

Advanced Research Infrastructure for Experimentation in Genomics (ARIES): A concept of a bioinformatics framework for the analysis of genomic data from zoonotic agents



Valeria Michelacci¹, Arnold Knijn¹, Massimiliano Orsini², Stefano Morabito¹

¹Istituto Superiore di Sanità, Rome, Italy

²Istituto Zooprofilattico Sperimentale dell'Abruzzo e Molise, Teramo, Italy

Contact: valeria.michelacci@iss.it; aries@iss.it



Introduction

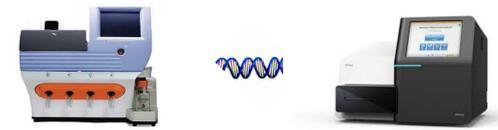
The European Reference Laboratory for *E. coli*, and the IT services of the Istituto Superiore di Sanità in Rome, have developed a web platform for the analysis of genomic data in the field of public health and food safety, with the aim of deploying a comprehensive bioinformatics approach to the study of food-borne zoonoses and infectious diseases at the human and animal interface. The web portal is termed ARIES (Advanced Research Infrastructure for Experimentation in genomics) and is powered by the Galaxy framework.

Aims

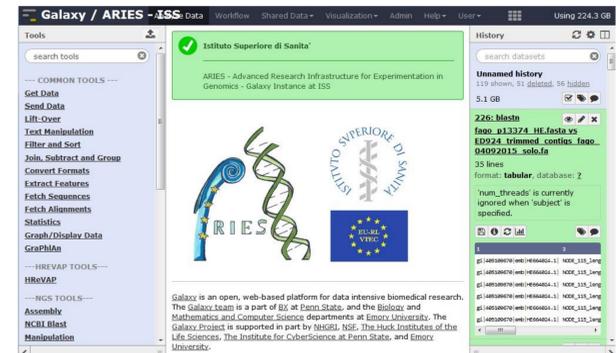
- Supporting scientists in genomic investigation of pathogens, by providing at the same time a unique pre-configured genomic environment characterised by reproducibility, transparency and easy data sharing as well as an open platform for implementation of new bioinformatics resources.
- Development of an Information System for the collection of genomic and epidemiological data to enable the Next Generation Sequencing (NGS)-based surveillance of infectious epidemics, foodborne outbreaks and diseases at the animal-human interface.
- Development of analytical pipelines enabling harmonized, real time multi-genome comparisons, to improve the detection of clusters of cases of infections and allowing the global bio-tracing of pathogens.
- Development of metagenomics models for the culture-independent detection and typing of pathogens and the study of their interactions with the microbiota in human and animal samples and in the vehicles of infections.

Concept

NGS Data production



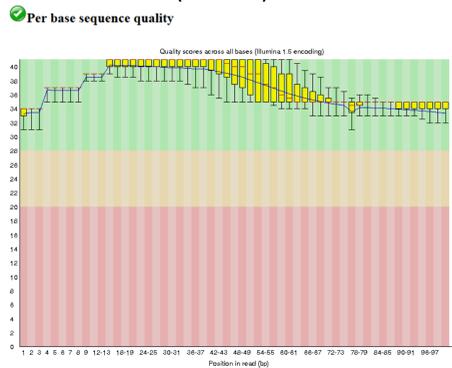
NGS Data Analysis



1. Choose the tool
2. Choose the data and set the parameters
3. Get the result

ARIES contains all the tools for basic and specialized genomics:

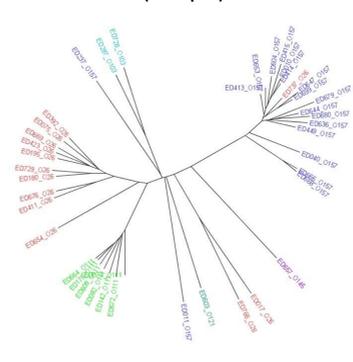
Quality check of the raw data (FastQC)



