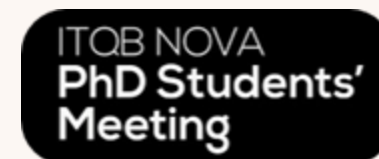# Computational Tools and Resources for the Life Sciences

**Gil Poiares-Oliveira**

INESC-ID | BioData.pt | ELIXIR Portugal
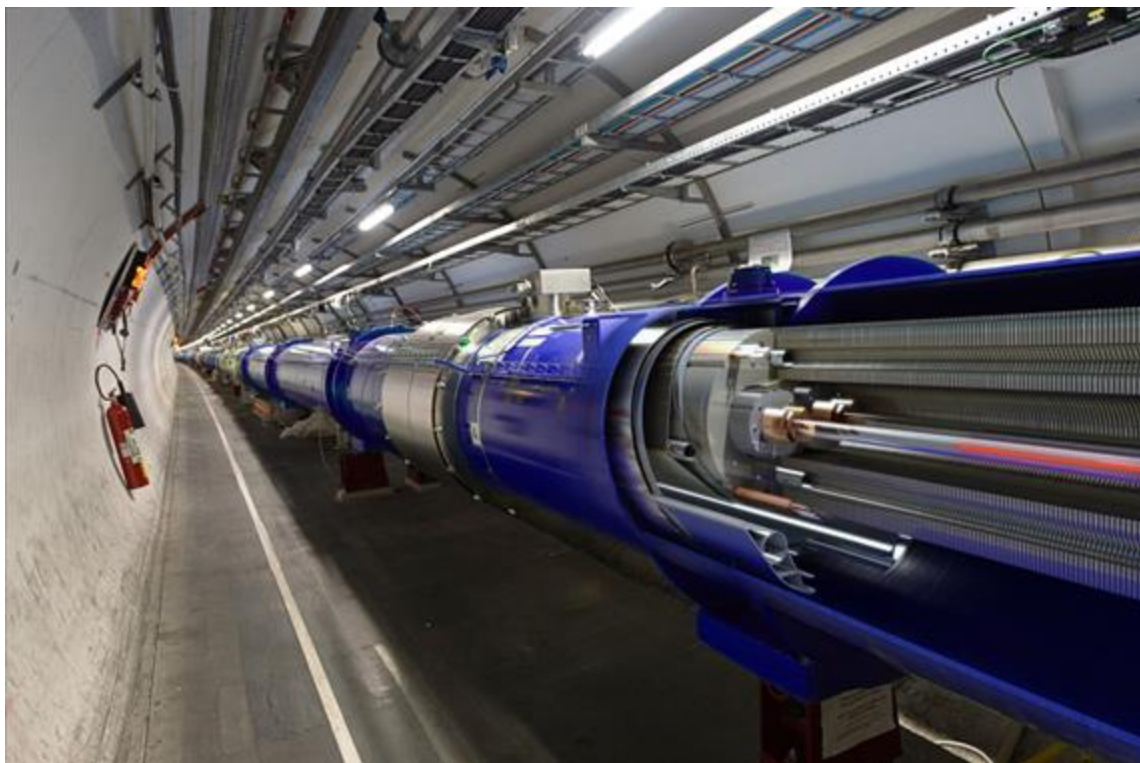
ITQB NOVA PhD Student Meeting • 2024-05-16
ITQB NOVA • Oeiras, Portugal

# Learning outcomes

- Understand the role of ELIXIR and BioData.pt in the European bioinformatics landscape

- Identify ELIXIR and BioData.pt services and projects like FEGA, Beacon and TeSS

- Perform basic genetic and phenoclinic queries using Beacon and the Beacon Network

