

# Systèmes d'information et analyses de données

mardi 05 avril, Paris

Isabelle Alic, Pascal Neveu, Cyril Pommier



**anr**®  
agence nationale  
de la recherche  
AU SERVICE DE LA SCIENCE  
ANR11-INBS-0012

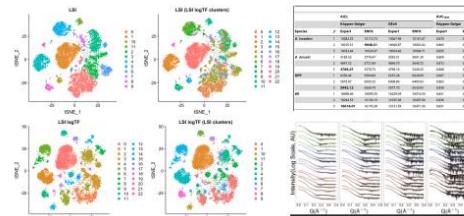
**INRAe**

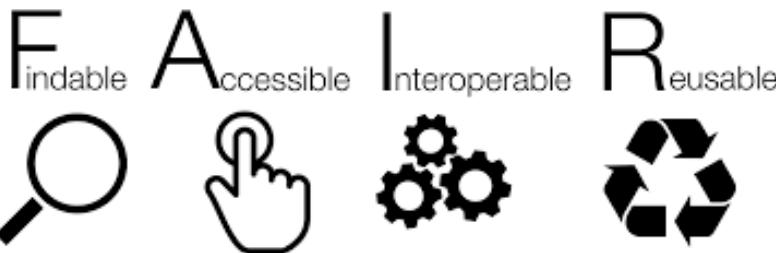
**ARVALIS**  
Institut du végétal

**Terres Inovia**  
l'agronomie en mouvement

# Phenomics Data Challenge

- **Experimentations in the domain of phenomics**
  - Expensive, require a lot of resources and often very hard
  - Cannot be reproduced
  - Huge and **very complex datasets**
- Strong needs of **transparency** and reproducibility
- **Give value to data:** re-analyses, meta-analyses and new analyses  
→ impossible without **advanced data management**





**Findable:** persistent ID, indexed in portals, standardized and relevant metadata

*Challenge: coordinated and sustainable data services*

**Accessible:** open and standardized protocols, license rights

*Challenge: cultural evolution of teams*

**Interoperable** (technology, syntax, semantic): shared standardized formats, vocabularies and formal languages for knowledge representation,

*Challenge: skill developement (interdisciplinary)*

**Reusable:** provenance, relevant metadata for understanding across disciplines,

*Challenge: new analysis methods*

❖ A need for new generation of Information Systems

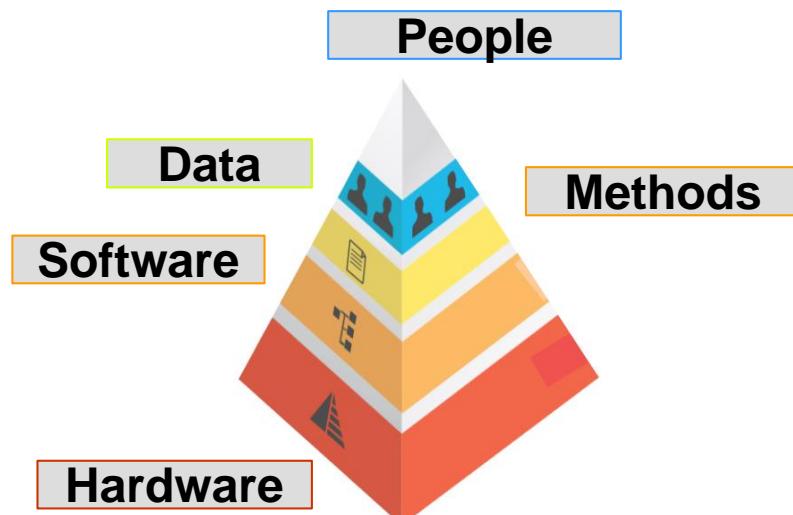


**PHIS:** Information System for phenomics based on **OpenSILEX**



**OpenSILEX: an open source software set**

- Methods, tools, components to implement information systems for experimental data in agriculture and environment
- Information System: Organized system for the collection, organisation, storage, exchange and treatment of information



# Data structuring

The **structuring of data** enables a computer system to perform store, retrieve, process data and implement good practices:

- Make **FAIR data**
- **Flexible**
- Ability to allow **understanding (and reproduce) data processing**
- Ability to enforce DMP and Open Science

## Structuring of data based on 2 key elements:

- **Identification and Naming convention**
  - Objects: plants, plots, experiments, sensors, events, etc.
  - Persistent, unambiguous, resolvable, globally unique
- **Semantic and tagging (based on ontology set)**
  - Controlled vocabulary
  - Formalized relationships between entities
  - Data annotation and enrichment (search engine friendly)

# Data Structuring: PHIS approach



## PHIS → Ontology driven

**Scientific object instances** (plant, plant organ, plot, etc.) are formalized  
Identified by **URI** standardized, unambiguous, shared, etc

**Events** (management, faults, meteo, etc) are formalized  
Identified by **URI**

**Variables, Observations, Factors, Documents, Devices, Softwares**  
are associated with these Objects and Events  
Identified by **URI**

**Organisation and linking of Objects and Events** → done with a controlled  
**semantic** (reference ontologies, vocabularies, thesaurus, taxonomies) and  
application Ontologies (**RDF, OWL, SKOS**)\*



## URI

- **Standardized** and easy integration in Web application
- **Unambiguous**
- **Actionable** (resolvable, dereferencable)

URI → generated by tools under responsibility of local coordinator

### URI of plant:

`<http://phenome.fr/arch/2017/c17000118>`

### URI of pot:

`<http://phenome.fr/arch/2013/pc13001542>`

### URI of cart:

`<http://phenome.fr/arch/2013/ct1300123>`

### URI of cabin:

`<http://phenome.fr/arch/2018/ac180015>`

### URI of camera:

`<http://phenome.fr/arch/2018/ac180019>`



### URI of image:

`<m3p:arch/2017/ic17002295855>`



URL identifies what exists on the Web

URI identifies, on the Web, what exists

IRI identifies, on the Web, in any language, what exists

**prefix m3p :** <<http://phenome.fr/arch>>

**URI of plant:**

< m3p:2017/c17000118 >



**URI of pot:**

< m3p:2013/pc13001542 >

**URI of cart:**

< m3p:2013/ct1300123 >

**URI of cabin:**

< m3p:2018/ac180015 >

**URI of camera:**

< m3p:2018/ac180019 >

**URI of image:**

< m3p:arch/2017/ic17002295855 >

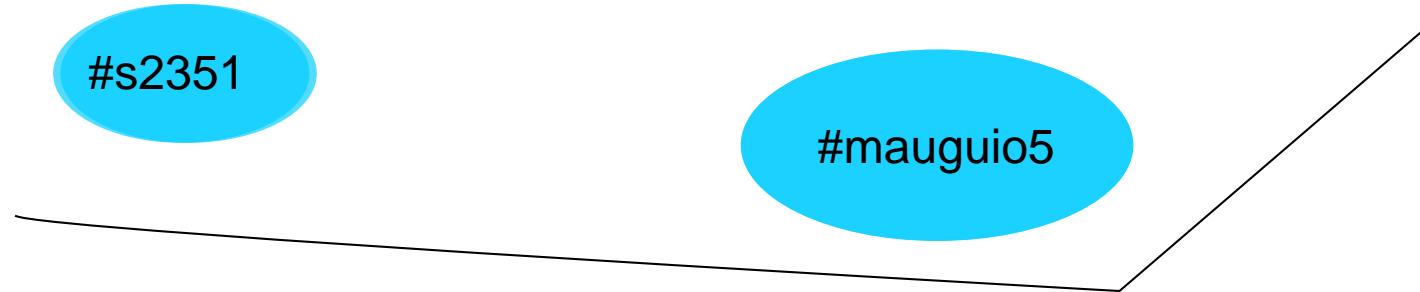


**Metadata / ontologies provide the meaning of data**  
→ Link data elements using a controlled and shared vocabulary  
and also **machine readable**



