# STEC Whole Genome Sequencing Project

**Eija Trees, Ph.D., D.V.M.**

Chief, PulseNet Next Generation Subtyping Methods Unit

16th Annual PulseNet Update Meeting
August 29th, 2012

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

CDC

# Objectives of the Presentation

- ❑ **Describe the STEC genome sequencing project**
- ❑ **Describe the sequence diversity seen in the STEC population**
- ❑ **Way forward**

# Outline

- ❑ **Background and Objectives of the Project**
- ❑ **Materials and Methods**
- ❑ **Results**
- ❑ **Conclusions**
- ❑ **Future Plans**
- ❑ **100K Pathogen Genome Sequencing Project**

# Background

- **CDC Bioinformatics Blue Ribbon Panel in June 2011**
- **2011 end-of-year  funding from CDC  OD, OID and NC**
  - Implemented a proof-of-concept bioinformatics core group through an outside contract
    - In 2012 focus on two pilot projects:
      - MicrobeNet: development of shared bioinformatic services and web platforms
      - Next Generation PulseNet: Whole genome sequencing, comparative genomics,  molecular epi

# Next Generation PulseNet Project Objectives

- **Perform whole genome sequencing, assembly and comparative analysis of 250 STEC strains**
  - Reference database
- **Understand diversity in the STEC population and determine correlation between WGS data, PFGE and epi**
- **Mine the resulting data for targets to predict strain type, virulence and antimicrobial resistance quickly and at high resolution**

# Next Generation PulseNet Project Objectives (cont'd)

❑ **Public health benefits:**

- Standardized workflows/SOPs for next generation sequencing across multiple platforms

- Reusable tools, techniques and software "pipelines" for a wide range of comparative genomics and molecular epidemiologic applications

- Builds capacity for rapid genome-level molecular epidemiology

# MATERIALS AND METHODS

# Strain Selection

❑ **Selection criteria should help us to understand variation:**

   1. within an outbreak or during a carrier state
   2. among epidemiologically unrelated isolates within a serotype and between serotypes

# Strain selection (cont'd)

❑ **Variation within an STEC outbreak or a carrier state**

- 10 isolates each from 5 outbreaks
  - O157:H7 (3), O111:NM (1), O145:NM (1)
  - Different types of outbreaks: clonal vs. polyclonal, short vs. long lasting, point source vs. vehicle never identified
- 11 isolates recovered from the same patient in the course of 2½ months

# Strain selection (cont'd)

❑ **Variation among epidemiologically unrelated STEC strains**

- 25 strains each from serogroups O26, O111, and O121
- 5 strains each from serogroups O45, O69, O91, O103, O118, O145
- 1 strain each from serogroups ranked in prevalence #11-22
- 1 representative each from past historical O157 (and non-O157) outbreaks (36 in total)
- Sporadic O157 isolates representing common PFGE patterns (23 in total)

❑ **Other *E. coli* pathotypes: ETEC (6), EPEC (3), EAggEC (3), EIEC (1)**

# Whole genome sequencing and assembly

❑ **Illumina HiSeq 2x100 bp; average coverage 400x**

- Pros: highly accurate reads with massive coverage
- Cons: reads short
- Assembled with CLC Bio; average of 200 contigs

❑ **Subset of 99 isolates sequenced with Pacific Biosciences SMRT (on-going)**

- Pros: longer reads (up to 10 kb)
- Cons: error rate high, multiple SMRT cells required per strain for adequate coverage
- Hybrid assembly with error correction using Illumina reads being optimized

# K-mer -based Approach to Identify SNPs

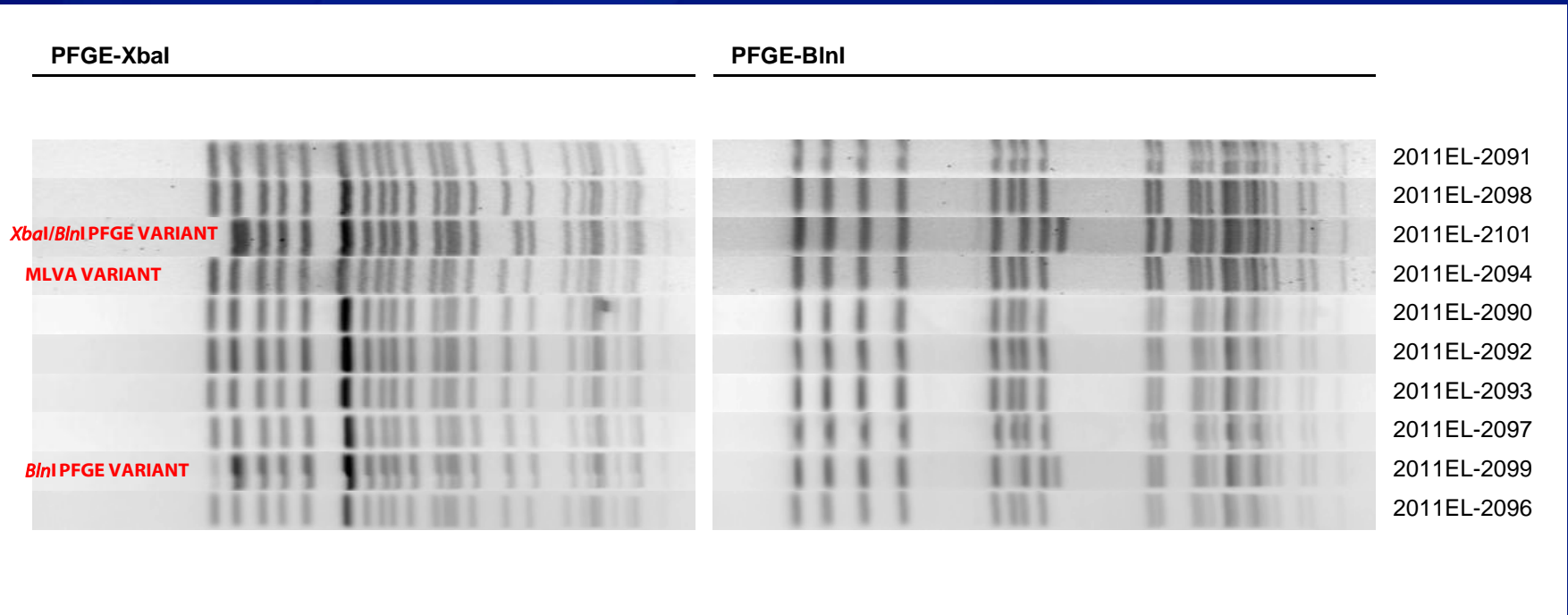| Suffix array to find all k-mers, where k=25 | Find candidate variant positions where k-mer exists, and central bases varies | Align k-mers with possible variations against target *E.coli* reference genome | Eliminate repetitive bases | Phylogenetic analysis (trees) based on variations |
|---|---|---|---|---|