- Create a bioinformatics analysis pipeline using Galaxy workflows

- Use ELIXIR TeSS to discover new training courses and materials

# What is a research infrastructure?



**Large Hadron Collider - CERN, Switzerland**

Image: © 2014 CERN. Dominguez, Daniel



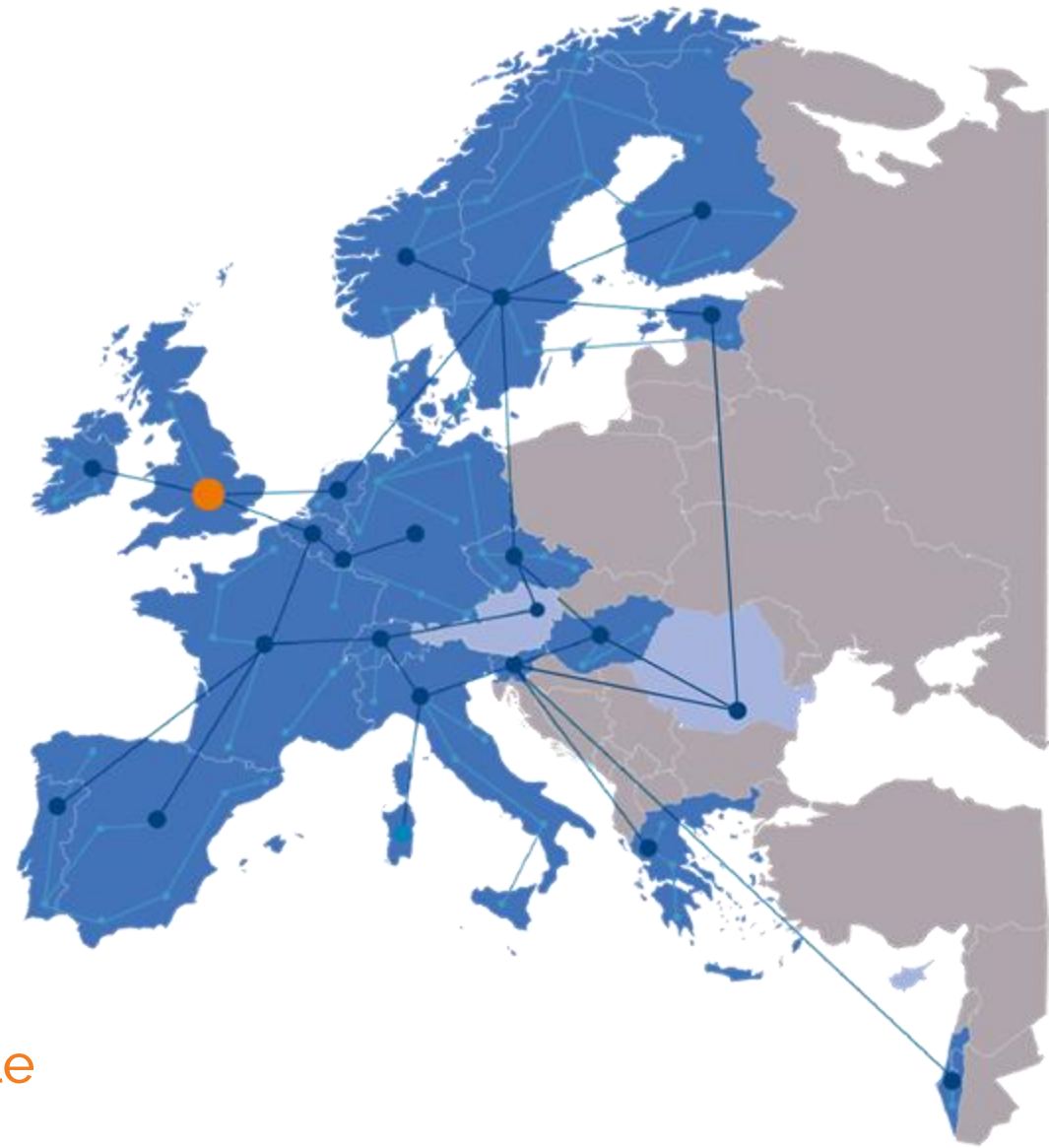**ELIXIR Nodes - ELIXIR, Europe**

Image: © 2024 ELIXIR

# ELIXIR Europe

ELIXIR is an intergovernmental organisation that brings together life science resources such as