#s2351

#mauguio5

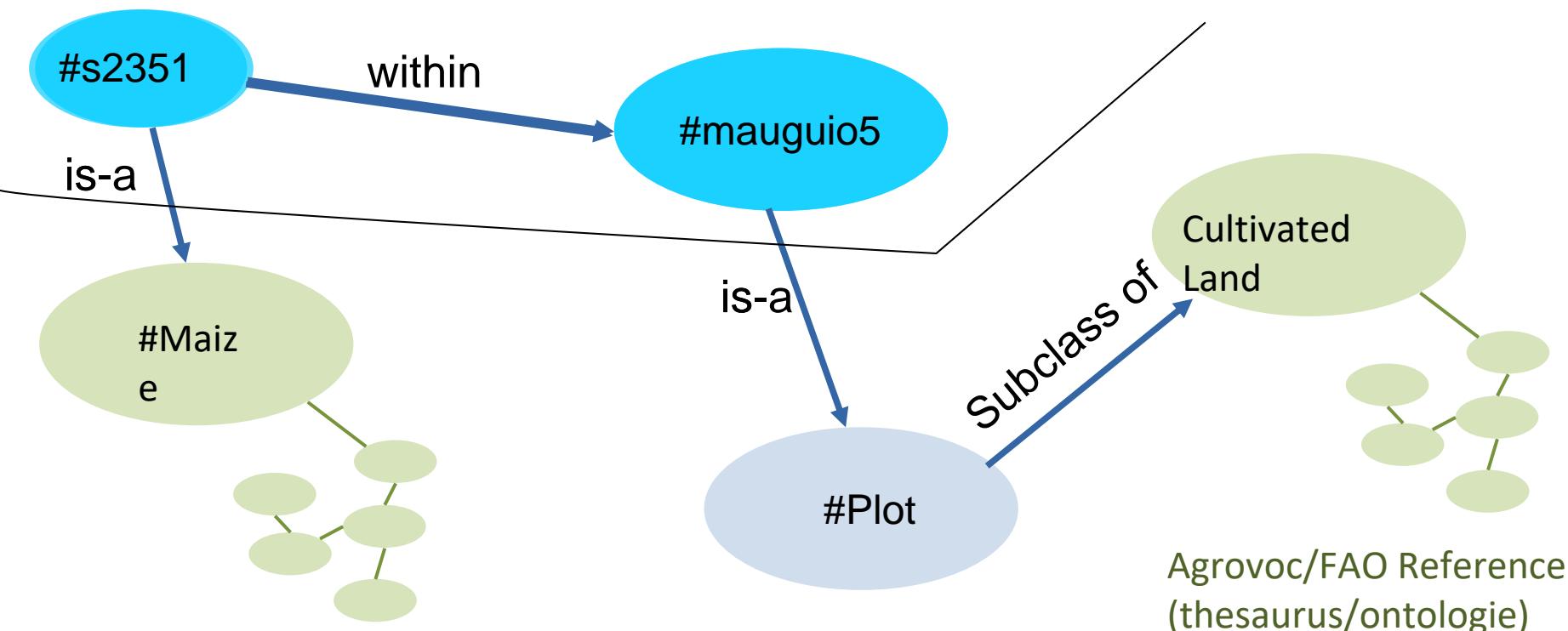




**Metadata / ontologies provide the meaning of data**

→ Link data elements using a controlled and shared vocabulary, and also machine readable

→ Data are structured in a graph that can be queried



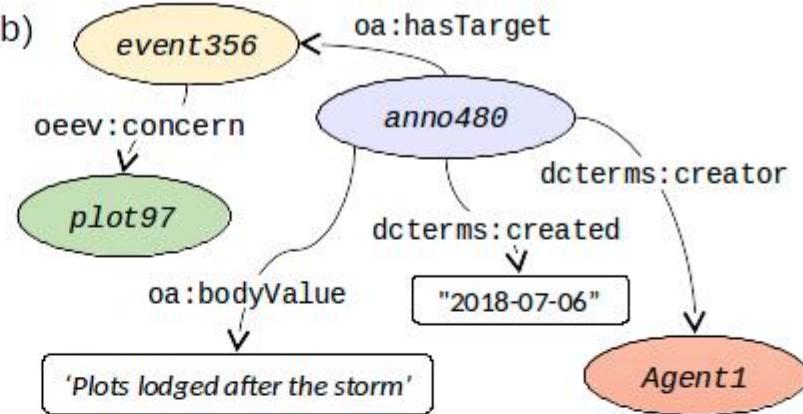
# Ontology driven Information System



(a)



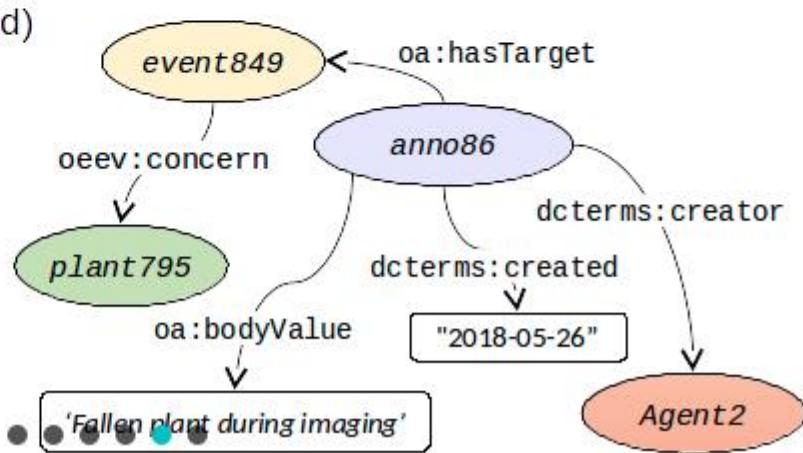
(b)

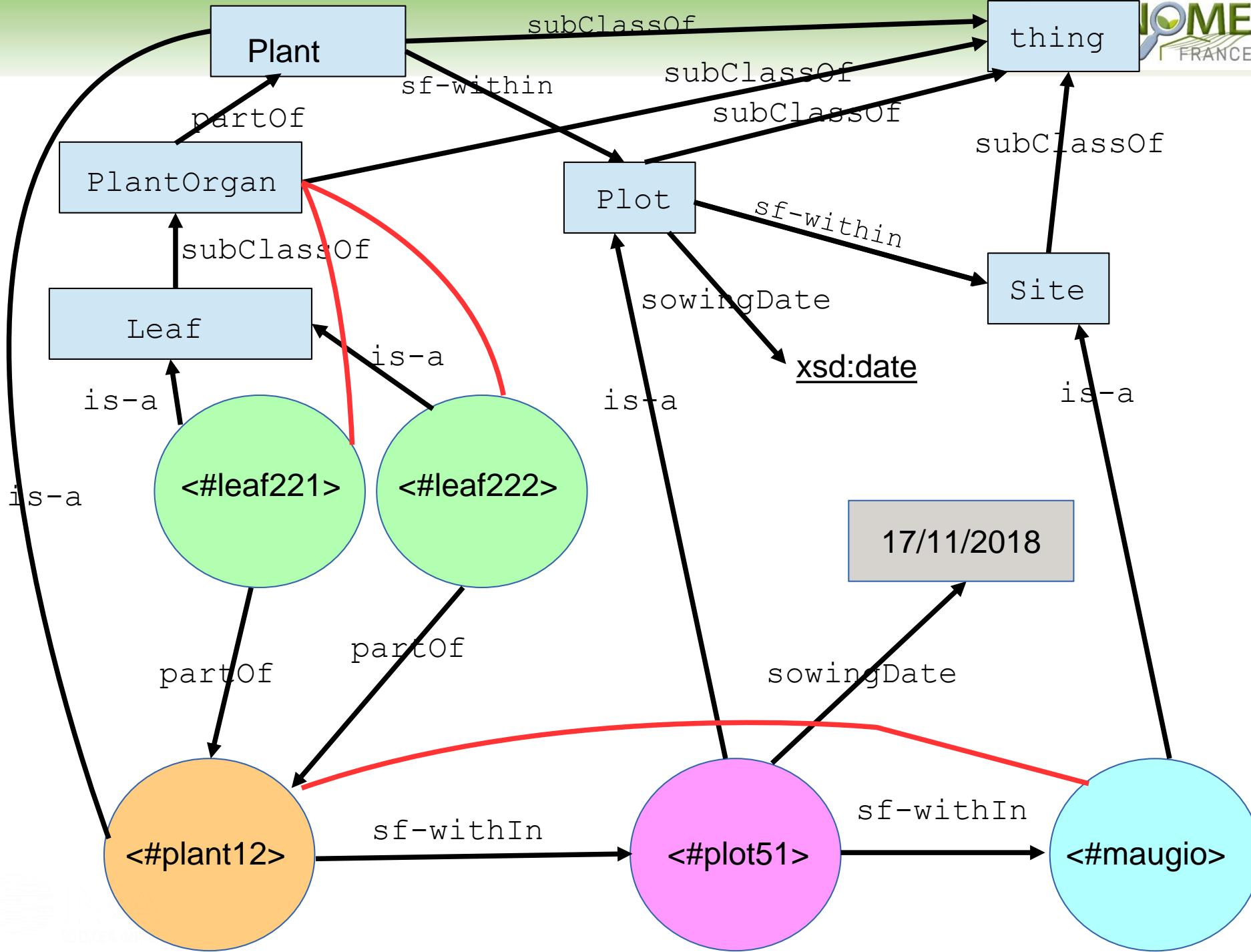


(c)



