# BioData.pt Talk

## The BioData.pt Plant Sciences Community

Célia Miguel (BioData.pt / FCUL / iBET),
Pedro Barros (BioData.pt / ITQB-NOVA),
Daniel Faria (BioData.pt / INESC-ID)

October 22nd, 2020

# Overview

- Introduction
  - ELIXIR
  - ELIXIR PT | Biodata.pt
  - The Biodata Plant Sciences community

- Showcase: The Cork Oak genome portal

- Showcase: Standards & Resources for FAIR Plant Data
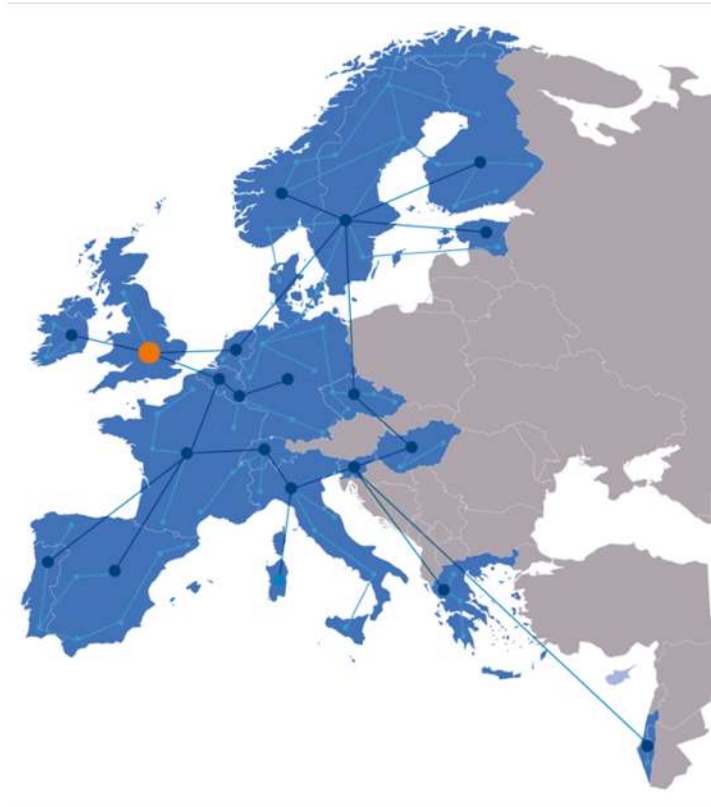
- Other Resources

- Conclusions

# ELIXIR

Intergovernmental organisation that brings together life science resources from across Europe (databases, software tools, training materials, standards and compute resources).

**Goal:** coordinate life science resources from across Europe so they form a single infrastructure

This makes it easier for scientists to:
- Find and share data
- Exchange expertise
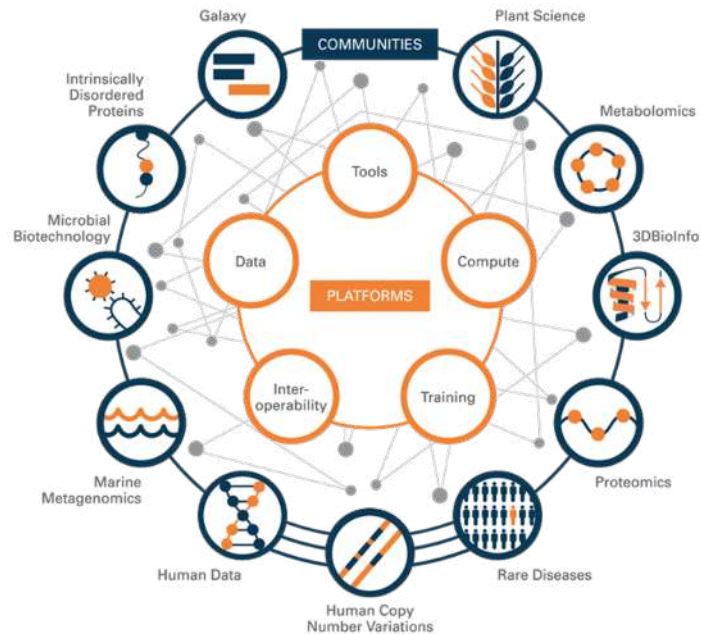- Agree on best practices in scientific research

