# Session X – Galaxy

**Miguel Juliá**

**BU-ISCIII**

**Unidades Comunes Científico Técnicas – SGSAFI-ISCIII**

04-15 Noviembre 2019, 2ª Edición
Programa Formación Continua, ISCIII

15/11/2019
Secuenciación de genomas bacterianos:
herramientas y aplicaciones
1

# Index

## Galaxy:

- Computing in Biosciences

- Change of Paradigm

- What is Galaxy

- Workflows

- The porject

- Galaxy training

Secuenciación de genomas  bacterianos:
herramientas y aplicaciones

# Computing in Biosciences

- Web-based platforms (i.e. Galaxy) and remote HPC

| Pros | Cons |
|---|---|
| No need to storage intermediate files | Your data is in someone else's computer<br>No backups or data management schemes |
| No need to install software<br>Partial control over installed software | No control over installed software, versions and future availability |
| Graphic interface | No control over hidden parameters |
| Analysis are partially reproducible | Quotas |

Secuenciación de genomas bacterianos: herramientas y aplicaciones

# Change of Paradigm I

## 1 sample

**Research only:** NGS was still a new thing, no applications 10 years ago

**Reproducibility is not needed:** Why would anyone reanalyse this?

**Storage is not an issue:** files of 1 sample fits everywhere in my HDD, maybe I will copy it in a CD-ROM

**Computing is simple:** no need to worry about resources or optimisation

## multiple samples

**Many applications:** research, clinical, industrial, forensic, military, …

**Reproducibility, scalability , portability and standardisation are required**

**Storage is challenging:** storage, indexation and backup required, privacy and legal standards

**Computing requires optimisation and lots of resources**

# Change of Paradigm II

- Nowadays scientific computing paradigm

| Pros | Cons |
|------|------|
| Data remains private Backups and data management schemes | High storage space Dedicated file systems Databases to index files |
| Control over software installed versions, open source programs | Many versions of the same software coexists |
| All parameters are available for the command | You have to understand all software variations |
| Analysis are reproducible and public | You have to publish and document your work |

Secuenciación de genomas bacterianos: herramientas y aplicaciones

# What is Galaxy I

**Data Intensive *analysis* for everyone**

- Versatile and reproducible workflows

- **Web** platform

- **Open source** under [Academic Free License](Academic Free License)

- Developed at Penn State, Johns Hopkins, OHSU and Cleveland Clinic with substantial outside contributions

Secuenciación de genomas bacterianos: herramientas y aplicaciones

# What is Galaxy II

- **Accessibility**
  - Users without programming experience can easily upload/retrieve data, run complex tools and workflows, and visualize data

