

Exposing the French agronomic resources as Linked Data

Pierre Larmande

Institute of Research for Development (IRD)
Head of data integration group at the Institute of
Computational Biology
pierre.larmande@ird.fr

Institut Français de Bioinformatique (PIA)

PIA



Institut Français de Bioinformatique (PIA)

Mission générale : fournir des **ressources** de base en bioinformatique à la communauté des sciences de la vie

Directeur : J-F Gibrat, DR INRA

Responsable cellule e-infrastructure : C. Blanchet, IR CNRS

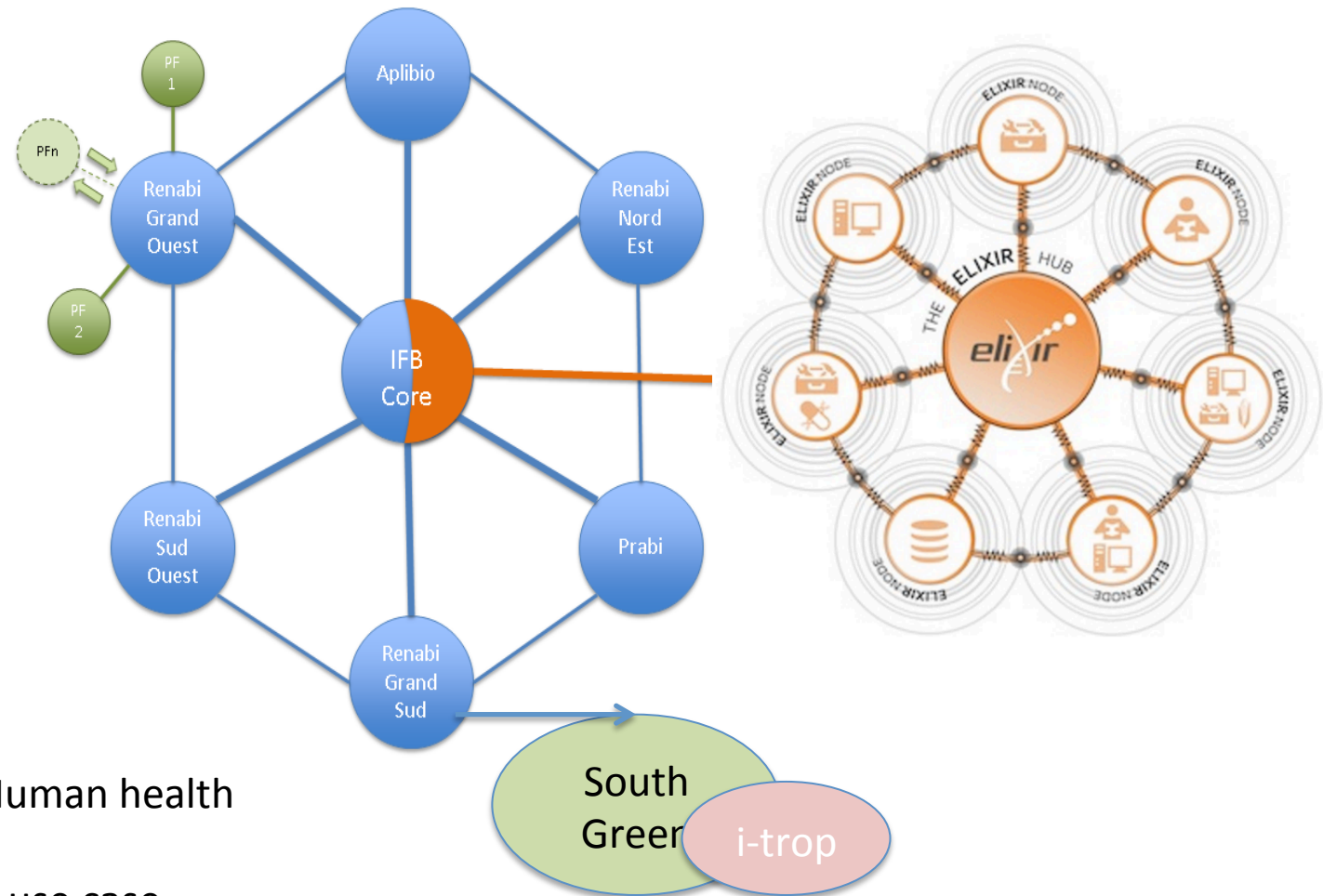
Infrastructure nationale de service en bioinformatique

- **Données** : Fournir un accès à des collections de données spécialisées à haute valeur ajoutée issues de l'expertise du laboratoire d'accueil
- **Outils** : Développer et mettre à disposition des outils et services en lignes pour analyser les données correspondant à l'expertise scientifique du laboratoire d'accueil
- **Appui** aux projets scientifiques et hébergement sur une infrastructure informatique
- **Infrastructure** : Mettre à disposition une infrastructure informatique dédiée à l'analyse des données des sciences du vivant (matériel, données, outils)
- **Formations**

cf. <http://france-bioinformatique.fr>

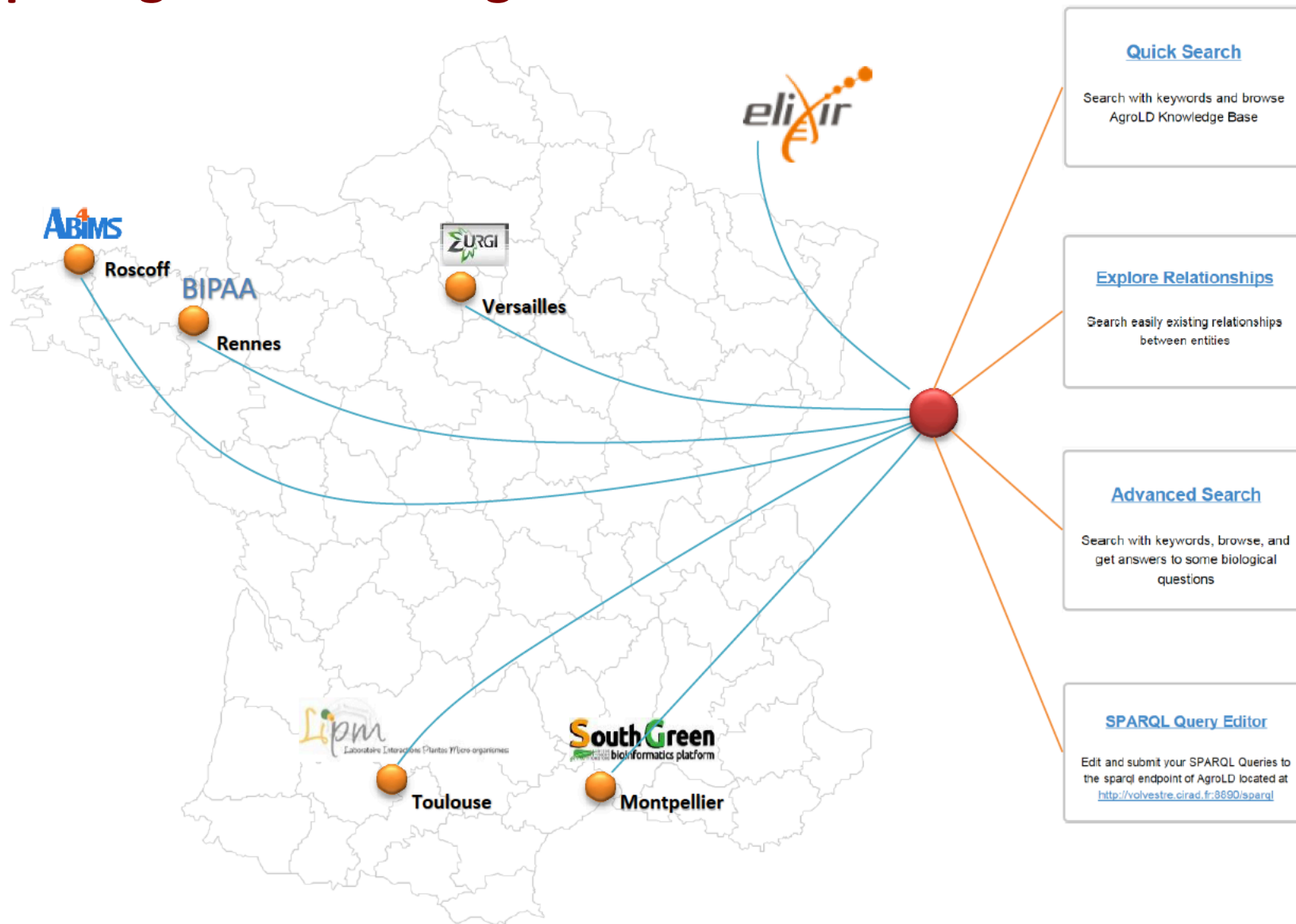
Elixir

ELIXIR has been awarded €19 million from the EU to accelerate the implementation of Europe's life-science data infrastructure over the next four years.



