



HAL
open science

Agronomic Linked Data (AgroLD): A knowledge-based system to enable integrative biology in agronomy

Aravind Venkatesan, Gildas Tagny Ngompe, Nordine El Hassouni, Imène Chentli, Valentin Guignon, Clement Jonquet, Manuel Ruiz, Pierre Larmande

► To cite this version:

Aravind Venkatesan, Gildas Tagny Ngompe, Nordine El Hassouni, Imène Chentli, Valentin Guignon, et al.. Agronomic Linked Data (AgroLD): A knowledge-based system to enable integrative biology in agronomy. PLoS ONE, 2018, 13 (11), pp.e0198270. 10.1371/journal.pone.0198270 . lirmm-01964772

HAL Id: lirmm-01964772

<https://hal-lirmm.ccsd.cnrs.fr/lirmm-01964772v1>

Submitted on 23 Dec 2018

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

RESEARCH ARTICLE

Agronomic Linked Data (AgroLD): A knowledge-based system to enable integrative biology in agronomy

Aravind Venkatesan^{1,2}, Gildas Tagny Ngompe^{1,2}, Nordine El Hassouni^{1,3,4}, Imene Chentli^{1,2}, Valentin Guignon^{4,5}, Clement Jonquet^{1,2}, Manuel Ruiz^{1,3,4,6}, Pierre Larmande^{1,2,4,7*}

1 Institut de Biologie Computationnelle (IBC), Univ. of Montpellier, Montpellier, France, **2** LIRMM, Univ. of Montpellier & CNRS, Montpellier, France, **3** UMR AGAP, CIRAD, Montpellier, France, **4** South Green Bioinformatics Platform, Montpellier, France, **5** Bioversity International, Montpellier, France, **6** AGAP, Univ. of Montpellier, CIRAD, INRA, INRIA, SupAgro, Montpellier, France, **7** DIADE, IRD, Univ. of Montpellier, Montpellier, France

* pierre.larmande@ird.fr



OPEN ACCESS

Citation: Venkatesan A, Tagny Ngompe G, Hassouni NE, Chentli I, Guignon V, Jonquet C, et al. (2018) Agronomic Linked Data (AgroLD): A knowledge-based system to enable integrative biology in agronomy. PLoS ONE 13(11): e0198270. <https://doi.org/10.1371/journal.pone.0198270>

Editor: Le Zhang, Sichuan University, CHINA

Received: May 14, 2018

Accepted: September 3, 2018

Published: November 30, 2018

Copyright: © 2018 Venkatesan et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: The data underlying this study have been uploaded to Zenodo and are accessible using the following DOI: <https://doi.org/10.5281/zenodo.1410742>.

Funding: This research was supported by the Computational Biology Institute of Montpellier (ANR-11-BINF-0002 - <http://www.agence-nationale-recherche.fr/ProjetIA-11-BINF-0002> - project: <http://www.ibc-montpellier.fr>), the Institut Francais de Bioinformatique (ANR-11-INBS-0013 - [## Abstract](http://www.agence-nationale-recherche.fr/ProjetIA-</p>
</div>
<div data-bbox=)

Recent advances in high-throughput technologies have resulted in a tremendous increase in the amount of omics data produced in plant science. This increase, in conjunction with the heterogeneity and variability of the data, presents a major challenge to adopt an integrative research approach. We are facing an urgent need to effectively integrate and assimilate complementary datasets to understand the biological system as a whole. The Semantic Web offers technologies for the integration of heterogeneous data and their transformation into explicit knowledge thanks to ontologies. We have developed the Agronomic Linked Data (AgroLD— www.agrold.org), a knowledge-based system relying on Semantic Web technologies and exploiting standard domain ontologies, to integrate data about plant species of high interest for the plant science community e.g., rice, wheat, arabidopsis. We present some integration results of the project, which initially focused on genomics, proteomics and phenomics. AgroLD is now an RDF (Resource Description Format) knowledge base of 100M triples created by annotating and integrating more than 50 datasets coming from 10 data sources—such as Gramene.org and TropGeneDB—with 10 ontologies—such as the Gene Ontology and Plant Trait Ontology. Our evaluation results show users appreciate the multiple query modes which support different use cases. AgroLD’s objective is to offer a domain specific knowledge platform to solve complex biological and agronomical questions related to the implication of genes/proteins in, for instances, plant disease resistance or high yield traits. We expect the resolution of these questions to facilitate the formulation of new scientific hypotheses to be validated with a knowledge-oriented approach.

11-INBS-0013 - project: <http://www.france-bioinformatique.fr>), the Labex Agro (ANR-10-LABX-001-01 - <http://www.agence-nationale-recherche.fr/ProjetIA-10-LABX-0001> - project: <http://www.agropolis-fondation.fr/>) all bypass of the French ANR Investissements d'Avenir program (<http://www.agence-nationale-recherche.fr/investissements-d-avenir>). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors have declared that no competing interests exist.

Introduction and background

Agronomy is a multi-disciplinary scientific discipline that includes research areas such as plant molecular biology, physiology and agro-ecology. Agronomic research aims to improve crop production and study the environmental impact on crops. Accordingly, researchers need to understand the implications and interactions of the various biological processes, by linking data at different scales (e.g., genomics, proteomics and phenomics). We are currently witnessing rapid advances in high throughput and information technologies that continue to drive a flood of data and analysis techniques within the domains mentioned above. However, much of these data or information are dispersed across different domain or model specific databases, varied formats and representations e.g., TAIR, GrainGenes and Gramene. Therefore, using these databases more effectively and adopting an integrative approach remains a major challenge.

Among the numerous research directions that the field of bioinformatics has taken, knowledge management has become a major area of research, focused on logically interlinking information and the representation of domain knowledge [1]. To this end, ontologies have become a cornerstone in the representation of biological and more recently agronomical knowledge [2]. Ontologies provide the necessary scaffold to represent and formalize biological concepts and their relationships. Currently, numerous applications exploit the advantages offered by biological ontologies such as: the Gene Ontology [3]—widely used to annotate genes and their products—Plant Ontology [4], Crop Ontology [5], Environment Ontology [6], to name a few. Ontologies have opened the space to various types of semantic applications [7,8] to data integration [9], and to decision support [10]. Semantic interoperability has been identified as a key issue for agronomy, and the use of ontologies declared a way to address it [11]. Furthermore, efficient knowledge management requires the adoption of effective data integration methodologies. This involves efficient semantic integration of the disparate data sources, making information machine-readable and interoperable. Accordingly, Semantic Web standards and technologies enforced by the W3C, and embracing Tim Berners-Lee's vision [12], offers a solution to facilitate integration and interoperability of highly diverse and distributed data resources. The Semantic Web technologies stack includes among others the following W3C Recommendations: the Resource Description Framework (RDF) [13] as a backbone language to describe resources with triples, RDF Schema (RDFS) [14] to build lightweight data schemas, Web Ontology Language (OWL) [15] to build semantically rich ontologies and the SPARQL Query Language (SPARQL) [16] to query RDF data. All of the previous languages rely on Unique Resource Identifiers (URIs) to define a resource and its components, enabling data interoperability across the Web. RDF describes a resource and its relationships/properties in the form of simple triples, i.e., *Subject-Predicate-Object* offering a very convenient framework for integrating data across multiple platforms assuming the platforms share some common vocabularies to describe their objects. These triples can be combined to construct large networks of information (also known as RDF graphs). A successfully implemented Semantic Web application allows scientists to pose very complex questions through a query or a set of queries that would return highly relevant answers to those questions, facilitating the formulation of research hypotheses [17,18].

There are other approaches to meet the current data integration challenges, e.g., data warehouses. For instance, Intermine [19] has developed a sophisticated application to accommodate the dynamic nature of biological data and simplify data integration. However, with integrative biology gaining popularity, it is necessary to preserve and share the semantics between the various datasets and make information machine interoperable, enabling large scale analyses of information available over the Web. The Semantic Web approach provides an added value, playing a complementary role to the traditional methods of data integration.