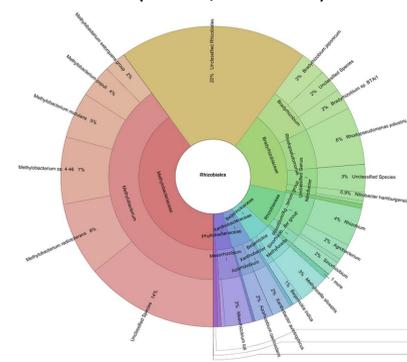
De novo assembly (SPADES, Edena, MIRA, Velvet)

```
>NODE_215_length_1206_cov_29.0156_ID_429
CCACCGAAAGTGGCGGGGATTCACACGATATATCTGTGAAATTAITAGGAT
SCAGATCTCTTGAATCATACACCGGATTTTCTTAAATGATGATGATGATG
TTTGGTTCCCTGAGGAAATTAAGGTTGATCATGCTGATGATGATGATGATG
TACTGAGACAGTATCTCATGATTAATCAAGATGGTGGATGGTCTGGGCTCA
TCATTTTGGATTAACCTACTAATTAATGCCCCCAAGAGGATGATTCATGTAI
AATGGTGGATCATATTCAGAAATAGGACATATCAATCTGACATGAAACTATGAC
GATGATATATGTCGACCGCGCTGTCGGGAGAGTTTCGGCTATGATATCTGAG
GATATGATGGGCGCATCTGATTAATGATTAATGAAATGATGATGATGATGATG
GGGGTACTGATTAATGAGCATCATGATGATTAATGATTAATGATTAATGATG
CTCAACTTGGTTAATCTGGGCGCTGATTAATGATTAATGATTAATGATTAAT
CATCAAAATAGTATTAATGATTAATGATTAATGATTAATGATTAATGATTAAT
GTATCACTAATCACTACTGCTCTCTGATTAATGATTAATGATTAATGATTAAT
ACGATATCGGTAAGTGGACGAGCAGATGGGTTACTGACGGGAGTGGATCAT
CTCCCTCCCGGTCATTCACGATCAACAGCTGACGCTGCGGATGATGATTTCC
AGCTCACTTACCTGACGAGCTGCTGCTGACGAGCCTGCTGCTGCTGCTGCTG
GTGAAATGATTAATGAGCGCAATGGCGCCACCGCTCTCTCTGCTGCTGCTG
ATTTACAGGATTAATTAATCAAGCGAGTGTGAGAGGCTGTACCTGAGAGCGCT
GATTAAGCCCTCTGCTGCGAGCGCTGAGCATTGAGAGGAGGATTAATGATG
TCGTTTCCCGGCAATACGCGCAACAGAGCTGCTGCTGCTGCTGCTGCTGCTG
CGAAGTTGATTCATGATTAATTAATGAGGCTCTCTGCTGCTGCTGCTGATTC
ATACTT
>NODE_212_length_1234_cov_52.9584_ID_423
TGATAGCTGCTGCTGCGCGCATTAACACAGTGTATTTGCGCTTCTCTGATTC
ACAACTGGCGAGCTTTTATTAATGATGTTGTGGTACGCGCATCCCGGAGCAAC
CCGATACCGCTTTCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTGCTG
```

Whole genome SNPs comparisons (Ksnpp3)



Metagenomics tools (Krona, Commet)



- Other tools:
- Mapping (BWA, Bowtie2)
 - Gene annotation (Prokka, Glimmer)
 - Alignment (BLAST+ suite) and counting...

At present ARIES features workflows specifically developed for the characterization of pathogenic *E. coli*, such as:

- identification of the presence of virulence genes
- *in silico* serotyping
- MLST typing

Gene	Score	Start	End
stx2B:40:AB232172:1f	100.00	264	264
gad:60:CP001671	100.00	599	1401
gad:61:FM180568	100.00	599	1401
gad:66:CP002167	100.00	687	1401
gad:67:CP001671	100.00	687	1401
ese:22:AB647391	100.00	2820	2820
espA:13:AF022336	100.00	579	579
tir:17:AB354737	100.00	678	1725
tir:17:AB354737	100.00	678	1707
gad:23:CP000970	100.00	164	1401
gad:59:EF547386	100.00	164	1401
gad:60:CP001671	100.00	164	1401
gad:61:FM180568	100.00	164	1401
gad:62:EF547387	100.00	164	1401
astA:4:AB042002	100.00	117	117
ese:17:AB647460	99.96	2820	2820

E. coli virulotyping workflow with example of results

Allele	Score	Count
FUMC4	100.00	469
ICD16	100.00	518
MDH9	100.00	452
GYRB12	100.00	460
RECA7	100.00	510
PUR47	100.00	478
ADK16	100.00	536

E. coli MLST alleles identifier with example of results

Possibility to easily develop and share new analytic pipelines

e.g. The pipeline for the **High Resolution Allelic Profiling of Shiga-toxin producing *E. coli*** was developed and is currently running on ARIES (see abstract from the corresponding presentation at the High throughput PCR pre-symposium workshop, October 7th)

CONCLUSIONS

- ARIES is currently optimized for the analysis of *E. coli* genomes, but the tools implemented can be adapted to the analysis of other pathogens upon request, by integrating reference databases
- The use of a common and easy-to-use platform will help the formation of a network of laboratories for the exchange of developed pipelines and results, driving towards the harmonization of NGS data analysis.
- The development of an integrated system for the production, collection and analysis of genomic data on agents of foodborne zoonoses and infectious diseases will be the ground for an enhanced real time monitoring and surveillance.

The platform is planned to open for external access in early 2016. Contact us at aries@iss.it to join us!

REFERENCES

- Goecks, J, Nekrutenko, A, Taylor, J and The Galaxy Team. Galaxy: a comprehensive approach for supporting accessible, reproducible, and transparent computational research in the life sciences. *Genome Biol.* 2010 Aug 25;11(8):R86.
- Blankenberg D, Von Kuster G, Coraor N, Ananda G, Lazarus R, Mangan M, Nekrutenko A, Taylor J. "Galaxy: a web-based genome analysis tool for experimentalists". *Current Protocols in Molecular Biology.* 2010 Jan; Chapter 19:Unit 19.10.1-21.
- Giardine B, Riemer C, Hardison RC, Burhans R, Elnitski L, Shah P, Zhang Y, Blankenberg D, Albert I, Taylor J, Miller W, Kent WJ, Nekrutenko A. "Galaxy: a platform for interactive large-scale genome analysis." *Genome Research.* 2005 Oct; 15(10):1451-5.