# Building an open, collaborative, online infrastructure for bioinformatics training
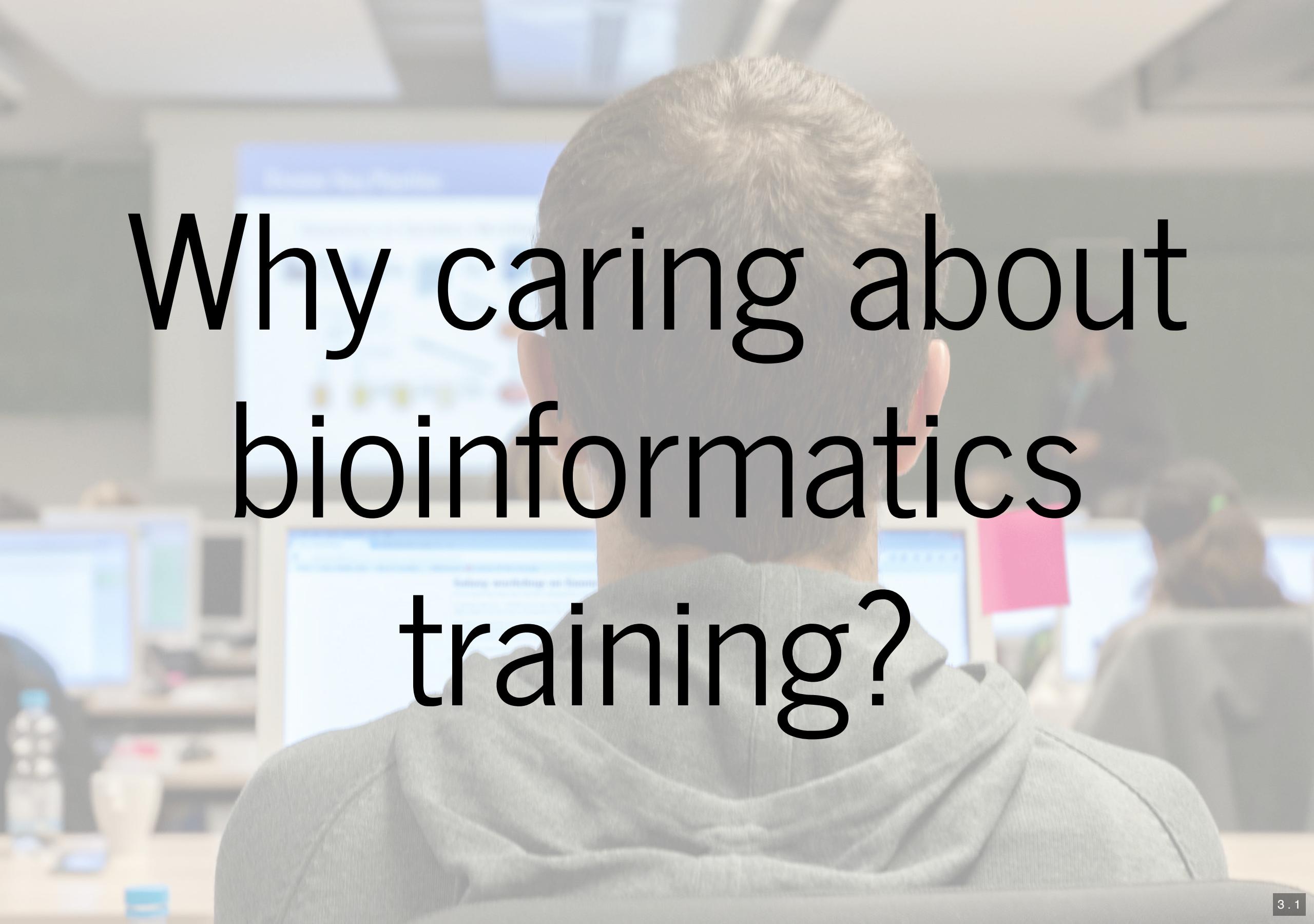