In the recent years, the biomedical community has strongly embraced the Semantic Web vision as demonstrated by a number of initiatives to provide ontologies [20,21] and use them for producing semantically rich data such as in Bio2RDF [22], OpenPHACTS [23], Linked Life Data [24], KUPKB [25], and the EBI RDF Platform [26]. In particular, OpenPHACTS serves as a good example of what can be achieved by using Semantic Web knowledge bases. The OpenPHACTS Explorer (<http://www.openphacts.org/open-phacts-discovery-platform/explorer>) provides use case driven tools that aid in browsing and visualizing the underlying knowledge represented in RDF which is very convenient for biologists.

Currently, there is a growing awareness within the agronomic domain towards efficient data interoperability and integration [2,27,28]. The need for an umbrella approach for providing uniform data is a widely-discussed topic. For instance, the Agriculture Data Interoperability Interest Group (<https://rd-alliance.org/groups/agriculture-data-interest-group-igad.html>) instituted by the Research Data Alliance (RDA) and agINFRA EU project (www.aginfra.eu) are initiatives that work on improving data standards and promoting data interoperability in agriculture. Moreover, the community has recently also started to adopt AgroPortal [11] as an vocabulary and ontology repository for agronomy—and related domains such as nutrition, plant sciences and biodiversity—that support browsing, searching and visualizing domain relevant ontologies, ontology alignments and creation of semantic annotations. While plant-centric ontologies are now being used to annotate data by various databases developers [2,5,28], unlike in the biomedical domain, the adoption of Semantic Web in agronomy is yet to be completely exploited. Given that agronomic studies involve multiple domains, publicly available knowledge bases such as EBI RDF, Linked Life Data and Bio2RDF serves only limited agronomical information. Hence, it is necessary to build on previous efforts and complete them to provide information compliant with Semantic Web principles within agronomic sciences. This adoption would certainly allow the homogenization of multi-scale information, thereby aiding in the discovery of new knowledge. Therefore, we have developed an RDF knowledge-based system, fully compliant with the Semantic Web vision, called Agronomic Linked Data (AgroLD—www.agrold.org) presented hereafter. The aim of our effort is to provide a portal (to discover) and an endpoint (to query) for integrated agronomic information and to aid domain experts in answering relevant biological questions.

The rest of the paper is organized as follows: in the next section, we describe the data sources integrated or used for the integration, the content and architecture of the knowledge-based system. In the following sections, we present the user interface with some examples queries, then we discuss about the contributions and the future directions.

Materials and methods

Information sources

AgroLD was conceived to accommodate molecular and phenotypic information available on various plant species (see Fig 1). The conceptual framework for the knowledge in AgroLD is based on well-established ontologies: GO, SO, PO, Plant Trait Ontology (TO) and Plant Environment Ontology (EO). Among these PO, TO and EO are currently developed by the Planteome project [29] (<http://planteome.org>). Furthermore, considering the scope of the effort, we decided to build AgroLD in phases. The current phase (phase I) covers information on genes, proteins, ontology associations, homology predictions, metabolic pathways, plant traits, and germplasm, relevant to the selected species. At this stage, we have incorporated the corresponding information from various databases, such as Gramene [30], UniprotKB [31], Gene Ontology Annotation [32], TropGeneDB [33], OryGenesDB [34], Oryza Tag Line [35], GreenPhylDB [36] and SNIPlay [37]. The selection of these data sources was considered based on

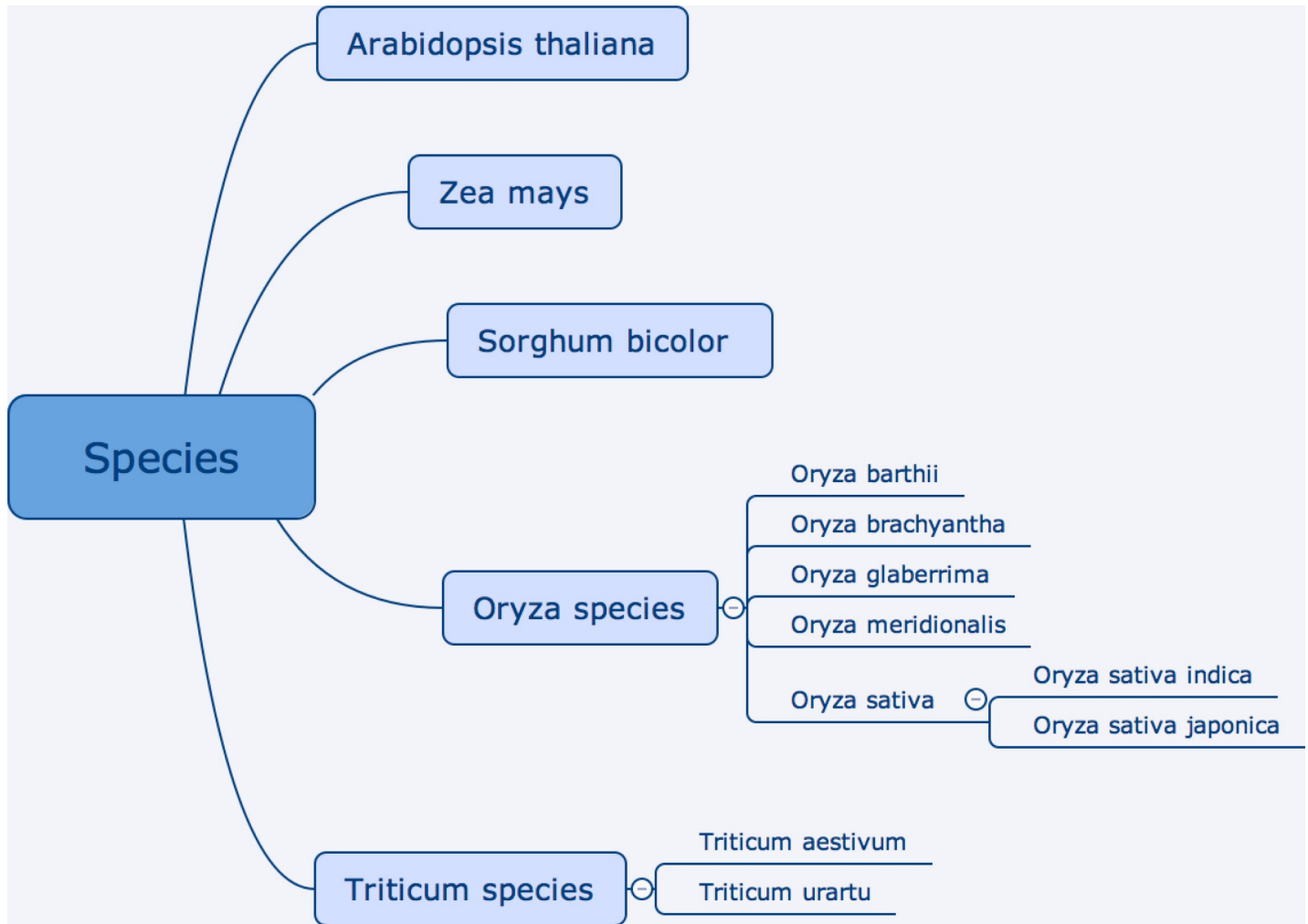


Fig 1. Current plant species included in AgroLD.

<https://doi.org/10.1371/journal.pone.0198270.g001>

popularity among domain experts such as GOA, Gramene, and complementary information hosted by the local research community, for instance, Oryza Tag Line and GreenPhylDB. Information on the integrated databases can be found in the documentation page (<http://www.agrold.org/documentation.jsp>). Table 1 provides a break-down of the data sources and the species covered.

