

---

# WHOLE GENOME SEQUENCING @ ISS

## EU RL for *E. coli*

Valeria Michelacci

Copenhagen, January 29th 2015



**Istituto Superiore di Sanità - Dip. Sanità Pubblica Veterinaria e Sicurezza Alimentare**  
**EU Reference Laboratory for *E.coli***



# Availability of Samples and WGS technology

---

**Strains collection**  $\cong$  2000 pathogenic *E. coli* strains

**Clinical samples**  $\cong$  150-200 specimens analysed/year (mainly BD, HC and HUS)

**Sequencing platforms @ ISS:**



2 x ION PGM



1 x MiSeq



1 x 454



Several ABI 3130 (Sanger)

# Sequencing throughput and IT

---

## Data production

### TODAY

70 strains fully sequenced

**TOMORROW** (projection to be achieved in the next two years)

200 strains /year + 20 metagenomics samples

## Data storage

TorrentServer, Central servers, firewalls, intrusion prevention systems, automated back-up services, storage units

### TODAY

15 Tb dedicated to sequence storage (backup and mirroring)

**TOMORROW** (projection to be achieved in the next two years)

60 Tb (local)



# Data analysis: Locally running softwares

---



- *de novo* assembly
- Alignment of sequences, production of VCF files, production of dendrograms
- MLST
- Search for interesting genes

**USER-FRIENDLY INTERFACE, Slow processing, RAM needed**



- *de novo* assembly
- Search for interesting genes
- Alignment of sequences, production of VCF files

**BUILT IN THE ION TORRENT TECHNOLOGY PACKAGE**

# Data analysis: web servers

---



- Species identification
- *de novo* assembly tools
- VirulenceFinder
- ResFinder
- MLST
- SNPs tree and newly developed NGS-driven phylogenetic tools

## **FREE, USER-FRIENDLY WEB INTERFACE**



- *de novo* assembly tools
- BLAST search of interesting genes
- Alignment of sequences, production of VCF files, production of dendrograms

## **OPEN SOURCE, USER-FRIENDLY WEB INTERFACE, OPEN FOR INTRODUCTION OF CUSTOMIZED TOOLS, ELECTION PLATFORM FOR DEVELOPING AND SHARING OF NEW TOOLS**



# Galaxy instance @ ISS

The screenshot displays the Galaxy web interface at galaxy.iss.it. The main workspace shows a phylogenetic tree visualization titled "170 Populations". The tree is rooted and shows various clades with bootstrap values indicated at the nodes. The left sidebar contains a "Tools" panel with categories like "Get Data", "Text Manipulation", "Statistics", "Genomics", "Phenotype Association", "Gene Annotation", "Mapping", "Assembly", "Blast", "Manipulation", "MetaPhlAn", "Phylogenetic Analysis", and "Workflows". The right sidebar shows a "History" panel with a list of datasets, including "TermoType MultiSample Table", "Log File", "Distance File", "Out Tree", "Outfile", "Summary File", "Extracted Regions", and "TermoType Tables".



Istituto Superiore di Sanità - Dip. Sanità Pubblica Veterinaria e Sicurezza Alimentare  
EU Reference Laboratory for *E.coli*



# Example of data analysis: O26 cluster of cases in 2013



22 Italian strains  
In house sequencing



9 strains  
GenBank data

- Virulence genes
- Multi-Locus Sequence Typing (MLST)
- Single Nucleotidic Polymorphisms (SNPs) based phylogenetic tree

13 Human isolates from  
Italy August 2013

**ST21**

**ST29**

***vtx1***

***vtx2***

***vtx-***

8

5

-

12

1

At least two different strains involved in the cases  
occurred in Italy in 2013

# SNPs analysis (Single Nucleotide Polymorphisms)



Epidemiologically related cases appear very different.  
Such a high sensitivity increases the risk for errors



