

Agronomic Linked Data (AgroLD): a Knowledge-based System to Enable Integrative Biology in Agronomy

1 **Aravind Venkatesan^{1,3}, Gildas Tagny^{1,3}, Nordine El Hassouni^{1,2,8}, Imene Chentli^{1,3}, Valentin**
2 **Guignon^{4,8}, Clement Jonquet^{1,3}, Manuel Ruiz^{1,2,8,9} and Pierre Larmande^{*1,3,7,8}**

3 ¹ Institut de Biologie Computationnelle (IBC), Univ. of Montpellier, Montpellier, France.

4 ² UMR AGAP, CIRAD, Montpellier, France.

5 ³ LIRMM, Univ. of Montpellier & CNRS, Montpellier, France.

6 ⁴ Bioversity International, Montpellier, France.

7 ⁷ DIADE, IRD, Univ. of Montpellier, Montpellier, France.

8 ⁸ South Green Bioinformatics Platform, Montpellier, France.

9 ⁹ AGAP, Univ. of Montpellier, CIRAD, INRA, INRIA, SupAgro, Montpellier, France

10 *** Corresponding author**

11 **Email:** pierre.larmande@ird.fr

Abstract

13 Recent advances in high-throughput technologies have resulted in a tremendous increase in the
14 amount of omics data produced in plant science. This increase, in conjunction with the
15 heterogeneity and variability of the data, presents a major challenge to adopt an integrative
16 research approach. We are facing an urgent need to effectively integrate and assimilate
17 complementary datasets to understand the biological system as a whole. The Semantic Web
18 offers technologies for the integration of heterogeneous data and their transformation into explicit
19 knowledge thanks to ontologies. We have developed the Agronomic Linked Data (AgroLD –
20 www.agrold.org), a knowledge-based system relying on Semantic Web technologies and
21 exploiting standard domain ontologies, to integrate data about plant species of high interest for
22 the plant science community e.g., rice, wheat, arabidopsis. We present some integration results
23 of the project, which initially focused on genomics, proteomics and phenomics. AgroLD is now

24 an RDF (Resource Description Format) knowledge base of 100M triples created by annotating
25 and integrating more than 50 datasets coming from 10 data sources –such as Gramene.org and
26 TropGeneDB– with 10 ontologies –such as the Gene Ontology and Plant Trait Ontology. Our
27 evaluation results show users appreciate the multiple query modes which support different use
28 cases. AgroLD’s objective is to offer a domain specific knowledge platform to solve complex
29 biological and agronomical questions related to the implication of genes/proteins in, for
30 instances, plant disease resistance or high yield traits. We expect the resolution of these questions
31 to facilitate the formulation of new scientific hypotheses to be validated with a knowledge-
32 oriented approach.

33 **Introduction and Background**

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

44 Among the numerous research directions that the field of bioinformatics has taken, knowledge
45 management has become a major area of research, focused on logically interlinking information and
46 the representation of domain knowledge [1]. To this end, ontologies have become a cornerstone in the
47 representation of biological and more recently agronomical knowledge [2]. Ontologies provide the

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

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

79 In the recent years, the biomedical community has strongly embraced the Semantic Web vision as
80 demonstrated by a number of initiatives to provide ontologies [20,21] and use them for producing
81 semantically rich data such as in Bio2RDF [22] , OpenPHACTS [23], Linked Life Data [24], KUPKB
82 [25] , and the EBI RDF Platform [26]. In particular, OpenPHACTS serves as a good example of what
83 can be achieved by using Semantic Web knowledge bases. The OpenPHACTS Explorer
84 (<http://www.openphacts.org/open-phacts-discovery-platform/explorer>) provides use case driven tools
85 that aid in browsing and visualizing the underlying knowledge represented in RDF which is very
86 convenient for biologists.

87 Currently, there is a growing awareness within the agronomic domain towards efficient data
88 interoperability and integration [2,27,28]. The need for an umbrella approach for providing uniform
89 data is a widely-discussed topic. For instance, the Agriculture Data Interoperability Interest Group
90 (<https://rd-alliance.org/groups/agriculture-data-interest-group-igad.html>) instituted by the Research
91 Data Alliance (RDA) and agINFRA EU project (www.aginfra.eu) are initiatives that work on
92 improving data standards and promoting data interoperability in agriculture. Moreover, the community
93 has recently also started to adopt AgroPortal [11] as an vocabulary and ontology repository for
94 agronomy –and related domains such as nutrition, plant sciences and biodiversity– that support
95 browsing, searching and visualizing domain relevant ontologies, ontology alignments and creation of

96 semantic annotations. While plant-centric ontologies are now being used to annotate data by various
97 databases developers [2,5,28], unlike in the biomedical domain, the adoption of Semantic Web in
98 agronomy is yet to be completely exploited. Given that agronomic studies involve multiple domains,
99 publicly available knowledge bases such as EBI RDF, Linked Life Data and Bio2RDF serves only
100 limited agronomical information. Hence, it is necessary to build on previous efforts and complete them
101 to provide information compliant with Semantic Web principles within agronomic sciences. This
102 adoption would certainly allow the homogenization of multi-scale information, thereby aiding in the
103 discovery of new knowledge. Therefore, we have developed an RDF knowledge-based system, fully
104 compliant with the Semantic Web vision, called Agronomic Linked Data (AgroLD – www.agrold.org)
105 presented hereafter. The aim of our effort is to provide a portal (to discover) and an endpoint (to query)
106 for integrated agronomic information and to aid domain experts in answering relevant biological
107 questions.

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

112 **Materials and Methods**

113 **Information sources**

114 AgroLD was conceived to accommodate molecular and phenotypic information available on various
115 plant species (see Fig 1). The conceptual framework for the knowledge in AgroLD is based on well-
116 established ontologies: GO, SO, PO, Plant Trait Ontology (TO) and Plant Environment Ontology (EO).
117 Among these PO, TO and EO are currently developed by the Planteome project [29]
118 (<http://planteome.org>). Furthermore, considering the scope of the effort, we decided to build AgroLD
119 in phases. The current phase (phase I) covers information on genes, proteins, ontology associations,
120 homology predictions, metabolic pathways, plant traits, and germplasm, relevant to the selected
121 species. At this stage, we have incorporated the corresponding information from various databases,
122 such as Gramene [30], UniprotKB [31], Gene Ontology Annotation [32], TropGeneDB [33],
123 OryGenesDB [34], Oryza Tag Line [35], GreenPhylDB [36] and SNIPlay [37]. The selection of these
124 data sources was considered based on popularity among domain experts such as GOA, Gramene, and
125 complementary information hosted by the local research community, for instance, Oryza Tag Line and
126 GreenPhylDB. Information on the integrated databases can be found in the documentation page
127 (<http://www.agrold.org/documentation.jsp>). Table 1 provides a break-down of the data sources and
128 the species covered.