Architecture

AgroLD relies on the RDF and SPARQL technologies for information modelling and retrieval. We use OpenLink Virtuoso (version 7.2) to store and access the RDF graphs. The data from the selected databases were parsed and converted into RDF using a semi-automated pipeline. The pipeline consists of several parsers to handle data in a variety of formats, such as the Gene Ontology Annotation File (GAF) [38], Generic File Format (GFF3) [39], HapMap [40] and Variant Call Format (VCF) [41]. Fig 2 shows the Extraction-Transform-Load (ETL) processes developed to transform in RDF various source data formats. The source code of the ETL workflow (<https://doi.org/10.5281/zenodo.1294660>) is available on GitHub (<https://github.com/SouthGreenPlatform/AgroLD>).

Table 1. Plant species and data sources in AgroLD.

Data sources	URLs	File format	#tuples	Crops	Ontologies used	#triples produced
GO associations	geneontology.org	GAF	1, 160K	R, W, A, M, S	GO, PO, TO, EO	6, 200K
Gramene	gramene.org	Custom flat file	1, 718K	R, W, M, A, S	GO, PO, TO, EO	4, 600K
UniprotKB	uniprot.org	Custom flat file	1, 400K	R, W, A, M, S	GO, PO	50, 000 K
OryGenesDB	orygenesdb.cirad.fr	GFF	1, 100K	R, S, A,	GO, SO	14, 800K
Oryza Tag Line	oryzatagline.cirad.fr	Custom flat file	22K	R	PO, TO, CO	300K
TropGeneDB	tropgenedb.cirad.fr	Custom flat file	2k	R	PO, TO, CO	20K
GreenPhylDB	greenphyl.org	Custom flat file	100K	R, A	GO, PO	700K
SNIplay	sniplay.southgreen.fr	HapMap, VCF	16K	R	GO	16, 000K
Q-TARO	Qtaro.abr.affrc.go.jp	Custom flat file	2K	R	PO,TO	20K
Oryzabase	shigen.nig.ac.jp/rice/oryzabase	Custom flat file	17K	R	GO,PO,TO	160K
TOTAL						92, 640K

The number of tuples gives an idea of the number of elements we have annotated from the data sources (e.g., 1160K Gene Ontology annotations). The crops & ontologies are referred as follows: R = rice, W = wheat, A = Arabidopsis, S = sorghum, M = maize, GO = Gene Ontology, PO = Plant Ontology, TO = Plant Trait Ontology, EO = Plant Environment Ontology, SO = Sequence Ontology, CO = Crop Ontology (specific trait ontologies).

<https://doi.org/10.1371/journal.pone.0198270.t001>

ETL process

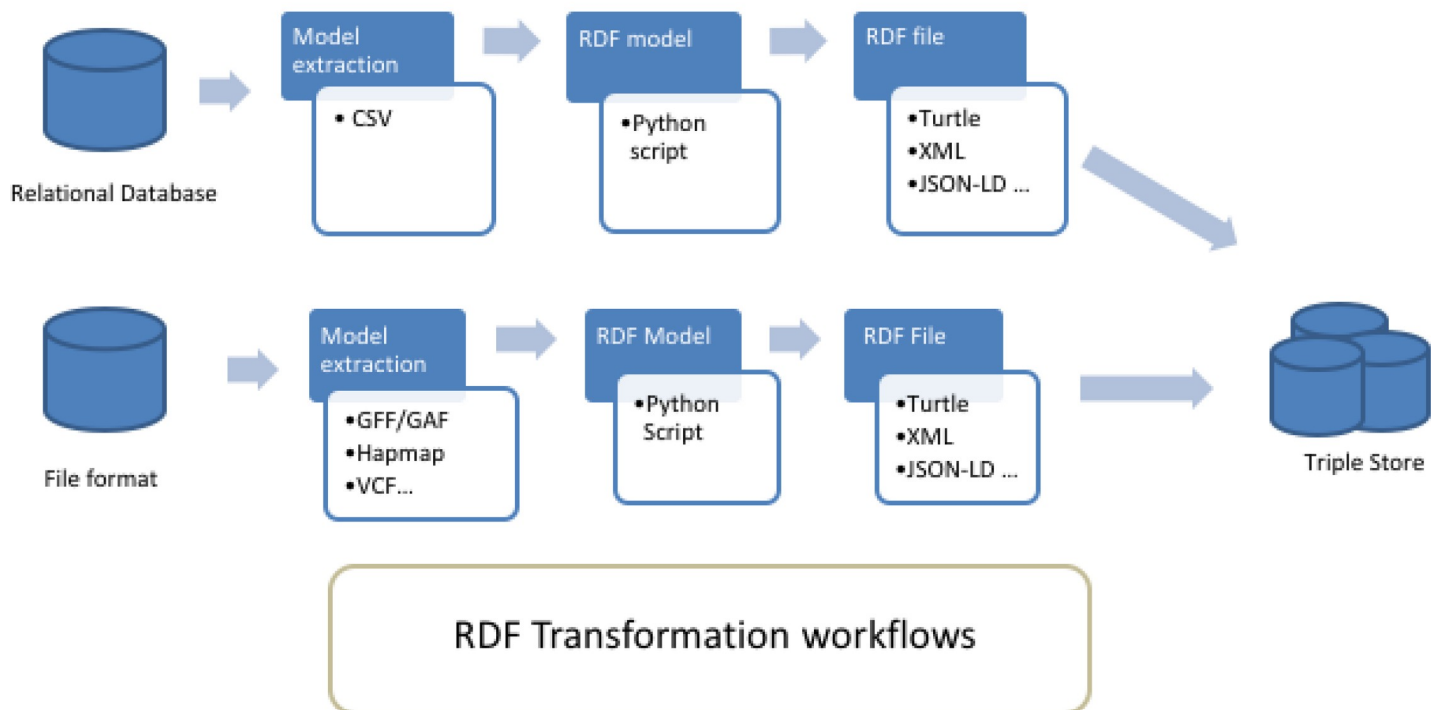


Fig 2. ETL workflow for the various datasets and data formats. The workflow shows two types of process: 1) from relational databases through a CSV file export: in that case, the transformation is tailored for the database model with some Python scripts converters. 2) from standards file formats: in that case, the transformation is generic with some Python packages used as converter tools. The workflow outputs can be produce in various type of RDF format such as turtle, JSON-LD, XML.

<https://doi.org/10.1371/journal.pone.0198270.g002>

For this phase, each dataset was downloaded from curated sources and was annotated with ontology terms URIs by reusing the ontology fields when provided by the original source. Additionally, we used the AgroPortal web service API to retrieve the URI corresponding to the taxon available for some data standards such as GFF. At the end of phase 1, early 2018, the AgroLD knowledge base contains around 100 million RDF triples created by converting more than 50 datasets from 10 data sources. Additionally, when available, we used some semantic annotation already present in the datasets such as, for instances, genes or traits annotated respectively with GO or TO identifiers. In that case, we produced additional properties with the corresponding ontologies thus adding 22% additional triples validated manually (see details in Table 1). The OWL versions of the candidate ontologies were directly loaded into the knowledge base but their triples are not counted in the total. We provided in the supplementary file S1 Table, a more comprehensive statistics analysis such as number of triples, classes, entities and properties for each graph stored in the knowledge base.

