**EU Reference Laboratory for *E. coli***
*Department of Veterinary Public Health and Food Safety*
*Unit of Foodborne Zoonoses*
**Istituto Superiore di Sanità**

# Basic Course on Bioinformatics tools

# for Next Generation Sequencing data mining

## 11-12 June, 2015

## SIDBAE Training Room

## (*Building 1, Floor B*)

## Istituto Superiore di Sanità

## Viale Regina Elena, 299 – Rome, Italy

# The trendy "Omics" approach I

- Is changing the labs language
- Computers are becoming more and more visible in the labs
- Increasing importance of Data Storage and Sharing
- Structural re-think of the Labs management (Storage and back-up facilities, fast internet, compression, encryption and other data protection measures)
- Computing hardware become quickly obsolete
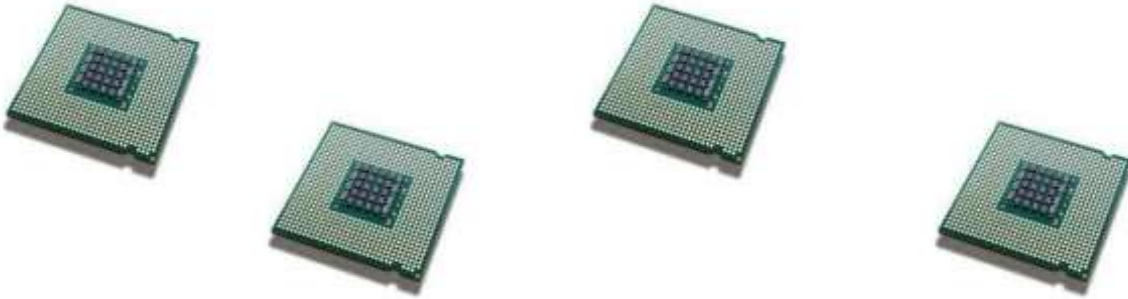
# High throughput Sequence Data

**High storage capacity**

**Table 4:** The comparison between PGM and MiSeq.

| | PGM | MiSeq |
|---|---|---|
| Output | 10 MB–100 MB | 120 MB–1.5 GB |
| Read length | ~200 bp | Up to 2 × 150 bp |
| Sequencing time | 2 hours for 1 × 200 bp | 3 hours for 1 × 36 single read<br>27 hours for 2 × 150 bp pair end read |
| Sample preparation time | 8 samples in parallel, less than 6 hrs | As fast as 2 hrs, with 15 minutes hand on time |
| Sequencing method | semiconductor technology with a simple sequencing chemistry | Sequencing by synthesis (SBS) |
| Potential for development | Various parameters (read length, cycle time, accuracy, etc.) | Limited factors, major concentrate in flowcell surface size, insert sizes, and how to pack cluster in tighter |
| Input amount | $\mu g$ | Ng (Nextera) |
| Data analysis | Off instrument | On instrument |

# Computationally intensive applications



The assembly of a bacterial genome (approx 600 MB) takes approx 30mins to 5 hours and completely occupies the computation capacity of a processor (or a core of a quad-cores processor)

The assembly of a metagenomics sample (up to more than 3 GB) may take days and in some cases it will not be assembled at all (regardless the number of cores available)

# The trendy "Omics" approach II

- Growth curve and colony morphology is going to be replaced by other concepts such as the Phred score or NG50

- Computing "omics" data requires expertise that is still scanty in the average NRL

- To get oriented in the blast of info proposing bioinformatics software might be a nightmare

# Data analysis: The Black Hole

# Data analysis: Software suites

- *de novo* assembly

- Alignment of sequences, production of VCF files, production of dendrograms

- MLST

- Search for interesting genes

**USER-(almost)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**
**IT ADMIN BY LifeTech**

# Data analysis: Specialized web servers



- Species identification

- *de novo* assembly tools

- VirulenceFinder

- ResFinder

- MLST

- SNPs tree and newly deleveloped NGS-driven philogenetic tools

- Other useful molecular microbiology/epi tools

**FREE, USER-FRIENDLY WEB INTERFACE, COMPLETELY CLOSED ENVIRONMENT, LIMTED POSSIBILITY TO INTERVENE FOR USERS**

# Data analysis: web servers for general "omics" analysis



- can be installed locally

- can run any commnd-line running scripts

- *de novo* assembly tools

- BLAST search of interesting genes

- Alignment of sequences, production of VCF files, production of dendrograms

- virtually unlimited possibilities……

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

# What we had in mind…

**Bring "omics" into the NRLs real life while keeping (almost) all those problems out**

**Build knowledge on these game-changing approaches in our network**

**Provide new analytical tools for *E. coli* detection and typing based on "omics"**

**Develop a flexible platform for Routine Work as well as Research**

**Keeping an eye on how the molecular surveillance will develop**

# ARIES: A shared workspace for intensive data analyses

# ARIES: Get started



*https://aries.iss.it*

# ARIES: Roadmap towards a Common analytical bioinformatics interface

**July 1, 2015: ARIES Opens to ISS Users**

- Accounts distributed
- New sections may be created
- Open to collaborations

**July 1, 2015:**

- More Molecular epidemiology tools
- Metagenomics

**October 31, 2015: ARIES publicly exposed (Beta)**

- Accounts available at National Level
- Accounts available for the *E. coli* network
- Stress test for the architecture

**October 31, 2015:**

- More tools for NGS data Mining
- More tools for NGS data Visualization

**January 1, 2016: ARIES on the World Wide Web (Beta)**

Accounts granted to international users
Requests addressed to the IT administrator

**January 2016 on…**

- Other "Omics"…..

# ARIES: Credits

**The Galaxy ARIES core group:**

- Stefano Morabito (EU-RL VTEC): stefano.morabito@iss.it

ARIES Scientific coordination, tools design, <u>contact person</u>

- Arnold Knijn (SIDBAE): arnold.knijn@iss.it

ARIES Administrator, Galaxy tools integration, <u>contact person</u>

- Valeria Michelacci (EU-RL VTEC): valeria.michelacci@iss.it

Tools design

- Massimiliano Orsini (IZS AM): m.orsini@izs.it

Tools design, code-writing