129 **Fig 1. Current plant species included in AgroLD.**

130 **Table 1. Plant species and data sources in AgroLD**

Data sources	URL s	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

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

135 Architecture

136 AgroLD relies on the RDF and SPARQL technologies for information modelling and retrieval. We use
 137 OpenLink Virtuoso (version 7.2) to store and access the RDF graphs. The data from the selected
 138 databases were parsed and converted into RDF using a semi-automated pipeline. The pipeline consists
 139 of several parsers to handle data in a variety of formats, such as the Gene Ontology Annotation File
 140 (GAF) [38], Generic File Format (GFF3) [39], HapMap [40] and Variant Call Format (VCF) [41]. Fig.
 141 2 shows the Extraction-Transform-Load (ETL) processes developed to transform in RDF various
 142 source data formats. The source code of the ETL workflow is available on GitHub¹².

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

148 For this phase, each dataset was downloaded from curated sources and was annotated with ontology
149 terms URIs by reusing the ontology fields when provided by the original source. Additionally, we used
150 the AgroPortal web service API to retrieve the URI corresponding to the taxon available for some data
151 standards such as GFF. At the end of phase 1, early 2018, the AgroLD knowledge base contains around
152 100 million RDF triples created by converting more than 50 datasets from 10 data sources.
153 Additionally, when available, we used some semantic annotation already present in the datasets such
154 as, for instances, genes or traits annotated respectively with GO or TO identifiers. In that case, we
155 produced additional properties with the corresponding ontologies thus adding 22% additional triples
156 validated manually (see details in Table 1). The OWL versions of the candidate ontologies were
157 directly loaded into the knowledge base but their triples are not counted in the total. We provided in
158 the supplementary file S1 Table, a more comprehensive statistics analysis such as number of triples,
159 classes, entities and properties for each graph stored in the knowledge base.

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

¹ <https://doi.org/10.5281/zenodo.1294660>

² <https://github.com/SouthGreenPlatform/AgroLD>

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

181 **Fig 3. Linking information in AgroLD.** The figure illustrates the linking of varies information in AgroLD.

182 To map the various data types and properties, we developed a lightweight schema (cf.
183 <https://github.com/SouthGreenPlatform/AgroLD>) that glues classes and properties identified in
184 AgroLD and the corresponding external ontologies. For instance, the class Protein
185 (<http://www.southgreen.fr/agrold/resource/Protein>) is mapped as *owl:equivalentClass* to class
186 polypeptide (http://purl.obolibrary.org/obo/SO_0000104) from SO. Similar mappings have been made
187 for properties, e.g., proteins/genes are linked to GO molecular function by the property
188 http://www.southgreen.fr/agrold/vocabulary/has_function, which is mapped as
189 *owl:equivalentProperty* to the corresponding Basic Formal Ontology (BFO) term
190 (http://purl.obolibrary.org/obo/BFO_0000085). When an equivalent property did not exist, we
191 mapped then to the closest upper level property using *rdfs:subPropertyOf* e.g., the property *has_trait*
192 (http://www.southgreen.fr/agrold/vocabulary/has_trait), links proteins to TO terms. It is mapped to a
193 more generic property, *causally related to* in the Relations Ontology [42]. For now, 55 mappings were
194 identified. Furthermore, mappings are both stored side by side with ontologies in AgroPortal, which
195 allows direct links between classes and instances of these classes in AgroLD. For example, the
196 following link will show the external mappings for SO:0000104 (polypeptide) stored in AgroPortal:
197 [http://agroportal.lirmm.fr/ontologies/SO/?p=classes&conceptid=http%3A%2F%2Fpurl.obolibrary.or](http://agroportal.lirmm.fr/ontologies/SO/?p=classes&conceptid=http%3A%2F%2Fpurl.obolibrary.org%2Fobo%2FSO_0000104&jump_to_nav=true#mappings)
198 [g%2Fobo%2FSO_0000104&jump_to_nav=true#mappings](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
199 resources (e.g., <http://www.southgreen.fr/agrold/page/biocyc.pathway/CALVIN-PWY>) are
200 dereferenced on a dedicated Pubby server [45]. For details on the graphs, URIs and properties, the
201 reader may refer to AgroLD’s documentation (<http://www.agrold.org/documentation.jsp>).

202 **User Interface**

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

- 204 • *Quick Search* (<http://www.agrold.org/quicksearch.jsp>), a faceted search plugin made available
205 by Virtuoso, that allows users to search by keywords and browse the AgroLD's content;
- 206 • *SPARQL Query Editor* (<http://www.agrold.org/sparqleditor.jsp>), that provides an interactive
207 environment to formulate SPARQL queries;
- 208 • *Explore Relationships* visualizer (<http://www.agrold.org/refinder.jsp>), which is an
209 implementation of RelFinder [46] that allows users to explore and visualize existing
210 relationships between entities;
- 211 • *Advanced Search* (<http://www.agrold.org/advancedSearch.jsp>), a query form providing entity
212 (e.g., gene) specific information retrieval.

213 Alternatively, some user management features have been implemented on the platform. Users have the
214 opportunity to save their search and results on a persistent history session attached to their own account.
215 Furthermore, they can manage search history by editing, deleting or re-running previous searches and
216 exporting results according several formats. In the future, we plan to develop some recommendation
217 features and sharing results between users. More detailed descriptions and figures of the different user
218 interfaces will be provided in the following section. Furthermore, other examples are shown in the User
219 Guide available in the supporting information S1 File.

220

221 **Results and Discussion**

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

230 We developed the SPARQL Query Editor based on the YASQE and YASR tools [47] and customized
231 it for our system. The SPARQL language is a powerful tool to mine and extract meaningful information
232 from the knowledge base. In the first example of the supplementary S3 file, we compare two queries
233 to answer the question: “Identify wheat proteins that are involved in root development.”. While the
234 first one (S3_Q1) using a simple search—which is a direct translation of SQL— with the corresponding
235 id (“GO_0048364”, “GO_2000280”) shows 73 entries, the second one (S3_Q2) using a property path
236 query (i.e., query the descending class hierarchy for a given trait ontology term) shows 137 entries,
237 thus more than 80% of additional results. In that case, the use of property path algorithm shows the
238 efficiency in retrieving a comprehensive answer. But the SPARQL language performs also very well
239 with complex queries such as: “Retrieve individuals which have positive SNP variant effect identified
240 for proteins associated with a QTL” available in S3_Q3. This type of query involves several datasets
241 and uses graph traversal property of SPARQL to perform the query.