The RDF graphs are named after the corresponding data sources (protein/qlt ontology annotations being the exception), sharing a common namespace: “<http://www.southgreen.fr/agrold/>”. The entities in the RDF graphs are linked by shared common URIs. As a design principle, we have used URI schemes made available by the sources (e.g., UniprotKB) or by Identifiers.org registry (<http://identifiers.org> - [42]). For instances, proteins from UnitProtKB are identified by the base URI: <http://purl.uniprot.org/uniprot/>; genes incorporated from Gramene/Ensembl plants are identified by the base URI: <http://identifiers.org/ensembl.plant/>. New URIs were minted when not provided by the sources or the by Identifiers.org such as TropGene and OryGenesDB; in such cases the URIs take the form [http://www.southgreen.fr/agrold/\[resource_namespace\]/\[identifier\]](http://www.southgreen.fr/agrold/[resource_namespace]/[identifier]). Furthermore, properties linking the entities took the form: [http://www.southgreen.fr/agrold/vocabulary/\[property\]](http://www.southgreen.fr/agrold/vocabulary/[property]). An outline of how the RDF graphs are linked is shown in Fig 3. About entity linking, we used the “key-based approach” which is the most common one. It combines the unique identifier/accession number of the entity shared with the community, with the URI basis pattern of the resource. Moreover, we also respected

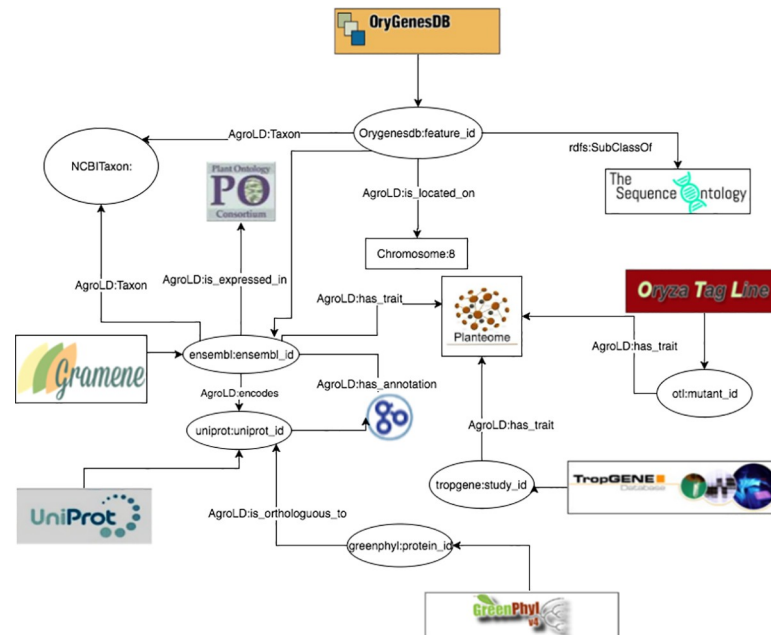


Fig 3. Linking information in AgroLD. The figure illustrates the linking of various information in AgroLD.

<https://doi.org/10.1371/journal.pone.0198270.g003>

the “common URI approach” which recommends to use the same URI pattern when the same accession number is used in different datasets. Therefore, defining the same URI for identical entities (represented by identifiers) in different datasets makes it possible to aggregate additional information for this entity. Additionally, we used cross-reference links (represented by identifiers from external datasets) by transforming them into URIs and linked the resource with the predicate “has_dbxref”. This greatly increases the number of outbound links, making AgroLD more integrated with other Linked Open Data. In the future, we will implement a “similarity-based approach” to identify correspondences between entities which have different URIs.

To map the various data types and properties, we developed a lightweight schema (cf. <https://github.com/SouthGreenPlatform/AgroLD>) that glues classes and properties identified in AgroLD and the corresponding external ontologies. For instance, the class Protein (<http://www.southgreen.fr/agrold/resource/Protein>) is mapped as *owl:equivalentClass* to class polypeptide (http://purl.obolibrary.org/obo/SO_0000104) from SO. Similar mappings have been made for properties, e.g., proteins/genes are linked to GO molecular function by the property http://www.southgreen.fr/agrold/vocabulary/has_function, which is mapped as *owl:equivalentProperty* to the corresponding Basic Formal Ontology (BFO) term (http://purl.obolibrary.org/obo/BFO_0000085). When an equivalent property did not exist, we mapped then to the closest upper level property using *rdfs:subPropertyOf* e.g., the property *has_trait* (http://www.southgreen.fr/agrold/vocabulary/has_trait), links proteins to TO terms. It is mapped to a more generic property, *causally related to* in the Relations Ontology [43]. For now, 55 mappings were identified. Furthermore, mappings are both stored side by side with ontologies in AgroPortal, which allows direct links between classes and instances of these classes in AgroLD. For example, the following link will show the external mappings for SO:0000104 (polypeptide) stored in AgroPortal: http://agroportal.lirmm.fr/ontologies/SO/?p=classes&conceptid=http%3A%2F%2Fpurl.obolibrary.org%2Fobo%2FSO_0000104&jump_to_nav=true#mappings. Additionally, classes, properties and resources (e.g., <http://www.southgreen.fr/agrold/page/biocyc.pathway/CALVIN-PWY>) are dereferenced on a dedicated Pubby server [44]. For details on the graphs, URIs and properties, the reader may refer to AgroLD’s documentation (<http://www.agrold.org/documentation.jsp>).

User interface

The AgroLD platform provides four entry points to access the knowledge base:

- *Quick Search* (<http://www.agrold.org/quicksearch.jsp>), a faceted search plugin made available by Virtuoso, that allows users to search by keywords and browse the AgroLD’s content;
- *SPARQL Query Editor* (<http://www.agrold.org/sparqleditor.jsp>), that provides an interactive environment to formulate SPARQL queries;
- *Explore Relationships* visualizer (<http://www.agrold.org/relfinder.jsp>), which is an implementation of RelFinder [45] that allows users to explore and visualize existing relationships between entities;
- *Advanced Search* (<http://www.agrold.org/advancedSearch.jsp>), a query form providing entity (e.g., gene) specific information retrieval.

Alternatively, some user management features have been implemented on the platform. Users have the opportunity to save their search and results on a persistent history session attached to their own account. Furthermore, they can manage search history by editing, deleting or re-running previous searches and exporting results according several formats. In the

future, we plan to develop some recommendation features and sharing results between users. More detailed descriptions and figures of the different user interfaces will be provided in the following section. Furthermore, other examples are shown in the User Guide available in the supporting information [S1 File](#).

Results and discussion

RDF knowledge bases are accessed via SPARQL endpoints and in certain cases equipped with faceted browser interfaces. Using SPARQL endpoints require a minimal knowledge of SPARQL, this may result in the resources not being exploited completely. Alternatively, faceted browser interfaces help the user in getting acquainted with information in the resource (e.g., retrieving a local neighborhood for a particular term), the presence non-textual details (e.g., URIs) in the results could be confusing. To this end, we attempted to lower the usability barrier by providing tools to explore the knowledge base. In this section, we demonstrate the complementary role of the *Advanced Search* and *Explore Relationships* query tools with that of the *SPARQL Query Editor*.

We developed the SPARQL Query Editor based on the YASQE and YASR tools [46] and customized it for our system. The SPARQL language is a powerful tool to mine and extract meaningful information from the knowledge base. In the first example of the supplementary [S3 File](#), we compare two queries to answer the question: “Identify wheat proteins that are involved in root development.”. While the first one (S3_Q1) using a simple search—which is a direct translation of SQL—with the corresponding id (“GO_0048364”, “GO_2000280”) shows 73 entries, the second one (S3_Q2) using a property path query (i.e., query the descending class hierarchy for a given trait ontology term) shows 137 entries, thus more than 80% of additional results. In that case, the use of property path algorithm shows the efficiency in retrieving a comprehensive answer. But the SPARQL language performs also very well with complex queries such as: “Retrieve individuals which have positive SNP variant effect identified for proteins associated with a QTL” available in S3_Q3. This type of query involves several datasets and uses graph traversal property of SPARQL to perform the query.

Because SPARQL is hard to handle for non-technical users, the *SPARQL Query Editor* includes a list of modularized example queries, customizable according to the users’ needs.