(d)





# Formalization of variables

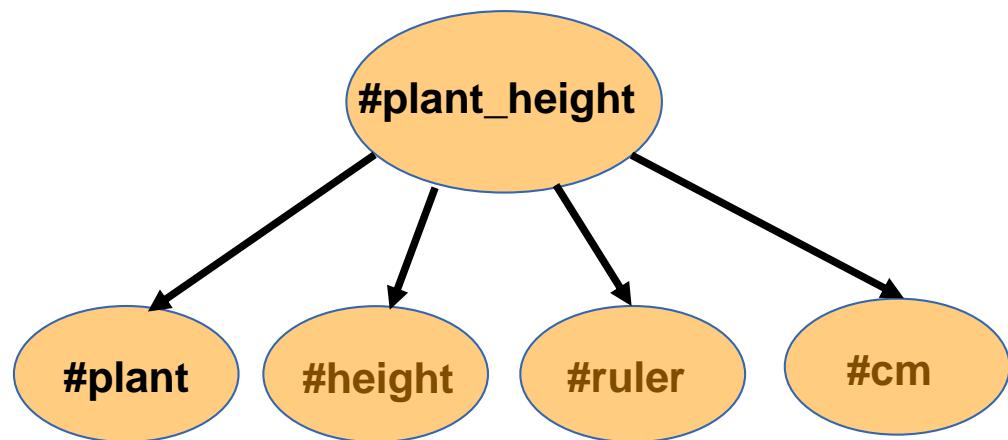


What do we recommended in PHIS

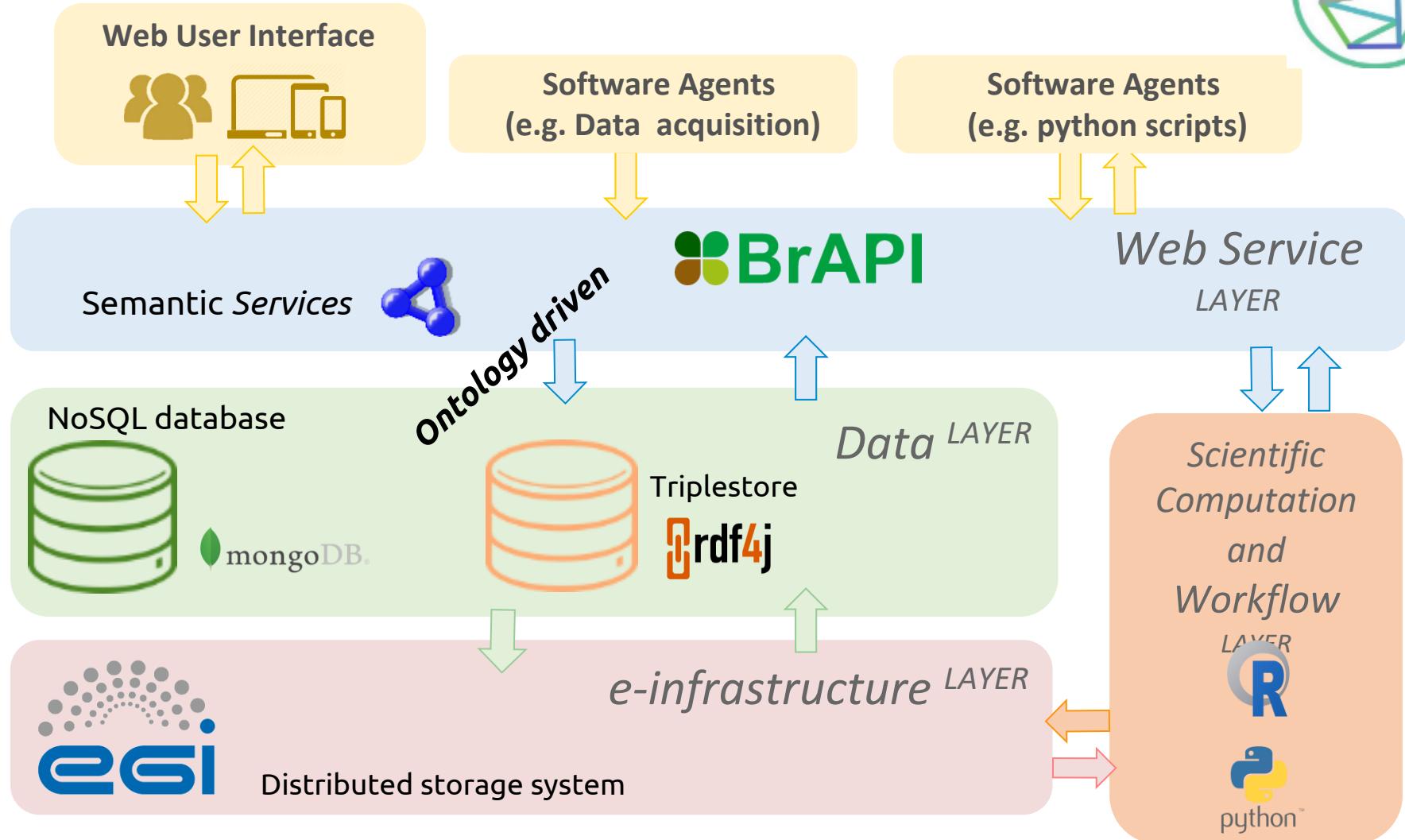
- Use URI (ID) for unambiguous variable (in global context)
- Actionable URI for an accessible description of variables  
Description can be **read by machine** and human
- Reuse existing variable description if available
- Use standardized representation schema for formalisation of new variable (and share it)

**Variable = Entity + Characteristic + Method + Scale**

**#plant\_height = URI**



# OpenSILEX - PHIS Architecture



## PHIS Web User Interfaces for the management of:

- Project information
- Experiment
- Facilities
- Devices
- Scientific objects
- Germplasms
- Experimental factors
- Data
- Data visualization
- Data provenance

The screenshot shows a web browser window with the PHIS logo at the top left. The page title is "MISTEA Server". The main content area displays a list of germplasm entries in a table format. The columns are labeled: Name, Type, and Species. The table includes entries for maize, Pearl millet, poplar, rice, sorghum, teosintes, upland cotton, accPoplar, B73, banana, barley, BC-seedlot-nonmais, bread wheat, CRAZI, and DKC4590. Each entry has a red trash icon and a blue edit icon. The sidebar on the left shows navigation links for Scientific Organization, Scientific Information, and Administration.