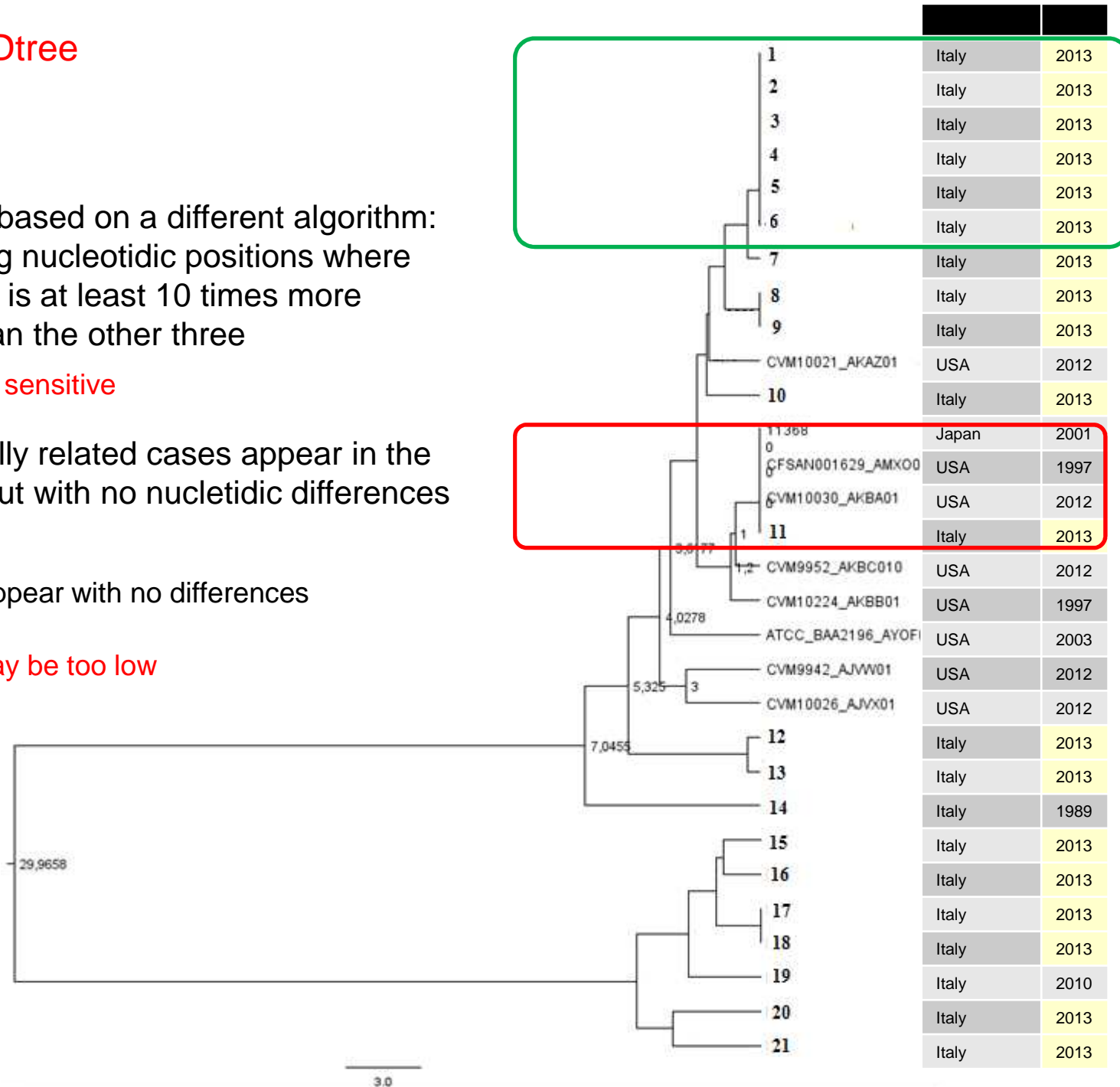
SNPs analysis based on a different algorithm:  
only considering nucleotidic positions where  
the assigned nt is at least 10 times more  
represented than the other three

More robust, less sensitive

Epidemiologically related cases appear in the  
same cluster, but with no nucleotidic differences  
at all

Very far strains appear with no differences

The sensitivity may be too low



# Issues to be addressed

---

- **Data production still needs to be streamlined**  
Reference laboratories only actively produce data as of today (6 NRLs in our network)
- **Cross-platform compatibility**  
Different platforms = different errors rates and types
- **Intrinsic **quality** of the sequence reads at the nucleotidic level**  
Filtering algorithms to be developed and harmonized
- **Refinement of existing tools for data analysis and development of new ones**  
Need for new approaches to typing
- **Need for education in bioinformatics**
- **Computationally intense data analysis**  
Accessibility of bioinformatic tools via open-source servers
- **Massive **data storage** and transfer**  
What data should be stored? Cloud storage?

# Problem solving

---