For the comparison, we consider a sample question: ‘*Retrieving genes that participate in Calvin cycle*’; (Q6 in the online list of modularized queries). As illustrated in [Fig 4](#), the user can run the query to retrieve the list of genes participating in the given pathway ([Fig 4A](#)). Additional information on a gene of interest can be retrieved by clicking on the URI. For example, clicking on AT1GI870 (<http://identifiers.org/ensembl.plant/AT1GI8270>) redirects the users to the gene information provided by Gramene/Ensembl Plants resource ([Fig 4B](#)). The query can be saved and the results can be downloaded in a variety of formats such as JSON, TSV, and RDF/XML. Additionally, user defined queries could also be uploaded.

The *Explore Relationships* tool is based on RelFinder visualization module. This tool aids in visualizing relationships between entities and searching entities by keyword when their URIs are ignored. However, the original version of RelFinder was developed (in ActionScript) and configured for DBpedia. We proposed a configuration and modification of the system suitable for AgroLD. The configuration mainly concerns the SPARQL access point, the properties to be considered for the search of entities and for the description of the resources. Furthermore, we have added some biological examples to guide users. In [Fig 5](#), the tool is used to search for genes involved in Calvin cycle by entering the name of the entities.

The *Advanced Search* query form is based on the REST API suite (<http://www.agrold.org/api-doc.jsp>), developed completely within the AgroLD project. The aim of this feature is to

Search > SPARQL Query Editor

Select a sample query and run it. The sample query could be used to modify the parameters accordingly. Alternatively, enter SPARQL code in the query box below.

Query Text

```

1 BASE <http://www.southgreen.fr/agrold/>
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 }
17 }
```

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

Query Patterns

1. Retrieve list of graphs ([select](#))
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: P0C127) ([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 ontological term "**homoacconitate 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:**P0C127** ([select](#))

Results

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

EnsemblPlants | BLAST | BioMart | Tools | Downloads | Documentation | Website help

Arabidopsis thaliana (TAIR10) | Location: 1:6,283,412-5,293,871 | Gene: AT1G18270

Gene: AT1G18270

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

Location: [Chromosome 1: 6,283,412-5,293,871](#) reverse strand.

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

Transcripts: [Show transcript table](#)

Fig 4. SPARQL query editor. Figure illustrates the execution of query Q6: (a) Q6 is one the examples queries on the top-right corner (highlighted in red). On executing the query, the results are rendered below the editor; (b) the user can look up specific genes of interest by clicking on the corresponding URI, which points to the original information source (in this case EsemblPlants).

<https://doi.org/10.1371/journal.pone.0198270.g004>

provide non-technical users with a tool to query the knowledge base while hiding the technical aspects of SPARQL query formulation. Fig 6 illustrates steps involved in retrieving information for Q6, using the query form:

- a. The user selects *Pathways* from the list of entities and enters the pathway of interest, in this case, Calvin cycle (Fig 6A);
- b. The list of genes involved in the pathway can be retrieved by selecting the pathway.

Furthermore, information on a gene of interest can be retrieved by selecting the specific gene (Fig 6B). For instance, clicking on AT1G1870 (Fig 6C) displays all the proteins the gene encodes and the pathways the gene participates in (apart from Calvin cycle). The RESTful API supports the query form and was developed for programmatic retrieval of entity specific knowledge represented in AgroLD. The current version of the API suite (ver. 1) can be used to retrieve gene and protein information, metabolic pathways, and proteins associated with ontological terms. This is achieved by querying entity by name or identifier.

User evaluation

AgroLD is being actively developed based on usability testing sessions conducted with domain experts, including doctoral students in biology, curators and senior researchers. Test sessions were designed to measure if:

Search > Explore

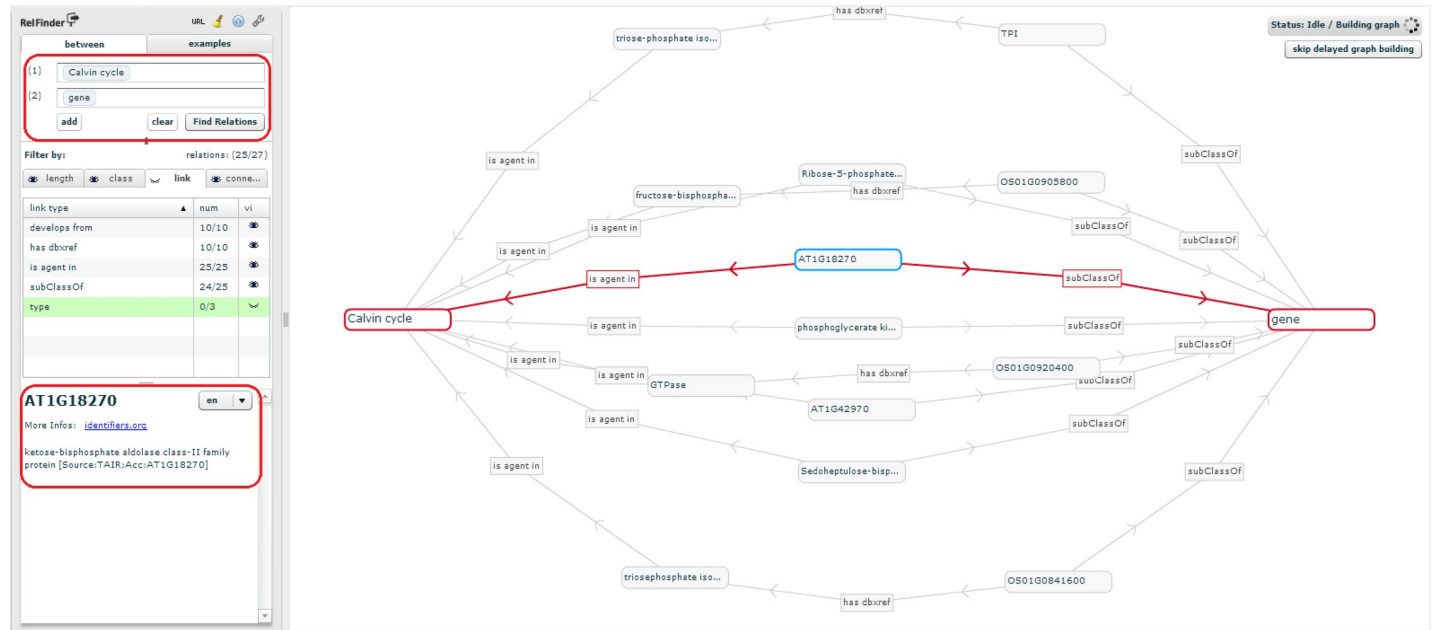


Fig 5. Exploring entity relationships in AgroLD. Figure illustrates differently the results obtained for Q6 using Explore Relationships tool. The results of Q6 can be visualized by entering the concepts (Calvin cycle and gene) in the left panel. On executing the query, all the genes involved in the chosen pathway are revealed. The visualized graph can be altered based on the user interest. Additionally, a gene could be selected (circled on the left) and further explored by clicking on the *More Info* link which directs the user to the information source.

<https://doi.org/10.1371/journal.pone.0198270.g005>

- Resources integrated in AgroLD are useful;
- AgroLD is easy to use.

For the evaluation of semantic search systems, Elbedweihy et al. [47] recommend a survey of users based on their experience with a few queries submitted to the system. We have used this approach to collect user opinions, comments and suggestions via a feedback form directly within the AgroLD web application. The form includes some questions from the "System Usability Scale" questionnaires [48] and other questions that we considered important. The three main criteria evaluated are:

1. Usability—ease to submit a query (number of attempts, time required) and presentation of the results;
2. Expressiveness—type of queries a user is able to formulate (e.g., keywords or more complex expressions);
3. Performance—speed, correctness and completeness of the results.

Recently, 20 participants were invited during 3 testing sessions, to search for concepts, genes, or pathways of their interests; and the online form was active (<http://agrold.org/survey.jsp>) to allow new feedbacks during the exploitation phase. Each question had 5 possible answers ranked from the highest to the lowest note (5 to 1). We reported the results of these sessions in *S2 File* as a supplementary document.