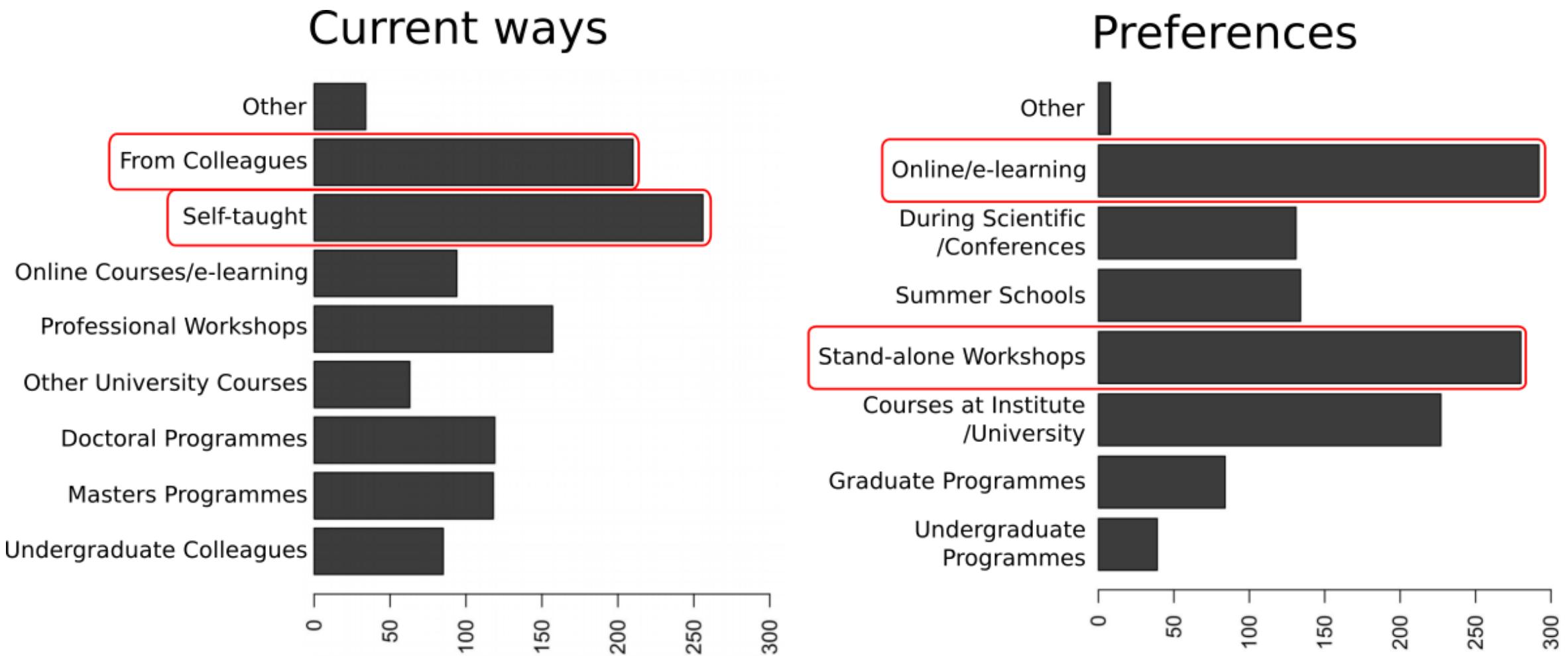**Bérénice Batut**

33rd TBI Winterseminar
Februar 2018

Why caring about bioinformatics training?