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

244 For the comparison, we consider a sample question: ‘*Retrieving genes that participate in Calvin cycle*’;
245 (Q6 in the online list of modularized queries). As illustrated in Fig 4, the user can run the query to
246 retrieve the list of genes participating in the given pathway (Fig 4a). Additional information on a gene
247 of interest can be retrieved by clicking on the URI. For example, clicking on AT1G1870
248 (<http://identifiers.org/ensembl.plant/AT1G18270>) redirects the users to the gene information provided
249 by Gramene/Ensembl Plants resource (Fig 4b). The query can be saved and the results can be
250 downloaded in a variety of formats such as JSON, TSV, and RDF/XML. Additionally, user defined
251 queries could also be uploaded.

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

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

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

269 The *Advanced Search* query form is based on the REST API suite (<http://www.agrold.org/api-doc.jsp>),
270 developed completely within the AgroLD project. The aim of this feature is to provide non-technical
271 users with a tool to query the knowledge base while hiding the technical aspects of SPARQL query
272 formulation. Fig 6 illustrates steps involved in retrieving information for Q6, using the query form:

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

276 Furthermore, information on a gene of interest can be retrieved by selecting the specific gene (Fig 6b).
277 For instance, clicking on AT1GI870 (Fig 6c) displays all the proteins the gene encodes and the
278 pathways the gene participates in (apart from Calvin cycle). The RESTful API supports the query form
279 and was developed for programmatic retrieval of entity specific knowledge represented in AgroLD.
280 The current version of the API suite (ver. 1) can be used to retrieve gene and protein information,
281 metabolic pathways, and proteins associated with ontological terms. This is achieved by querying
282 entity by name or identifier.

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

288

289 **User Evaluation**

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

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

295 For the evaluation of semantic search systems, Elbedweihy et al. [48] recommend a survey of users
296 based on their experience with a few queries submitted to the system. We have used this approach to
297 collect user opinions, comments and suggestions via a feedback form directly within the AgroLD web
298 application. The form includes some questions from the "System Usability Scale" questionnaires [49]
299 and other questions that we considered important. The three main criteria evaluated are:

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

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

310 Globally, participants found the platform useful and easy to use. Overall, the idea of data navigation
311 and traversal through knowledge graphs was well received. However, many of them needed help with
312 some features. The general observation is that testing users ranked *Advanced Search* first then *Quick*
313 *Search* after. We explain this by the display output that looks friendlier for Advanced Search. *Quick*
314 *Search* won votes for usability and performance despite several comments to improve the ranking and
315 presentation of results (4 user's comments). *Advanced and Explore search* got average scores but good
316 comments on the capability of discovering unexpected results (e.g., nearest neighbour entities in the
317 graph for the Explore Search and additional results from external Web services for Advanced Search).
318 With no surprise, evaluation results show the *SPARQL Query Editor* is the most difficult to handle.
319 We mitigate this by offering examples of query pattern to help users handle query formulation. In the
320 future, we will improve the examples by offering a large spectrum of search type which will follow the
321 new phase of data integration. Furthermore, we will provide links to some SPARQL tutorials in the
322 documentation. These user feedbacks reinforced the need for knowledge bases such as AgroLD,
323 wherein users could retrieve information across various data types and sources. This knowledge
324 discovery is supported by the use of shared URI schemes and domain ontologies. The testing sessions
325 also helped us to identify areas for further improvement. Plus, we received suggestions on improving
326 the AgroLD's coverage with more data types such as gene expression data, and protein-protein
327 interactions. Considering, linked data and Semantic Web are still not widely adopted in agronomy,
328 increasing AgroLD's coverage will be an incremental process engaging our user community. This
329 situation is expected to improve with new community efforts such as the Agrisemantics RDA Working
330 Group (<https://rd-alliance.org/groups/agrisemantics-wg.html>), which role is to reinforce the adoption
331 of semantic technologies in the agri-food domain. We may also mention the AgBioData consortium
332 (<https://www.agbiodata.org>, [2]) which promotes the FAIR (Findable, Accessible, Interoperable and
333 Reusable) data principles [50] within agricultural research.

334 Furthermore, we observed that although the information integrated in AgroLD came from curated
335 sources, scientists often prefer to validate these knowledge statements against assertions made in
336 scientific articles. Currently, we have implemented an external Web Services as part of the *Advanced*
337 *Search Form* to automatically search for publications related to a protein or gene of interest in PubMed
338 Central and aggregates them within the result of the AgroLD query. However, this feature does not
339 provide detailed (sentence level) assertions described in those publications. This is an area that requires
340 further work. With the recent developments towards making text mined (sentence level) annotations
341 available as RDF [51], query federation can be explored to retrieve entity specific assertions. This
342 would serve as an additional provenance layer.

343 **Limits and Perspectives**

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

- 348 • The coverage content has to be extended to a larger number of biological entities (e.g., miRNA,
349 mRNA) or interaction between them (e.g., co-expression, regulation and interaction networks)
350 in order to capture a broad view of the molecular interactions.
- 351 • We have observed many information remains hidden in RDF literal contents such as biological
352 entities or relationship between them. This information is poorly annotated (i.e., plain text not
353 formally expressed) and new research methods to identify biological entities and reconstruct
354 their relations further allowing the discovery of relevant links between related resources are
355 required.

- 356 • The explosion of data in agronomy forces database providers to augment the frequency of their
357 releases. The survey shows a growing interest of using up to date information from the original
358 sources. This have to be taken into account for the updating process in AgroLD.
- 359 • The user interfaces show some limitations to manage responses with large number of results,
360 e.g., to filter and rank them with precision score.
- 361 These limitations identified in the current version of AgroLD will be improved in the following
362 versions. We will focus on the following areas:
- 363 • User Interface: we plan to explore features offered by Elastic search tool
364 (<https://www.elastic.co>), to enabling *Quick Search* retrieving more textual information and
365 hiding the technical details. Further, we will improve the performance and expand the API suite
366 to cover other entities represented in AgroLD (e.g., genomic annotation and homology
367 information).
- 368 • Content: integrate information on gene expression such as IC4R [52], Gene Expression Atlas
369 [53], on gene regulatory networks such as RiceNetDB [54] and explore linking text-mined
370 annotations from publications. Support molecular interaction networks per species and also
371 allow knowledge transfer between species.
- 372 • Knowledge discovery: explore methods to aid generating hypotheses by retrieving implicit
373 knowledge, e.g., inference rules, automatic data linking, entity recognition, text mining,
374 automatic semantic annotations.
- 375 • Data provenance: develop a provenance and annotation model. Set up a validation process to
376 allow users validating computed facts such as semantic annotations automatically produced and
377 attached to a biological entity.

- 378 • Updates: To keep AgroLD updated with the latest available data, by processing regular data
379 updates and potentially re-building the entire repository from scratch every 12 months.³
380 Additionally, we plan to fully automate the current ETL workflow.