- Rare disease
- Genomics for Human health
- Plant use case
- Marine biology use case

Exposing the French agronomic resources as Linked Data





The Plant bioinformatics node

(33 FTEs)



Genetics and genomics resources for
plants and crop parasites

(INRA)



Genomic resource for southern and
mediterranean plants.

(CIRAD, INRA, IRD)

Contributors (data, tools, and expertise)



Resources for plants, symbionts and
pathogens

(INRA, CNRS)



Marine biology analysis


(CNRS, UPMC)



Arthropods for Agroecosystems

(INRA)

Services



analysis

genomics annotation

software hosting

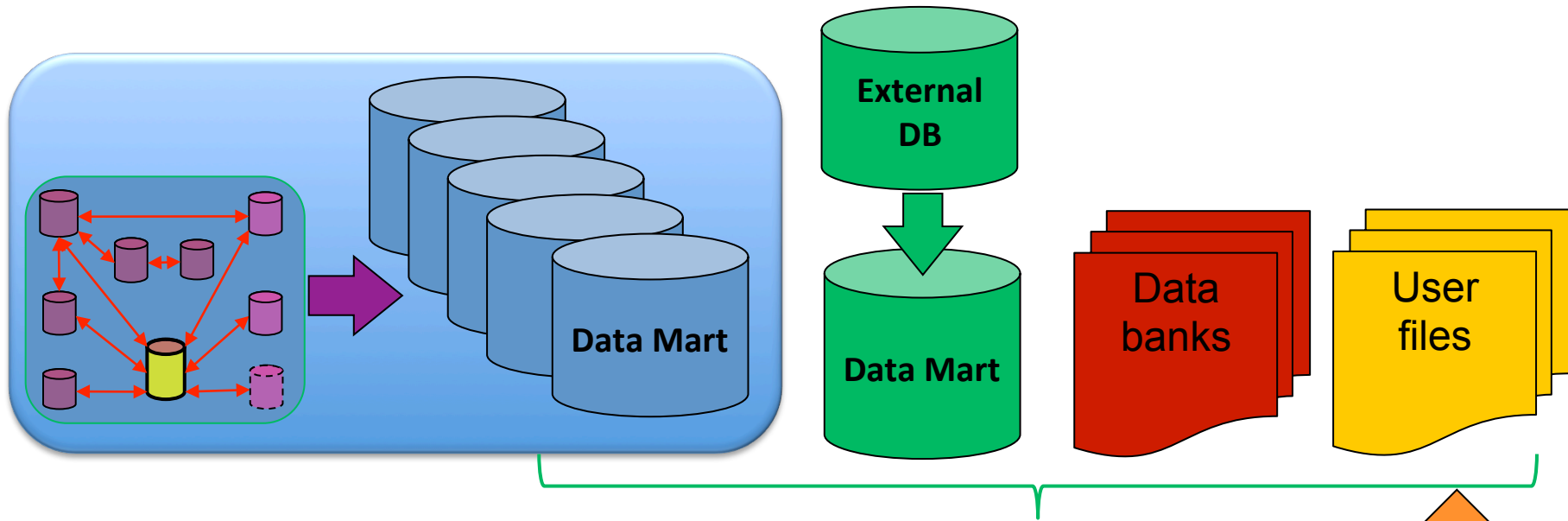
data repository

data integration

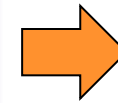
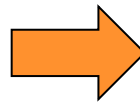
database design

software engineering

Data Integration / data mining Workbench



Developer



User

Resources

urgi.versailles.inra.fr

www.southgreen.fr

The screenshot shows the URGi website interface. At the top, there's a header with the URGi logo and the text 'PLANT AND FUNGI DATA INTEGRATION'. Below this is a navigation menu with categories like Platform, Research, Projects, Data, Tools, and Species. A search bar is also present. The main content area features a 'WHAT'S NEW?' section with news items dated 13 Mar 2015 and 12 Feb 2015. There are also sections for 'EVENT & PUBLICATIONS' and 'RESEARCH'.

The screenshot shows the South Green bioinformatics platform website. The header features the 'South Green bioinformatics platform' logo. Below the header is a navigation menu with options like Home, Tools, Databases, Developer tools, Species, Projects, Teaching, and Platform. The main content area displays a grid of logos for various bioinformatics tools and databases, including Galaxy, The Banana Genome Hub, Coffee Genome Hub, GreenPhyl, SNIPlay, TropGENE, DryGenesDB, Oryza Tag Line, EURGEN, GNPannot, EST, HaploPhylo, and GInDIVERSITY.



The screenshot shows the GnpIS website interface. The header reads 'GnpIS GENETIC AND GENOMIC INFORMATION SYSTEM'. Below the header is a search bar and a navigation menu. The main content area features a central graphic of a green flower with a white center, surrounded by various data categories like Genomes, Taxons, Sequences, Genetic maps, Polymorphisms, Phenotypes, Association, Genetic resources, and Transcriptomic. There are also 'ADVANCED TOOLS' like BIOMART, GALAXY, and INTERMine.



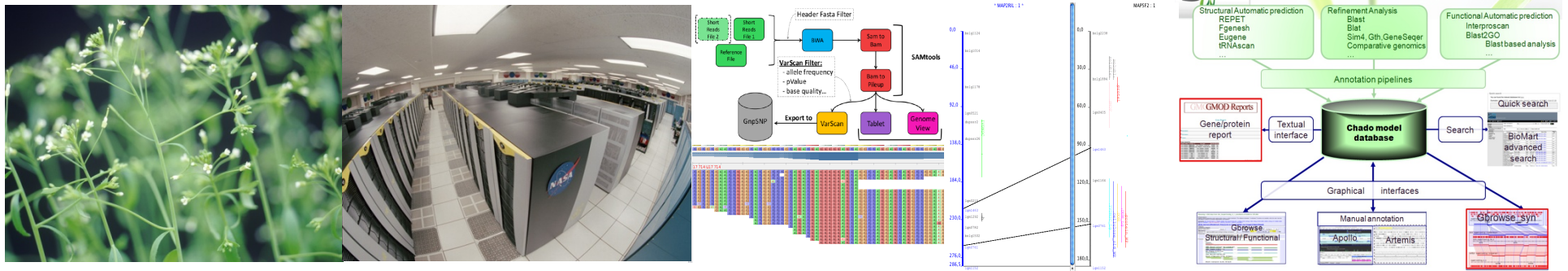
The screenshot shows the TropGENE Database website. The header includes logos for CIRAD, TropGENE Database, and South Green bioinformatics platform. Below the header is a navigation menu with options like Presentation, Species, Updates, Staff, and Links. The main content area describes the TropGENE database as a resource for managing genomic, genetic, and phenotypic information about tropical crops. It lists several modules currently online, including BANANA, COCOA, COCONUT, COFFEE, COTTON, OIL PALM, RICE, RUBBER TREE, and SUGARCANE.

From H. Quesneville

IFB technical project

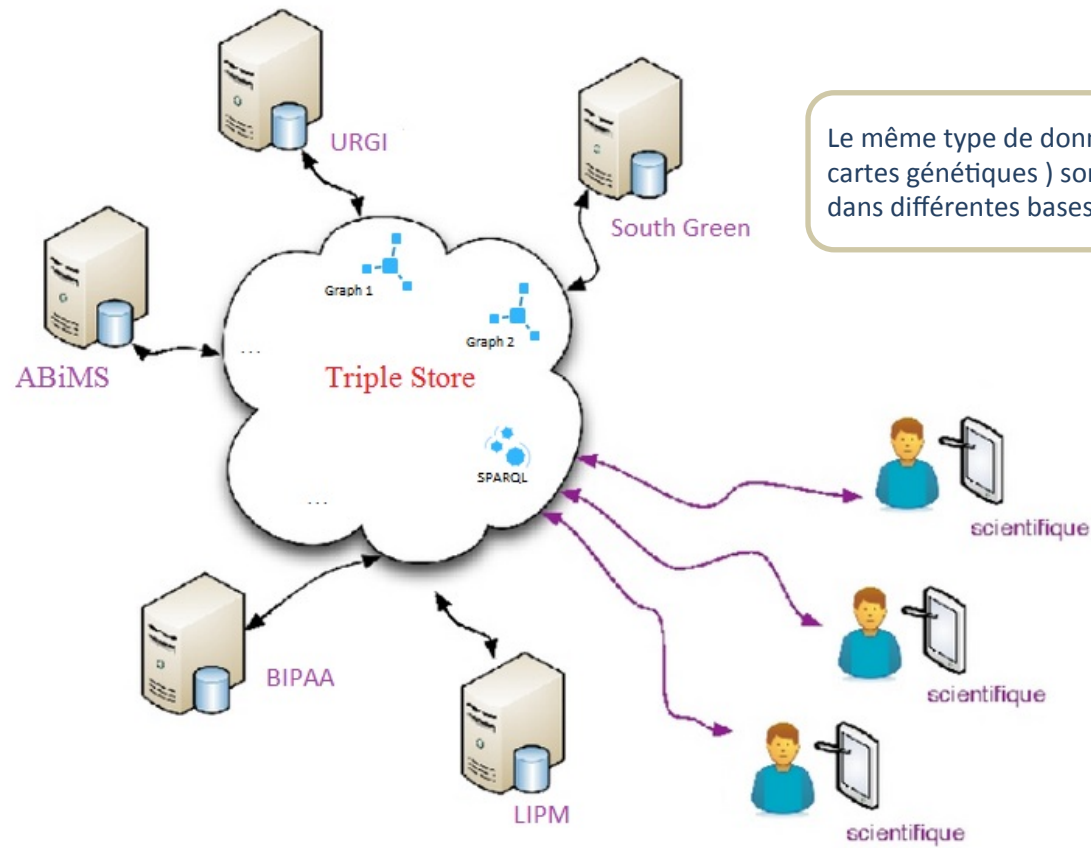
- WP1: Develop a RDF-based semantic interoperability between the Plant bioinformatics node databases
- WP2: Develop an intuitive Google-like search portal
- WP3: Setup the “IT-core hub” workflows under the Galaxy