- databases
- software tools
- training resources
- interoperability resources
- compute resources
- data management support

The goal of ELIXIR is to coordinate bioinformatics resources from across Europe so they form a single infrastructure.

elixir-europe.org

# Meet BioData.pt



▶ Portuguese distributed infrastructure for life & health data, and home of ELIXIR Portugal

▶ We offer **training**, **tools**, and **repositories** to the life sciences community and **connect** them with an **international network**

# BioData.pt Services

**Training**

Bioinformatics and RDM

**Bioinformatics analysis**

RNA-Seq, DE, multiomics

**Computing**

servers, tool hosting

**Mentorship**
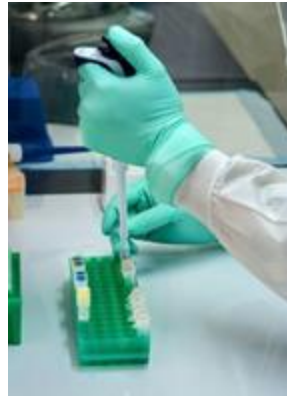
1-on-1 sessions

BioData.pt | elixir PORTUGAL

# Our associates



▶ If you're affiliated with one of our associates, you can benefit from our services

# How does a patient DNA sample get into a computer?



Collection

DNA isolation & library prep

Sequencing

FASTQ

Assembly

Reference sequence

Long read sequencing

Short read sequencing

Reconstructed genome

BioData.pt | elixir PORTUGAL

# Example FASTQ file



1. @ sample ID
2. Nucleotide seq.
3. +
4. Quality score

# 1st Challenge: Getting the data

What we know:

- First recorded symptoms in Bergen, Norway

- Some carriers are asymptomatic, genetic factors may have influence

- University of Bergen has patients' sequencing data

- We are researchers at ITQB NOVA

**How can we help from here?**

**Federated EGA, a resource for discovery and access of human data across national borders**
https://www.youtube.com/watch?v=bRFHyI6hnNk

**User journey**

1. User **registration** in EGA

2. User approval as **Submitter**

3. Creation of Data Access Committee (**DAC**) and associated Data Access Policies.

4. **Upload** of a real genome file into the local server.

5. Submission of **metadata** by user.

6. **Request** to access to the genome, DAC approval and rejection

7. Genome **download** and decryption.

▶ A Beacon allows an anonymous queries on genomic datasets.
e.g. Check if a genome contains a genome with a given base at a particular coordinate

beacon-project.io

🇵🇹 beacon.biodata.pt

▶ Query all ELIXIR's beacons

beacon-network.elixir-europe.org

# Challenge 1

1. How many people in the Portuguese Beacon have an adenine instead of guanine at position 69848 of chromosome 1?

2. How many people have asthma?

3. How many of them are women?

4. How many of the female asthma patients are of African ethnicity?

5. Which country's beacon has a patient with a mitochondrial mutation at the 10th nucleotide which switches a thymine to a cytosine

BioData.pt | elixir PORTUGAL

# Step 2: Analysing the data

Here's what a sample pipeline looks like:

1. **Trim** adapters and low quality reads

2. Perform **quality control** on the data

3. **Map reads** against reference genome

4. **Analyse** according to biological question

   e.g. variant calling - use variant detector to compare

   against reference genome

Tools:

e.g. Trimmomatic, seqtk_trimq

e.g. FastQC

e.g. STAR, HISAT2, Bowtie2, TopHat2

e.g. FreeBayes

# Minimum requirements:

1. A powerful enough computer...

- Trimmomatic
- FASTQC
- STAR
- HISAT2
- Bowtie2
- TopHat2
- FreeBays
- *That cool tool you've just discovered on WoS*