381 **Conclusion**

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

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

396 AgroLD is one of the first RDF linked open data knowledge-based system in the agronomic domain.
397 It demonstrates a first step toward adopting the Semantic Web technologies to facilitate research by
398 integrating numerous heterogeneous data and transforming them into explicitly knowledge thanks to

³ Processing regular data update is a hard issue has 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.

399 ontologies. We expect AgroLD will facilitate the formulation of new scientific hypotheses to be
400 validated with its knowledge-oriented approach.

401 **Funding**

402 This research was supported by the Computational Biology Institute of Montpellier (ANR-11-BINF-
403 0002), the Institut Francais de Bioinformatique (ANR-11-INBS-0013), the Labex Agro (ANR-10-
404 LABX-001-01) all bypass of the French ANR *Investissements d'Avenir* program.

405 **Authors' contributions**

406 AV designed and implemented the AgroLD project and wrote the manuscript. GT designed and
407 implemented the API and the website. NEH contributed to the integration of data and set up of the
408 RDF store. IC tested and formulated biological queries. VG contributed to the integration of data. CJ
409 reviewed the manuscript. MR helped conceive the AgroLD project and reviewed the manuscript. PL
410 conceived, designed, implemented the AgroLD project and wrote the manuscript. All the authors
411 approved the final manuscript.

412 **Acknowledgments**

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

418 **References**

- 419 1. Goble C., and Stevens R. State of the nation in data integration for bioinformatics. *J Biomed*
420 *Inform.* Elsevier; 2008;41: 687–693. doi:10.1016/j.jbi.2008.01.008
- 421 2. Harper L., Campbell J., Cannon E.K., Jung S., Main D., Poelchau M., Walls R., Andorf C.,
422 Arnaud E., Berardini T., Birkett C., Cannon S., Carson J., Condon B., Cooper L., Dunn N.,

- 423 Farmer A., Ficklin S., Grant D., et al. AgBioData Consortium Recommendations for
424 Sustainable Genomics and Genetics Databases for Agriculture. Database. 2018; 1–7.
- 425 3. Ashburner M., Ball C.A., Blake J.A., Botstein D., Butler H., Cherry J.M., Davis A.P., Dolinski
426 K., Dwight S.S., Eppig J.T., Harris M.A., Hill D.P., Issel-Tarver L., Kasarskis A., Lewis S.,
427 Matese J.C., Richardson J.E., Ringwald M., Rubin G.M., et al. Gene ontology: tool for the
428 unification of biology. The Gene Ontology Consortium. *Nat Genet.* 2000;25: 25–29.
429 doi:10.1038/75556
- 430 4. Cooper L., Walls R.L., Elser J., Gandolfo M.A., Stevenson D.W., Smith B., Preece J., Athreya
431 B., Mungall C.J., Rensing S., Hiss M., Lang D., Reski R., Berardini T.Z., Li D., Huala E.,
432 Schaeffer M., Menda N., Arnaud E., et al. The plant ontology as a tool for comparative plant
433 anatomy and genomic analyses. *Plant Cell Physiol.* 2013;54: e1. doi:10.1093/pcp/pcs163
- 434 5. Shrestha R., Matteis L., Skofic M., Portugal A., McLaren G., Hyman G., and Arnaud E.
435 Bridging the phenotypic and genetic data useful for integrated breeding through a data
436 annotation using the Crop Ontology developed by the crop communities of practice. *Front*
437 *Physiol.* 2012;3: 326. doi:10.3389/fphys.2012.00326
- 438 6. Buttigieg P.L., Morrison N., Smith B., Mungall C.J., Lewis S.E., and ENVO Consortium. The
439 environment ontology: contextualising biological and biomedical entities. *J Biomed*
440 *Semantics.* 2013;4: 43. doi:10.1186/2041-1480-4-43
- 441 7. Walls R.L., Deck J., Guralnick R., Baskauf S., Beaman R., Blum S., Bowers S., Buttigieg
442 P.L., Davies N., Endresen D., Gandolfo M.A., Hanner R., Janning A., Krishtalka L.,
443 Matsunaga A., Midford P., Morrison N., Ó Tuama É., Schildhauer M., et al. Semantics in
444 support of biodiversity knowledge discovery: an introduction to the biological collections
445 ontology and related ontologies. *PLoS One.* 2014;9: e89606.
446 doi:10.1371/journal.pone.0089606
- 447 8. Oellrich A., Walls R.L., Cannon E.K., Cannon S.B., Cooper L., Gardiner J., Gkoutos G. V,
448 Harper L., He M., Hoehndorf R., Jaiswal P., Kalberer S.R., Lloyd J.P., Meinke D., Menda N.,
449 Moore L., Nelson R.T., Pujar A., Lawrence C.J., et al. An ontology approach to comparative
450 phenomics in plants. *Plant Methods.* 2015;11: 10. doi:10.1186/s13007-015-0053-y
- 451 9. Wang Y., Wang Y., Wang J., Yuan Y., and Zhang Z. An ontology-based approach to
452 integration of hilly citrus production knowledge. *Comput Electron Agric.* Elsevier; 2015;113:
453 24–43. doi:10.1016/J.COMPAG.2015.01.009
- 454 10. Lousteau-Cazalet C., Barakat A., Belaud J.-P., Buche P., Busset G., Charnomordic B.,
455 Dervaux S., Destercke S., Dibie J., Sablayrolles C., and Vialle C. A decision support system
456 for eco-efficient biorefinery process comparison using a semantic approach. *Comput Electron*
457 *Agric.* Elsevier; 2016;127: 351–367. doi:10.1016/J.COMPAG.2016.06.020
- 458 11. Jonquet C., Toulet A., Arnaud E., Aubin S., Dzalé Yeumo E., Emonet V., Graybeal J., Laporte
459 M.A., Musen M.A., Pesce V., and Larmande P. AgroPortal: A vocabulary and ontology
460 repository for agronomy. *Comput Electron Agric.* 2018;144: 126–143.
461 doi:10.1016/j.compag.2017.10.012
- 462 12. Berners-lee T., Hendler J., and Lassila O. The Semantic Web. *Sci Am.* 2001;284: 35–43.