WP1 : RDF "store" permettant d'interconnecter sémantiquement l'ensemble des bases de données végétales des plateformes bioinformatiques IFB



IFB plant node: (URGI), South Green, IFB Grand Ouest (ABiMS), GenOuest (BIPAA) et LIPM.

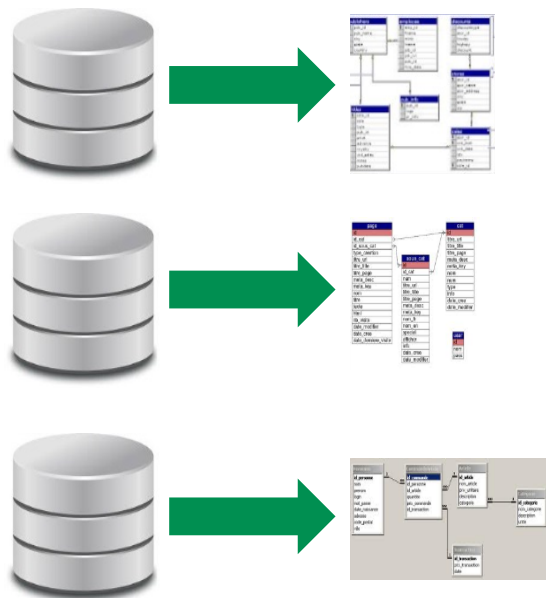
Develop a web semantic interoperability



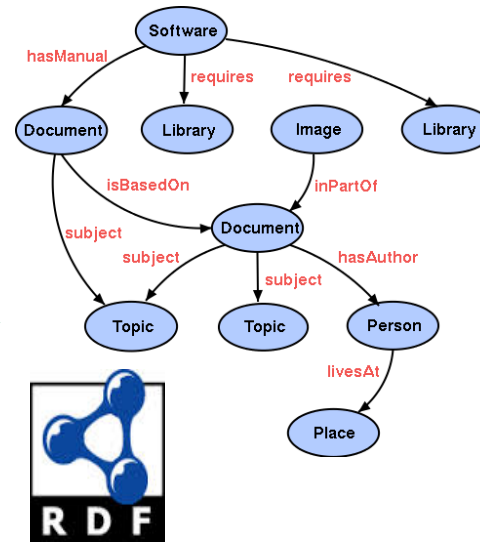
Le même type de données (par exemple SNP , collections génétiques , cartes génétiques) sont souvent représentées différemment dans différentes bases de données.

Développement d'une application Web interrogeant les différentes sources de données.

Develop a web semantic interoperability



Ontology based annotation of database schemas (GO, PO, TO, CO, ...)



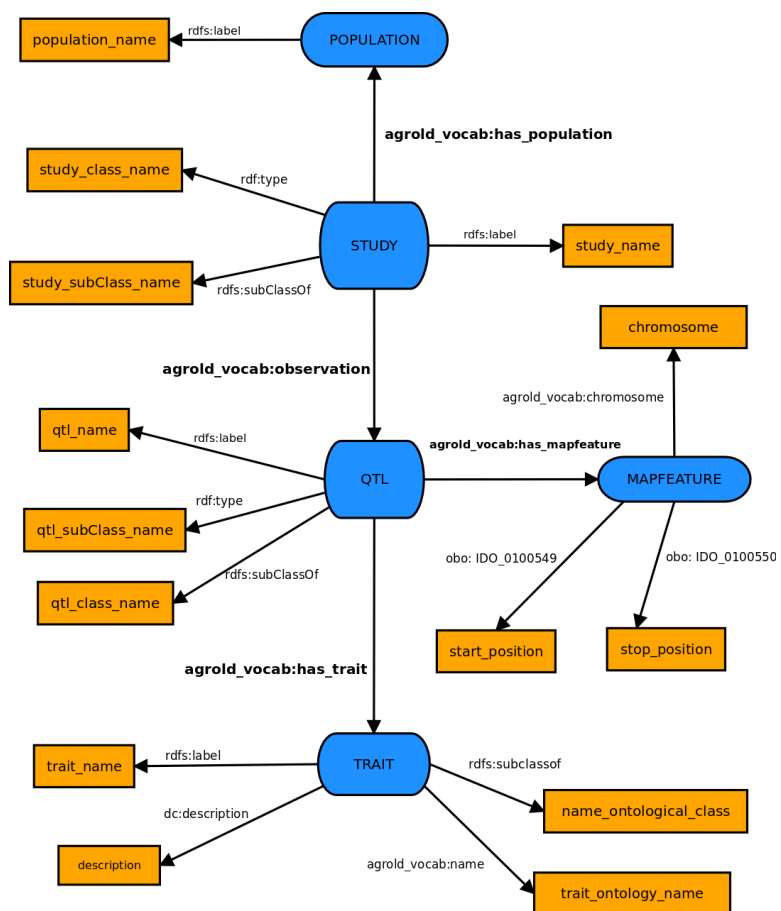
RDF modeling of the databases schemas

RDF triple store storage



Query of dispersed data for data integration through Semantic web services (SPARQL, web user interface)

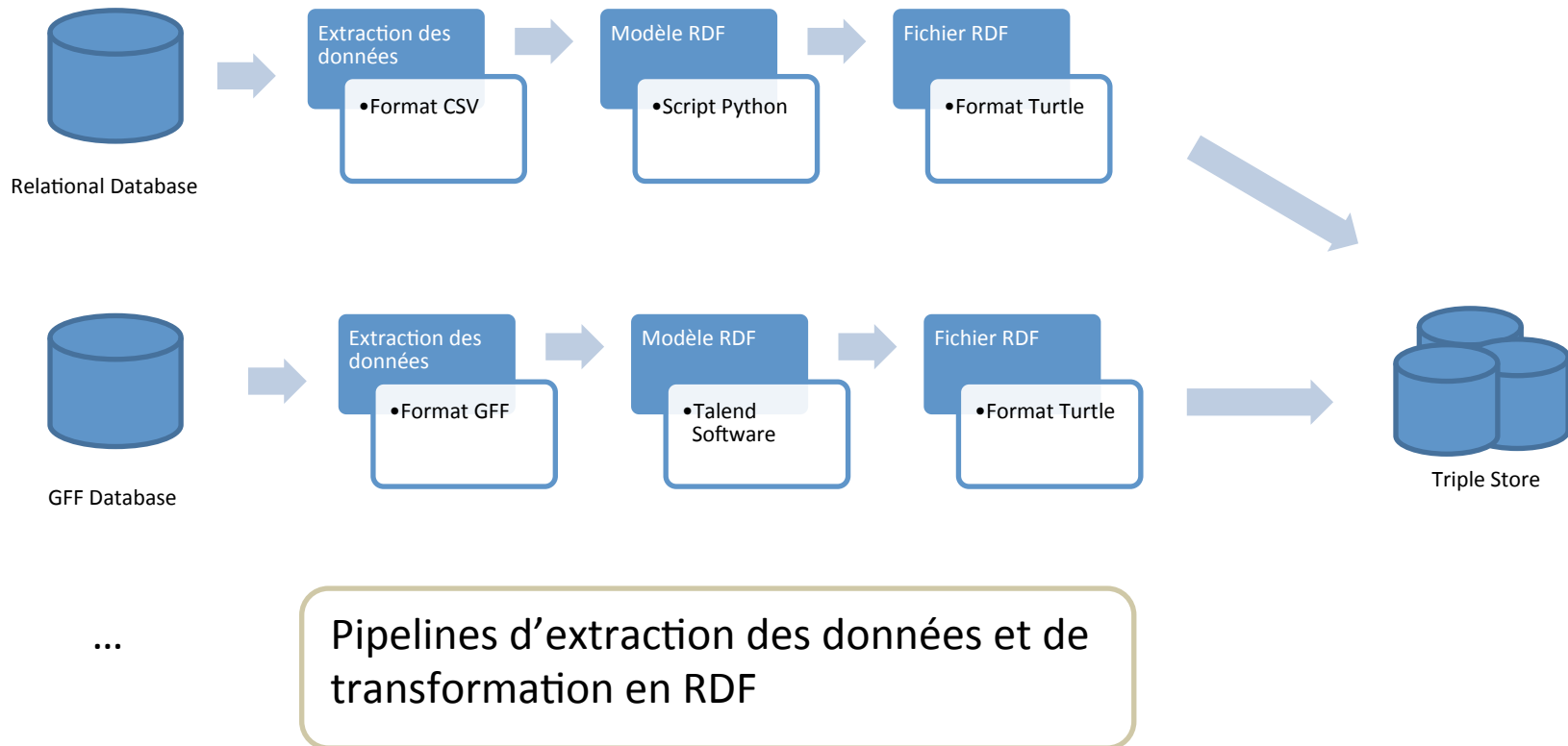
Le Modèle RDF présente l'organisation générale des données.