12 bases     13th     12 bases

TTGATAGGGCAA**A**AGCGCCGATTTT

- When k-mer 25 is used, the 13th base position is used for determining if there is a variation

- No prior sequence assembly or multiple sequence alignment required
- Fast completion of analysis
- Phylogenies less affected by regions that have undergone strong selection, deletions or horizontal gene transfer
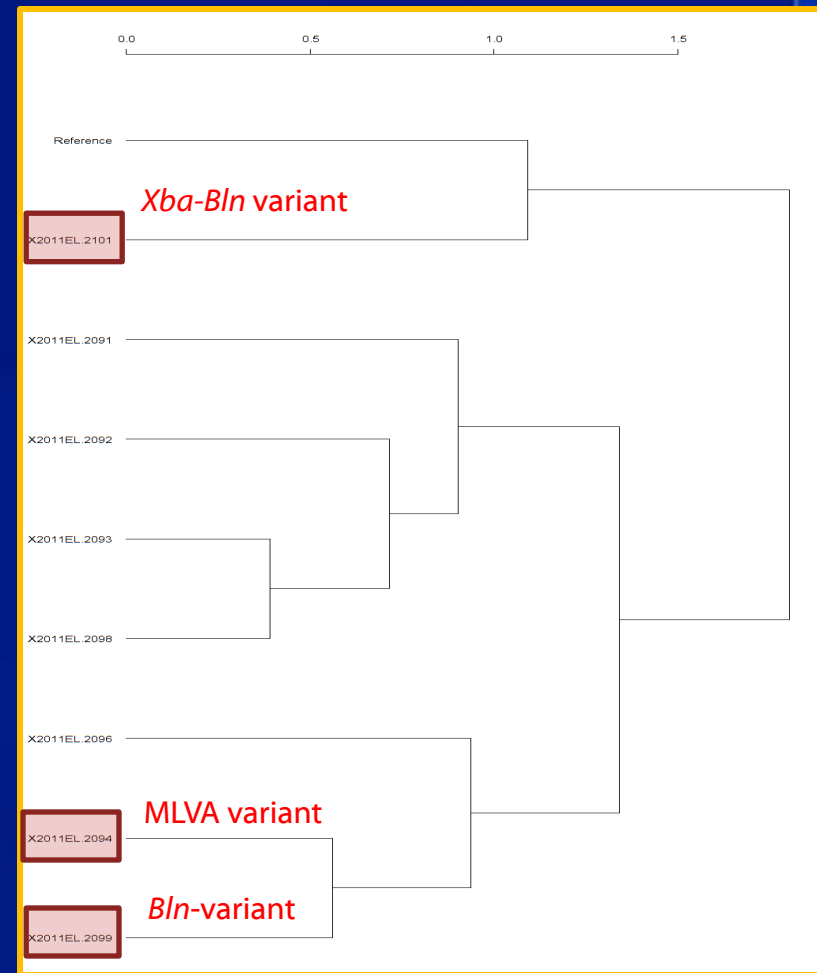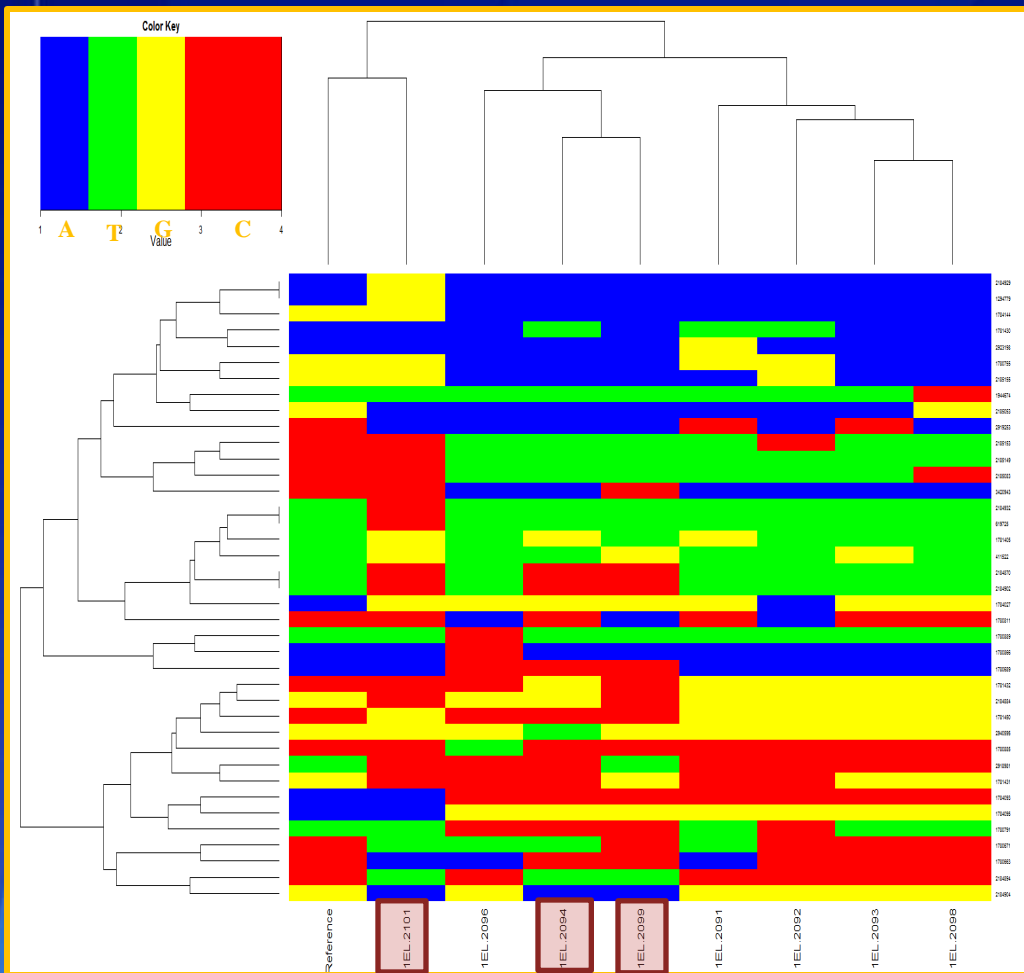- Dependent on high quality data
- Reference genome required