Globally, participants found the platform useful and easy to use. Overall, the idea of data navigation and traversal through knowledge graphs was well received. However, many of them needed help with some features. The general observation is that testing users ranked *Advanced*

Search > 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'

Pathway Calvin cycle Search

Search pathway with keyword " Calvin cycle "

Id	Name	URI
1 CALVIN-PWY (display)	Calvin cycle	http://www.southgreen.fr/agrold/biocyc.pathway/CALVIN-PWY (in Sparql)

Showing 1 to 1 of 1 entries

a)

PATHWAY : CALVIN-PWY / Calvin cycle

URI: http://www.southgreen.fr/agrold/biocyc.pathway/CALVIN-PWY

Participating genes Next page>>

geneid	gene_name	taxon	taxon_name	URI
1 AT1G18270 (display)	fructose-bisphosphate aldolase	http://purl.obolibrary.org/obo/NCBITaxon_3702 (in Sparql)	Arabidopsis thaliana	http://identifiers.org/ensembl.plant/AT1G18270 (in Sparql)
2 AT1G42970 (display)	glyceraldehyde-3-phosphate dehydrogenase	http://purl.obolibrary.org/obo/NCBITaxon_3702 (in Sparql)	Arabidopsis thaliana	http://identifiers.org/ensembl.plant/AT1G42970 (in Sparql)

b)

GENE : AT1G18270 / fructose-bisphosphate aldolase

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

URI: http://identifiers.org/ensembl.plant/AT1G18270

encodes proteins ±

Pathways ±

c)

Fig 6. Advanced search query form: Figure demonstrates the steps involved in retrieving the results for Q6 using the Advanced Search query form: (a) query Q6 can be executed by selecting the type of entity (Pathways–highlighted in red) to search and entering the name of the entity (Calvin cycle). The API then displays the matched results; (b) Clicking on the result displays the genes participating in Calvin cycle; (c) selecting a gene of interest displays more information pertaining to that gene, for instance, encoding proteins and pathways this selected gene participates in.

<https://doi.org/10.1371/journal.pone.0198270.g006>

Search first then Quick Search after. We explain this by the display output that looks friendlier for Advanced Search. Quick Search won votes for usability and performance despite several comments to improve the ranking and presentation of results (4 user’s comments). Advanced and Explore search got average scores but good comments on the capability of discovering unexpected results (e.g., nearest neighbour entities in the graph for the Explore Search and additional results from external Web services for Advanced Search). With no surprise, evaluation results show the SPARQL Query Editor is the most difficult to handle. We mitigate this by offering examples of query pattern to help users handle query formulation. In the future, we will improve the examples by offering a large spectrum of search type which will follow the new phase of data integration. Furthermore, we will provide links to some SPARQL tutorials in the documentation. These user feedbacks reinforced the need for knowledge bases such as AgroLD, wherein users could retrieve information across various data types and sources. This knowledge discovery is supported by the use of shared URI schemes and domain ontologies. The testing sessions also helped us to identify areas for further improvement. Plus, we received

suggestions on improving the AgroLD's coverage with more data types such as gene expression data, and protein-protein interactions. Considering, linked data and Semantic Web are still not widely adopted in agronomy, increasing AgroLD's coverage will be an incremental process engaging our user community. This situation is expected to improve with new community efforts such as the Agrisemantics RDA Working Group (<https://rd-alliance.org/groups/agrisemantics-wg.html>), which role is to reinforce the adoption of semantic technologies in the agri-food domain. We may also mention the AgBioData consortium (<https://www.agbiodata.org>, [2]) which promotes the FAIR (Findable, Accessible, Interoperable and Reusable) data principles [49] within agricultural research.

Furthermore, we observed that although the information integrated in AgroLD came from curated sources, scientists often prefer to validate these knowledge statements against assertions made in scientific articles. Currently, we have implemented an external Web Services as part of the *Advanced Search Form* to automatically search for publications related to a protein or gene of interest in PubMed Central and aggregates them within the result of the AgroLD query. However, this feature does not provide detailed (sentence level) assertions described in those publications. This is an area that requires further work. With the recent developments towards making text mined (sentence level) annotations available as RDF [50], query federation can be explored to retrieve entity specific assertions. This would serve as an additional provenance layer.

Limits and perspectives

With the achievement of the first phase of AgroLD, many plant scientists can benefit from the interoperability of the data, but user feedback reveals some limitations and challenges on the current version of AgroLD. In order to achieve the expectations of the scientists for the use of Semantic Web technologies in agronomy, a number of issues need to be addressed:

- The coverage content has to be extended to a larger number of biological entities (e.g., miRNA, mRNA) or interaction between them (e.g., co-expression, regulation and interaction networks) in order to capture a broad view of the molecular interactions.
- We have observed many information remains hidden in RDF literal contents such as biological entities or relationship between them. This information is poorly annotated (i.e., plain text not formally expressed) and new research methods to identify biological entities and reconstruct their relations further allowing the discovery of relevant links between related resources are required.
- The explosion of data in agronomy forces database providers to augment the frequency of their releases. The survey shows a growing interest of using up to date information from the original sources. This have to be taken into account for the updating process in AgroLD.
- The user interfaces show some limitations to manage responses with large number of results, e.g., to filter and rank them with precision score.

These limitations identified in the current version of AgroLD will be improved in the following versions. We will focus on the following areas:

- User Interface: we plan to explore features offered by Elastic search tool (<https://www.elastic.co>), to enabling *Quick Search* retrieving more textual information and hiding the technical details. Further, we will improve the performance and expand the API suite to cover other entities represented in AgroLD (e.g., genomic annotation and homology information).
- Content: integrate information on gene expression such as IC4R [51], Gene Expression Atlas [52], on gene regulatory networks such as RiceNetDB [53] and explore linking text-

mined annotations from publications. Support molecular interaction networks per species and also allow knowledge transfer between species.

- **Knowledge discovery:** explore methods to aid generating hypotheses by retrieving implicit knowledge, e.g., inference rules, automatic data linking, entity recognition, text mining, automatic semantic annotations.
- **Data provenance:** develop a provenance and annotation model. Set up a validation process to allow users validating computed facts such as semantic annotations automatically produced and attached to a biological entity.
- **Updates:** To keep AgroLD updated with the latest available data, by processing regular data updates and potentially re-building the entire repository from scratch every 12 months. Processing regular data update is a hard issue as the original databases do not always provide an automatic way to obtain the differential data between releases. From experience, we know that regularly rebuilding the entire knowledge base is for us a good alternative to avoid dealing with data diffs. Additionally, we plan to fully automate the current ETL workflow.

Conclusion

Data in the agronomic domain are highly heterogeneous and dispersed. For agronomic researchers to make informed decisions in their daily work it is critical to integrate information at different scales. Current traditional information systems are not able to exploit such data (i.e., genes, proteins, metabolic pathways, plant traits, and phenotypes), in efficient way. To this end, the application of Semantic Web, initiated in the biomedical domain, provides a good example to follow by capitalizing on previous experiences and addressing weaknesses.

To further build on this line of research in agronomy, we have developed AgroLD. We have demonstrated the advantages of AgroLD in data integration over multiple data sources using plant domain ontologies and Semantic Web technologies. To date, AgroLD contains 100M of triples created by transforming more than 50 datasets coming from 10 data and annotating with 10 ontologies. The impact of AgroLD is expected to grow with an increase in coverage (with respect to the species and the data sources) and user inputs. For instance, when user feedback and implementation of inference rules are put within a context that supports searching and recommendations, then we have the beginnings of a platform that can support automated hypotheses generation.

AgroLD is one of the first RDF linked open data knowledge-based system in the agronomic domain. It demonstrates a first step toward adopting the Semantic Web technologies to facilitate research by integrating numerous heterogeneous data and transforming them into explicitly knowledge thanks to ontologies. We expect AgroLD will facilitate the formulation of new scientific hypotheses to be validated with its knowledge-oriented approach.