Un Graphe nommé est un concept clé de l'architecture du Web sémantique dans lequel un ensemble de **ressources déclarées** sont **identifiées en utilisant une URI**.

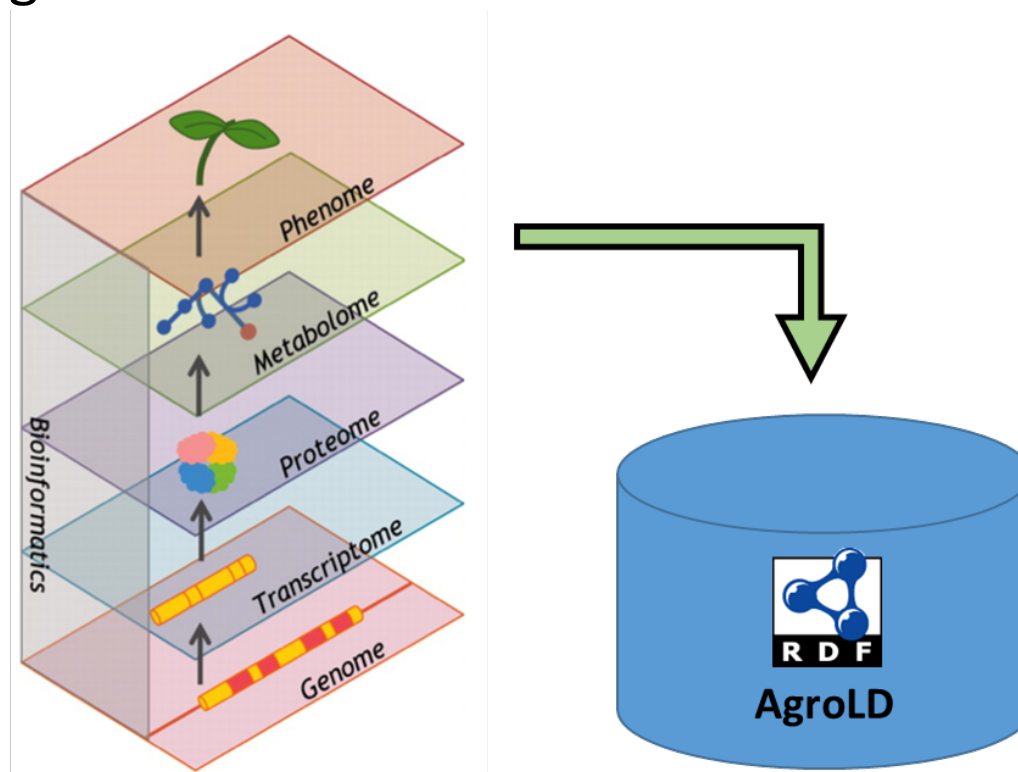
Exemple d'URI:
<http://www.southgreen.fr/tropgene.trait/35>

Démarche



Agronomic Linked Data (AgroLD)

- RDF knowledge base that integrates data from a variety of plant resources.
- Integrate information at different levels.



AgroLD

(www.agrold.org)

– AgroLD is developed in phases:

– Phase I: includes information on:

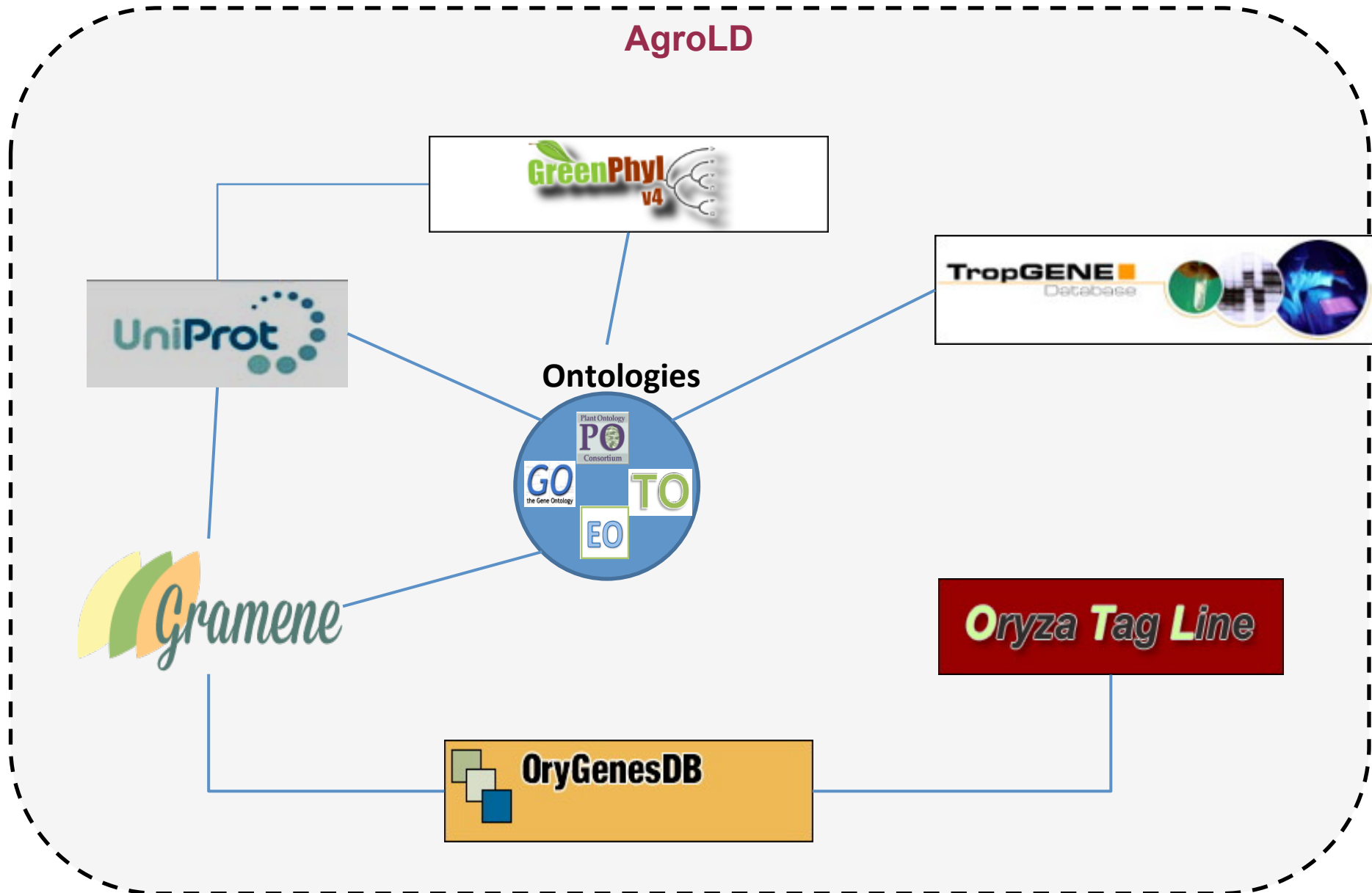
- **Arabidopsis thaliana**
- **Wheat (Triticum spp.)**
 - Triticum aestivum
 - Triticum urartu
- **Rice (Oryza spp.)**
 - Oryza barthi
 - Oryza brachyantha
 - Oryza Sativa
 - Oryza glaberimma
- **Sorghum (Sorghum bicolor)**
- **Maize/Corn (Zea mays)**