# RESULTS

# Probing Diversity within an Outbreak – O157:H7 in MN Daycare (2008)



**PFGE-XbaI**

**PFGE-BlnI**

*XbaI*/*BlnI* PFGE VARIANT

MLVA VARIANT

*BlnI* PFGE VARIANT

2011EL-2091
2011EL-2098
2011EL-2101
2011EL-2094
2011EL-2090
2011EL-2092
2011EL-2093
2011EL-2097
2011EL-2099
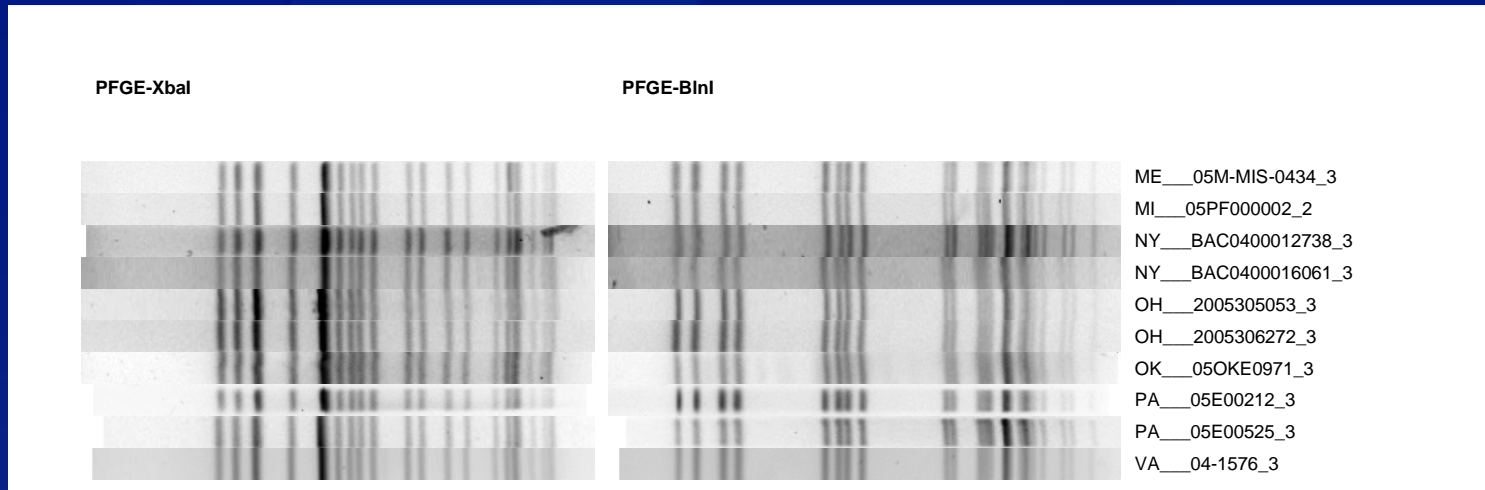2011EL-2096

- **Relatively clonal short lasting point source outbreak**

# Probing Diversity within a Clonal Point Source O157 Outbreak – MN Daycare



**Reference Sample:** 2011EL-2090_O157:H7_EXHX01.0589_EXHA26.1376_main MLVA
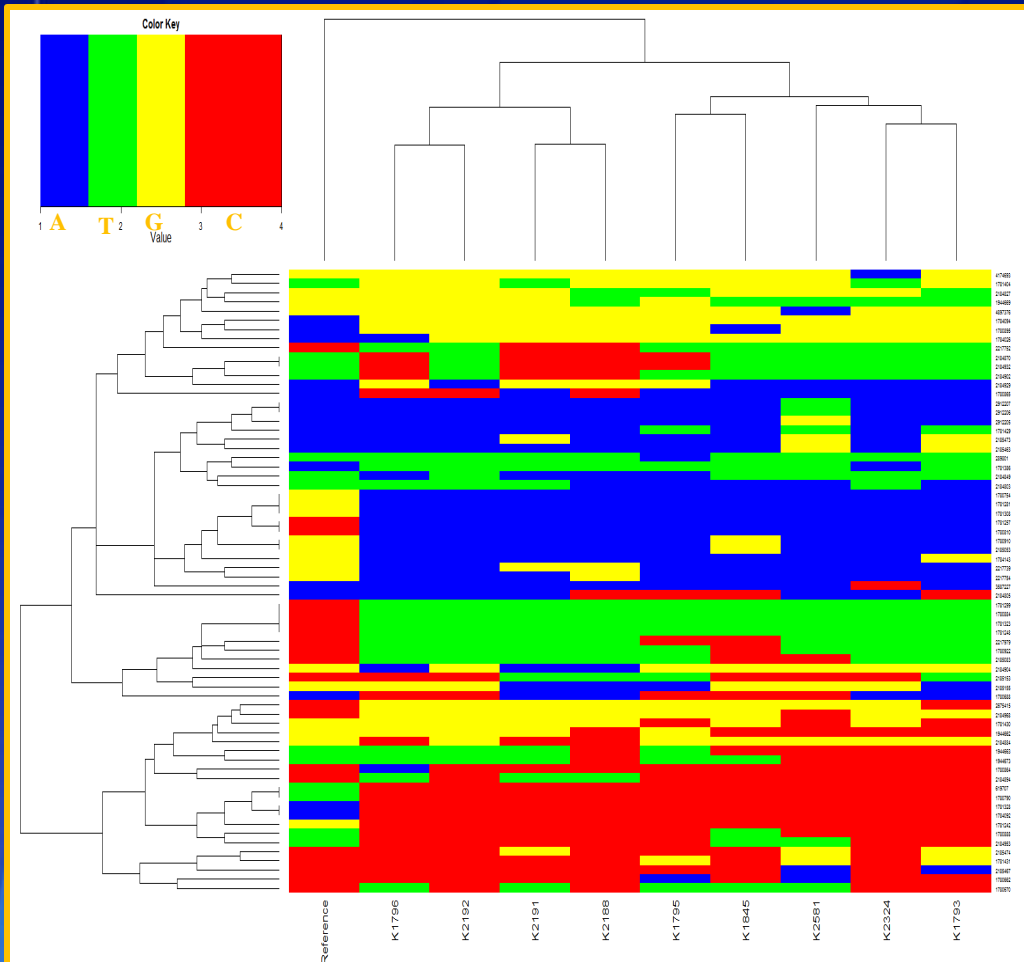- The total number of SNPs found across all the samples from this outbreak at 100% frequency is **39**.