Name	Type	Species	Action
maize		Species	
Pearl millet		Species	
poplar		Species	
rice		Species	
sorghum		Species	
teosintes		Species	
upland cotton		Species	
accPoplar	Accession	poplar	
B73	Variety	maize	
banana		Species	
barley		Species	
BC-seedlot-nonmais	Seed Lot	banana	
bread wheat		Species	
CRAZI	Variety	maize	
DKC4590	Variety	maize	



- ✓ Allows management of huge and complex data
- ✓ Enables and facilitates cloud computing (data center, EGI)
  - distributed computing, distributed storage, backup
- ✓ Manage semantics (ontologies, standardized vocabularies)
- ✓ Flexible design
- ✓ Open technologies, Web APIs and portal interoperability
- ✓ Provenance and reproducibility for data processing
- ✓ Over 10 instances of PHIS for various installations (field and greenhouse)
- ✓ Distributed PHIS instances: >1 petabytes of data
- ✓ Other implementations of OpenSilex: WEIS, SunAGRI, Sixtine
- ✓ Open Software - support and development (MISTEA team)

# Conclusion

- ✓ A new generation of information systems (e.g. PHIS) is needed
- ✓ Giving value to complex data requires structuring according to FAIR principles
- ✓ A better formalization of concepts (using ontologies) and data is required for interdisciplinary research
- ✓ Advanced data management makes data available for AI and data analytics
- ✓ Support and training available



## ➤ PHIS demonstration

- <http://phis.inra.fr/>
- Research paper:  
<https://nph.onlinelibrary.wiley.com/doi/full/10.1111/nph.15385>

## ➤ How to contribute to OpenSILEX?

- Github repository: <https://github.com/OpenSILEX/>
- Developer documentation: <https://opensilex.github.io/docs-community-dev/>

## ➤ User documentation of the version in development:

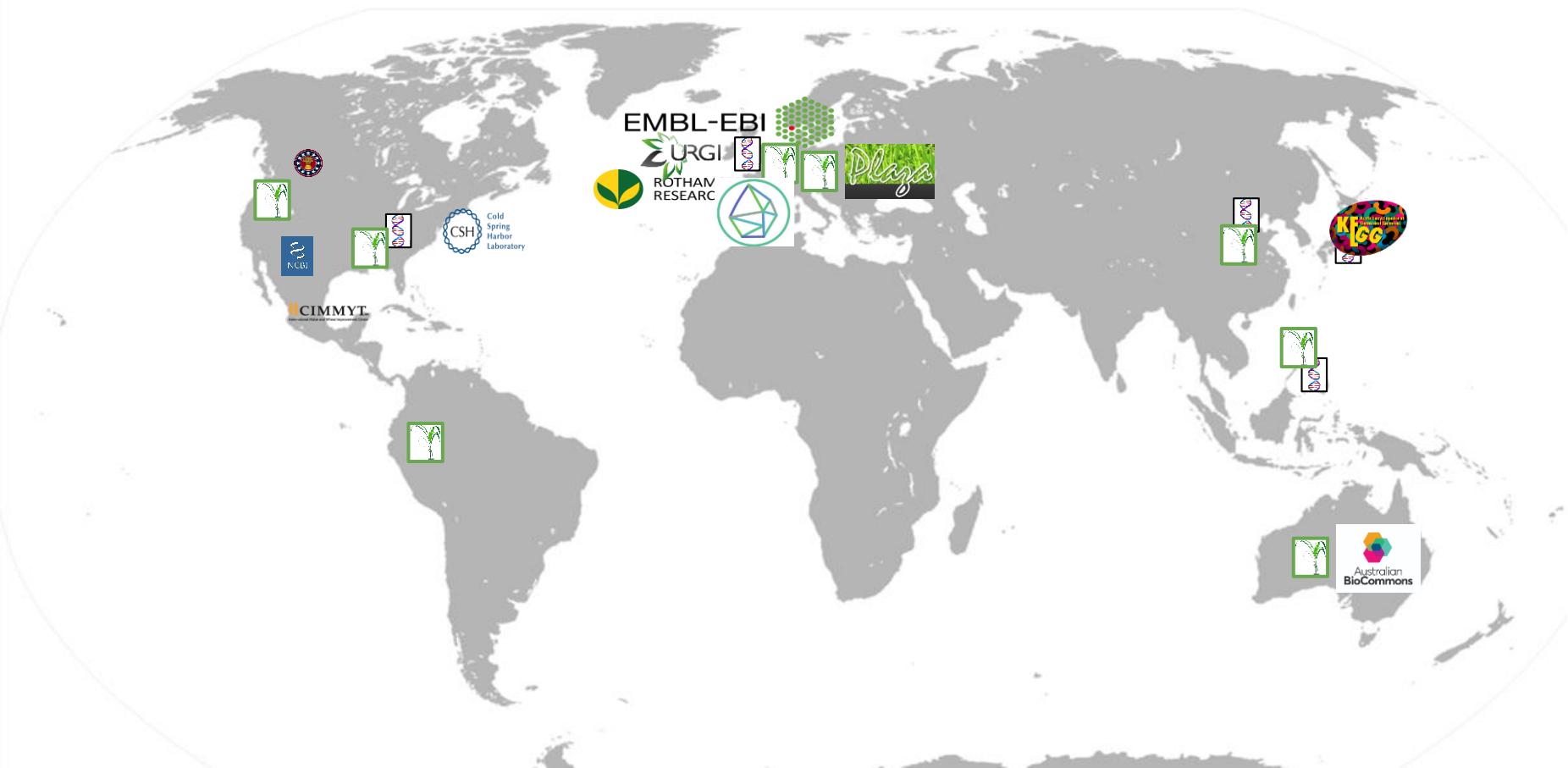
- <https://opensilex.github.io/phis-docs-community/>