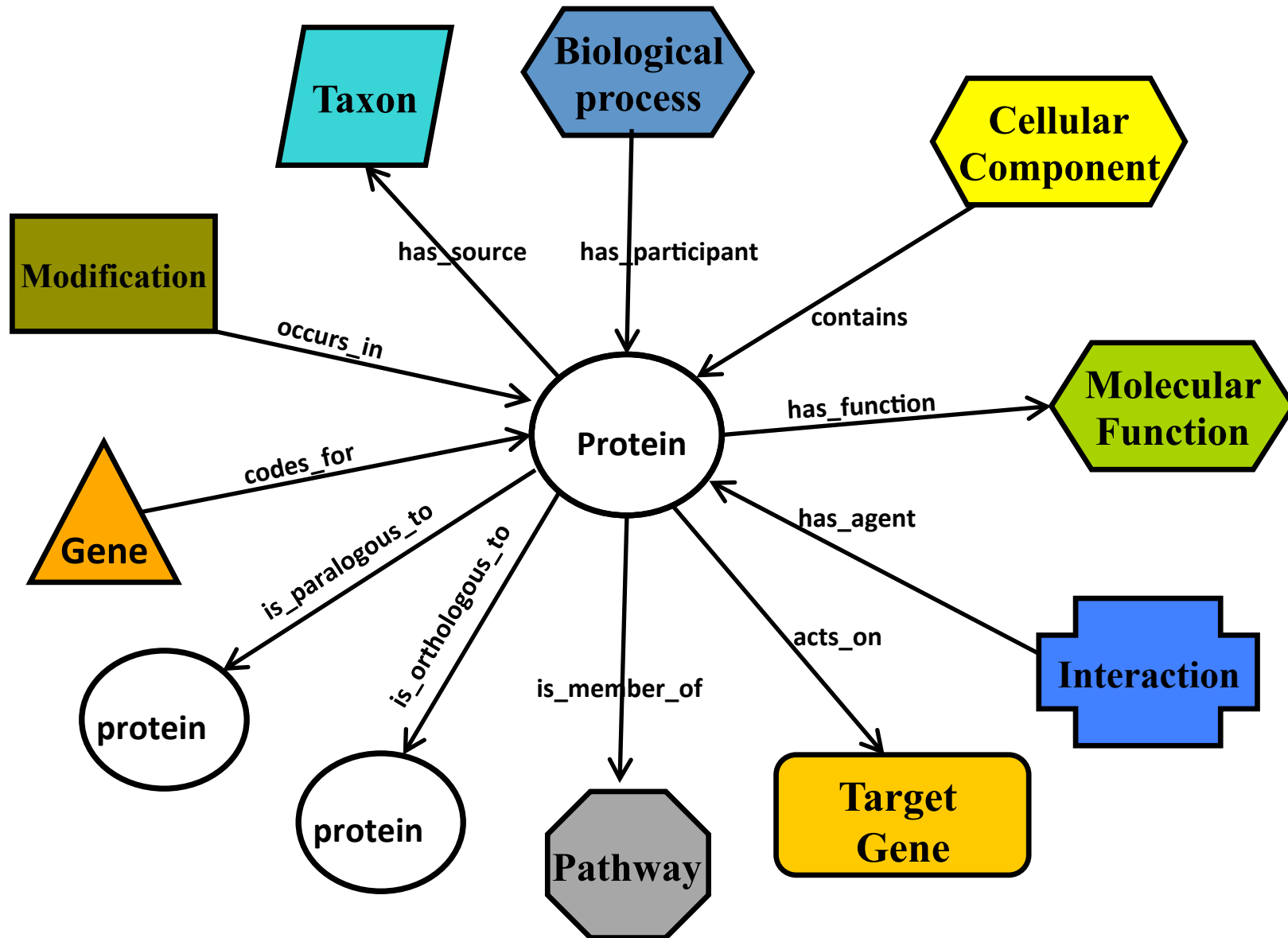
Information in AgroLD

- Integrates information from:
 - **Ontologies:** Gene Ontology (GO), Sequence Ontology (SO), Plant Ontology (PO), Plant Trait Ontology (TO), Plant Environment Ontology (EO), NCBI Taxonomy
 - **Information sources:**
 - **Ontology association:** GOA, Gramene (TO, PO and EO)
 - **Gene/Protein information:** OryGenesDB, Gramene, UniPort
 - **QTL information:** TropGeneDB, Gramene
 - **Pathway information:** Gramene - Cyc
 - **Phenotype information:** Oryza Tag Line
 - **Homology prediction:** GreenPhylDB

Knowledge in AgroLD



Knowledge representation in AgroLD



Agronomic Linked Data (AgroLD)

[Home](#)[Search](#)[Documentation](#)[About](#)[Please send us your feedback!](#)

The Agronomic Linked Data (AgroLD) Project

At the Institute of Computational Biology (IBC), we are involved in developing methods to aid data integration and knowledge management within the plant biology domain to improve information accessibility of heterogeneous data. Among others, a solution for the data integration challenges is offered by the Semantic Web technologies. The semantic web has emerged as one of the most promising solutions for high scale integration of distributed resources. This is made possible by a stack of technologies such as the Resource Description Framework (RDF), RDF Schema (RDFS), Web Ontology Language (OWL) and the SPARQL Query Language (SPARQL) proposed by the World Wide Web Consortium (W3C). RDF forms the basis of the stack allows modeling information as a directed graph composed of triples that can be queried using SPARQL.

AgroLD is a RDF knowledge base that consists of data integrated from a variety of plant resources and ontologies. The aim of the Agronomic Linked Data (AgroLD) project is to provide a portal for bioinformatics and domain experts to exploit the homogenized data models towards efficiently generating research hypotheses.

[Quick Search](#)

Search with keywords and browse
AgroLD Knowledge Base

[Advanced Search](#)

Search with keywords, browse, and
get answers to some biological
questions

[Explore Relationships](#)

Search easily existing relationships
between entities

[SPARQL Query Editor](#)

Edit and submit your SPARQL Queries to
the sparql endpoint of AgroLD located at
<http://volvestre.cirad.fr:8890/sparql>



© AgroLD 2015

www.agrold.org



Search and browse AgroLD

Search examples: ontological concepts - 'plant height' or 'regulation of gene expression'; gene names - 'GRP2' or 'TCP12'.

Search Text

Entity	Title	Named Graph	
http://www.identifie...gramene.qtl/AQFS209	PHTT	http://www.southgree...old/qtl.annotations	plant height.
http://www.identifie...gramene.qtl/AQBK027	BNL6.32	http://www.southgree...old/qtl.annotations	plant height.
http://www.identifie...gramene.qtl/AQCN007	BNL6.32	http://www.southgree...old/qtl.annotations	plant height.
http://www.southgree...notype/ALGB06_otl_6		http://www.southgreen.fr/agrold/otl	... dark green leaves; wide flag leaves; compact plant 1 plant.
http://www.southgree...notype/AEVD11_otl_2		http://www.southgreen.fr/agrold/otl	Decreased height 30 ; late flowering; compact plant ; erect leaves 1 plant.
http://www.identifie...gramene.qtl/AQEA058	BNL6.32	http://www.southgree...old/qtl.annotations	plant height.
http://www.identifie...gramene.qtl/AQEA071	BNL6.32	http://www.southgree...old/qtl.annotations	plant height.
http://www.identifie...ramene.qtl/AQFS1461	PHTT	http://www.southgree...old/qtl.annotations	plant height.
http://www.identifie...gramene.qtl/AQFS183	PHTT	http://www.southgree...old/qtl.annotations	plant height.