# ELIXIR Organisation

ELIXIR coordinates activities through at least one of the five areas of activities called **Platforms**:

- Compute
- Data
- Interoperability
- Tools
- Training

These Platforms are driven by eleven ELIXIR **Communities** which develop standards, services, and training within their life science domains.

# ELIXIR Members

- ELIXIR members host Nodes, which represent centres of excellence in bioinformatics.

- In total there are more than 700 researchers from over 220 research institutes.

- ELIXIR activities are coordinated by the ELIXIR Hub, based at the Wellcome Genome Campus, UK.



**ELIXIR Members**

Belgium, Czech Republic, Denmark, EMBL, Estonia, Finland, France, Germany, Greece, Hungary, Ireland, Italy, Luxembourg, Israel, Netherlands, Norway, Slovenia, Portugal, Spain, Sweden, Switzerland, United Kingdom

**ELIXIR Observers**

Cyprus

# ELIXIR Portugal | BioData.pt

## Consortium Partners

# ELIXIR Portugal | BioData.pt

## Organization



Strategically advises direction

Composed by heads of all member institutions

Manages the day-to-day activities

Life science domains (e.g. agro-food and forestry, sea, health)

Represent areas of activity transversal to communities (e.g computing, training)

Provided by BioData.pt to the (inter)national research community

Scientific Advisory Board

Management Board

Board of Directors

Communities

Platforms

Services

# The BioData.pt Plant Sciences Community

**"Plant" Partners**

# The BioData.pt Plant Sciences Community

## Context

- Became active during the ELIXIR-EXCELERATE project, where we co-led the Plant Sciences use-case

- **Goal**: support the publication and access to phenotypic and genotypic data, annotated according to established standards (<u>woody plants</u> as the main target)

  - Critical for addressing some of the major global challenges:
    - Sustainable supply of food and non-food materials
    - Competitive life-sciences industry sector
    - Environmental protection

# The BioData.pt Plant Sciences Community

## Motivation: Plant Breeding (GxExP) remains challenging

- Advances can be obtained from the integration of genomic/genotyping data with diverse types of phenotyping data

- Systematic study of phenotypes on a genome-wide scale, and its association with genomic information under different environmental conditions



Djande et al 2020 DOI 10.3390/agronomy10060831

- Genomics/genotyping and phenotyping datasets are growing in number and size ⇒ gaps in genotype-phenotype associations

# The BioData.pt Plant Sciences Community

## Motivation: How to Make Plant Data FAIR?



Phenomics NL

- Problems with plant data:

  - Heterogeneous (different settings, types of data...)

  - Complex and diverse experimental designs

  - Dispersed (no comprehensive public archive)

  - Poorly annotated (weaknesses in standards)

- FAIR data principles (Findable, Accessible, Interoperable and Reusable) are key to enable knowledge discovery

Célia Miguel, Pedro Barros, Daniel Faria

BioData.pt Talk

# The BioData.pt Plant Sciences Community

## Objectives & Activities

Develop/recommend standards and ontologies to enable FAIR plant phenotyping data

Develop/implement repositories for plant phenotyping and genomic data

Develop/implement user-friendly interfaces for data deposition and retrieval

Annotate and provide curated plant data sets

Develop/provide tools for plant data analysis

Provide training on plant data management

Engage with industry to exchange and apply knowledge

Icons made by inipagistudio, phatplus, Eucalyp and Freepik

# The BioData.pt Plant Sciences Community

**Species of Interest**

- Cork oak (*Q. suber*)

- Maritime Pine (*P. pinaster*)

- Eucalyptus (*E. globulus*)

- Olive tree (*O. europaea*)

- Grapevine (*V. vinifera*)

- Rice (*O. sativa*)

- ...