# Global Information systems

Dispersed data

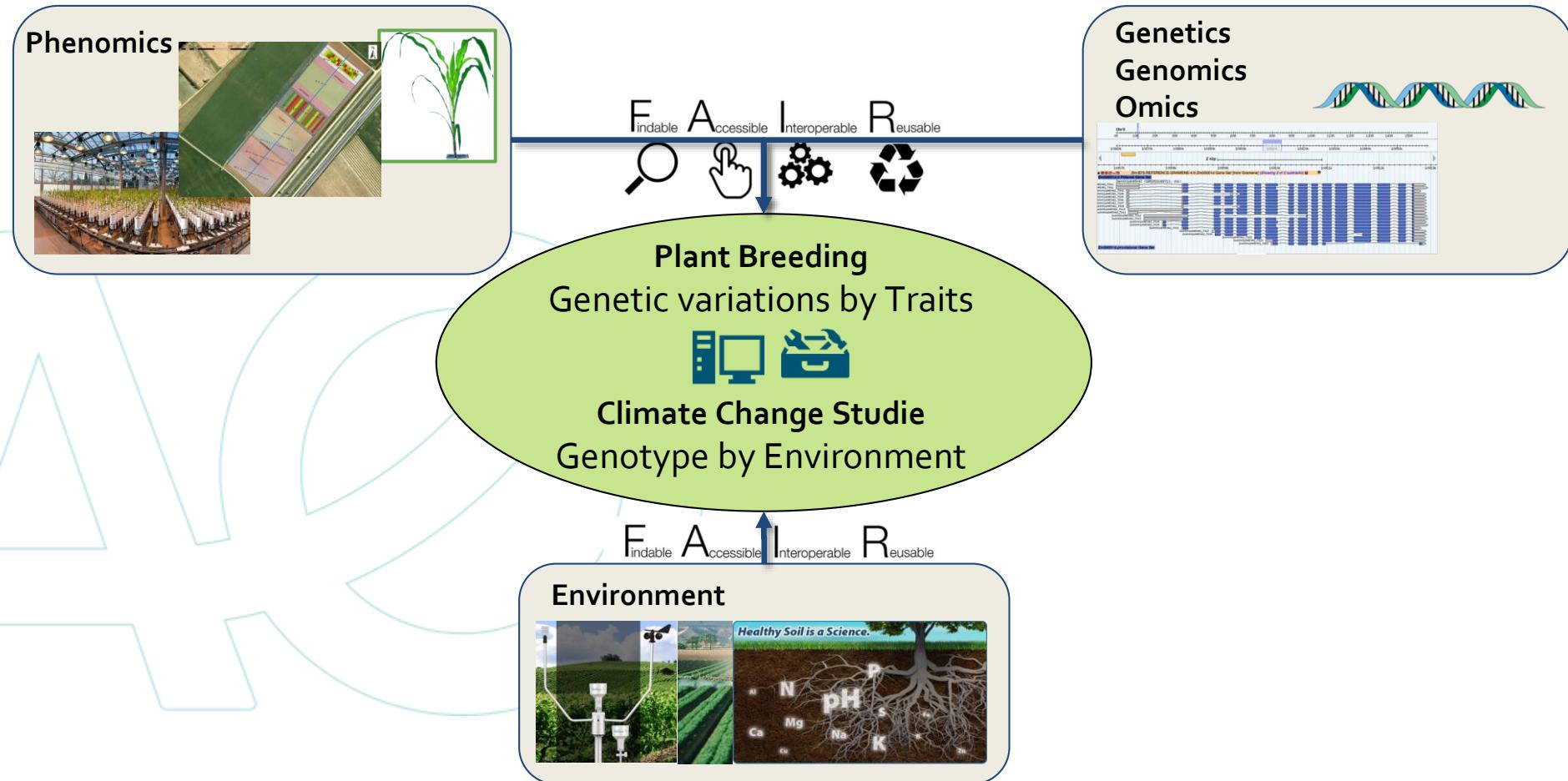
Heterogenous data

Dedicated repositories & Archives



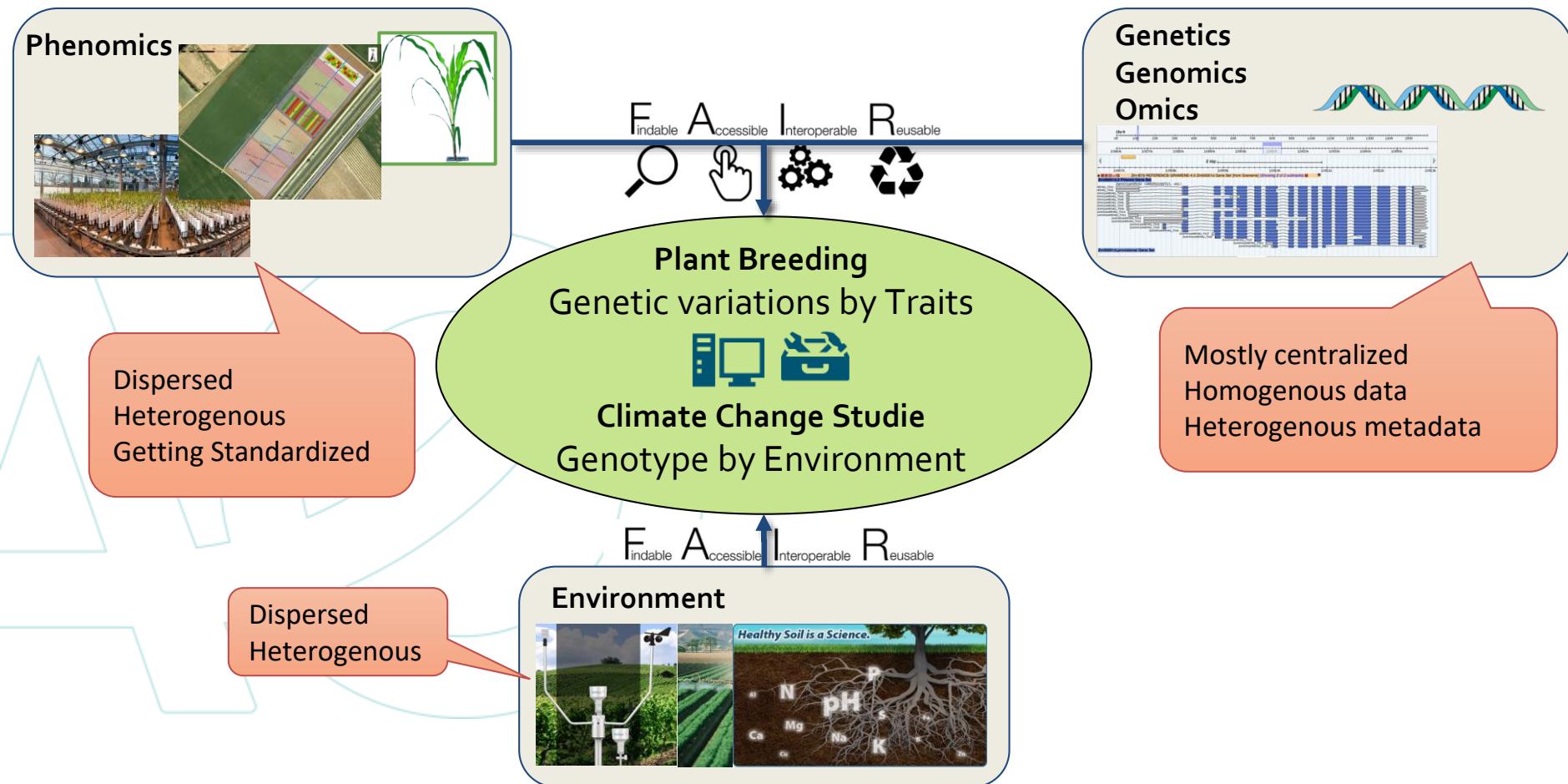
# Data standardisation and exchange

Environment / Phenome / Genomic / \*omic / Genetic



# PLANT use case

Environment / Phenome / Genomic / \*omic / Genetic



# Phenotyping data life cycle

« Raw » data, pheno/env measures, variables

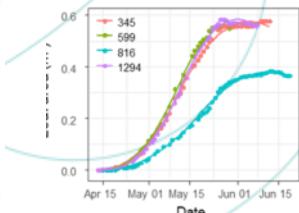
Soisson s x high N	Charger x high N	Charger x low input
Soisson s x low input	Soisson s x low input	Soisson s x high N
Charger x high N	Soisson s x high N	Soisson s x high N

« computed » data, reduced, indicators

Soisson s x high N	Charger x high N
Soisson s x low input	Charger x low input

Derivation, Reduction

Genotype	Treatment	N input	Date	Rep	Fusariose
Soissons	low input	15,32253129	15/11/2011	1	5
Soissons	low input	15,31430556	16/11/2011	2	7



Genotype	Treatment	Fusariose
Soissons	low input	6

661300270 Ardon	2004 45.645632645603683	12/01/2004 284.3
661300270 Ardon	2005	
661300444 Ardon	2004 38.96112577281653	12/01/2004 228.8
661300444 Ardon	2005	
661300312 Cavallermaggiore	2004 52.4	01/01/2004 249.9
661300312 Cavallermaggiore	2005	
661300371 Cavallermaggiore	2004 45.74	01/01/2004 230.2
661300371 Cavallermaggiore	2005	
661300487 Cavallermaggiore	2004 72.52	01/01/2004 309.8
661300487 Cavallermaggiore	2005	
661300585 Cavallermaggiore	2004 71.73999999999995	01/01/2004 305.7
661300585 Cavallermaggiore	2005	
661300468 Headley	2004 45.27	01/01/2004
661300468 Headley	2005	
661300469 Headley	2004 70.93000000000007	01/01/2004
661300469 Headley	2005	
661300533 Headley	2004 57.67	01/01/2004 258.8
661300533 Headley	2005	

# Plant Phenotyping Life cycle

## Raw data long term conservation

### Data acquisition

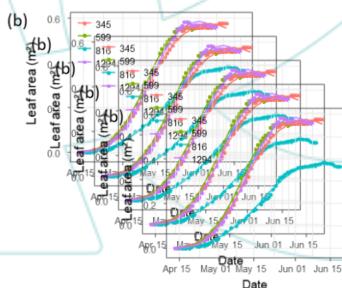
- **VARIABLES**
- Plant/microplot level
- Traceability
- Raw measures
- Data Cleaning
- Platform IS (Emphasis IS, PHIS, ...)
- Analysis Reproducibility
- Provenance

### Data computation

- **INDICATORS**
- Statistical integration
- Genotype level (mostly)
- New computation for each scientific question
- One raw dataset → many computed datasets

### Data publication

- One Data Publication by datasets.
- **Platform IS**
  - Phenomic, plant level
- **FAIR Data Repositories**
  - Reduced



Genotype	traitement	Fusariose
Soisson	low input	5
Soisson	high N	7
Charger	low input	1
Charger	high N	2

Variety charger is resistant to fusariose under intensiv cultural practice

# Plant Phenotyping Life cycle

## Raw data long term conservation

### Data acquisition

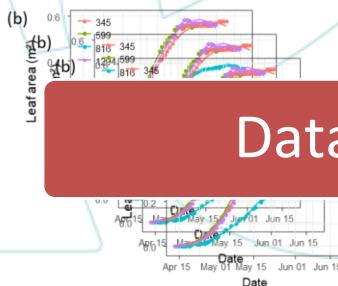
- **VARIABLES**
- Plant/microplot level
- Traceability
- Raw measures
- Data Cleaning
- Platform IS (Emphasis IS, PHIS, ...)
- Analysis Reproducibility
- Provenance

### Data computation

- **INDICATORS**
- Statistical integration
- Genotype level (mostly)
- New computation for each scientific question
- One raw dataset → many computed datasets