OpenLinks Faceted search

String matching on literal + ranking based on occurrence number

Sparql query editor

```

1 BASE <http://www.southgreen.fr/agroid/>
2 PREFIX rdf:<http://www.w3.org/1999/02/22-rdf-syntax-ns#>
3 PREFIX rdfs:<http://www.w3.org/2000/01/rdf-schema#>
4 PREFIX obo:<http://purl.obolibrary.org/obo/>
5 PREFIX uniprot:<http://purl.uniprot.org/uniprot/>
6 PREFIX vocab:<vocabulary/>
7 PREFIX graph:<gramene.cyc>
8 PREFIX pathway:<biocyc.pathway/CALVIN-PWY>
9
10 SELECT DISTINCT ?gene ?name ?taxon_name
11 WHERE {
12 GRAPH graph: {
13 ?gene vocab:is_agent_in pathway:.
14 ?gene rdfs:label ?name.
15 ?gene vocab:taxon ?taxon_name.
16 }

```

Execution timeout: 20000 milliseconds (values less than 1000 are ignored) Results Format: RDF/XML Download Results

Filename to Save As: query.sparql Save Query Choose File No file chosen Load Selected Query File

2. Search terms by label ([select](#))
3. List relation types in a given graph ([select](#))
4. Retrieve the local neighbourhood of Oryza sativa japonica protein: **IAA16** - Auxin-responsive protein (UniProt accession:POC127) ([select](#))
5. Identify Wheat proteins that are involved in root development. ([select](#))
6. Retrieve genes that participate in a given pathway: **Calvin cycle** ([select](#))
7. Retrieve Proteins associated with a given QTL: **DTHD** (days to heading) ([select](#))
8. Get the ID corresponding to the ontology term "**homoaconitate hydratase activity**" ([select](#))
9. Get the name of the ontological element that has the ID "**GO:0003824**" ([select](#))
10. Get the level **4** ancestor of **GO:0004409** ([select](#))
11. Get the level **2** descendance of **GO:0003824** ([select](#))
12. Get protein ids associated with the ontological id **GO:0003824** ([select](#))
13. Get QTL ids associated with the ontological id **EO:0007403** ([select](#))
14. Describe **uniprot:POC127** ([select](#))

Results

Raw Response Table Pivot Table

Search: Show 50 entries

gene	name	taxon_name
http://identifiers.org/ensembl.plant/AT1G18270	fructose-bisphosphate aldolase	obo:NCBITaxon_3702
http://identifiers.org/ensembl.plant/AT1G42970	glyceraldehyde-3-phosphate dehydrogenase	obo:NCBITaxon_3702
http://identifiers.org/ensembl.plant/AT1G43670	fructose-1,6-bisphosphatase	obo:NCBITaxon_3702

EnsemblPlants | BLAST | BioMart | Tools | Downloads | Documentation | Website help | Login/Register

Arabidopsis thaliana (TAIR10) | Location: 1:6,283,412-6,293,871 | Gene: AT1G18270

Search Ensembl Plants...

Gene-based displays

- Summary
- Splice variants
- Transcript comparison
- Supporting evidence
- Gene alleles
- Sequence
- Secondary Structure
- Gene families
- External references
- Regulation
- Literature
- Ontology
- GO: biological process
- GO: molecular function
- PO: plant structure development

Gene: AT1G18270

Description: ketose-bisphosphate aldolase class-II family protein [Source:TAIR;Acc:AT1G18270]

Location: [Chromosome 1: 6,283,412-6,293,871](#) reverse strand.

About this gene: This gene has 3 transcripts ([splice variants](#)), [37 orthologues](#) and [6 paralogues](#).

Transcripts: [Show transcript table](#)

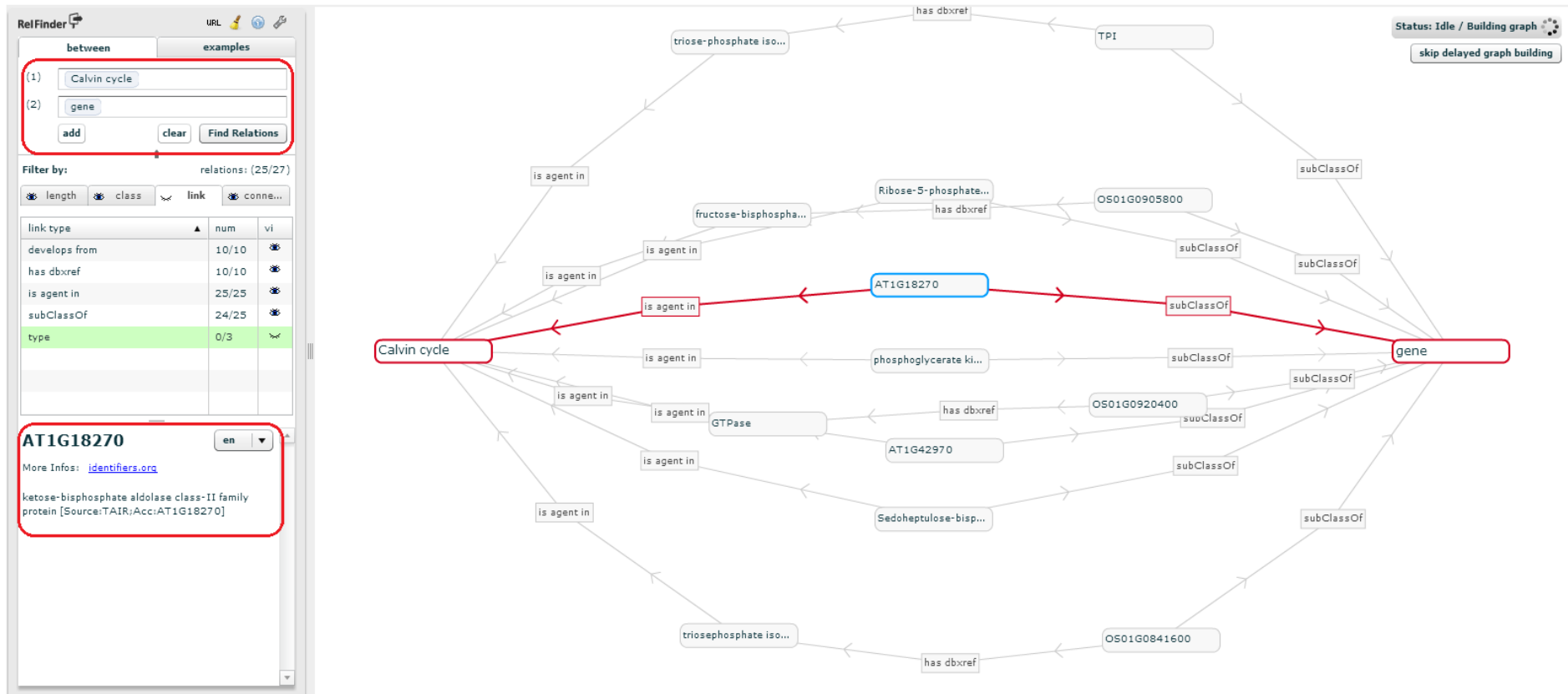
Summary

Gene type: Protein coding

Annotation Method: Gene annotation by [TAIR](#) through a process of automatic and manual curation.

Visualisation of queries

Search > Explore



Based on RelFinder [Heim et al, 2009]

Web Services API

AgroLD API 1.0 Interactive Documentation

This page provides information about the RESTful API (version 1.0) of AgroLD.

More documentations

Find out more about Swagger

<http://swagger.io>

[Contact the developer](#)

[IBC Montpellier](#)

gene : Services over genes

Show/Hide | List Operations | Expand Operations

POST	/genes{format}	Retrieve complete URI and description of all genes from AgroLD in JSON format
POST	/genes/publications/byId	Get publications of a gene
POST	/genes/byKeyword{format}	Retrieve genes with the URI or the name or the description containing the given keyword
POST	/genes/encodingProtein{format}	Get URIs, ids, and name of genes encoding a protein given its ID
POST	/genes/participatingInPathway{format}	Get URIs, ids, and name of genes participating in a pathway given its ID

graphs : General Services

Show/Hide | List Operations | Expand Operations

ontologies : Services over ontologies used in AgroLD

Show/Hide | List Operations | Expand Operations

pathway : Services over pathways

Show/Hide | List Operations | Expand Operations

protein : Services over proteins

Show/Hide | List Operations | Expand Operations

qtl : Services over QTLs

Show/Hide | List Operations | Expand Operations

Advanced form-based search

Search examples: ontological concepts - 'plant height' or 'regulation of gene expression'; gene names - 'GRP2' or 'TCP2'.