- 463 13. W3C. Resource Description Framework (RDF): Concepts and Abstract Syntax [Internet].
464 [cited 3 Apr 2010]. Available: <http://www.w3.org/TR/2004/REC-rdf-concepts-20040210/>
- 465 14. W3C. RDF Schema 1.1 [Internet]. [cited 27 Apr 2018]. Available:
466 <https://www.w3.org/TR/rdf-schema/>
- 467 15. W3C. OWL 2 Web Ontology Language Structural Specification and Functional-Style Syntax
468 [Internet]. [cited 3 Apr 2010]. Available: [http://www.w3.org/TR/2009/REC-owl2-syntax-
469 20091027/](http://www.w3.org/TR/2009/REC-owl2-syntax-20091027/)
- 470 16. The W3C SPARQL Working Group. SPARQL 1.1 Overview [Internet]. [cited 15 Apr 2013].
471 Available: <http://www.w3.org/TR/sparql11-overview/>
- 472 17. Luciano J.S., Andersson B., Batchelor C., Bodenreider O., Clark T., Denney C.K., Domarew
473 C., Gambet T., Harland L., Jentsch A., Kashyap V., Kos P., Kozlovsky J., Lebo T., Marshall
474 S.M., McCusker J.P., McGuinness D.L., Ogbuji C., Pichler E., et al. The Translational
475 Medicine Ontology and Knowledge Base: driving personalized medicine by bridging the gap
476 between bench and bedside. *J Biomed Semantics*. 2011;2 Suppl 2: S1. doi:10.1186/2041-
477 1480-2-S2-S1
- 478 18. Venkatesan A., Tripathi S., Sanz de Galdeano A., Blondé W., Lægreid A., Mironov V., and
479 Kuiper M. Finding gene regulatory network candidates using the gene expression knowledge
480 base. *BMC Bioinformatics*. 2014;15: 386. doi:10.1186/s12859-014-0386-y
- 481 19. Smith R.N., Aleksic J., Butano D., Carr A., Contrino S., Hu F., Lyne M., Lyne R., Kalderimis
482 A., Rutherford K., Stepan R., Sullivan J., Wakeling M., Watkins X., and Micklem G.
483 InterMine: a flexible data warehouse system for the integration and analysis of heterogeneous
484 biological data. *Bioinformatics*. Oxford University Press; 2012;28: 3163–5.
485 doi:10.1093/bioinformatics/bts577
- 486 20. Smith B., Ashburner M., Rosse C., Bard J., Bug W., Ceusters W., Goldberg L.J., Eilbeck K.,
487 Ireland A., Mungall C.J., Leontis N., Rocca-Serra P., Ruttenberg A., Sansone S.-A.,
488 Scheuermann R.H., Shah N., Whetzel P.L., Lewis S., and Lewis S. The OBO Foundry:
489 coordinated evolution of ontologies to support biomedical data integration. *Nat Biotechnol*.
490 Nature Publishing Group; 2007;25: 1251–1255. doi:10.1038/nbt1346
- 491 21. Noy N.F., Shah N.H., Whetzel P.L., Dai B., Dorf M., Griffith N., Jonquet C., Rubin D.L.,
492 Storey M.-A., Chute C.G., and Musen M.A. BioPortal: ontologies and integrated data
493 resources at the click of a mouse. *Nucleic Acids Res*. 2009;37: W170-173.
494 doi:10.1093/nar/gkp440
- 495 22. Belleau F., Nolin M.-A., Tourigny N., Rigault P., and Morissette J. Bio2RDF: towards a
496 mashup to build bioinformatics knowledge systems. *J Biomed Inform*. Elsevier; 2008;41: 706–
497 716. doi:10.1016/j.jbi.2008.03.004
- 498 23. Williams A.J., Harland L., Groth P., Pettifer S., Chichester C., Willighagen E.L., Evelo C.T.,
499 Blomberg N., Ecker G., Goble C., and Mons B. Open PHACTS: Semantic interoperability for
500 drug discovery. *Drug Discovery Today*. 2012. pp. 1188–1198.
501 doi:10.1016/j.drudis.2012.05.016

- 502 24. Momtchev V., Peychev D., Primov T., and Georgiev G. Expanding the Pathway and
503 Interaction Knowledge in Linked Life Data. International Semantic Web Challenge. 2009.
- 504 25. Jupp S., Klein J., Schanstra J., and Stevens R. Developing a kidney and urinary pathway
505 knowledge base. J Biomed Semantics. 2011;2 Suppl 2: S7. doi:10.1186/2041-1480-2-S2-S7
- 506 26. Jupp S., Malone J., Bolleman J., Brandizi M., Davies M., Garcia L., Gaulton A., Gehant S.,
507 Laibe C., Redaschi N., Wimalaratne S.M., Martin M., Le Novère N., Parkinson H., Birney E.,
508 and Jenkinson A.M. The EBI RDF platform: linked open data for the life sciences.
509 Bioinformatics. 2014; 1–2. doi:10.1093/bioinformatics/btt765
- 510 27. Venkatesan A., El Hassouni N., Phillipe F., Pommier C., Quesneville H., Ruiz M., and
511 Larmande P. Towards efficient data integration and knowledge management in the Agronomic
512 domain. APIA'15: premiere Conference Applications Pratiques de l'Intelligence Artificielle.
513 2015.
- 514 28. Leonelli S., Davey R.P., Arnaud E., Parry G., and Bastow R. Data management and best
515 practice for plant science. Nat Publ Gr. Macmillan Publishers Limited; 2017;3: 1–4.
516 doi:10.1038/nplants.2017.86
- 517 29. Cooper L., Meier A., Laporte M.A., Elser J.L., Mungall C., Sinn B.T., Cavaliere D., Carbon
518 S., Dunn N.A., Smith B., Qu B., Preece J., Zhang E., Todorovic S., Gkoutos G., Doonan J.H.,
519 Stevenson D.W., Arnaud E., and Jaiswal P. The Planteome database: An integrated resource
520 for reference ontologies, plant genomics and phenomics. Nucleic Acids Res. 2018;
521 doi:10.1093/nar/gkx1152
- 522 30. Monaco M.K., Stein J., Naithani S., Wei S., Dharmawardhana P., Kumari S., Amarasinghe V.,
523 Youens-Clark K., Thomason J., Preece J., Pasternak S., Olson A., Jiao Y., Lu Z., Bolser D.,
524 Kerhornou A., Staines D., Walts B., Wu G., et al. Gramene 2013: Comparative plant genomics
525 resources. Nucleic Acids Res. 2014;42. doi:10.1093/nar/gkt1110
- 526 31. Magrane M., and Consortium U.P. UniProt Knowledgebase: A hub of integrated protein data.
527 Database. 2011;2011. doi:10.1093/database/bar009
- 528 32. Barrell D., Dimmer E., Huntley R.P., Binns D., O'Donovan C., and Apweiler R. The GOA
529 database in 2009 - An integrated Gene Ontology Annotation resource. Nucleic Acids Res.
530 2009;37. doi:10.1093/nar/gkn803
- 531 33. Hamelin C., Sempere G., Jouffe V., and Ruiz M. TropGeneDB, the multi-tropical crop
532 information system updated and extended. Nucleic Acids Res. 2013;41.
533 doi:10.1093/nar/gks1105
- 534 34. Droc G., Ruiz M., Larmande P., Pereira A., Piffanelli P., Morel J.B., Dievart A., Courtois B.,
535 Guiderdoni E., and Périn C. OryGenesDB: a database for rice reverse genetics. Nucleic Acids
536 Res. 2006;34: D736-40. doi:10.1093/nar/gkj012
- 537 35. Larmande P., Gay C., Lorieux M., Périn C., Bouniol M., Droc G., Sallaud C., Perez P.,
538 Barnola I., Biderre-petit C., Martin J., Morel J.B., Johnson A.A.T., Bourgis F., Ghesquière A.,
539 Ruiz M., Courtois B., and Guiderdoni E. Oryza Tag Line, a phenotypic mutant database for
540 the Génoplatte rice insertion line library. Nucleic Acids Res. 2008;36: 1022–1027.