### Data publication

- One Data Publication by datasets.
- **Platform IS**
  - Phenomic, plant level
- **FAIR Data Repositories**
  - Reduced



Data

Genotype	traitement	Fusariose
Soisson	low input	5
Soisson	high N	7
Charger	low input	1
Charger	high N	2

Knowledge

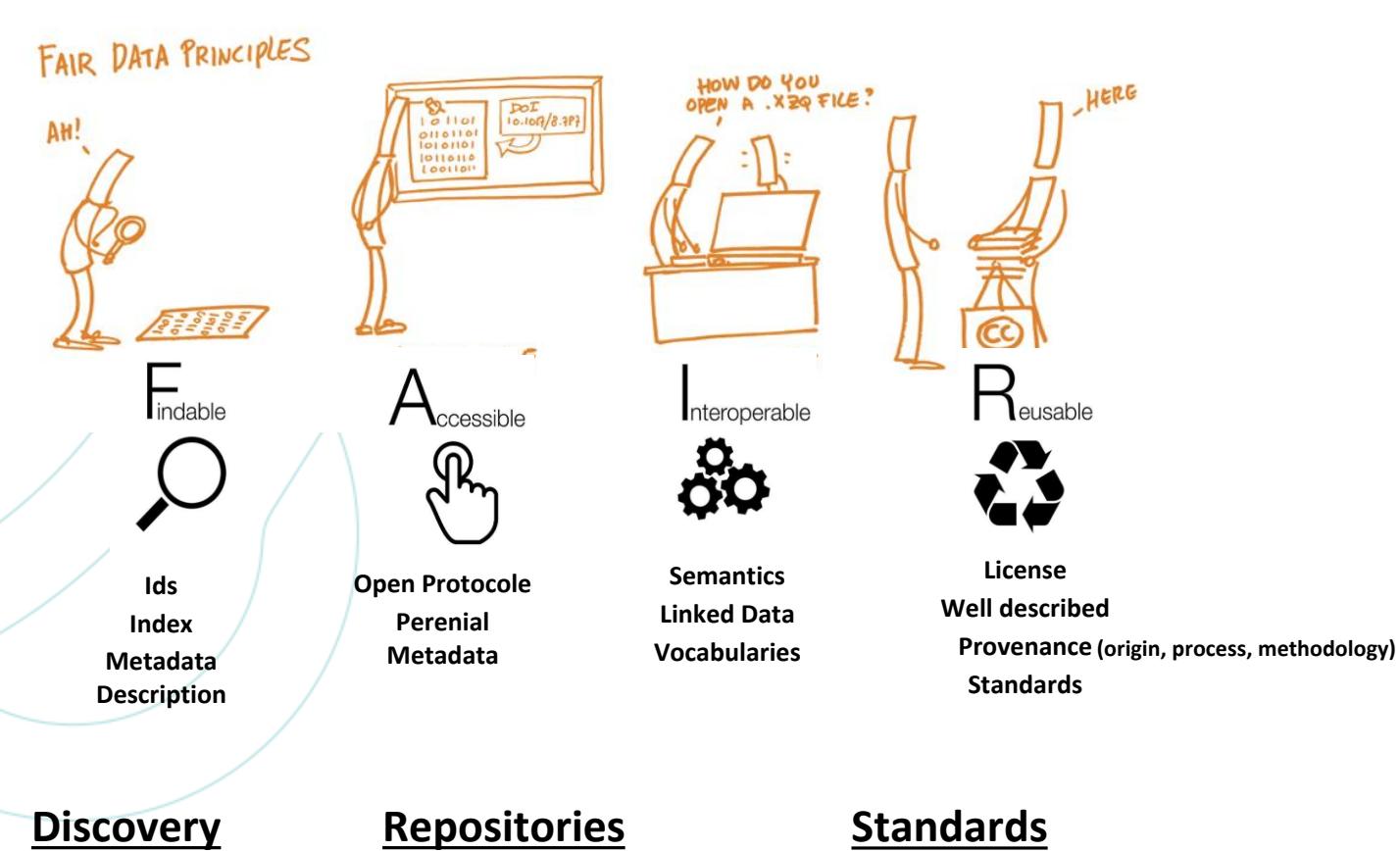
Variety charger  
to  
under

intensive  
cultural  
practice

# Open science through FAIR data principles

Sustainable data access over decades

Wilkinson et al., *The FAIR Guiding Principles for scientific data management and stewardship*.  
*Scientific Data* 3 (2016)



# FAIR data repositories

<https://projet-recherchedatagv.ouvrirlascience.fr/>



## Découvrir Recherche Data Gouv

Accueil Pourquoi Recherche Data Gouv ? Le projet Les acteurs du projet Ateliers de la donnée Actualités

La plateforme nationale fédérée des données de la recherche

LANCEMENT Printemps 2022

ACCOMPAGNER LES ÉQUIPES DE RECHERCHE

DÉPOSER & PUBLIER DES DONNÉES DE RECHERCHE

DÉCOUVRIR LES DONNÉES DE RECHERCHE

À propos de Recherche Data Gouv



Université de Lille



# FAIR data repositories

<https://projet-recherchedatagv.ouvrirlascience.fr/>

<https://data.inrae.fr/>

Portail Data INRAE > INRAE > Omics Dataverse > URGI Plant and Fungi Dataverse >



## A multi-site experiment in a network of European fields for assessing the maize yield response to environmental scenarios

Version 2.0



Millet, Emilie J.; Pommier, Cyril; Buy, Mélanie; Nagel, Axel; Kruijer, Willem; Welz-Bolduan, Therese; Lopez, Jeremy; Richard, Cécile; Racz, Ferenc; Tanzi, Franco; Spitkot, Tamas; Canè, Maria-Angela; Negro, Sandra S.; Coupel-Ledru, Aude; Nicolas, Stéphane D.; Palaffre, Carine; Bauland, Cyril; Praud, Sébastien; Ranc, Nicolas; Presterl, Thomas; Bedo, Zoltan; Tuberosa, Roberto; Usadel, Björn; Charcosset, Alain; van Eeuwijk, Fred A.; Draye, Xavier; Tardieu, François; Welcker, Claude, 2019, "A multi-site experiment in a network of European fields for assessing the maize yield response to environmental scenarios", <https://doi.org/10.15454/IASSTN>, Portail Data INRAE, V2, UNF:6:zF9w0A2f+MHeW7maeeXJWA== [fileUNF]

### Citer le dataset ▾

Pour en apprendre davantage sur le sujet, consulter le document [Data Citation Standards \[en\]](#).

### Description ⓘ

This dataset comes from the European Union project DROPS (DR panel of 256 maize hybrids was grown with two water regimes (irrigation in 2012 and 2013, respectively, spread along a climatic transect from one site in Chile in 2013. This resulted in 29 experiments defined by

### Modalités d'accès au dataset ▾

Contact

Partager

### Statistiques d'utilisation sur les datasets



a scenario but largely differ between scenarios.

Sciences; Agricultural Sciences

g, Genotype-by-environment interaction, Multi-env

f Yield in Europe: Allelic Effects Vary with Drought

ude Welcker, Willem Kruijer, Sandra Negro, Aude C

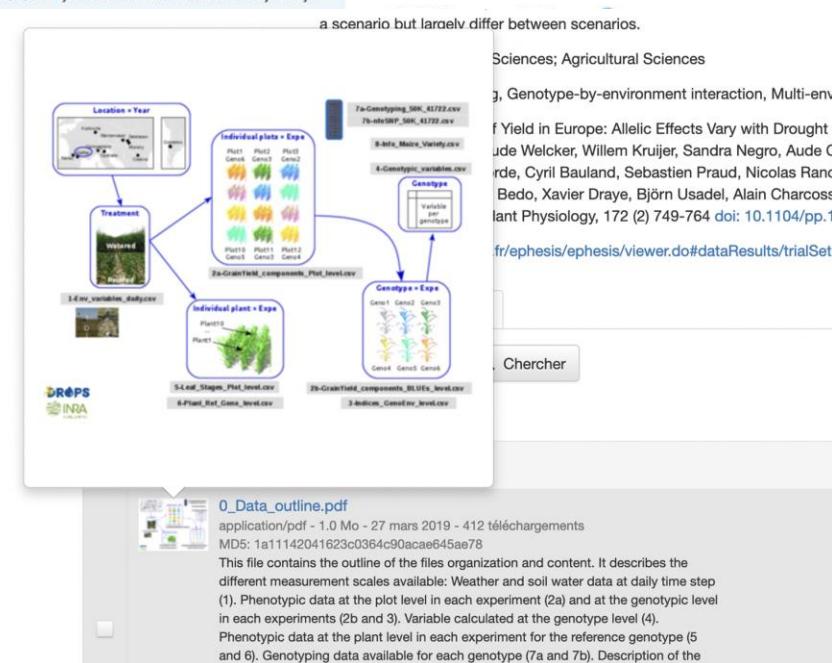
arde, Cyril Bauland, Sébastien Praud, Nicolas Ranc