# Need for bioinformatic training

*Bioinformatics has become too central to biology
to be left to specialist bioinformaticians*

- Explosion of data to analyze
- Access to computational power
- Thousand of possible tools for specialized analyses

# An increasing demand
# for learning bioinformatics
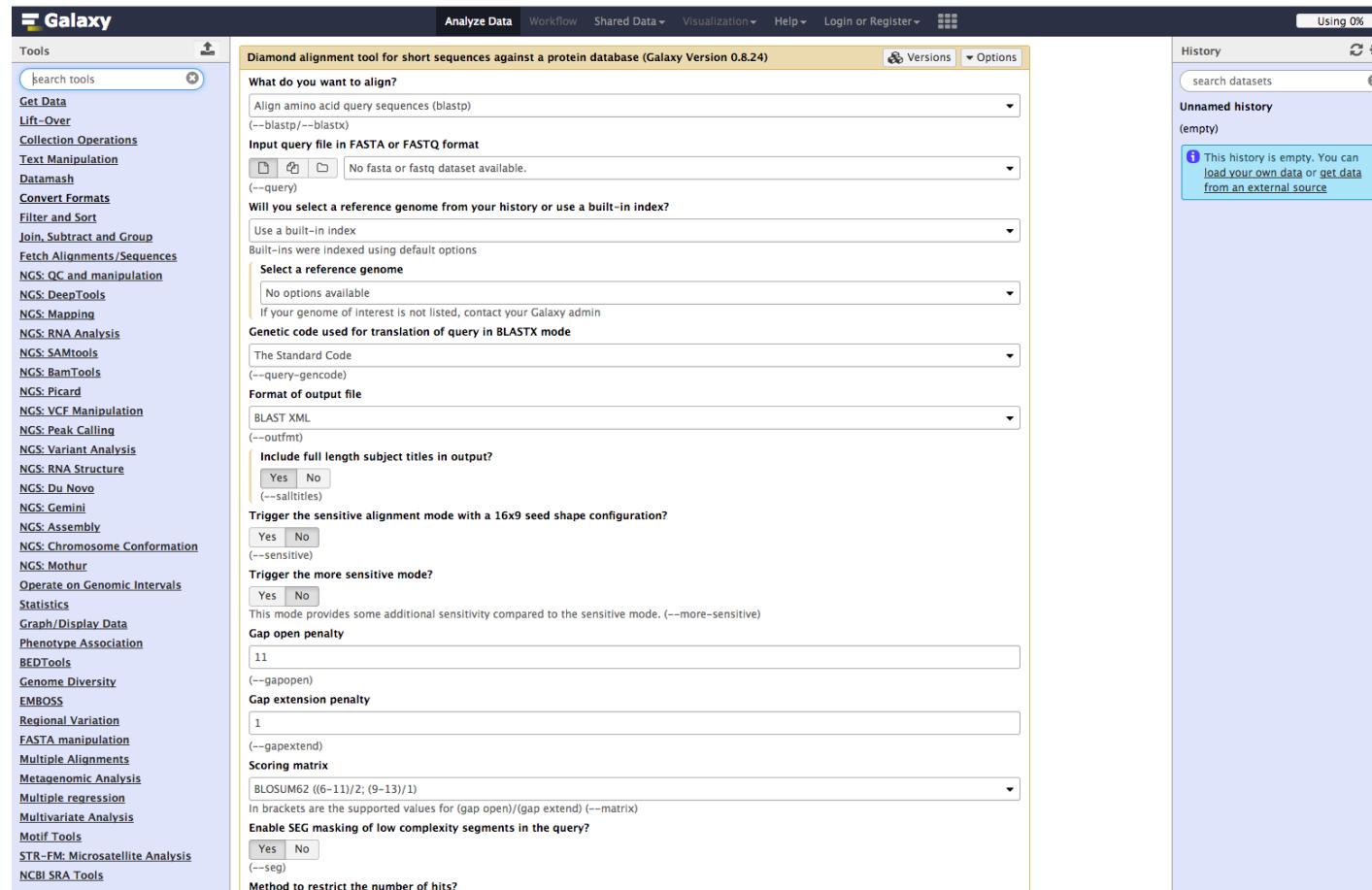


Graphs of Brazas et al, 2017

Galaxy
a great solution !