# Probing Diversity within Clonal Long Lasting O157:H7 Outbreak (2004-2005)



- 110 2-enzyme matches in 20 states within 15 months
- Rare PFGE type
- An indistinguishable new MLVA pattern
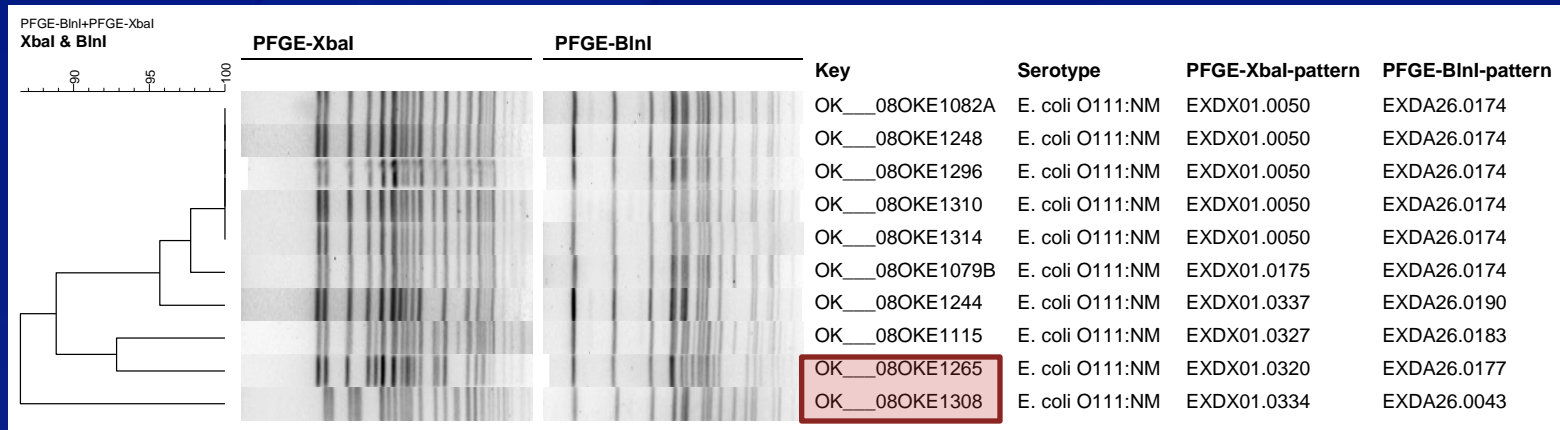- Vehicle never identified

# Probing Diversity within a Clonal Long Lasting O157 Outbreak – Not resolved



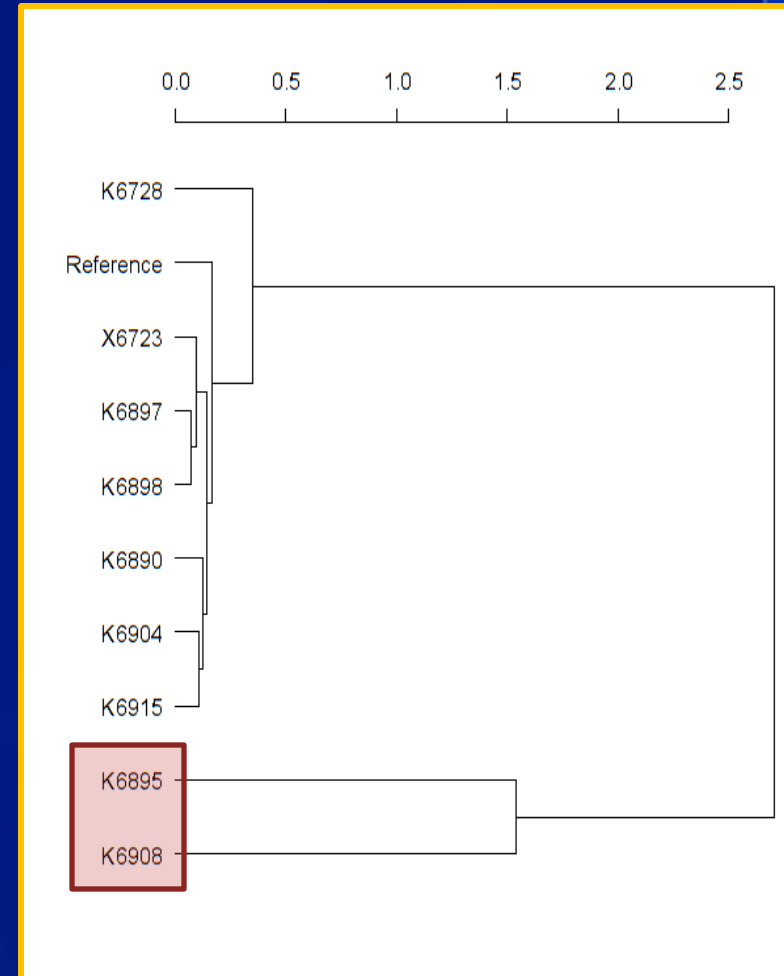**Reference Sample: K1792_O157:H7_EXHX01.0086_EXHA26.0576**
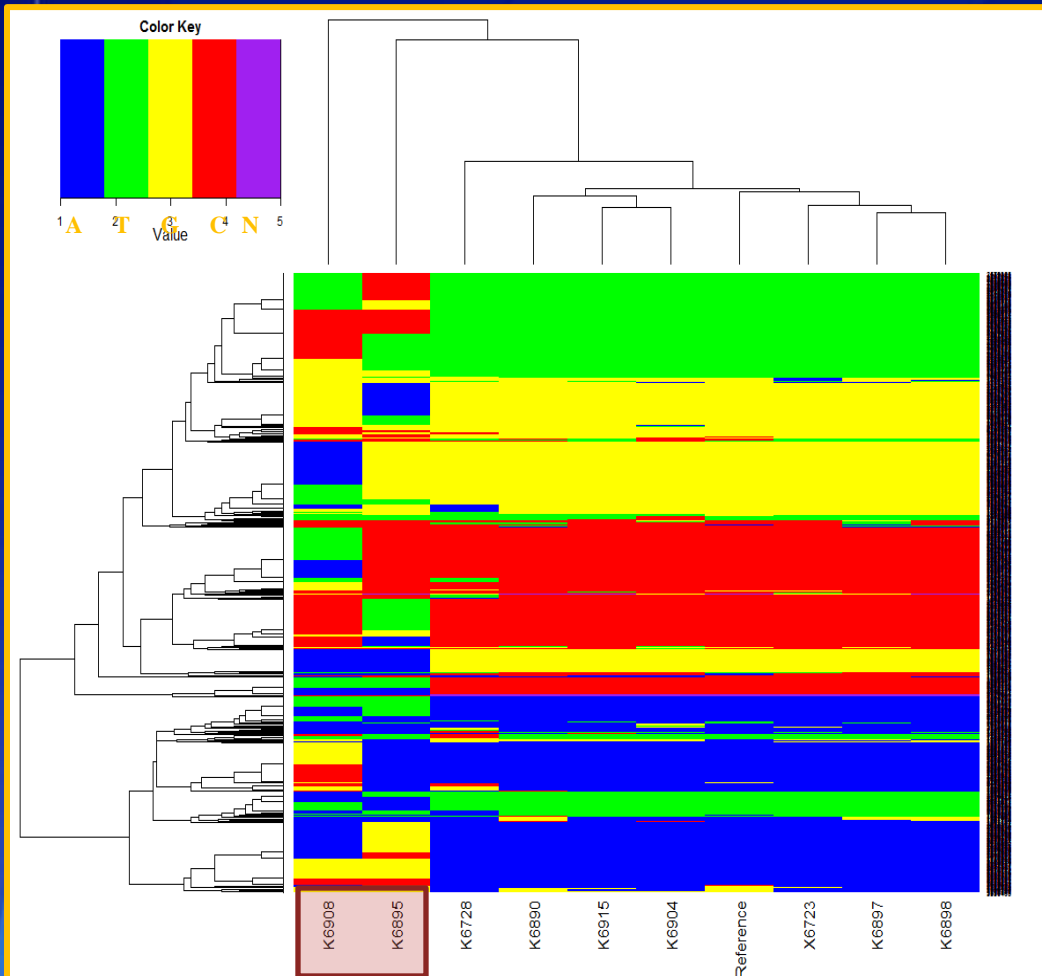- The total number of SNPs found across all the samples from this outbreak at 100% frequency is 68