Bedo, Xavier Draye, Björn Usadel, Alain Charcosse

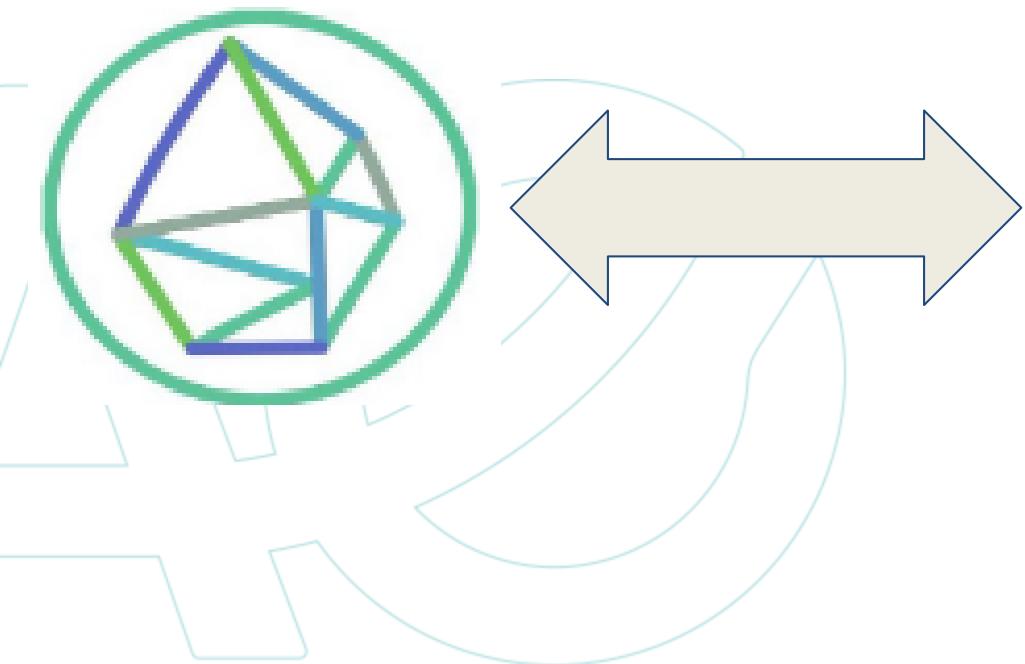
Plant Physiology, 172 (2) 749-764 doi: 10.1101/44

fr/ephesis/ephesis/viewer.do#dataResults/trialSett

Chercher



# FAIR data repositories



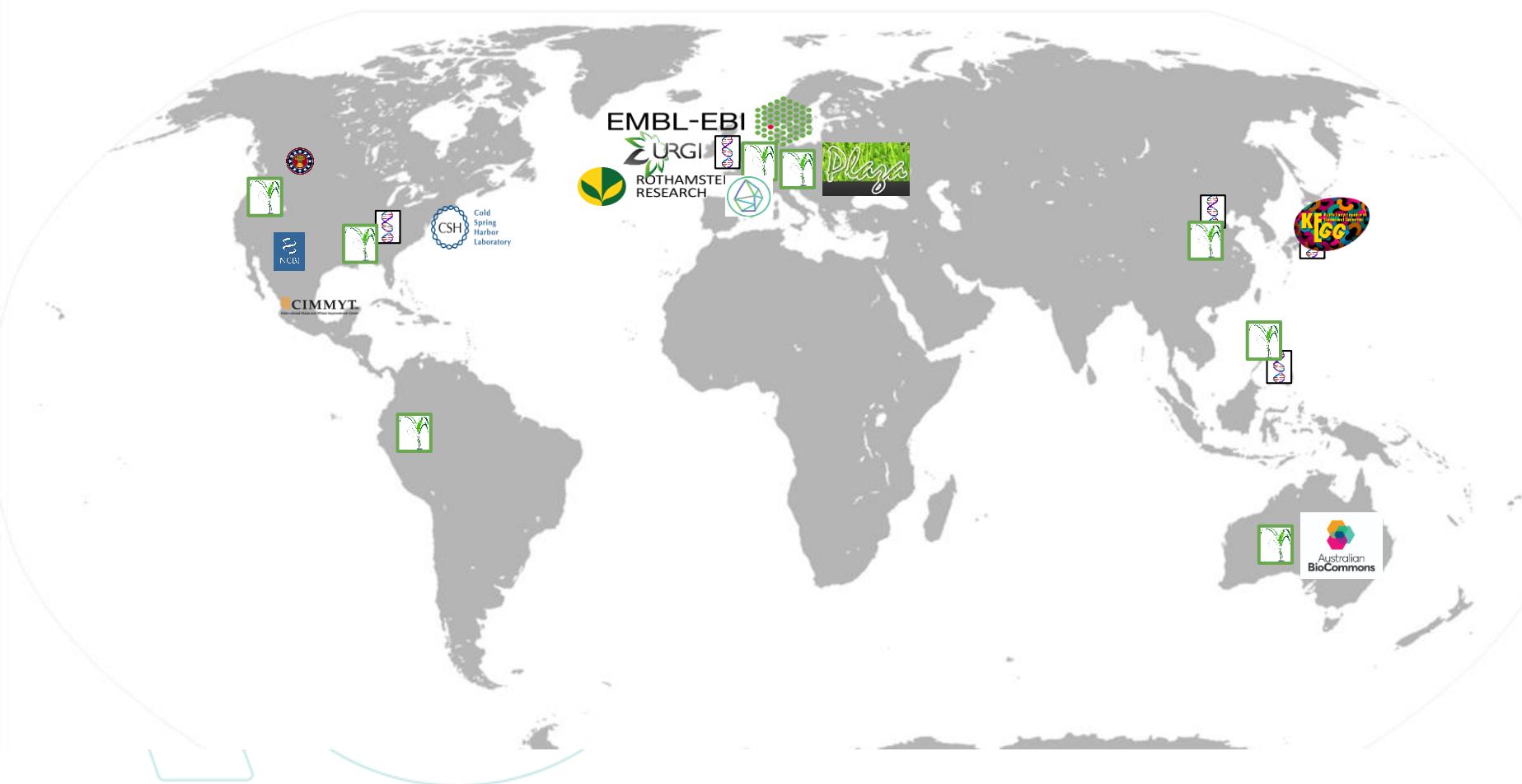
data **INRAe**

# Global Data discovery portal

Dispersed data

Heterogenous data

Dedicated repositories & Archives



# Global Data discovery portal

Dispersed data

Heterogenous data

Dedicated repositories & Archives

yield

Search

Results 1 to 20 of 156

**10.3389/fpls.2018.00529 - OpenMinTeD@GnplS**

Bibliography Triticum Triticum aestivum

Global QTL Analysis Identifies Genomic Regions on Chromosomes 4A and 4B Harboring Stable Loci for Yield-Related Traits Across Different Environments in Wheat (Triticum aestivum L.). 2018 Global QTL Analysis Identifies Genomic Regions on Chromosomes 4A and ... (expand)

**10.1186/s12864-019-6005-6 - OpenMinTeD@GnplS**

Bibliography Triticum Triticum aestivum

Genome-wide association study reveals new loci for yield-related traits in Sichuan wheat germplasm under stripe rust stress. 2019 Genome-wide association study reveals new loci for yield-related traits in Sichuan wheat germplasm under stripe rust stress Ba ... (expand)

**10.1186/s12863-019-0785-1 - OpenMinTeD@GnplS**

Bibliography Triticum

Accuracy of genomic selection for grain yield and agronomic traits in soft red winter wheat. 2019 Accuracy of genomic selection for grain yield and agronomic traits in soft red winter wheat Background Genomic selection has the potential to increase genetic ... (expand)

Species (21)

Filter on Species...