# Computational knowledge: Not required!



- Web interface for numerous bioinformatics tools
- Scalable
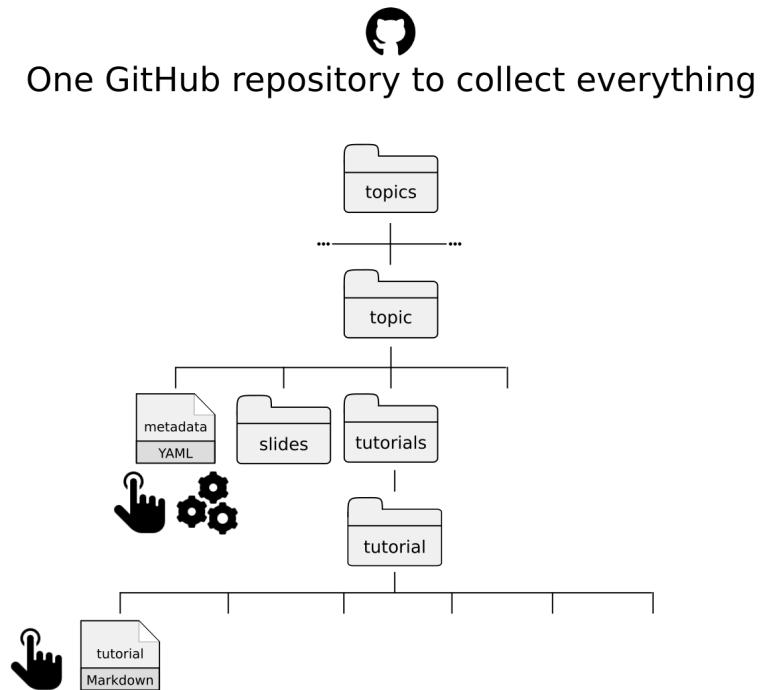- No issue with computer configuration during training

Building a new **open**, **collaborative** and **FAIR** model for bioinformatics training

# Requirements

- Easy to use
- Support for effective training for
  - Individual users
  - Instructors
- Definition of technological infrastructure
- Limited redundancy

# An open, collaborative, FAIR, online infrastructure



One GitHub repository to collect everything

# Separation between content and format



Markdown



User-friendly HTML

# An open, collaborative, FAIR, online infrastructure



Separation between content and form

One GitHub repository to collect everything

Findable
Accessible
Interoperable
Reusable

# Findable material via TeSS



https://tess.elixir-europe.org/

# An open, collaborative, FAIR, online infrastructure



One GitHub repository to collect everything

**Separation between content and form**

**Population of external bioinformatic training resources**

Findable

Accessible

Interoperable

Reusable

# Storage of training datasets on Zenodo

# An open, collaborative, FAIR, online infrastructure

One GitHub repository to collect everything

Separation between content and form

Population of external bioinformatic training resources

Findable, accessible, persistent and citable data

topics

topic

metadata
YAML

slides

tutorials

tutorial

tutorial
Markdown

data lib
YAML

Input file

Findable

Accessible

Interoperable

Reusable

# An open, collaborative, FAIR, online infrastructure



Separation between content and form

One GitHub repository to collect everything

Self-training boxes

Findable, accessible, persistent and citable data

Population of external bioinformatic training resources

Annotation & configuration

Findable

Accessible

Interoperable

Reusable

# An open, collaborative, FAIR, online infrastructure

One GitHub repository to collect everything

Separation between content and form

Population of external bioinformatic training resources

Self-training boxes

Findable, accessible, persistent and citable data

Annotation & configuration

An open development process

Branch creation

Pull Request opening

Discussion, review and improvement

Pull Request merging

Automatic testing of website generation, links, images, ..

Findable

Accessible

Interoperable

Reusable

# Galaxy Training materials



http://training.galaxyproject.org

# Thank you!



An open, collaborative, FAIR, online infrastructure

🌐 training.galaxyproject.org

github.com/galaxyproject/training-material