# Plant Sciences Community Showcase

## **CorkOakDB: The Cork Oak Genome Portal**

# CorkOakDB: The Cork Oak Genome Portal

## History

2014                                      2018         2019         2020

- Comprehensive database of cork oak transcriptome obtained by cDNA sequencing (ESTs)

- No knowledge of genome structure

- Incomplete gene coding sequences

www.corkoakdb.org

# CorkOakDB: The Cork Oak Genome Portal

## History



2014                                    2018          2019          2020

- Reference genome sequenced

- Improved gene structure annotation

- New transcriptomic datasets publicly available

# CorkOakDB: The Cork Oak Genome Portal

## Goal

- Create an integrated repository dedicated to cork oak 'omics'

- Development of tools for data visualization and retrieval for core genomics analyses

- Become a reference hub for research, aggregating all available genomic and transcriptomic (...) data



**CorkOakDB**

*Node Service*

# CorkOakDB: The Cork Oak Genome Portal

## Data

- Genome sequence and structural annotation (gene, exon, intron, CDS)

- Gene expression (retrieved from publicly available studies)

- Curated metadata

## Structure

- Web portal based on the Tripal framework (Drupal CMS + Chado database)



Node Service

# CorkOakDB: Genome Browser

## Tools

- **Gene search and sequence retrieval**

- Genome visualization (JBrowse)

- Homology search (Blast)

# CorkOakDB: Genome Browser

## Tools

- Gene search and sequence retrieval

- **Genome visualization (JBrowse)**

- Homology search (Blast)

# CorkOakDB: Genome Browser

## Tools

- Gene search and sequence retrieval

- Genome visualization (JBrowse)

- **Homology search (Blast)**



**BLAST Results**

Download: Alignment, Tab-Delimited, XML, GFF3

Query Information: /tmp/2020Apr07_164625_query.fasta
Search Target: Cork Oak Proteins
Submission Date: Tue, 07/04/2020 – 16:46
BLAST Command executed: blastp –max_target_seqs 500 –evalue 0.001 –word_size 3 –gapopen 11 –gapextend 1 –matrix BLOSUM62

Number of Results: 217

**Resulting BLAST hits**

The following table summarizes the results of your BLAST. Click on a *triangle* on the left to see the alignment and a visualization of the hit, and click the *target name* to get more information about the target hit.