Data type

Bibliography [151]

None [5]

Ontology annotation (30)

Filter on Ontology annotat

Expand search ?

Database

OpenMinTeD@GnplS [151]

GEBIS [5]

# FAIDARE: Global Data discovery portal

<https://urgi.versailles.inrae.fr/faidare/>

The diagram illustrates the FAIDARE architecture. At the top, a world map serves as a background. Overlaid on it are several rounded rectangular boxes representing different data sources and services:

- miappe**: Phenomics (EMPHASIS, PIPPA), Genetics (Genomics, EMBL-EBI)
- BrAPI**: A central hub connecting the miappe modules and other services.
- Dataverse**: data INRAE
- Historical Datadiscovery**: WHEAT INITIATIVE
- Literature**: Europe PMC
- Bioschemas sources**: A green hexagonal icon.

A large blue arrow points downwards from the BrAPI hub towards a search interface. This interface includes:

- URGI** and **More...** dropdown menus.
- A search bar containing the word **yield**.
- Species (21)**: Filter on Species...
- Data type**: Bibliography [151], None [5]
- Ontology annotation (30)**: Filter on Ontology annotation, Expand search
- Database**: Searchable fields for URL, Title, DOI, ID, Author, Subject, and Abstract.

To the right of the search interface, two detailed windows show ontology variable selection and a results table:

- Ontology variable selection**: Shows a tree view of categories like Woody Plant Ontology, Biochemical, Morphological, Other, and Budflush, with specific variables like Budflush, BF\_score\_BL, and BS\_date selected.
- Results table**: A table with columns for Identifier, Name, Description, Entity, Attribute, Class, and Status. It shows entries for Bud date protocol, Bud date, Calendar day, and Documentation links.

On the far right, three descriptive text blocks are aligned vertically:

- Full text** + **Fine criteria** + **Link back**

**EUROPEAN ELIXIR FRANCE** is visible at the bottom left.

# Data standards for FAIR

## Semantic

- Description of the data
- Controlled vocabularies: term name and definitions
- Ontologies: semantic links between Crop Ontology for agricultural data

*Biologist driven*



Persistent Unique Identifiers  
URI, gene ID, accessions ID, Trait ID, DOI,...

## Structure

- Formatting and Organizing the data
- Data Models
- Standards : CSV, VCF, GFF, MIAPPE ([www.miappe.org](http://www.miappe.org)) , etc...
- *Biologist & Computer scientist driven*



## Technical

- Data integration and sharing
- Interoperability : tools and systems
  - GA4GH
  - Breeding API [www.brapi.org](http://www.brapi.org)
- *Computer scientist driven*



# Phenotype Structure Standard

## Minimal Information About Plant Phenotyping Experiment

[www.miappe.org](http://www.miappe.org)

Many stakeholders

Elixir, Emphasis, Bioversity, North American PPI

Open Community:

Request for comments

Github Feature requests

Mailing lists

Meetings & Workgroups

Crops and woody plants



line #	MIAPPE Check list	Definition	Example	Format	Cardinality
DM-1	<b>Investigation</b>	Investigations are research projects with defined aims. They can exist at various scales (for example, they could encompass a grant-funded programme of work, the various components comprising a peer-reviewed publication, or a single experiment).			1 per MIAPPE submission
DM-2	<b>Investigation unique ID</b>	Identifier comprising the unique name of the institution/database hosting the submission of the investigation data, and the accession number of the investigation in that institution.	EBI12345678	Unique identifier	0..1
	<b>Investigation title</b>	Human-readable string summarising the investigation.	Adaptation of Maize to Temperate Climates: Mid-Density Genome-Wide Association Genetics and Diversity Patterns Reveal Key Genetic Resources with Free text (short)		1
<b>Environment</b>					
ENV-1	<b>Environment parameters</b>	Definition	Example	Format	
ENV-2					
ENV-3					
ENV-4	<b>Air temperature</b>	List of hourly air temperature throughout the experiment.	22 °C	Numeric	
ENV-5	<b>Organ temperature</b>	List of hourly organ temperatures throughout the experiment	18 °C	Numeric	
<b>Experimental Factors</b>					
TR-1	<b>Factor type</b>	Definition	Example factor values	Format	
TR-2					
TR-3	<b>Seasonal environment</b>	A plant treatment (EO:0001001) involving an exposure to a given conditions of regional seasons.	Spring season; dry season	Plant Environment Ontology 'EO_0007038'	
TR-4	<b>Air treatment regime</b>	The treatment involving an exposure to wind/air with varying degree of temperature, which may depend on the study type or the regional environment.	28/25°C ( Day/Night )	Plant Environment Ontology 'EO_0007161'	
TR-5	<b>Soil temperature regime</b>	A physical plant treatment (EO:0007316) involving an exposure to varying degree of temperature, which may depend on regional environment.	27/25°C ( Day/Night )	Plant Environment Ontology 'EO_0007161'	

Papoutsoglou *et al.* (2020) Enabling reusability and interoperability of plant phenomic datasets with MIAPPE 1.1. *New Phytol*, 227:260-273; <https://doi.org/10.1111/nph.16544>

# Phenotype Technical Standard, MIAPPE Implementations



## Ontology, OWL Implementation

<https://github.com/MIAPPE/MIAPPE-ontology>

<http://agroportal.lirmm.fr/ontologies/PPEO>

Data model representation

Formal concepts and constraints

## File Archive

ISA Tab: data + metadata

RO Crate studies

## Web Services

Breeding API

International collaboration

Standard Open Web Service API

Information Exchange, Main target: Breeding

Excellence in Breeding platform (CGIAR)



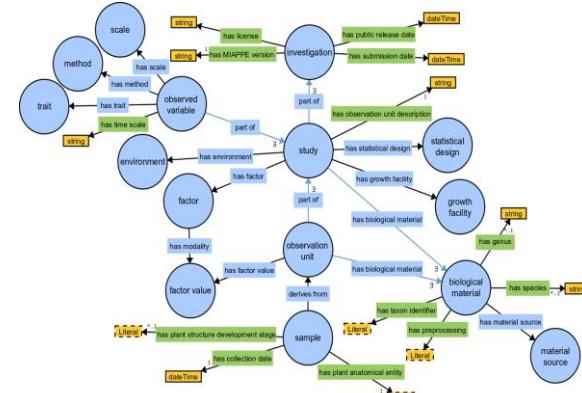
## Adoption

PhenospeX

ELIXIR

EMPHASIS

European Genbanks (EURISCO)



# Take Home Message

Research Data Reuse, including by yourself, Public or not

- → Complex
- → Relies on FAIR data principles and in particular
  - Data Management
    - Information systems
    - Data Management Plans
  - Data repositories



- Support and Guidances  
<https://rdmkit.elixir-europe.org/>

**RDMkit**

# Acknowledgments



## Elixir Plant community & platforms

Beier S., Gruden C., Pommier C., Coppens F., Scholz U.,  
Lange M., Contreras B., Adam Blondin AF., Faria D.,  
Chavez I., Miguel C., Droedsbek B., Finkers R.,  
Papoutsoglou E., Olster R., Ramsak Z., ...



## Emphasis

Tardieu F., Usadel B., Alic I., Arend D., Junker A.,  
Poorter H., Neveu P., Pierushka R., Shur U...  
And many more!



EMPHASIS

EMBL-EBI



## Breeding API

Selby P., Mueller L., Robbins K.,  
Backlund JE., ... ,  
And many more!

## Crop Ontology

Arnaud E., Laporte MA., ...



## MIAPPE community

ELIXIR Plant Community,  
Krajewsky P., Cwiek H., Tardieu F., Usadel B., Arend D., Arnaud E.,  
Junker A., King G., Laporte MA., Poorter H., Reif J., Rocca-Serra P.,  
Sansone SA., Kersey P.,  
And many more!

<https://www.phenome-emphasis.fr/>

## H2020 AGENT

N. Stein (IPK, coord), P. Kersey (RBGK), M. Alaux (INRAE), S. Weise (IPK), C. Pommier (INRAE), M. Lange (IPK), R. Finkers (WUR), J. Destin (INRAE)