QTL ID: 'AQAA003' ; protein name: 'TBP1'

✓ --- Select a type* ---

- Gene
- Protein
- QTL
- Pathway
- Ontology

Search

Search **protein** with keyword " **TBP1** "



Search: Show 30 entries

Id	Name	Description	URI
1 (display)	Q9LL45 TBP1	Telomere-binding protein 1	http://purl.uniprot.org/uniprot/Q9LL45 (in Sparql)
2 (display)	Q9LL45 TBP1_ORYSJ	Telomere-binding protein 1	http://purl.uniprot.org/uniprot/Q9LL45 (in Sparql)
3 (display)	AT3G13445.1 "AT3G13445.1"^^xsd:string	" Symbols: TBP1, TFIID-1 TATA binding protein 1 chr3:4380317-4381869 FORWARD LENGTH=200"^^xsd:string	http://www.southgreen.fr/agrold/greenphyl.sequence/AT3G13445.1 (in Sparql)
4 (display)	P46465 TBP1	26S protease regulatory subunit 6A homolog	http://purl.uniprot.org/uniprot/P46465 (in Sparql)

Showing 1 to 4 of 4 entries

Results are combined with external services

PROTEIN : Q9LL45 / TBP1

Telomere-binding protein 1

URI: <http://purl.uniprot.org/uniprot/Q9LL45>

is encoded by



Search: Show entries

Id	Name	Description	URI
1 OS02G0817800 (display)	TBP1	Telomere-binding protein 1 [Source:UniProtKB/Swiss-Prot;Acc:Q9LL45]	http://identifiers.org/ensembl.plant/OS02G0817800 (in Sparql)

Showing 1 to 1 of 1 entries

QTL associations \pm

Ontology associations \pm

Publication

1. Yu EY, Kim SE, Kim JH, Ko JH, Cho MH, Chung IK., " **Sequence-specific DNA recognition by the Myb-like domain of plant telomeric protein RTBP1.** ", *J Biol Chem*, 2000
More at: <http://www.ncbi.nlm.nih.gov/pubmed/10811811>
2. International Rice Genome Sequencing Project., " **The map-based sequence of the rice genome.** ", *Nature*, 2005
More at: <http://www.ncbi.nlm.nih.gov/pubmed/16100779>
3. ice Annotation Project, Tanaka T, Antonio BA, Kikuchi S, Matsumoto T, Nagamura Y, Numa H, Sakai H, ..., " **The Rice Annotation Project Database (RAP-DB): 2008 update.** ", *Nucleic Acids Res*, 2008
More at: <http://www.ncbi.nlm.nih.gov/pubmed/18089549>
4. Hong JP, Byun MY, Koo DH, An K, Bang JW, Chung IK, An G, Kim WT., " **Suppression of RICE TELOMERE BINDING PROTEIN 1 results in severe and gradual developmental defects accompanied by genome instability in rice.** ", *Plant Cell*, 2007
More at: <http://www.ncbi.nlm.nih.gov/pubmed/17586654>
5. Ko S, Yu EY, Shin J, Yoo HH, Tanaka T, Kim WT, Cho HS, Lee W, Chung IK., " **Solution structure of the DNA binding domain of rice telomere binding protein RTBP1.** ", *Biochemistry*, 2009
More at: <http://www.ncbi.nlm.nih.gov/pubmed/19152316>



Galaxy Wrapper available for AgroLD

Workflow Canvas | test workflow

```
graph LR; A[Input dataset] --> B[Cut]; B --> C[Line/Word/Character count];
```

testWorkflowResult
3 shown
138.3 KB

3: Line/Word/Character count on data 2

1 line
format: **tabular**, database: ?

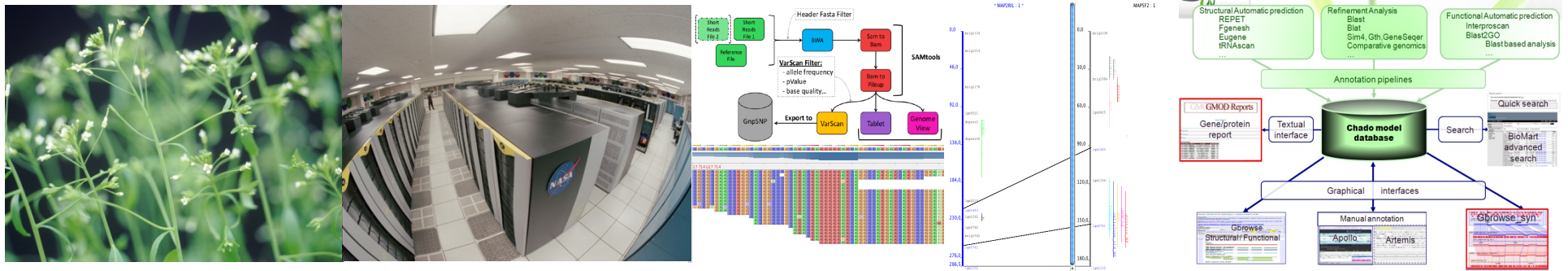
1
548

2: testResult.tsv

548 lines
format: **tabular**, database: ?

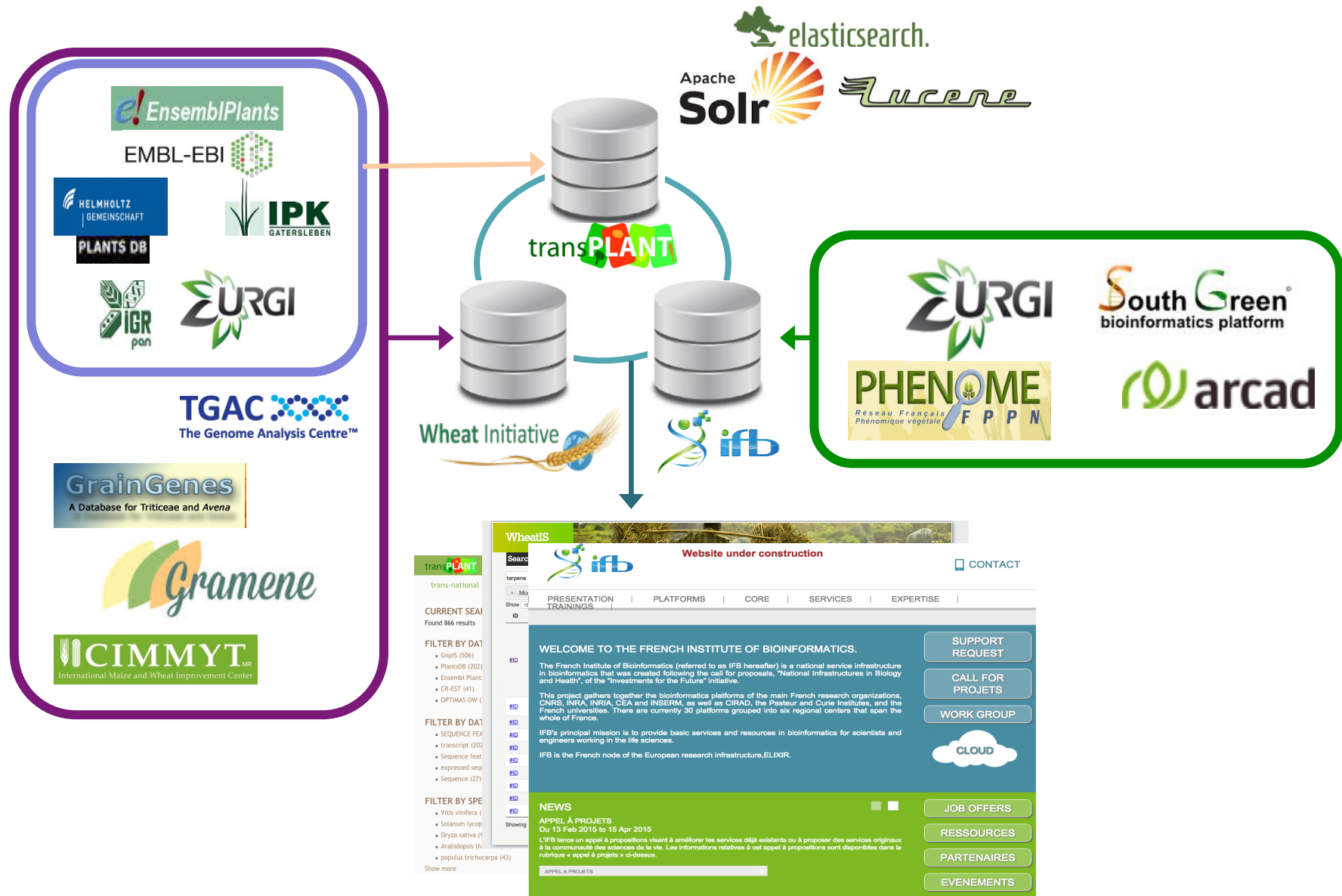
1	2	3
"geneId"	"gene_name"	"taxon_name"
"AT1G18270"	"fructose-bisphosphate aldolase"	"Arabidopsis thaliana"
"AT1G42970"	"glyceraldehyde-3-phosphate dehydrogenase"	"Arabidopsis thaliana"
"AT1G43670"	"fructose-1,6-bisphosphatase"	"Arabidopsis thaliana"