- 541 doi:10.1093/nar/gkm762
- 542 36. Conte M.G., Gaillard S., Lanau N., Rouard M., and Périn C. GreenPhylDB: a database for
543 plant comparative genomics. *Nucleic Acids Res.* 2008;36: D991-998. doi:10.1093/nar/gkm934
- 544 37. Dereeper A., Homa F., Andres G., Sempere G., Sarah G., Hueber Y., Dufayard J.-F., and Ruiz
545 M. SNIPlay3: a web-based application for exploration and large scale analyses of genomic
546 variations. *Nucleic Acids Res.* 2015;43: W295-300. doi:10.1093/nar/gkv351
- 547 38. The Gene Ontology Consortium. Gene Annotation File (GAF) specification [Internet]. [cited
548 20 Mar 2018]. Available: <http://geneontology.org/page/go-annotation-file-format-20>
- 549 39. Sequence Ontology consortium. GFF3 Specification [Internet].
- 550 40. Gibbs R.A., Belmont J.W., Hardenbol P., Willis T.D., Yu F., Zhang H., Zeng C., Matsuda I.,
551 Fukushima Y., Macer D.R., Suda E., Stein L.D., Cunningham F., Kanani A., Thorisson G.A.,
552 Chakravarti A., Chen P.E., Cutler D.J., Kashuk C.S., et al. The International HapMap Project.
553 *Nature.* 2003;426: 789–796. doi:10.1038/nature02168
- 554 41. Danecek P., Auton A., Abecasis G., Albers C.A., Banks E., DePristo M.A., Handsaker R.E.,
555 Lunter G., Marth G.T., Sherry S.T., McVean G., and Durbin R. The variant call format and
556 VCFtools. *Bioinformatics.* 2011;27: 2156–8. doi:10.1093/bioinformatics/btr330
- 557 42. Juty N., Le Novère N., and Laibe C. Identifiers.org and MIRIAM Registry: community
558 resources to provide persistent identification. *Nucleic Acids Res.* 2012;40: D580-6.
559 doi:10.1093/nar/gkr1097
- 560 43. Manel A., Zohra B., and Konstantin T. A survey on web data linking. *Ingénierie des systèmes*
561 *d'information.* 2016;21: 11–29. doi:10.3166/isi.21.5-6.11-29
- 562 44. Smith B., Ceusters W., Klagges B., Kohler J., Kumar A., Lomax J., Mungall C., Neuhaus F.,
563 Rector A., and Rosse C. Relations in biomedical ontologies. *Genome Biol.* 2005;6: R46.
- 564 45. Cyganiak R. (National U. of I., and Bizer C. Pubby - A Linked Data Frontend for SPARQL
565 Endpoints. 2008; Available: <http://wifo5-03.informatik.uni-mannheim.de/pubby/>
- 566 46. Heim P., Hellmann S., Lehmann J., Lohmann S., and Stegemann T. RelFinder: Revealing
567 relationships in RDF knowledge bases. *Lecture Notes in Computer Science (including*
568 *subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics).* 2009.
569 pp. 182–187. doi:10.1007/978-3-642-10543-2_21
- 570 47. Rietveld L., and Hoekstra R. The YASGUI Family of SPARQL Clients. *Semant Web J.*
571 2015;0: 1–10.
- 572 48. Elbedweihy K., Wrigley S.N., Ciravegna F., Reinhard D., and Bernstein A. Evaluating
573 semantic search systems to identify future directions of research. *The Semantic Web: ESWC*
574 *2012 Satellite Events.* Springer; 2012. pp. 148–162.
- 575 49. Brooke J. SUS-A quick and dirty usability scale. *Usability Eval Ind.* London; 1996;189: 4–7.
- 576 50. Wilkinson M.D., Dumontier M., Aalbersberg I.J., Appleton G., Axton M., Baak A., Blomberg

- 577 N., Boiten J.-W., da Silva Santos L.B., Bourne P.E., Bouwman J., Brookes A.J., Clark T.,
578 Crosas M., Dillo I., Dumon O., Edmunds S., Evelo C.T., Finkers R., et al. The FAIR Guiding
579 Principles for scientific data management and stewardship. *Sci Data*. 2016;3.
580 doi:10.1038/sdata.2016.18
- 581 51. Venkatesan A., Kim J.-H., Talo F., Ide-Smith M., Gobeill J., Carter J., Batista-Navarro R.,
582 Ananiadou S., Ruch P., McEntyre J., Venkatesan A., Kim J.-H., Talo F., Ide-Smith M.,
583 Gobeill J., Carter J., Batista-Navarro R., Ananiadou S., Ruch P., et al. SciLite: a platform for
584 displaying text-mined annotations as a means to link research articles with biological data.
585 *Wellcome Open Res*. 2016;1: 25. doi:10.12688/wellcomeopenres.10210.1
- 586 52. IC4R Project Consortium, Hao L., Zhang H., Zhang Z., Hu S., and Xue Y. Information
587 Commons for Rice (IC4R). *Nucleic Acids Res*. 2016;44: D1172–D1180.
588 doi:10.1093/nar/gkv1141
- 589 53. Petryszak R., Keays M., Tang Y.A., Fonseca N.A., Barrera E., Burdett T., Füllgrabe A.,
590 Fuentes A.M.P., Jupp S., Koskinen S., Mannion O., Huerta L., Megy K., Snow C., Williams
591 E., Barzine M., Hastings E., Weisser H., Wright J., et al. Expression Atlas update - An
592 integrated database of gene and protein expression in humans, animals and plants. *Nucleic
593 Acids Res*. 2016;44: D746–D752. doi:10.1093/nar/gkv1045
- 594 54. Lee T., Oh T., Yang S., Shin J., Hwang S., Kim C.Y., Kim H., Shim H., Shim J.E., Ronald
595 P.C., and Lee I. RiceNet v2: An improved network prioritization server for rice genes. *Nucleic
596 Acids Res*. 2015;43: W122–W127. doi:10.1093/nar/gkv253

597

598 **Supporting information**

599 **S1 File. AgroLD User Guide.** This document shows how to use the various features of the platform.