- Basic Unix command line knowledge
- Introductory programming skills[*]
- Familiarity with Python/R[*]
- Knowledge of Pandas/matplotlib/...[*]

[*]additionally, depending on the tool

or

- A bioinformatician

or

# Galaxy
PROJECT

BioData.pt | elixir PORTUGAL

# Galaxy Project

▶ Galaxy is a free, open-source system for analyzing data, authoring workflows, training and education, publishing tools, managing infrastructure, and more.

galaxyproject.org

# Galaxy & the Pulsar Network

## Self-hosted



Your institute's server

Galaxy software

## Pulsar Network



usegalaxy.eu

# Challenge 2: build a workflow in Galaxy

Get FASTQ from EGA
EGAF00004859455

FASTQ →

Trim for quality (seqtk_trimfq)

→ QC (FastQC)

FASTQ →

Map against human genome (HISAT2)

BioData.pt | elixir PORTUGAL

# Compare with mine

BOD workflow

# Step 3: Learning more



▶ Find training courses and materials from ELIXIR Nodes and institutes

tess.elixir-europe.org

BioData.pt | elixir PORTUGAL

# Challenge 3: Diving into TeSS

**Think about some bioinformatics/data science skill or tool that you've been meaning to learn**

phylogenomics

Bash

Nextflow

multiomics

RNA-Seq

Docker

Git

PCA

R

differential expression

metagenomics

disease modelling

Python

Data Management Plans

data visualisation

machine learning

metabolic modelling

protein structure

FAIR principles

Linux

protein function prediction

BioData.pt | eliXir PORTUGAL

# Challenge 3: Diving into TeSS

**Try to answer the following questions about your chosen topic using TeSS**

1.  How many training events are scheduled on that topic?
2.  How many training events, in total are registered on that topic?
3.  How many materials could you find on that topic?
4.  How many materials, on all topics, are from the Portugal node?
5.  How many materials, on all topics, are targeted at PhD Students?

# There's many more...

**bio.tools**

bio.tools helps you find and select bionformatics software and connect it in workflows.

**BioContainers**

Search a repository of containerised software that you can build into workflows.

**WorkflowHub**

A registry for sharing and publishing scientific computational workflows.

**FAIRsharing.org**

FAIRsharing.org allows you to search for databases and data policies by aspects such as domain, species and country.

**TeSS**

Search for training courses, webinars, training materials and workflows in TeSS, ELIXIR's training portal.

**Overview of good data management practices**

The Research Data Management Kit (RDMkit) guides you through the whole data management life cycle and includes advice specific to your domain, your role and your country.

**Step-by-step instructions**

The FAIR Cookbook contains step-by-step recipes to accomplish specific data management tasks and to make your data FAIR (Findable, Accessible, Interoperable, Reusable).

**Data management plan wizard**

The Data Stewardship Wizard (DSW) is an online tool that guides researchers and data stewards through their data management planning.

**LS Login**

The Life Science Login enables researchers to use their home organisation credentials or community or other identities o sign in and access data and services they need.

**Bioconda**

Bioconda lets you install thousands of software packages related to biomedical research using the conda package manager.

BioData.pt | elixir PORTUGAL

# Take home messages

- ELIXIR is the European life sciences data research infrastructure

- BioData.pt is the national life and health data research infrastructure, and hosts the Portuguese ELIXIR node

- Both provide tools and services to support life sciences researchers, e.g. FEGA, Beacon, Galaxy, TeSS

- You can count on us to help you navigate through the complexities of biological data :)

BioData.pt     eliXir PORTUGAL

# Thanks!

This presentation will be available at

github.com/BioData-PT/computational-tools-resources

**Gil Poiares-Oliveira, MSc**

**gpo@biodata.pt**

INESC-ID | BioData.pt | ELIXIR Portugal

BioData.pt
*Living data*

biodata.pt | info@biodata.pt | (+351) 937 990 500