Supporting information

S1 File. AgroLD user guide. This document shows how to use the various features of the platform.

(PDF)

S2 File. Report of the online survey. Report of 3 sessions evaluating the AgroLD user interfaces.

(PDF)

S3 File. Examples of SPARQL queries. Example of SPARQL queries showing the benefits of property path algorithm, and complex queries.
(PDF)

S1 Table. AgroLD graph statistics.
(PDF)

Acknowledgments

Authors thank the technical staffs of the South Green Bioinformatics platform for their support. Authors thank the providers of databases listed in Fig 1, who kindly gave access to their publicly datasets. Authors thank the expert biologists and bioinformaticians who contributed to the testing sessions and helped us to improve the content of the system and the user interface. Authors specially thank Dr. Patrick Valduriez and Dr. Eric Rivals for their supports and advises in this project.

Author Contributions

Conceptualization: Aravind Venkatesan, Pierre Larmande.

Data curation: Aravind Venkatesan, Pierre Larmande.

Formal analysis: Aravind Venkatesan.

Funding acquisition: Manuel Ruiz, Pierre Larmande.

Investigation: Aravind Venkatesan.

Methodology: Aravind Venkatesan.

Project administration: Pierre Larmande.

Resources: Aravind Venkatesan.

Software: Aravind Venkatesan, Gildas Tagny Ngompe, Nordine El Hassouni, Imene Chentli, Valentin Guignon, Pierre Larmande.

Supervision: Pierre Larmande.

Validation: Aravind Venkatesan.

Writing – original draft: Aravind Venkatesan.

Writing – review & editing: Clement Jonquet, Manuel Ruiz, Pierre Larmande.

References

1. Goble C, Stevens R. State of the nation in data integration for bioinformatics. *J Biomed Inform.* Elsevier; 2008; 41: 687–693. <https://doi.org/10.1016/j.jbi.2008.01.008> PMID: 18358788
2. Harper L, Campbell J, Cannon EK, Jung S, Main D, Poelchau M, et al. AgBioData Consortium Recommendations for Sustainable Genomics and Genetics Databases for Agriculture. *Database.* 2018; 1–7.
3. Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, et al. Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nat Genet.* 2000; 25: 25–29. <https://doi.org/10.1038/75556> PMID: 10802651
4. Cooper L, Walls RL, Elser J, Gandolfo MA, Stevenson DW, Smith B, et al. The plant ontology as a tool for comparative plant anatomy and genomic analyses. *Plant Cell Physiol.* 2013; 54: e1. <https://doi.org/10.1093/pcp/pcs163> PMID: 23220694
5. Shrestha R, Matteis L, Skofic M, Portugal A, McLaren G, Hyman G, et al. Bridging the phenotypic and genetic data useful for integrated breeding through a data annotation using the Crop Ontology

- developed by the crop communities of practice. *Front Physiol.* 2012; 3: 326. <https://doi.org/10.3389/fphys.2012.00326> PMID: 22934074
6. Buttigieg PL, Morrison N, Smith B, Mungall CJ, Lewis SE, ENVO Consortium. The environment ontology: contextualising biological and biomedical entities. *J Biomed Semantics.* 2013; 4: 43. <https://doi.org/10.1186/2041-1480-4-43> PMID: 24330602
 7. Walls RL, Deck J, Guralnick R, Baskauf S, Beaman R, Blum S, et al. Semantics in support of biodiversity knowledge discovery: an introduction to the biological collections ontology and related ontologies. *PLoS One.* 2014; 9: e89606. <https://doi.org/10.1371/journal.pone.0089606> PMID: 24595056
 8. Oellrich A, Walls RL, Cannon EK, Cannon SB, Cooper L, Gardiner J, et al. An ontology approach to comparative phenomics in plants. *Plant Methods.* 2015; 11: 10. <https://doi.org/10.1186/s13007-015-0053-y> PMID: 25774204
 9. Wang Y, Wang Y, Wang J, Yuan Y, Zhang Z. An ontology-based approach to integration of hilly citrus production knowledge. *Comput Electron Agric. Elsevier;* 2015; 113: 24–43. <https://doi.org/10.1016/J.COMPAG.2015.01.009>
 10. Lousteau-Cazalet C, Barakat A, Belaud J-P, Buche P, Busset G, Charnomordic B, et al. A decision support system for eco-efficient biorefinery process comparison using a semantic approach. *Comput Electron Agric. Elsevier;* 2016; 127: 351–367. <https://doi.org/10.1016/J.COMPAG.2016.06.020>
 11. Jonquet C, Toulet A, Arnaud E, Aubin S, Dzalé Yeumo E, Emonet V, et al. AgroPortal: A vocabulary and ontology repository for agronomy. *Comput Electron Agric.* 2018; 144: 126–143. <https://doi.org/10.1016/j.compag.2017.10.012>
 12. Berners-lee T, Hendler J, Lassila O. The Semantic Web. *Sci Am.* 2001; 284: 35–43.
 13. W3C. Resource Description Framework (RDF): Concepts and Abstract Syntax [Internet]. [cited 3 Apr 2010]. Available: <http://www.w3.org/TR/2004/REC-rdf-concepts-20040210/>
 14. W3C. RDF Schema 1.1 [Internet]. [cited 27 Apr 2018]. Available: <https://www.w3.org/TR/rdf-schema/>
 15. W3C. OWL 2 Web Ontology Language Structural Specification and Functional-Style Syntax [Internet]. [cited 3 Apr 2010]. Available: <http://www.w3.org/TR/2009/REC-owl2-syntax-20091027/>
 16. The W3C SPARQL Working Group. SPARQL 1.1 Overview [Internet]. [cited 15 Apr 2013]. Available: <http://www.w3.org/TR/sparql11-overview/>
 17. Luciano JS, Andersson B, Batchelor C, Bodenreider O, Clark T, Denney CK, et al. The Translational Medicine Ontology and Knowledge Base: driving personalized medicine by bridging the gap between bench and bedside. *J Biomed Semantics.* 2011; 2 Suppl 2: S1. <https://doi.org/10.1186/2041-1480-2-S2-S1> PMID: 21624155
 18. Venkatesan A, Tripathi S, Sanz de Galdeano A, Blondé W, Læg Reid A, Mironov V, et al. Finding gene regulatory network candidates using the gene expression knowledge base. *BMC Bioinformatics.* 2014; 15: 386. <https://doi.org/10.1186/s12859-014-0386-y> PMID: 25490885
 19. Smith RN, Aleksic J, Butano D, Carr A, Contrino S, Hu F, et al. InterMine: a flexible data warehouse system for the integration and analysis of heterogeneous biological data. *Bioinformatics.* Oxford University Press; 2012; 28: 3163–5. <https://doi.org/10.1093/bioinformatics/bts577> PMID: 23023984
 20. Smith B, Ashburner M, Rosse C, Bard J, Bug W, Ceusters W, et al. The OBO Foundry: coordinated evolution of ontologies to support biomedical data integration. *Nat Biotechnol.* Nature Publishing Group; 2007; 25: 1251–1255. <https://doi.org/10.1038/nbt1346> PMID: 17989687
 21. Noy NF, Shah NH, Whetzel PL, Dai B, Dorf M, Griffith N, et al. BioPortal: ontologies and integrated data resources at the click of a mouse. *Nucleic Acids Res.* 2009; 37: W170–173. <https://doi.org/10.1093/nar/gkp440> PMID: 19483092
 22. Belleau F, Nolin M-A, Tourigny N, Rigault P, Morissette J. Bio2RDF: towards a mashup to build bioinformatics knowledge systems. *J Biomed Inform. Elsevier;* 2008; 41: 706–716. <https://doi.org/10.1016/j.jbi.2008.03.004> PMID: 18472304
 23. Williams AJ, Harland L, Groth P, Pettifer S, Chichester C, Willighagen EL, et al. Open PHACTS: Semantic interoperability for drug discovery. *Drug Discovery Today.* 2012. pp. 1188–1198. <https://doi.org/10.1016/j.drudis.2012.05.016> PMID: 22683805
 24. Momtchev V, Peychev D, Primov T, Georgiev G. Expanding the Pathway and Interaction Knowledge in Linked Life Data. *International Semantic Web Challenge.* 2009.
 25. Jupp S, Klein J, Schanstra J, Stevens R. Developing a kidney and urinary pathway knowledge base. *J Biomed Semantics.* 2011; 2 Suppl 2: S7. <https://doi.org/10.1186/2041-1480-2-S2-S7> PMID: 21624162
 26. Jupp S, Malone J, Bolleman J, Brandizi M, Davies M, Garcia L, et al. The EBI RDF platform: linked open data for the life sciences. *Bioinformatics.* 2014; 1–2. <https://doi.org/10.1093/bioinformatics/btt765> PMID: 24413672