# Probing Diversity within Polyclonal Outbreak – O111:NM in OK Restaurant (2008)



| Key | Serotype | PFGE-XbaI-pattern | PFGE-BlnI-pattern |
|---|---|---|---|
| OK____08OKE1082A | E. coli O111:NM | EXDX01.0050 | EXDA26.0174 |
| OK____08OKE1248 | E. coli O111:NM | EXDX01.0050 | EXDA26.0174 |
| OK____08OKE1296 | E. coli O111:NM | EXDX01.0050 | EXDA26.0174 |
| OK____08OKE1310 | E. coli O111:NM | EXDX01.0050 | EXDA26.0174 |
| OK____08OKE1314 | E. coli O111:NM | EXDX01.0050 | EXDA26.0174 |
| OK____08OKE1079B | E. coli O111:NM | EXDX01.0175 | EXDA26.0174 |
| OK____08OKE1244 | E. coli O111:NM | EXDX01.0337 | EXDA26.0190 |
| OK____08OKE1115 | E. coli O111:NM | EXDX01.0327 | EXDA26.0183 |
| OK____08OKE1265 | E. coli O111:NM | EXDX01.0320 | EXDA26.0177 |
| OK____08OKE1308 | E. coli O111:NM | EXDX01.0334 | EXDA26.0043 |

- 1 *Xba* variant, 4 *Xba-Bln* variants
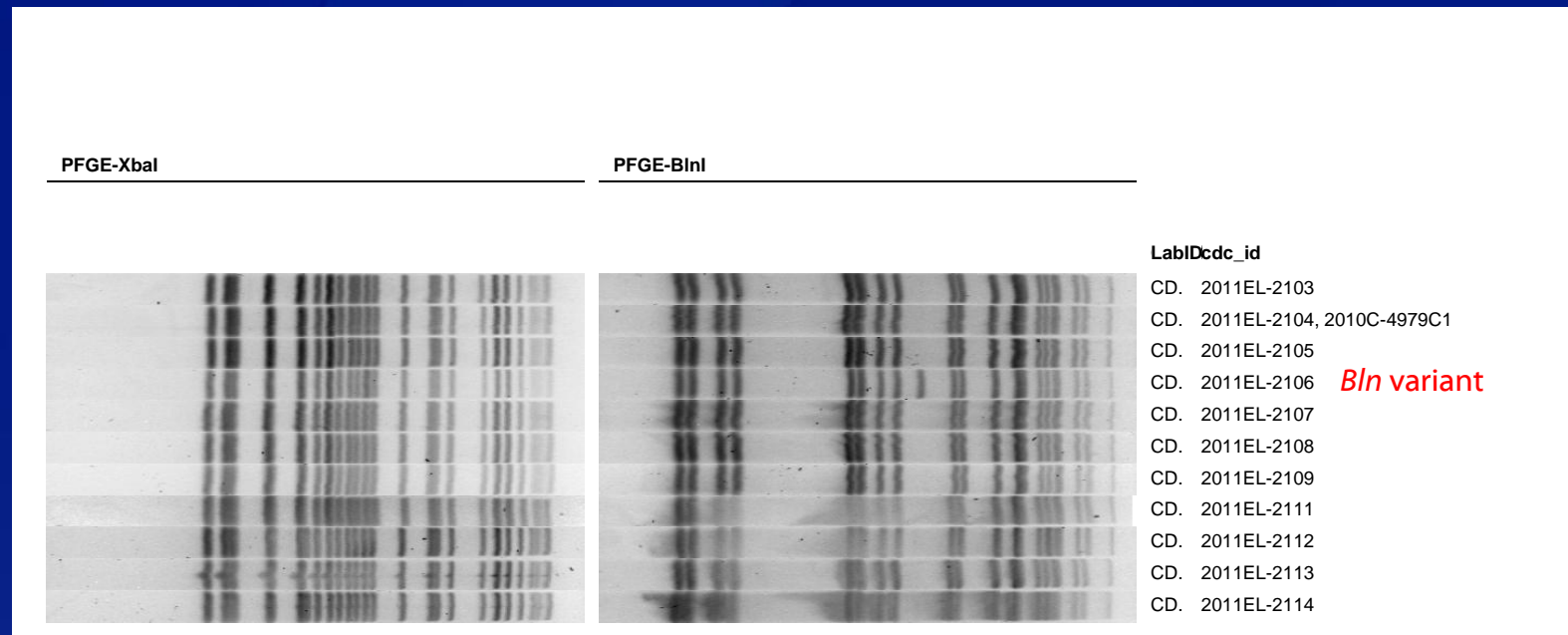- Vehicle not identified