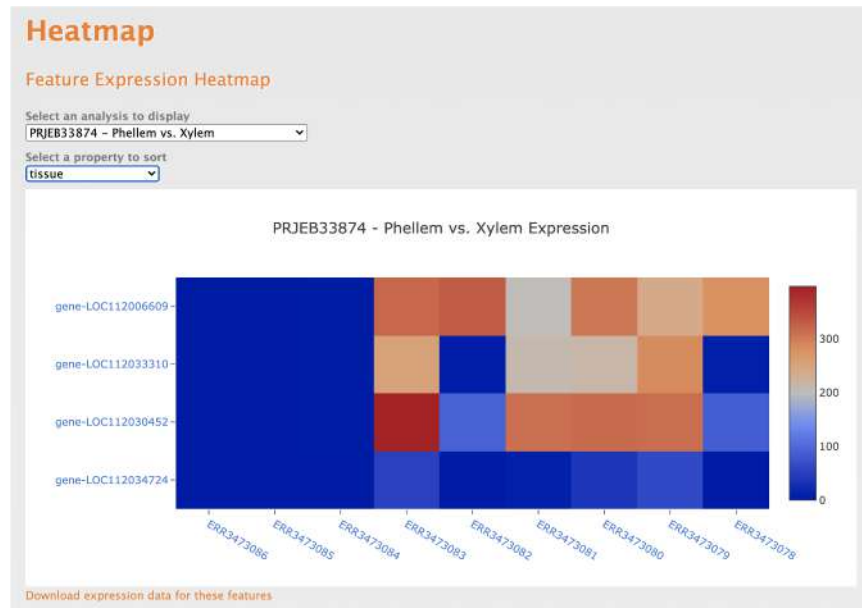| # | Query Name (Click for alignment & visualization) | Target Name | E-Value |
|---|---|---|---|
| ▼ 1 | AT4G17785|MYB39|SUBERMAN | ref|XP_023923310.1| | 5.76552E-86 |
| ▼ 2 | AT4G17785|MYB39|SUBERMAN | ref|XP_023894664.1| | 3.16396E-77 |
| ▼ 3 | AT4G17785|MYB39|SUBERMAN | ref|XP_023918909.1| | 6.30974E-74 |
| ▼ 4 | AT4G17785|MYB39|SUBERMAN | ref|XP_023921866.1| | 1.09019E-72 |
| ▼ 5 | AT4G17785|MYB39|SUBERMAN | ref|XP_023921864.1| | 1.09019E-72 |
| ▼ 6 | AT4G17785|MYB39|SUBERMAN | ref|XP_023877052.1| | 1.07984E-70 |
| ▼ 7 | AT4G17785|MYB39|SUBERMAN | ref|XP_023893622.1| | 2.21112E-69 |
| ▼ 8 | AT4G17785|MYB39|SUBERMAN | ref|XP_023893623.1| | 3.39408E-69 |
| ▼ 9 | AT4G17785|MYB39|SUBERMAN | ref|XP_023921729.1| | 9.60762E-69 |
| ▼ 10 | AT4G17785|MYB39|SUBERMAN | ref|XP_023886483.1| | 9.56498E-68 |

# CorkOakDB: Gene Expression

## Data

- EST sequencing data

- Transcriptomic data for different **tissues**, **developmental stages** and **growth conditions**

- 65 RNA-seq libraries from 15 studies publicly available

# CorkOakDB: Gene Expression

## Visualization

- RNA-seq reads aligned with the reference genome to **estimate gene expression**

- Heatmap for **comparative analysis** of multiple genes

# CorkOakDB: Next Steps



www.corkoakdb.org

- Linking ESTs IDs (old portal) with gene IDs (new portal)

- Integration of curated data from published results (functional validation)

- Integration of other types of data (SNPs, phenotype, …)

