

The ELIXIR infrastructure for plant phenotyping-genotyping data management

Anne-Françoise Adam-Blondon

▶ To cite this version:

Anne-Françoise Adam-Blondon. The ELIXIR infrastructure for plant phenotyping-genotyping data management. Workshop "Plant Phenotyping-Genotyping Data Management", ITQB; ELIXIR-PT, Nov 2019, Oiras, Portugal. hal-03316380

HAL Id: hal-03316380 https://hal.inrae.fr/hal-03316380

Submitted on 6 Aug 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License



The ELIXIR infrastructure for plant phenotyping-genotyping data management

A-F Adam-Blondon, INRA, ELIXIR-FR



Context





Predictive biology: Genomic x Environment x Phenomic



From C. Miguel and C. Pommier

3

Reusing and integrating the data : e.g. modelling the impact of climate change using plant phenology

Global plant phenology data portal: <u>www.plantphenology.org</u> Pan European PEP725 Plant phenology database: http://www.pep725.eu/



Phenology data : different sources, different accuracy in terms of identification of the species, scoring methods, record formats

- Modelers of the impact of climate change
- Geneticists, Breeders
- Genbank managers
- Experimental station managers
- Civil (economic) society: e.g. vintage dates, cherry blooming date...

Europe has been continuously reinforcing its policy for facilitating open access to data



"Facilitating access to results encourages the re-use of research outputs and supports Open Science. This is essential for Europe's ability to enhance its economic performance and improve the capacity to compete through knowledge. [...] Results of publicly-funded research can therefore be disseminated more broadly and faster, to the benefit of researchers, innovative industry and citizens.



Data management plans are evaluated in EU projects proposals

This policy is implemented via the European Open Science Cloud (EOSC)

- The European infrastructures in Life Science should contribute to EOSC
- E.g. European Infrastructure of Bioinformatics for Life-sciences: ELIXIR (23 countries)



From N. Blomberg

ELIXIR's objectives

In 2023: Continent-scale, standards-based infrastructure for accessing and analysing life-science data



Contacts Plants: Cyril Pommier (FR), Astrid Junker (DE), Kristina Gruden (SL)

From N. Blomberg



From Adam-Blondon et al (2016) HortRes, 3:16056

Insertion of data repositories in federations of information systems



Featuring federations of information systems

- A network of (stable/sustainable) nodes
- A central portal offering services (e.g. search data)



Backbone of good practices enabling such infrastructures



Wilkinson et al (2016) SCIENTIFIC DATA, 3:160018, DOI: 10.1038/sdata.2016.18

Examples of federations



Three "use cases" of federations

- The Wheat Initiative (G20 Initiative) and its Wheat Information System Expert Working group (<u>www.wheatis.org</u>). Also supported by the Research Data Alliance.
- The European Infrastructure for Multi-scale Plant Phenomics and Simulation (EMPHASIS) and its information system (<u>https://emphasis.plant-</u> <u>phenotyping.eu/e-Infrastructure</u>). In the frame of a strong collaboration with ELIXIR
- The data federation of the ELIXIR Plant Science community (<u>https://elixir-europe.org/communities/plant-</u> <u>sciences</u>)







https://emphasis.plant-phenotyping.eu/e-Infrastructure



Developing a federation of FAIR plant data repositories

Development of guidelines: e.g. www. wheatis.org



RESEARCH DATA ALLIANCE

Dzale-Yeumo et al (2017) F1000Research, 6 : 1843

Registries of standards and guidelines



Example of needed resources for a federation of phenotyping data

Semantic

- Description of the data
- Controlled vocabularies: term name and definitions
- Ontologies: semantic links between terms
- Biologist driven



Structure

- Formatting and Organizing the data
- Data Models
- Standards : VCF, GFF, MIAPPE (<u>www.miappe.org</u>) , etc...



Biologist & Computer scientist driven

Technical

- Data integration and sharing
- Interoperability : tools and systems
 - GA4GH 🌔
 - Breeding API <u>www.brapi.org</u>
 BrAPI
- Computer scientist driven



Development of a metadata standard for phenotyping experiments



- MIAPPE: Minimum Information About Phenotyping Experiment
- www.miappe.org
- Steering committee Emphasis, Elixir CGIARs

Last release MIAPPE v1.1 (Jan. 2019). Major improvements:

- Extension to accommodate woody plants as an additional use-case.
- Specification of a data model for easier implementation in various formats and automatic validation.
- Improved compatibility with <u>ISA-Tools</u> and <u>Breeding API (BrAPI)</u>.
- Provision of clear definitions and examples for all fields.



