# Plant biology open data interoperability in the big data era

Anne-Francoise Adam-Blondon, Sophie S. Durand, Erik Kimmel,
Raphaël-Gauthier R.-G. Flores, Cyril Pommier, Michael M. Alaux, Delphine
Steinbach, Hadi Quesneville

# Plant biology: open data interoperability in the big data era

**A-F Adam-Blondon**, S Durand, E Kimmel, R Flores, C Pommier, M Alaux, D Steinbach, H Quesneville

# A bioinformatic unit for crops and pathogens



**URGI DBs « GnpIS »**

**Phenotypes**
- GxE
- QTL maps
- GWAS, GS

**Genetics**
- Genetic maps
- Genetic markers
- Genetic resources

**Genome**
- Annotation
- Transcriptome
- SNP / Structural variant

AFAQ
ISO 9001
VERSION 2000

# URGI is a node of the french network of bio-informatics facilities (IFB-ReNaBi)

# Challenges

# Necessity to connect data stored into different information systems

- Because the volumes are becoming too big for one information system (Ex: NGS)

- Because it is impossible to store all data in a single data model (Ex: phenotyping)

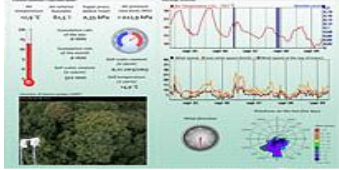- Because data relevant for a scientific question may be stored in different databases dedicated to other purposes

INRA
SCIENCE & IMPACT

# Necessity to organize and query heterogeneous data collected in different laboratories/context

Different data structures <-> different initial question
Potentially different experimental protocoles

**Climate, environment**

**Physical measurements, sensors…**



**Metabolites, proteins, genomic data…**

Development of guidelines, ontologies and standards by the **community of data producers/researchers**

**Bibliography, human sciences…**

**Post-harvest**

# Consequences on information systems

Single
database

DB
Interoperability

Information
system

Towards distributed systems

Distributed
information
system

# Work in progress

# Towards distributed information systems

**WheatIS:** the information system of the International Wheat Initiative (coord. H. Lucas): (chair: H. Quesneville)

Google-like query tool allowing to retrieve information in the databases of the transnational **TransPLANT** infrastructure (coord P. Kersey, EBI)
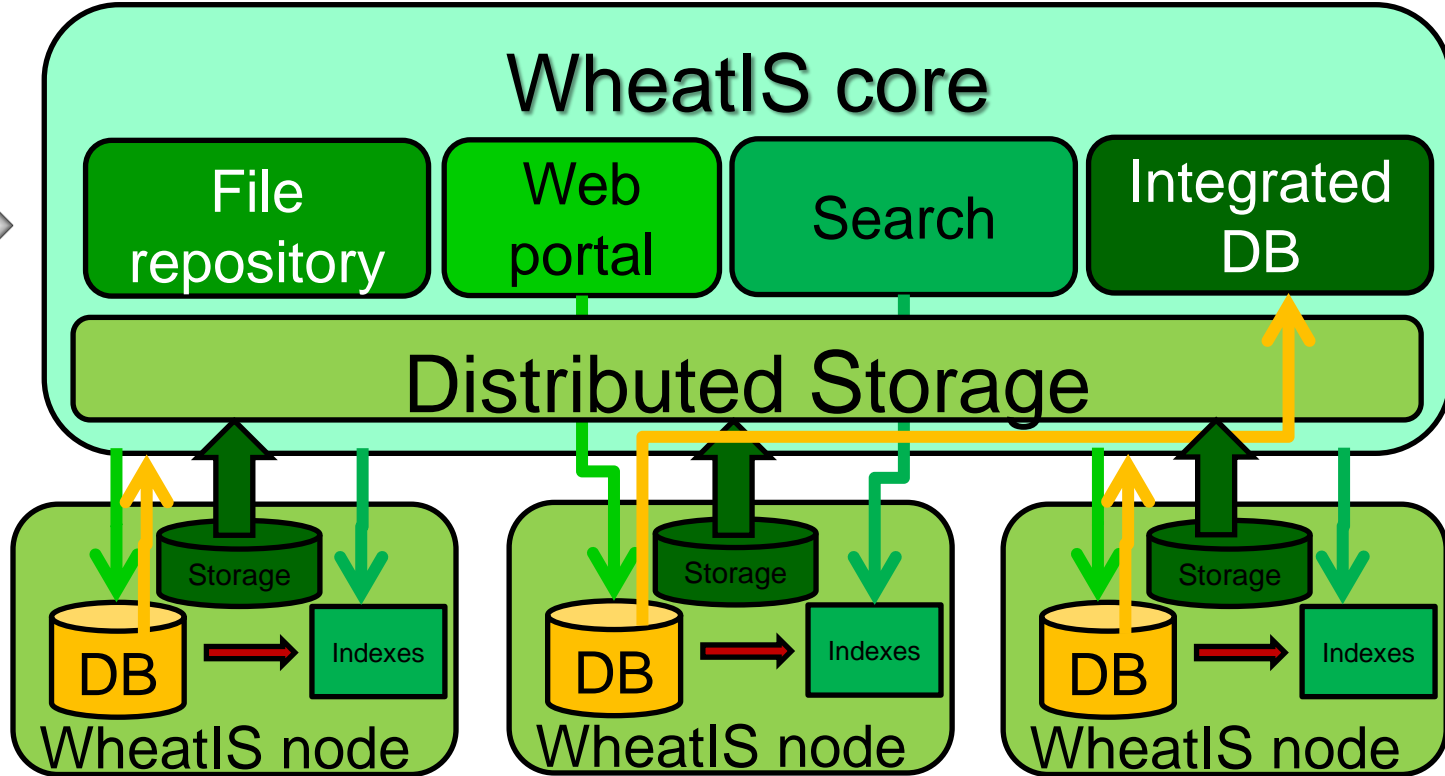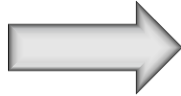
Information system for French Plant Phenotyping Network (**Phenome**, coord F. Tardieu)