- **Reproducibility**
  - Galaxy captures information so that any user can understand and repeat a complete computational analysis
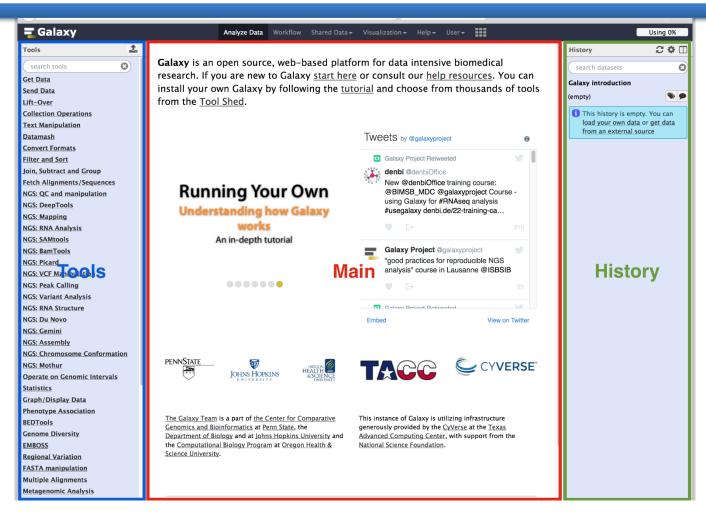
- **Transparency**
  - Users can share or publish their analyses (histories, workflows, visualizations)
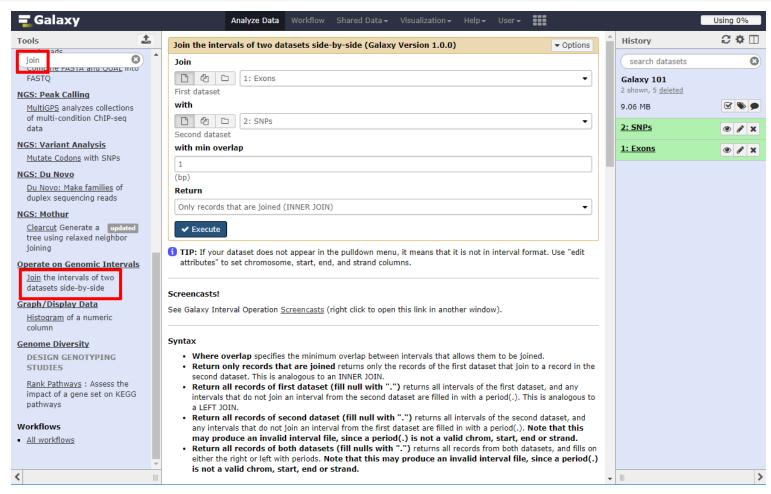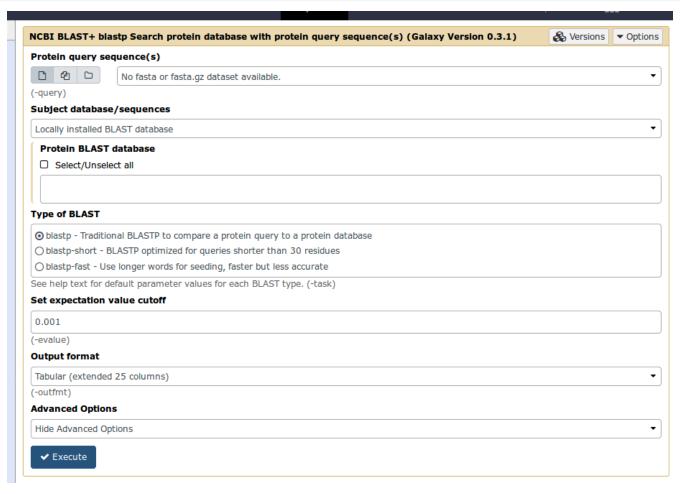  - Pages: online Methods for your paper

Secuenciación de genomas  bacterianos: herramientas y aplicaciones

# What is Galaxy III

Secuenciación de genomas bacterianos: herramientas y aplicaciones

# What is Galaxy IV

Secuenciación de genomas bacterianos: herramientas y aplicaciones

# What is Galaxy V

Secuenciación de genomas bacterianos:
herramientas y aplicaciones

# What is Galaxy VI

Secuenciación de genomas bacterianos: herramientas y aplicaciones
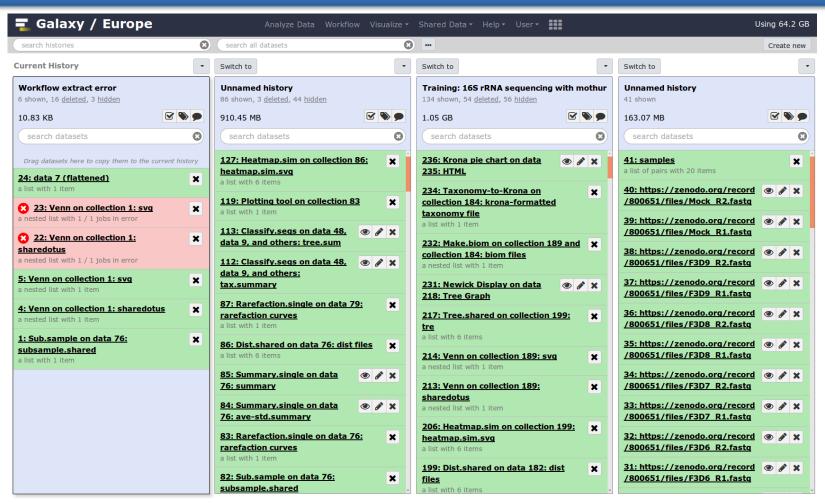
# What is Galaxy VII

- Location of all analyses
  - collects all datasets produced by tools
  - collects all operations performed on the data

- For each dataset (the heart of Galaxy's reproducibility), the history tracks
  - name, format, size, creation time, datatype-specific metadata
  - tool id, version, inputs, parameters
  - standard output (stdout) and error (stderr)
  - state (waiting, running, success, failed)
  - hidden, deleted, purged