600 **S1 Table. AgroLD graph statistics.**

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

602 **S3 File. Examples of SPARQL queries.** Example of SPARQL queries showing the benefits of
603 property path algorithm, and complex queries.

604

Fig 1

[Click here to download Figure Fig1.eps](#)

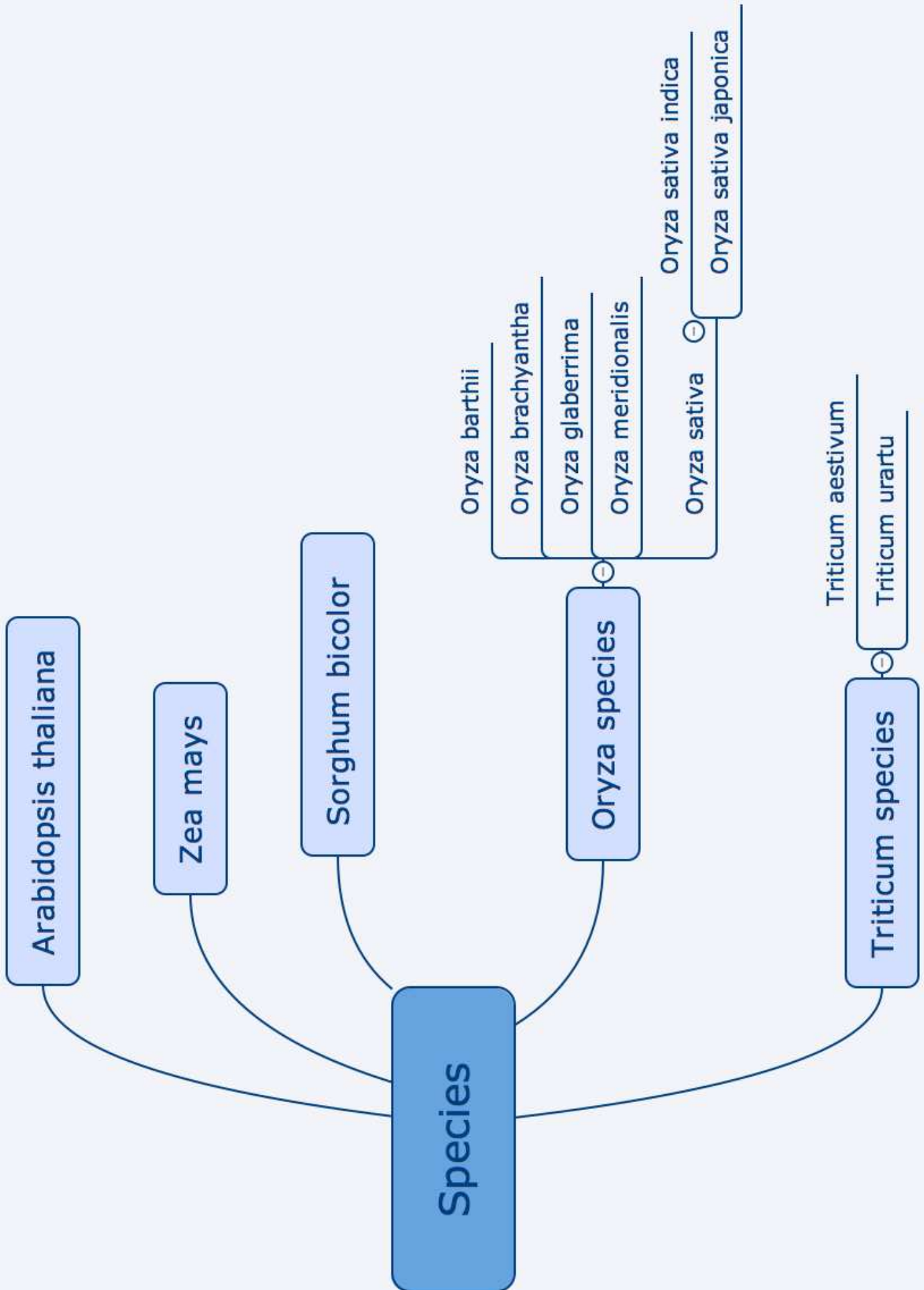


Fig 2

[Click here to download Figure Fig2.eps](#)

