 1013 | y=0.18x+27 | 0.58 |
| 254 | y=0.79x+1.2 | 0.98 |
| 94 | y=0.78x+1.7 | 0.96 |

Katz et al 2017, *Lyve-SET*, Frontiers in Microbiology.

# Which algorithm should you use?

| | Kmer-based | wgMLST | hqSNP |
|---|---|---|---|
| Diversity | ✓✓ | ✓ | ✗✗ |
| Outbreak-level resolution | ✗ | ✓✓ | ✓✓ |
| Further genomic information | ✗ | ✓ | ✓ |
| Minimal upfront effort | ✓ | ✗✗ | ✓ |
| Fast | ✓✓ | ✓✓ | ✗ |
| Easy to use for anyone | ✗ | ✓ | ✗ |

# Fun examples

More information can be found in the virtual lab talk given on Jan 7, 2019: https://youtu.be/YPnU63Le53Y?t=1234

# Multistate outbreak of farmstead cheeses

# How to read a phylogeny



Scale bar

Direction of evolution

0.01

Outgroup1 ] **taxon**

Outgroup2 ] **taxon**

Root
(LCA of
all taxa)

Vertical distance is
irrelevant

100

taxon3

taxon4

**taxa**

taxon5

100

taxon6

taxon7

Percent
confidence in
hypothetical
ancestor

Hypothetical last common
ancestor (**LCA**) of taxa 3-7
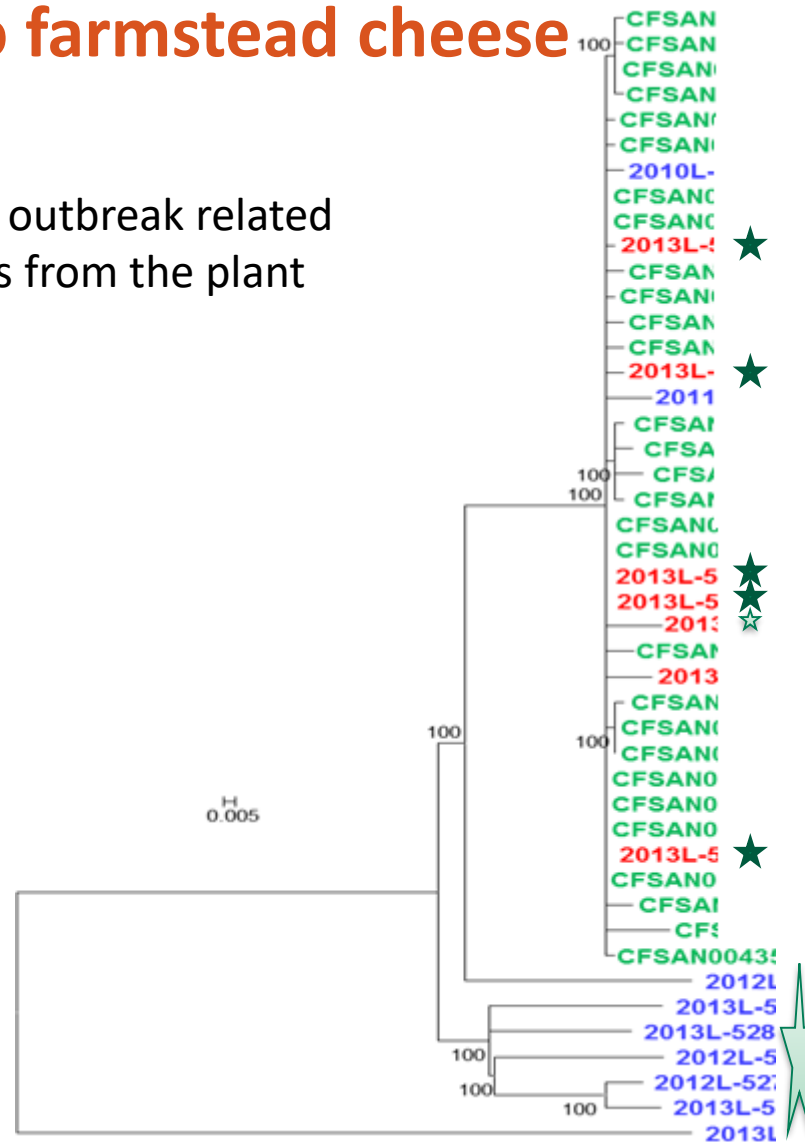
# 2013 outbreak linked to farmstead cheese

Red= epi-related clinical isolates

Blue= retrospective clinical cases or not outbreak related

Green= historical environmental isolates from the plant

★ Exposure

☆ No Exposure

# Phylogenetically related outbreak of unknown etiology, December 2013

# In conclusion

- **WGS provides high resolution**

- **We have many tools for differing levels of resolution**

- **We can and have used it on outbreak investigations**

# Questions?

Lee Katz
gzu2@cdc.gov

lskatz
github.com/lskatz

**College of Agricultural & Environmental Sciences**
*Center for Food Safety*
**UNIVERSITY OF GEORGIA**

National Center for Emerging and Zoonotic Infectious Diseases
Division of Foodborne, Waterborne, and Environmental Diseases