# Plant Sciences Community Showcase

## Standards & Resources for FAIR Plant Data

# The FAIR Data Principles

## Findability

- Persistent identifiers
- Rich metadata
- Searchable repository

## Accessibility

- Access protocol
- Authentication & authorization

## Interoperability

- Knowledge representation language
- Controlled vocabularies

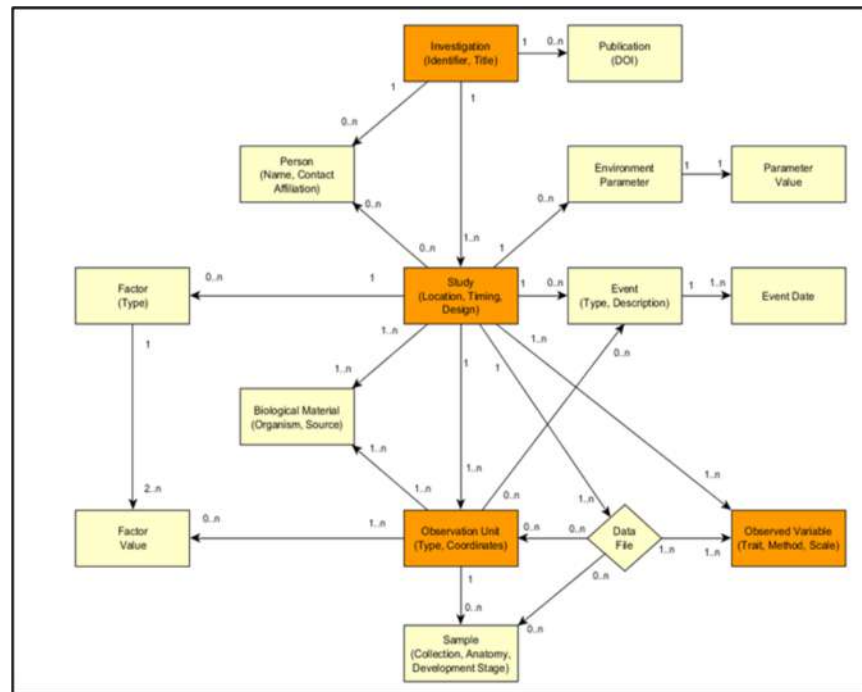## Reusability

- Rich metadata, following community standards
- License and provenance

# Plant Phenotyping Data FAIRness?

**Findability**

- ? Persistent identifiers
- ✓ Rich metadata
- ✓ Searchable repository

**Accessibility**

- ✓ Access protocol
- ? Authentication & authorization

**Interoperability**

- ✓ Knowledge representation language
- ✓ Controlled vocabularies

**Reusability**

- ✓ Rich metadata, following community standards
- ✓ License and provenance

isatools

Crop Ontology
for agricultural data

miappe

BrAPI

# Updating MIAPPE

## Minimum Information About a Plant Phenotyping Experiment 1.1

- Scope extension: woody plants (e.g. identification through GPS)

- Revised structure to match ISA

- Enriched examples

- Data model specification (a standard is more than a flat list)

# Updating MIAPPE

## Minimum Information About a Plant Phenotyping Experiment 1.1

- Explicit ontology recommendations

- Mapping to BrAPI

- ISA-Tab templates

- OWL encoding (PPEO) to enable RDF

- JSON-schema version (in progress)

# Developing Ontologies for Plant Data Annotation

## Crop Ontologies

- Describe traits/features and methods for measuring them in specific plant species (MIAPPE observed variables)
  - Woody Plant Ontology
  - Rice Ontology

## Plant Experimental Assay Ontology

*Node Service*

- Describes experimental procedures and pipelines in Plant (Molecular) Biology

# Updating BrAPI

## Breeding API 2.0

- Full coverage of mandatory MIAPPE fields (based on mapping)

- BrAPI to ISA-Tab exporting

- BrAPI to RDF or JSON-LD exporting (using PPEO as context)

# Setting up a national BrAPI end-point

## Breeding API

- BrAPI web service implements 13 BrAPI calls

- Underlying SQL database reverse-engineered from BrAPI specifications (shifting to RDF)

- Currently includes manually curated datasets on cork oak (Q. suber), rice (O. sativa) and J. curcas

# Setting up a federated data lookup service

**FAIR Data-finder for Agronomic REsearch**

# Enabling Web Page Findability with BioSchemas

- BioSchemas:
  - Extension of schema.org for the life science domain
  - Mark-up of web pages to enable findability

- Goals:
  - Map BioSchemas to MIAPPE fields critical for dataset findability
  - Mark-up CorkOakDB pages with BioSchemas
  - Setup FAIDARE to ingest BioSchemas mark-up

# Improving MIAPPE Usability

## User Interfaces for MIAPPE Dataset Submission

- ISA-Tab template for ISA-Tools, not the most user-friendly

- Templates for popular data management platforms:
  - Dataverse (ready for use)
  - SEEK 4 Science (nearly finalized)

- Stand-alone OWL-based interface (in progress)

# Providing Training on MIAPPE

## Plant Data Management Workshops

- Goal: Teach participants to describe a dataset according to MIAPPE v1.1

- Methodology: Lectures on introductory data management and MIAPPE, followed by a hands-on group exercise on a mock dataset, using a prepared MIAPPE template

- Past Events: Oeiras, Nov-2019; Paris, Feb-2020

- Future Events: TBD, virtual!