# Probing Diversity within Polyclonal O111 Outbreak – OK Restaurant



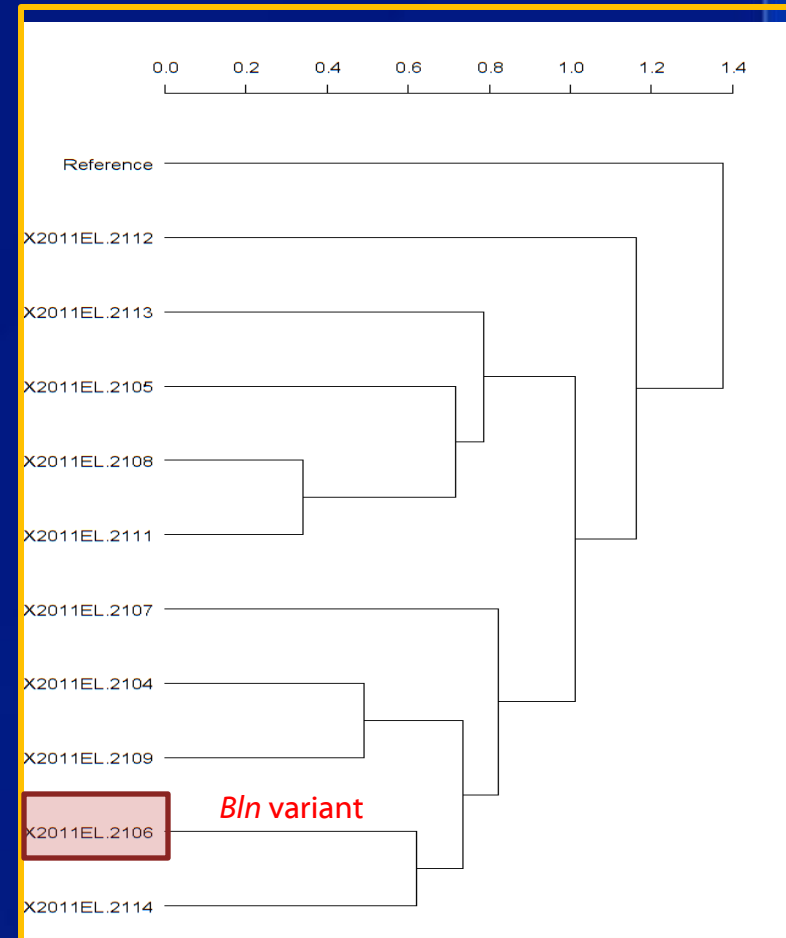**Reference Sample:** **K6722_O111:IsolatetoCDC_EXDX01.0050_EXDA26.0174**
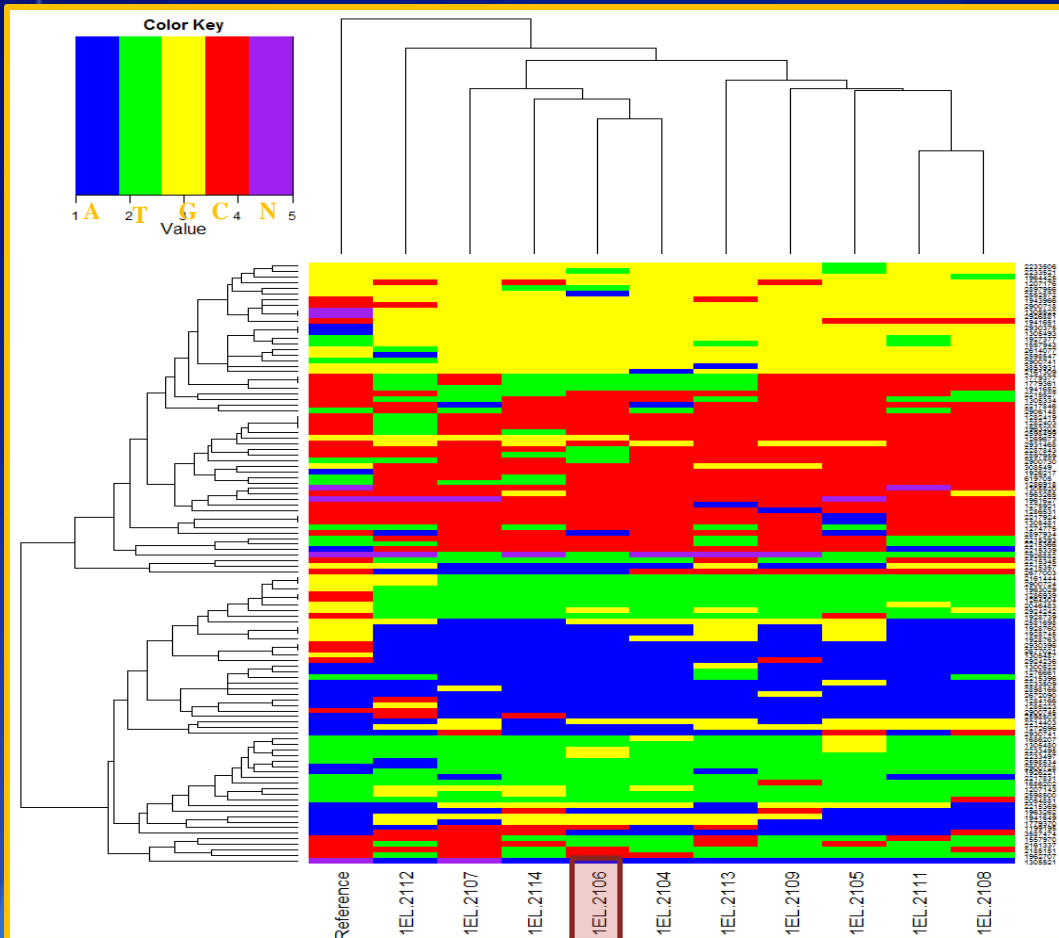- The total number of SNPs found across all the samples from this outbreak at 100% frequency is **947**

# Probing Diversity During a Carrier State



- **Collected from the same person within a period of 2½ months**
- **An indistinguishable MLVA pattern**