WP2: Develop an intuitive Google-like search portal

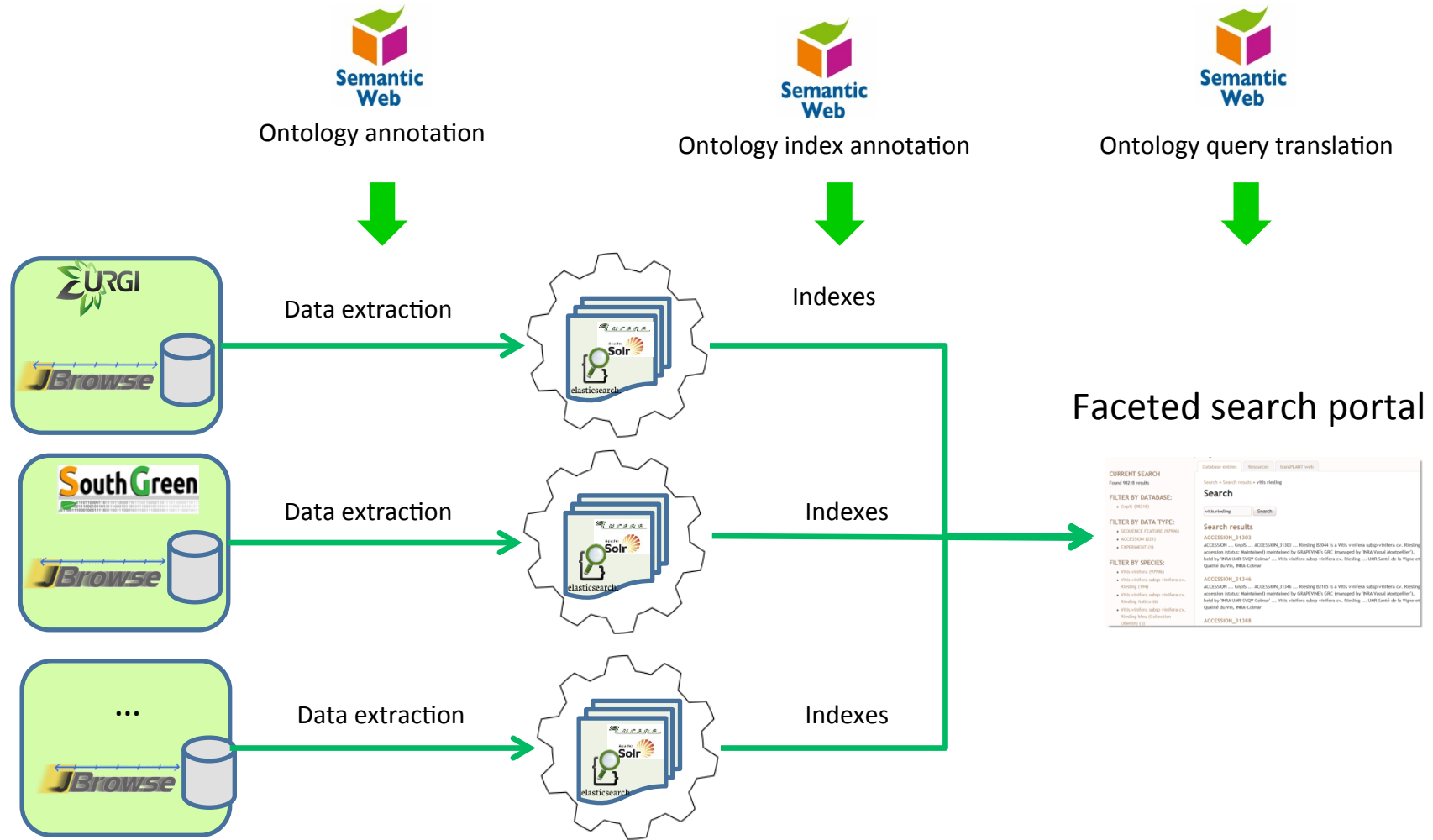


- *Use the very fast distributed indexation system developed in international projects: WheatIS, transplant*
- *Inject Semantic web capabilities in this system*

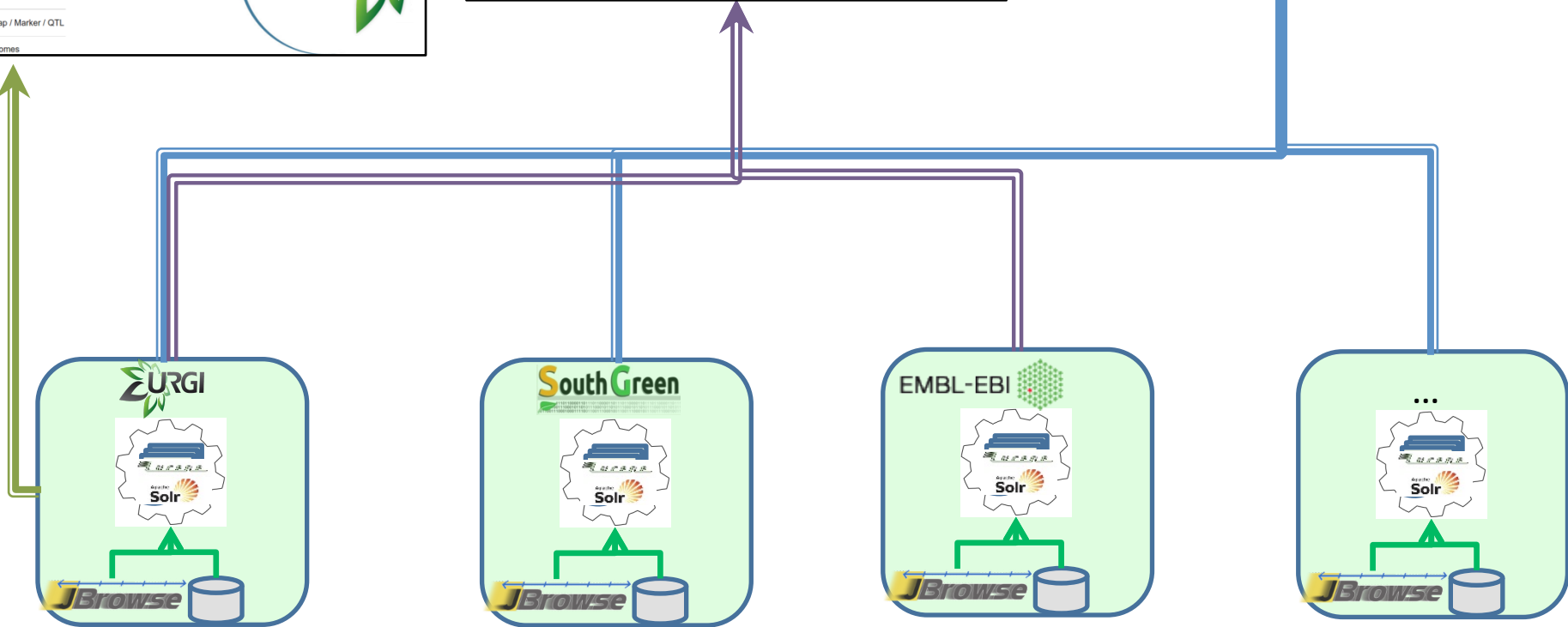
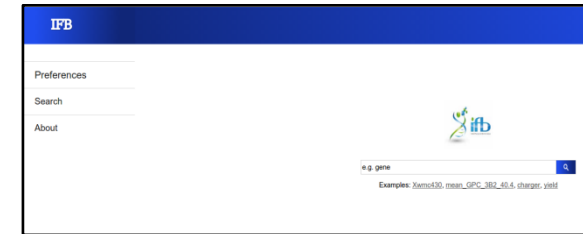
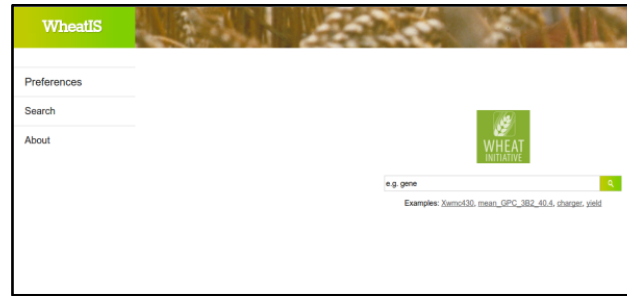
Full text search of distributed databases



Data discovery



WP2 Search portal



Future directions

- Phase II: having both wider and deeper coverage to promote comparative analysis
 - Include varied data types – gene expression data, protein-protein interaction, Transcription factor- target gene
- Developing methods to aid the process of hypotheses generation - e.g. inference rules.
- Query translation – natural language query translation.
- Engage with biologists to mobilise ‘user-pull’:
 - Develop real world use cases – studying the molecular mechanism of panicle differentiation in rice

Acknowledgements



**Elizabeth Arnaud,
Leo Valette,
Marie-Angelique Laporte,
Julian Pietragalla**



**Hadi Quesneville
Esther Dzalé-Yeumo,
Cyril Pommier
Florian Philippe**



**Pierre Larmande
Alexis Dereeper**

Contact: pierre.larmande@ird.fr



**Manuel Ruiz,
Guilhem Sempere,
Nordine El Hassouni**



**Aravind Venkatesan,
Gildas Tagny
Luyen Le Ngoc
Imene Chently**

**Clement Jonquet
Konstantin Todorov**