Adoption of the Crop Ontology format for the description of the phenotyping variables (www.cropontology.org)

Variable = trait + method + scale

Examples

- Woody Plant Ontology
- Wheat INRA Phenotyping Ontology v1.3 beeing merged with the Wheat Ontology developed by the CYMMIT
- Grape Ontology (v2)
- Also used to describe environment variables
- Strong curation efforts still needed: documentation of the methods, standardize the vocabulary (traits: entity and quality), ...
- Environments facilitating the development and curation of ontologies are needed



Registries of identifiers for key objects

- DOI for plant accessions (following the FAO recommandations)
- Biosample : identifiers for samples derived from accessions
- Crop Ontology : identifiers for phenotyping variables
- DOI associated to phenotyping trial sets (and data papers)

Developing services across federations: e.g. search and access data



WheatIS Generic Data Discovery Tool



Spannagl et al. 2016, https://doi.org/10.3835/plantgenome2015.06.0038











WheatIS data discovery: https://urgi.versailles.inra.fr/wheatis/



Alaux et al. Genome Biology 2018, https://doi.org/10.1186/s13059-018-1491-4

transPLANT data model



elițir

Perspective: map the data model to Bioschema terms (schema.org) to enhance its web findability





WheatIS data discovery tool: evolution

WheatIS nodes



WheatIS data discovery tool

Wheat@URGI WheatIS	Wheat Initiative		
URGI IWGSC@GnpIS [18 566 139] GnpIS [92 214] OpenMinTeD [3 398] WheatIS File Repository [6] EBI Ensembl Plants [2 122 980]	WheatIS Wheat Information System		1×
IPK CR-EST [199 220]	Examples: yield,	fhb	Search
GEBIS [50 875]			

MetaCrop [177] Gramene Gramene [229 789]

South Green AaroLD [137 060]

Wheat Pangenome [167 167]

The Triticeae Toolbox [138 441]

Rothamsted Research

KNetMiner [110 775] GrainGenes GrainGenes [15 827] Wheat Gene Catalog at

Komugi [3 043] PGSB

IPGPAS PlantPhenoDB [3]

CrowsNest [13 324] CIMMYT

CIMMYT Dspace [981]

CIMMYT dataverse [1]

UWA

Τ3

Open software, very generic, that can (and is) adapted to any type of federation:

e.g. the federation of information systems for of the french infrastructure of genetic resources for research in agriculture, AgroBRC-RARe (https://urgi.versailles.inra.fr/rare/)

Light metadata based on the data model and no constraint on vocabularies

No data integration

Towards data integration in federations





Breeding API initiative

http://www.brapi.org/

The Breeding API (BrAPI) Project is an **international effort** to create a RESTful specifications for a web service that enables interoperability among plant breeding databases.





BrAPI Breeding API initiative

http://www.brapi.org/

- Development of a standard API :
 - Calls for plant material aligned with MCPD
 - Calls for phenotyping experiments aligned with MIAPPE and of a supporting data model
- Next steps: develop the same type work on the calls for genotyping data (coordination with GA4GH)

Selby et al. Bioinformatics (2019), doi.org/10.1093/bioinformatics/btz190



Enabling improvements of data services by an international community



Anyone can develop a service that will « consume » the data of the federation exposed with the **BrAPI** standard: « machines understand your data »



ELIXIR Plant Data Search Service: FAIDARE

- Open source software
- Links
- Senomic Centralised repository
 - Phenotype Distributed repositories





Sources	Germplasm Trait	Reset all	~
URGI GnplS (81,335) EBI European Nucleotide Archive (44,975) CIRAD TropGENE (722) VIB PIPPA (692) BET BioData (67)	Crops (common name, species, genus, subtaxa &	Search crops	
	synonyms) Germplasm list (papel, collection & population)	Search germplasm lists	
Types Germplasm (94,589) Genotyping Study (32,210) Phenotyping Study (992)	Accession (accession name, number & synonyms)	Search germplasm accession	
	,		



elixir





FAIDARE : Web Interface

URGI - Data providers - More... -

FAIDARE

Sources

Types

Populus ×

Search cr

Search ge

Search g

EBI European Nucleotide Archive (6)

Germplasm (10,101)

Genotyping Study (6)

Phenotyping Study (5)



URGI • Data providers • More... •

FAIR Data-finder for Agron

Sources

URGI GnpIS (5)

Types

Germplasm (10,101) Genotyping Study (6)

Germplasm Trait Crops (common name, species, genus, subtaxa & synonyms)

Germplasm list (panel, collection & population)

Accession

(accession name, number & synonyms)

Results:



(URGI GnpIS data source link)

Description: "bacterial canker resistance test of ma 2000-04-01 (seasons: 2002).