# Probing Diversity during a Carrier State –Serial Isolates from the Same Person
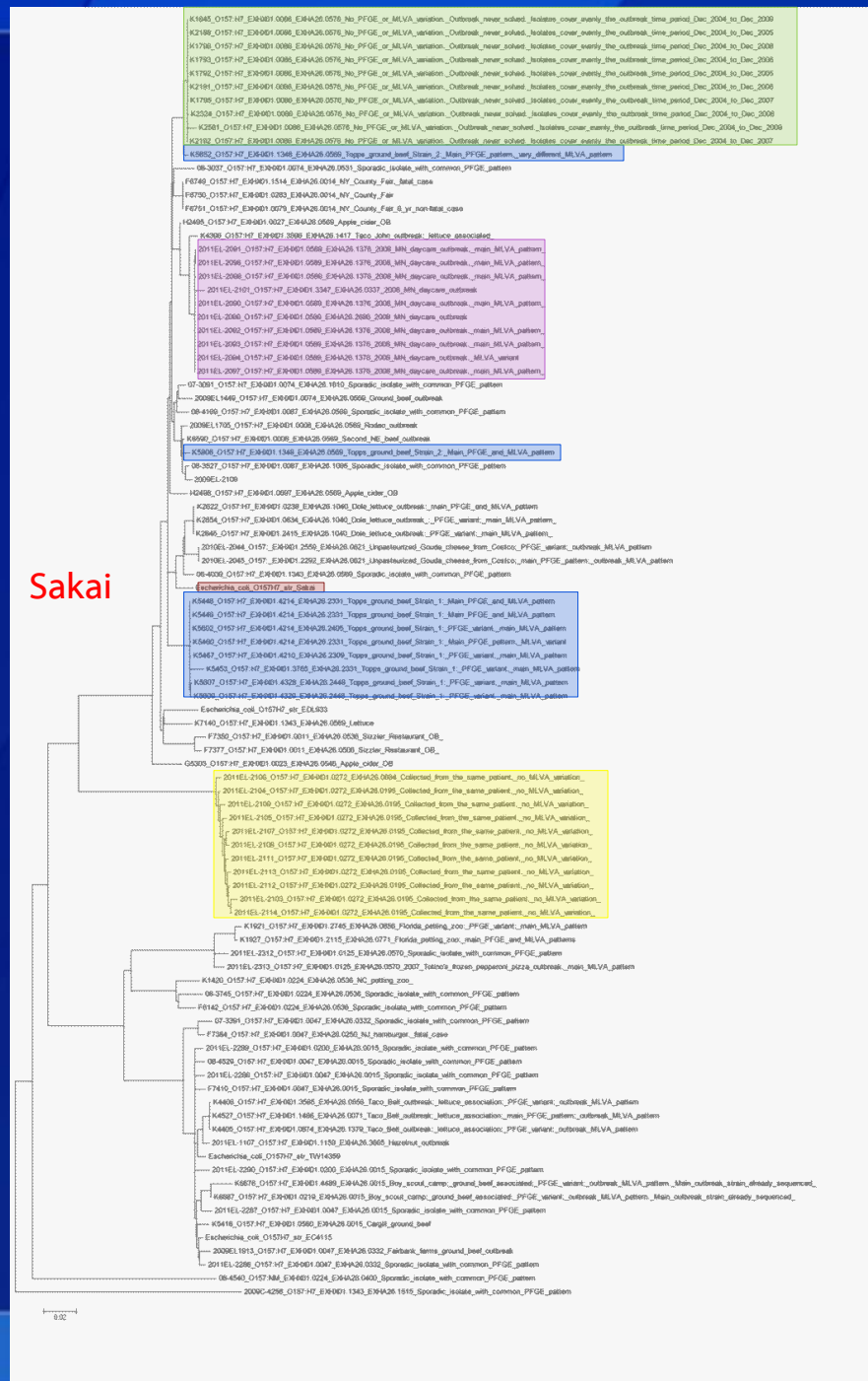


**Reference Sample: 2011EL-2103_O157:H7_EXHX01.0272_EXHA26.0195**
- The total number of SNPs found across all the samples from this "outbreak" at 100% frequency is 108

**K-mer –based Phylogenetic Tree for STEC O157**

93 to 2101 SNPs compared to the reference

Sakai

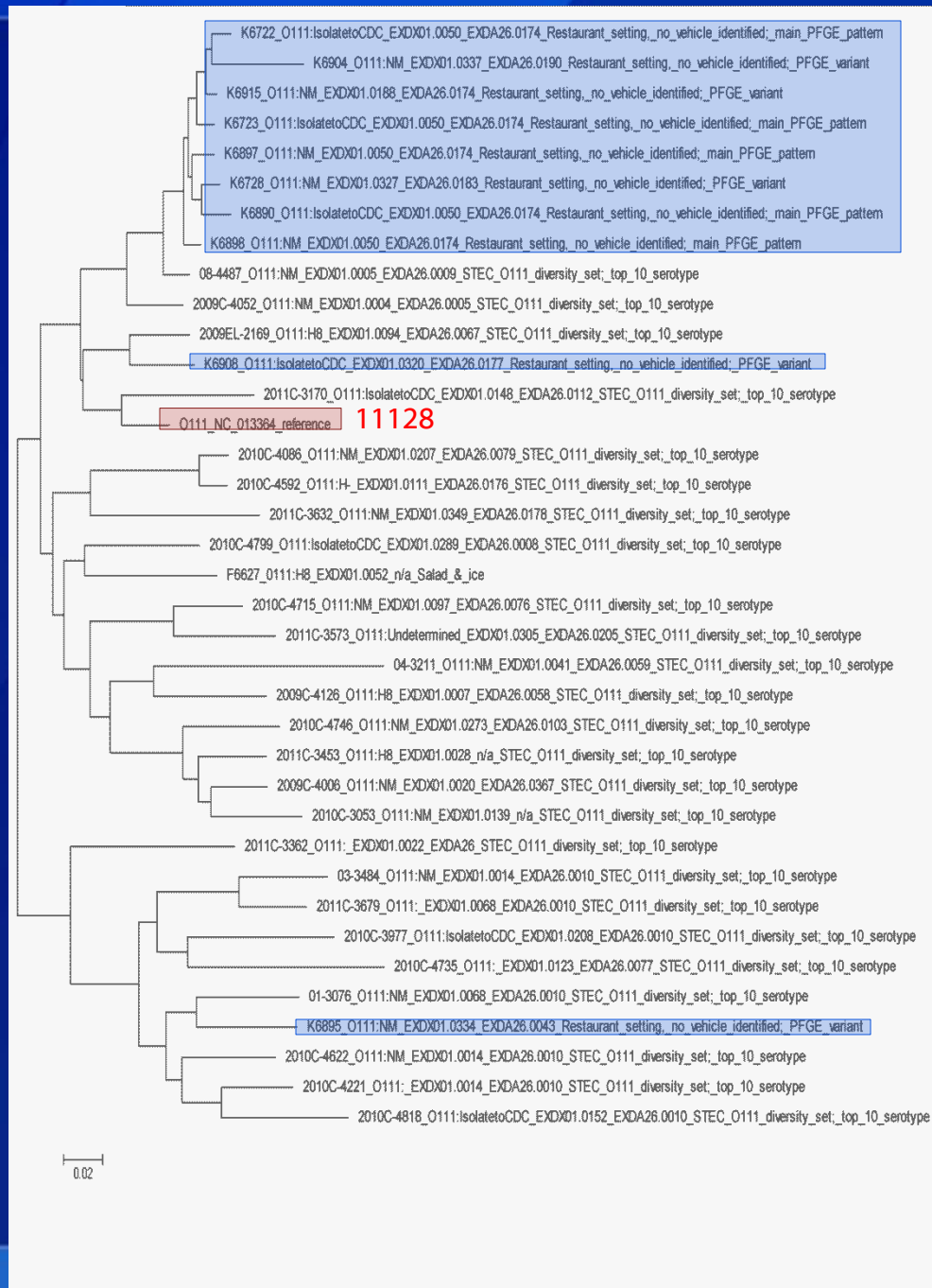Outbreak 1 (Not solved, 2004-05)