Building a portal for the french crop **germplasm collections** (ARCAD-FEDER, J-L Pham coord)

# WheatIS architecture

# Definition of standards

**Survey** of existing standards:  **(1) data, (2) ontologies, (3) meta-data**

"**Cookbook**": how to produce easily shareable, reusable and interoperable "wheat data"

Identification of **end-users** categories and **WheatIS nodes**

**Challenge : adoption of the recommendations by the community**
- simple ontologies
- good balance between genericity and necessary specifity
- alignment with other international initiatives
- tools to help users

# Ontologies / Thesaurus

# Full text queries of distributed databases



GnpIS
URGI

e! EnsemblPlants
EMBL-EBI

MetaCrop
IPK GATERSLEBEN

CR-EST
IPK GATERSLEBEN

GEBIS
IPK GATERSLEBEN

HELMHOLTZ GEMEINSCHAFT
PLANTS DB

PolapgenDB
IGR pan

http://www.transplantdb.eu/

User web interface

Google like query

Google like list of results

Lucene indexes

Apache Solr
Lucene
elasticsearch.

Wheat Initiative

EMBL-EBI

# Perspectives

# Develop a web semantic interoperability between the plant databases of the French Elixir node



**Ontology based annotation of database schemes**

**RDF triple store modeling of the databases schemes integrated with existing ontologies**

**Semantic web services**

# Summary

Challenge: the query of high volumes of heterogeneous and distributed data

$\Rightarrow$ Federation of information systems

- ❖ At the national and european level
- ❖ through noSQL technologies (SolR, ElasticSearch…)
- ❖ Web semantic layer

# Aknowledgements

**URGI team**



**arcad**

Jean-Louis Pham
Christophe Jenny
Felix Homa

**EPPN** European Plant Phenotyping Network

**PHENOME** Réseau Français Phénomique végétale F P P N

**Wheat Initiative**

Hélène Lucas
WheatIS Expert Working Group
   Mario Caccamo
   Dave Edwards
   Gerard Lazo…

P. Bento
L. Couderc
C. Gageat
A. Keliet
T. Letellier

M. Loaec
A. Ménard
C. Michotey
N. Mohellibi
C. Viseux
S. Meilo

François Tardieu
Pascal Neveu
Jacques LeGouis
Eric Duchêne

**RDA** RESEARCH DATA ALLIANCE

Esther Dzalé Kaboré
Sophie Aubin

**transPLANT**

Paul Kersey
Dan Bolser…

**INRA** SCIENCE & IMPACT

# WheatIS timeline

**Step 1: Network building**

## Definition of standards

- Define standards, nomenclature, formats.
- Meta-data exchange

WheatIS
=
A web platform to exchange data

**Step 2: Integrated portal**

## Search of data

- DBs federation
- Google-like search

WheatIS
=
A portal to access a network of DBs

**Step 3: Integrated DB**

## Integration of data

- In one place
- Focused on relevant data sets
- Consolidated and consistent

WheatIS
=
A integrated DB