# Engaging with the Industry

## Data Producers: The Navigator Company

- A major player in the international pulp and paper market and one of Portugal's strongest brands on the world stage

- Massive amounts of data on eucalyptus breeding and genetics:

  - Genotypic and phenotypic data on over 300,000 specimens across a range of sites and covering up to 4 generations of pedigree

# Engaging with the Industry

**Data Producers: The Navigator Company**

Goals:

- Process, annotate and ingest pilot datasets into our BrAPI end-point

- Share knowledge and demonstrate the value of plant data FAIRification (e.g. integration with external datasets)

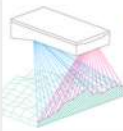- Automate the process so the company can submit its datasets to BrAPI in bulk
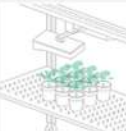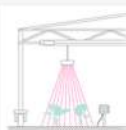
# Engaging with the Industry

## Technology Producers: Phenospex

- Provide state-of-the-art automated plant phenotyping solutions to researchers across the world

- Software team interested in ensuring their information system is MIAPPE compliant and can export data via BrAPI calls

- An ideal partnership to ensure researchers produce FAIR data!

# Putting it All Together

## One Standard to Rule Them All

- Whichever data management tool or data submission platform you choose for plant phenotyping data shall be MIAPPE compliant

- Metadata from all tools and platforms shall be interchangeable

- Genotyping and genomic data shall comply with MIAPPE specifications of Biological Materials to enable integration

- Bioschemas shall map to upper layers of MIAPPE

# Putting it All Together

## The ELIXIR RDM Toolkit – Plant Domain

- Document tools and resources for plant FAIR data management

- Address FAQs and common needs

- Prepare Data Management Plan templates for the plant domain, referencing the resources developed by the community and listed in the RDM Toolkit

# Other Resources

# Conclusions

# Plant sRNA Portal

## miRPursuit

- Automated workflow for downstream analysis of gene expression data and prediction of sRNA target genes

## sRNA Database

- Repository of sRNA either annotated using miRPursuit or publicly available

**Node Service**



**Under Development!**
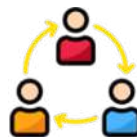
# Conclusions & Perspectives

Engagement and adoption of resources by the community to promote FAIR data and provide easier ways to analyse and gain new insights from available data…

Training sessions on the usage of resources, plant data management, training materials…

Icons made by Freepik

Collaboration with other ELIXIR nodes - European projects (Converge,…), Implementation Studies, ELIXIR Knowledge Exchange and Staff Exchange Programs,…

Interaction with EMPHASIS, Crop Ontology, Bioversity,…

# Thanks!

The BioData.pt Plant Sciences Community

Célia Miguel, FCUL / iBET
Margarida Oliveira, ITQB-NOVA
Nelson Saibo, ITQB-NOVA
Pedro Barros, ITQB-NOVA
Marcos Ramos, CEBAL
Inês Chaves, ITQB-NOVA
Daniel Faria, INESC-ID / IGC
Bruno Costa, FCUL / INESC-ID
Marta Silva, ITQB-NOVA
Filippo Bergeretti, ITQB-NOVA
André Cordeiro, ITQB-NOVA

Past collaborators:
Daniel Sobral
Cirenia Arias-Baldrich

All the elixir Plant Sciences Community

Especially our close collaborators:
Cyril Pommier, ELIXIR-FR
Anne-Françoise Adam-Blondon, ELIXIR-FR
Celia Michotey, ELIXIR-FR
Richard Finkers, ELIXIR-NL
Evangelia Papoutsoglou, ELIXIR-NL
Frederik Coppens, ELIXIR-BE
Paul Kersey, previously ELIXIR-EBI