Outbreak 2 (Daycare, 2008)

Outbreak 3 (Topps ground beef, 2007)
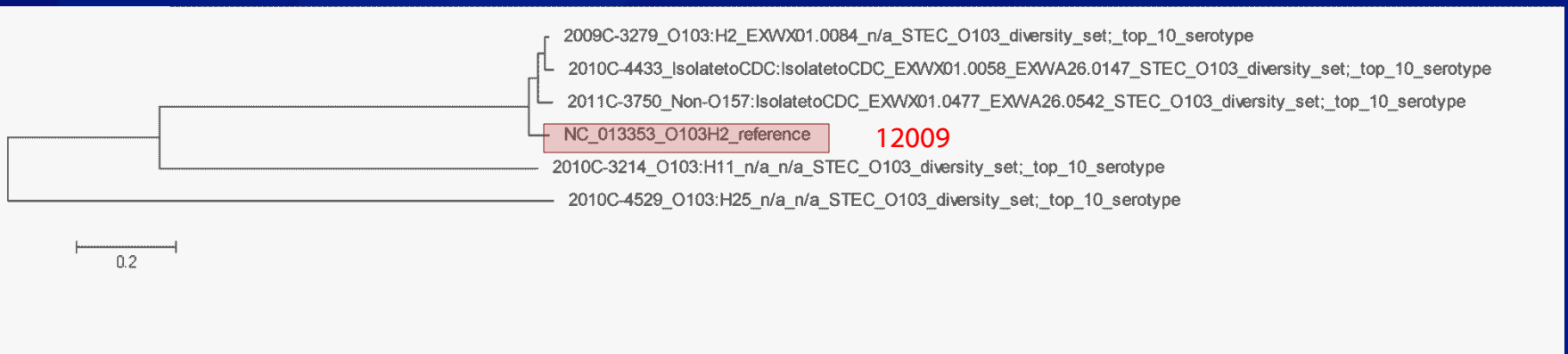
Long term carrier patient

**K-mer -based Phylogenetic Tree for STEC O111**

195 to 1229 SNPs compared to reference

Restaurant outbreak, 2008

# K-mer -based Phylogenetic Tree for STEC O103



2009C-3279_O103:H2_EXWX01.0084_n/a_STEC_O103_diversity_set;_top_10_serotype

2010C-4433_IsolatetoCDC:IsolatetoCDC_EXWX01.0058_EXWA26.0147_STEC_O103_diversity_set;_top_10_serotype

2011C-3750_Non-O157:IsolatetoCDC_EXWX01.0477_EXWA26.0542_STEC_O103_diversity_set;_top_10_serotype

NC_013353_O103H2_reference    12009

2010C-3214_O103:H11_n/a_n/a_STEC_O103_diversity_set;_top_10_serotype

2010C-4529_O103:H25_n/a_n/a_STEC_O103_diversity_set;_top_10_serotype

0.2

1789 to 20,094 SNPs
compared to the reference

# Conclusions Based on the Illumina Data

- **K-mer –based clustering appears to have good correlation with epidemiological and PFGE data**
- **Variation within a (relatively) clonal outbreak $< 100$ SNPs**
- **The variation between unrelated strains of the same serotype hundreds to ~ 2000 SNPs**
- **The variation between serotypes > 10,000 SNPs**

# Future Plans

- **Improve the assemblies using PacBio sequences**
    - Make the draft genomes publicly available
- **Additional sequencing to close reference genomes**
- **Define a cross serotype gene set to be used to determine strain subtype, virulence and antimicrobial resistance profiles ('super' MLST)**
    - Presence / absence
    - SNPs
- **Add prospective strains from outbreaks to the database**

# Questions to be Answered

❑ **Which clustering method will work for surveillance?**

1. K-mer –based clustering as a primary tool with gene-based 'super' MLST as a secondary tool OR

2. Gene-based 'super' MLST as a primary tool

❑ **How to standardize data across the platforms?**

▪ Different error profiles

❑ **How to get an actionable report out from raw reads?**

# 100K Pathogen Genome Sequencing Project

❑ **Dilemma: the safety and security of the world food supply is hindered by lack of food-related genomes**

❑ **Partnership between BGI and UC Davis**

- ▪ http://100kgenome.vetmed.ucdavis.edu

❑ **Goal: sequence 100,000 foodborne pathogen isolates within 5 years**

- ▪ Sequences will be made publicly available

❑ **Steering committee**

- ▪ BGI
- ▪ UC Davis School of Veterinary Medicine
- ▪ FDA
- ▪ CDC
- ▪ Agilent Technologies, Inc.
- ▪ Mars, Inc.

# 100K Pathogen Genome Sequencing Project (cont'd)

## Isolate priority list

| Tier 1 | Tier 2 | Tier 3 |
|---|---|---|
| *Salmonella* | *Yersinia* | Toxigenic bacilli |
| *Campylobacter* | *Shigella* | Norovirus |
| *E. coli* | *Clostridium* | Hepatitis |
| *Vibrio*<br>*Listeria* | *Enterococcus*<br>*Cronobacter* | Rotovirus |

- Seeking isolates from diverse types of food, distinct environments, diverse regions of the world and longitudinal isolate collections

# Acknowledgements

# Thank You for Your Attention! Questions?

**For more information please contact Centers for Disease Control and Prevention**

1600 Clifton Road NE, Atlanta, GA  30333
Telephone: 1-800-CDC-INFO (232-4636)/TTY: 1-888-232-6348
E-mail:  cdcinfo@cdc.gov     Web:  http://www.cdc.gov

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.

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

CDC