Secuenciación de genomas  bacterianos: herramientas y aplicaciones

# What is Galaxy VIII

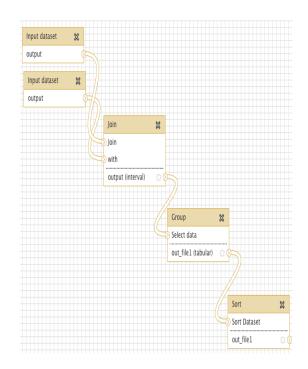Secuenciación de genomas bacterianos: herramientas y aplicaciones
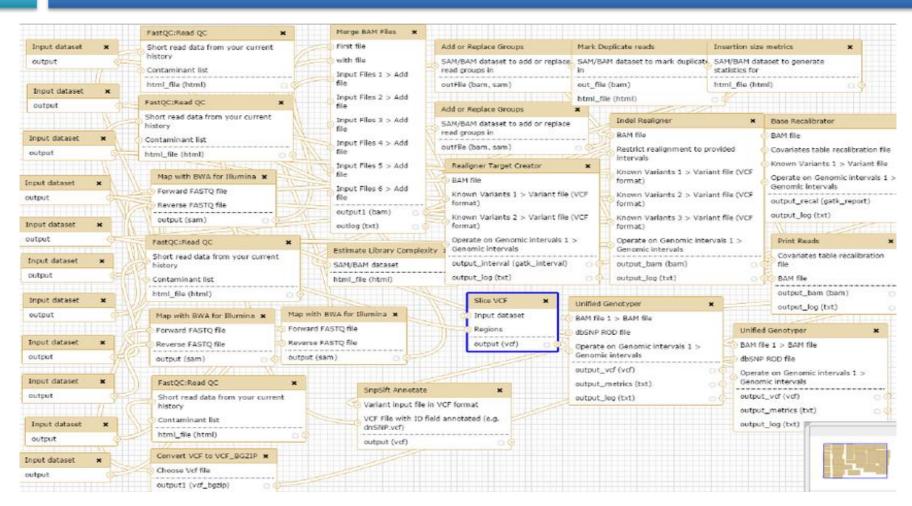
# Workflows I

- Bioinformatic analyses invariably involve shepherding files through a series of transformations, called a **pipeline** or a **workflow.**

- These transformations are done by executable **command line software** written for Unix-compatible operating systems.

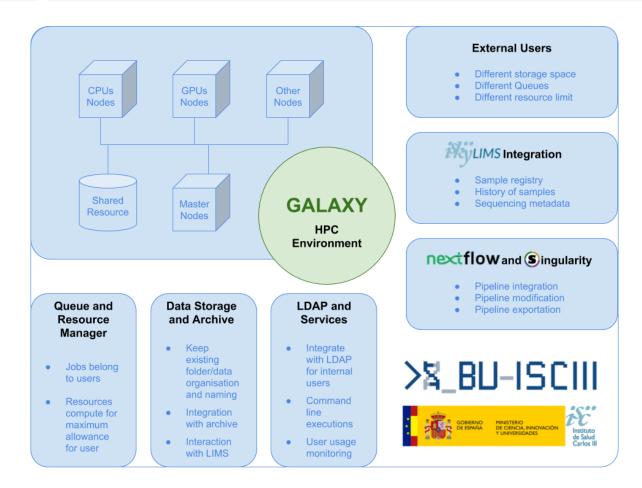- They need to be **reproducible, easy to maintain, portable and scalable.**

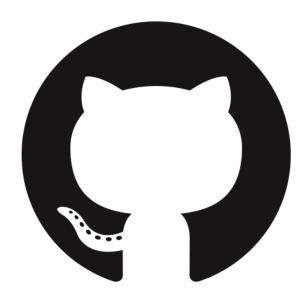Secuenciación de genomas bacterianos: herramientas y aplicaciones

# Workflows II

Secuenciación de genomas bacterianos: herramientas y aplicaciones

# The Project



- Clinical data storage

- Hospitals

- Patient oriented research

- Training

Secuenciación de genomas bacterianos: herramientas y aplicaciones

# Thanks for your attention!

And this is only the tip of the iceberg...
Check this if you wanna know what's really going under the hood:



**https://github.com/BU-ISCIII**

Secuenciación de genomas bacterianos:
herramientas y aplicaciones