(URGI GnpIS data source link)

Description: "clonal test of mapping pedigree 0504 (seasons: 2004, 2005). this study is part of the POPY



(URGI GnpIS data source link)

Description: "clonal test of mapping pedigree 0504



Germplasm	Trait	Reset a
Crops (common name, species, genus, subtaxa & synonyms)		Populus ×
		Search crops
Germplasm list (panel, collection & p	population)	Search germplasm lists
Accession (accession name, nur	nber & synonyms)	Search germplasm accession

Results:

From 1 to 10 over 6 documents

elixii

Genotyping Study EBI European Nucleotide Archive Identifying the genetic bases of plant traits and community composition in Populus tremuloides

(EBI European Nucleotide Archive data source link)

Description: "Identifying the genetic bases of plant traits and community composition in Populus tremuloides" is a Genotyping study.

Genotyping Study EBI European Nucleotide Archive Populus euphratica and Populus pruinosa Raw sequence reads

(EBI European Nucleotide Archive data source link)

Description: "Populus euphratica and Populus pruinosa Raw sequence reads" is a Genotyping study.

Genotyping Study EBI European Nucleotide Archive MiRNA-targets regulate adventitious rooting in Populus using degradome sequencing

L MONIME O MELICENSE

(EBI European Nucleotide Archive data source link)

Genotypes	661300230	661300230	661300230	Populus x generosa
	661300232	661300232	661300232	Populus x depense



How do we (easily) get the data FAIR into the data federation?





Example of five large crop centered French national projects on sugar beet, maize, wheat, pea and rapeseed (2012-2020)

AKER, AMAIZING, BREEDWHEAT, PEAMUST, RAPSODYN Transverse actions in relation with data management

- GnpIS for the integration of heterogeneous public and private data according to the data management plans described in the consortium agreements
- Optimisation of the developments of GnpIS between projects



 Public private partnership to develop a suite of tools aiming at facilitating the insertion and integration of partner's data in GnpIS



Nb: these projects started before we knew about the FAIR principles and at the very start of the Wheat Initiative

Publishing FAIR data in GnpIS

Enforcement of good practices for data traceability in the consortia

- Plant material identified under the responsibility of the genebanks
- Development of crop ontologies for phenotypic variables (e.g. <u>https://www.cropontology.org/ontology/CO_333/Beet%20Ontology</u>)

	Genetic Resources			Phenotyping			Geno	typing	Association		
	Public	Priv	ate	Public		Privat	te	Public	Private	Public	Private
Brassica	981	18		5							
Beta	10783			5							
Zea	1861	869		20		3		1	4	3	
Pisum	1015	872		0		86			3		
Wheat	10448	281	1	821		57		1	22	43	1970
		Sequen		ce	Genetic		tic	Мар	Genome browse		owser
		Public	Priv	vate	Pub	olic	Pri	ivate	Public	Priv	vate
Brassica		4	2		2				1		
Beta											
Zea		8	8		18		21		4	3	
Pisum		5	2				12		1		
Wheat		7	27		29		1		6		

Publishing FAIR data in GnpIS



Active data stewardship helps a lot

Publishing FAIR data in GnpIS

Enforcement of good practices for data traceability in the consortia

3001

873 H

 Association of DOI to datasets using data.inra before integration in GnpIS (e.g. DROPS/AMAIZING dataset)



Improve/have a control on the FAIRness of the data sets that are deposited

Enabling GnpIS insertion in data federation

Web Interface: get MIAPPE compliant ISATab data

BrAPI: MIAPPE compliant programatic access

Data submission to GnpIS compliant with **MIAPPE** requirements



38

Conclusions



5 * Open Data



Courtesy of Daniel Jacob

Progressing towards FAIR and Open Data requires a multidisciplinary cooperation :

- Biologists
- Bioinformaticians
- specialists of ontologies/semantics

Data life cycle in (plant) biology research activities



Challenge: provide operational environments facilitating data management all along its life cycle



elizír

Priority 2020-2024 : ELIXIR-Converge H2020 project.

Use case Plant coordinated by ELIXIR-PT

ELIXIR Plant community





Thank you!

www.elixir-europe.org



Acknowledgements



- H. Quesneville C. Pommier
- M. Alaux
- M. Buy
- D. Charruaud
- G. Cornut
- J. Destin
- S. Diagne

S. Durand B. El-Houdaigui R. Flores C. Guerche E. Kimmel *T. Letelllier* C. Michotey N. Mohellibi National and International infrastructures/initiatives







Consortium





Other national and international projects





🖉 Pea MUST