27. Venkatesan A, El Hassouni N, Phillippe F, Pommier C, Quesneville H, Ruiz M, et al. Towards efficient data integration and knowledge management in the Agronomic domain. *APIA'15: premiere Conference Applications Pratiques de l'Intelligence Artificielle*. 2015.
28. Leonelli S, Davey RP, Arnaud E, Parry G, Bastow R. Data management and best practice for plant science. *Nat Publ Gr. Macmillan Publishers Limited*; 2017; 3: 1–4. <https://doi.org/10.1038/nplants.2017.86> PMID: 28585570
29. Cooper L, Meier A, Laporte MA, Elser JL, Mungall C, Sinn BT, et al. The Planteome database: An integrated resource for reference ontologies, plant genomics and phenomics. *Nucleic Acids Res*. 2018; <https://doi.org/10.1093/nar/gkx1152> PMID: 29186578
30. Monaco MK, Stein J, Naithani S, Wei S, Dharmawardhana P, Kumari S, et al. Gramene 2013: Comparative plant genomics resources. *Nucleic Acids Res*. 2014; 42. <https://doi.org/10.1093/nar/gkt1110> PMID: 24217918
31. Magrane M, Consortium UP. UniProt Knowledgebase: A hub of integrated protein data. *Database*. 2011;2011. <https://doi.org/10.1093/database/bar009> PMID: 21447597
32. Barrell D, Dimmer E, Huntley RP, Binns D, O'Donovan C, Apweiler R. The GOA database in 2009—An integrated Gene Ontology Annotation resource. *Nucleic Acids Res*. 2009; 37. <https://doi.org/10.1093/nar/gkn803> PMID: 18957448
33. Hamelin C, Sempere G, Jouffe V, Ruiz M. TropGeneDB, the multi-tropical crop information system updated and extended. *Nucleic Acids Res*. 2013; 41. <https://doi.org/10.1093/nar/gks1105> PMID: 23161680
34. Droc G, Ruiz M, Larmande P, Pereira A, Piffanelli P, Morel JB, et al. OryGenesDB: a database for rice reverse genetics. *Nucleic Acids Res*. 2006; 34: D736–40. <https://doi.org/10.1093/nar/gkj012> PMID: 16381969
35. Larmande P, Gay C, Lorieux M, Périn C, Bouniol M, Droc G, et al. Oryza Tag Line, a phenotypic mutant database for the Génoplante rice insertion line library. *Nucleic Acids Res*. 2008; 36: 1022–1027. <https://doi.org/10.1093/nar/gkm762> PMID: 17947330
36. Conte MG, Gaillard S, Lanau N, Rouard M, Périn C. GreenPhyIDB: a database for plant comparative genomics. *Nucleic Acids Res*. 2008; 36: D991–998. <https://doi.org/10.1093/nar/gkm934> PMID: 17986457
37. Dereeper A, Homa F, Andres G, Sempere G, Sarah G, Hueber Y, et al. SNIPlay3: a web-based application for exploration and large scale analyses of genomic variations. *Nucleic Acids Res*. 2015; 43: W295–300. <https://doi.org/10.1093/nar/gkv351> PMID: 26040700
38. The Gene Ontology Consortium. Gene Annotation File (GAF) specification [Internet]. [cited 20 Mar 2018]. Available: <http://geneontology.org/page/go-annotation-file-format-20>
39. Sequence Ontology consortium. GFF3 Specification [Internet].
40. Gibbs RA, Belmont JW, Hardenbol P, Willis TD, Yu F, Zhang H, et al. The International HapMap Project. *Nature*. 2003; 426: 789–796. <https://doi.org/10.1038/nature02168> PMID: 14685227
41. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, et al. The variant call format and VCFtools. *Bioinformatics*. 2011; 27: 2156–8. <https://doi.org/10.1093/bioinformatics/btr330> PMID: 21653522
42. Juty N, Le Novère N, Laibe C. Identifiers.org and MIRIAM Registry: community resources to provide persistent identification. *Nucleic Acids Res*. 2012; 40: D580–6. <https://doi.org/10.1093/nar/gkr1097> PMID: 22140103
43. Smith B, Ceusters W, Klagges B, Kohler J, Kumar A, Lomax J, et al. Relations in biomedical ontologies. *Genome Biol*. 2005; 6: R46. <https://doi.org/10.1186/gb-2005-6-5-r46> PMID: 15892874
44. Cyganiak R (National U of I, Bizer C. Pubby—A Linked Data Frontend for SPARQL Endpoints. 2008; Available: <http://wifo5-03.informatik.uni-mannheim.de/pubby/>
45. Heim P, Hellmann S, Lehmann J, Lohmann S, Stegemann T. RelFinder: Revealing relationships in RDF knowledge bases. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*. 2009. pp. 182–187. https://doi.org/10.1007/978-3-642-10543-2_21
46. Rietveld L, Hoekstra R. The YASGUI Family of SPARQL Clients. *Semant Web J*. 2015; 0: 1–10.
47. Elbedweihy K, Wrigley SN, Ciravegna F, Reinhard D, Bernstein A. Evaluating semantic search systems to identify future directions of research. *The Semantic Web: ESWC 2012 Satellite Events*. Springer; 2012. pp. 148–162.
48. Brooke J. SUS-A quick and dirty usability scale. *Usability Eval Ind*. London; 1996; 189: 4–7.
49. Wilkinson MD, Dumontier M, Aalbersberg IJ, Appleton G, Axton M, Baak A, et al. The FAIR Guiding Principles for scientific data management and stewardship. *Sci Data*. 2016; 3. <https://doi.org/10.1038/sdata.2016.18> PMID: 26978244

50. Venkatesan A, Kim J-H, Talo F, Ide-Smith M, Gobeill J, Carter J, et al. SciLite: a platform for displaying text-mined annotations as a means to link research articles with biological data. *Wellcome Open Res.* 2016; 1: 25. <https://doi.org/10.12688/wellcomeopenres.10210.2> PMID: 28948232
51. IC4R Project Consortium, Hao L, Zhang H, Zhang Z, Hu S, Xue Y. Information Commons for Rice (IC4R). *Nucleic Acids Res.* 2016; 44: D1172–D1180. <https://doi.org/10.1093/nar/gkv1141> PMID: 26519466
52. Petryszak R, Keays M, Tang YA, Fonseca NA, Barrera E, Burdett T, et al. Expression Atlas update—An integrated database of gene and protein expression in humans, animals and plants. *Nucleic Acids Res.* 2016; 44: D746–D752. <https://doi.org/10.1093/nar/gkv1045> PMID: 26481351
53. Lee T, Oh T, Yang S, Shin J, Hwang S, Kim CY, et al. RiceNet v2: An improved network prioritization server for rice genes. *Nucleic Acids Res.* 2015; 43: W122–W127. <https://doi.org/10.1093/nar/gkv253> PMID: 25813048