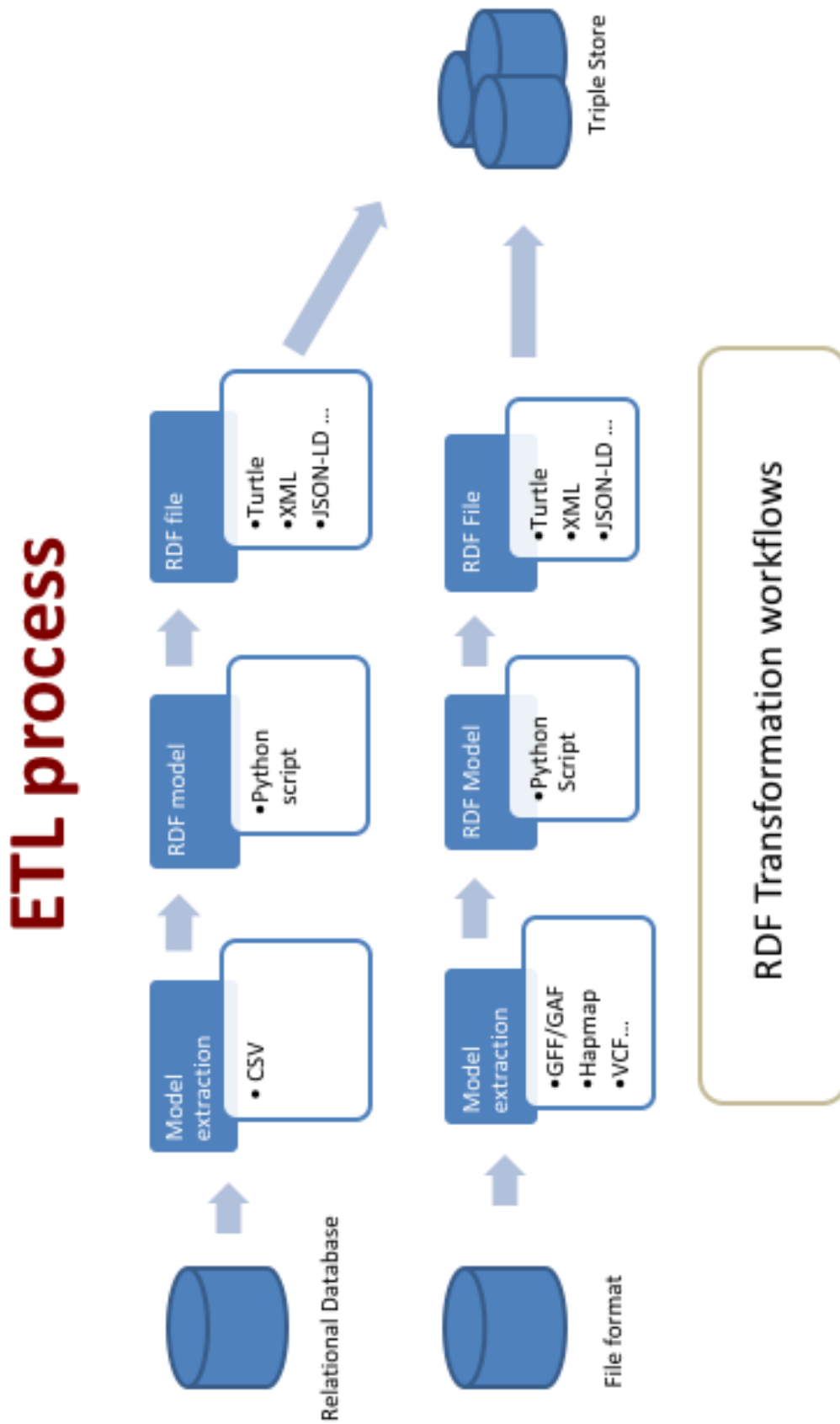


Fig 3

[Click here to download Figure Fig3.eps](#)

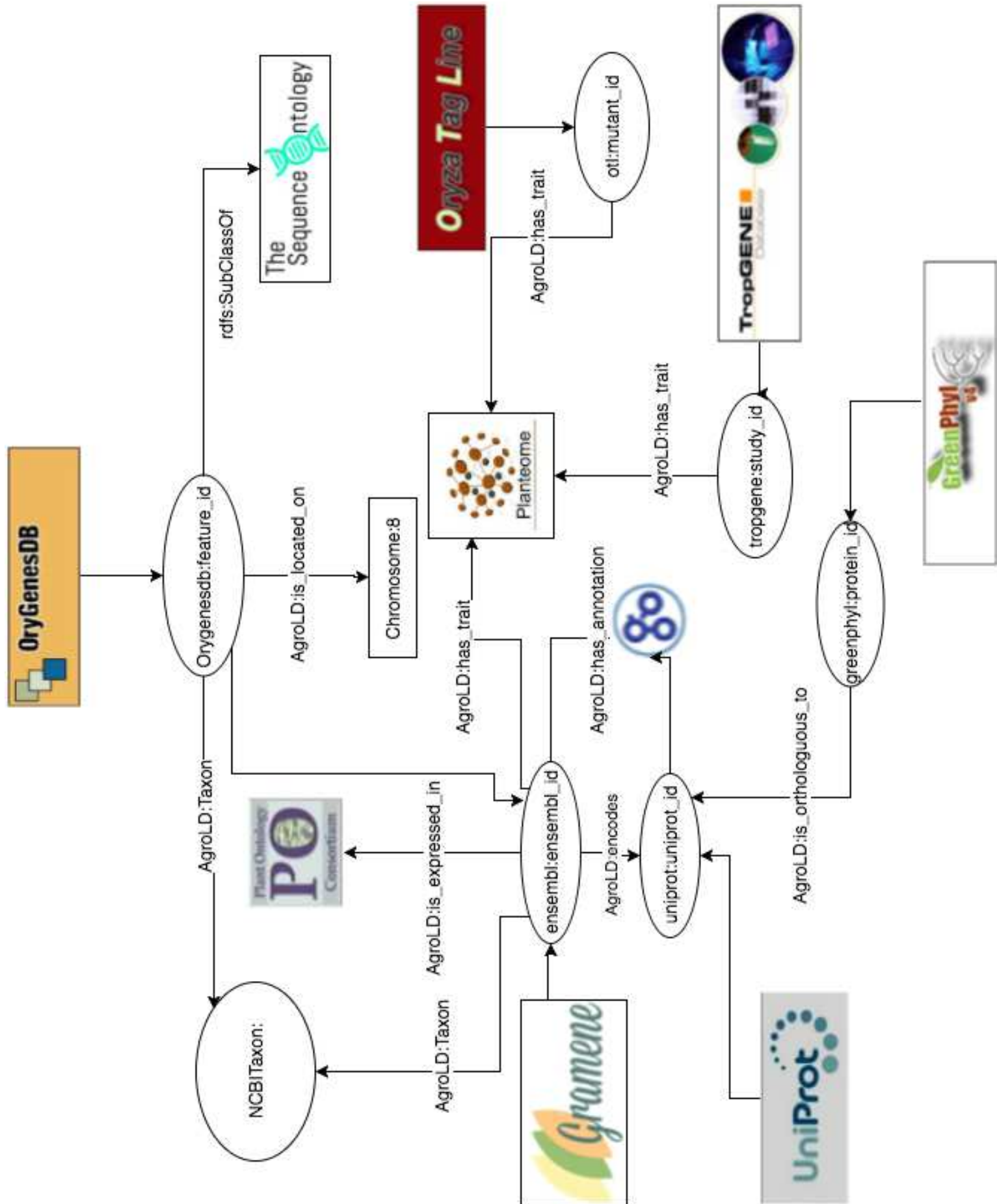


Fig 4

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.

KEYBOARD COMMANDS

Results Format RDF/XML

Download Results

```

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)

Filename to Save As: query.sparql

Choose File

No file chosen

Load Selected Query File

a)

Watch how!

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: 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

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

EnsemblPlants

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

Gene: AT1G18270

ketose-bisphosphate aldolase class-II family protein [Source:TAIR;Acc:AT1G18270]
Chromosome 1: 6,283,412-5,293,871, reverse strand
This gene has 3 transcripts (splice variants), 37 orthologues and 6 paralogues.

Summary

Show transcript table

Splice variants

Transcript comparison
Supporting evidence
Gene alleles
Sequence
Secondary Structure
Gene families
External references

Description

ketose-bisphosphate aldolase class-II family protein [Source:TAIR;Acc:AT1G18270]
Chromosome 1: 6,283,412-5,293,871, reverse strand
This gene has 3 transcripts (splice variants), 37 orthologues and 6 paralogues.

Location

Chromosome 1: 6,283,412-5,293,871, reverse strand
This gene has 3 transcripts (splice variants), 37 orthologues and 6 paralogues.

About this gene

Transcripts

Transcripts

Show transcript table

Raw Response

Table
Pivot Table

Search

Show 50
entries

Login/Regis...

Search Ensembl Plants...

b)

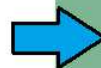


Fig 5

[Click here to download Figure Fig5.eps](#)

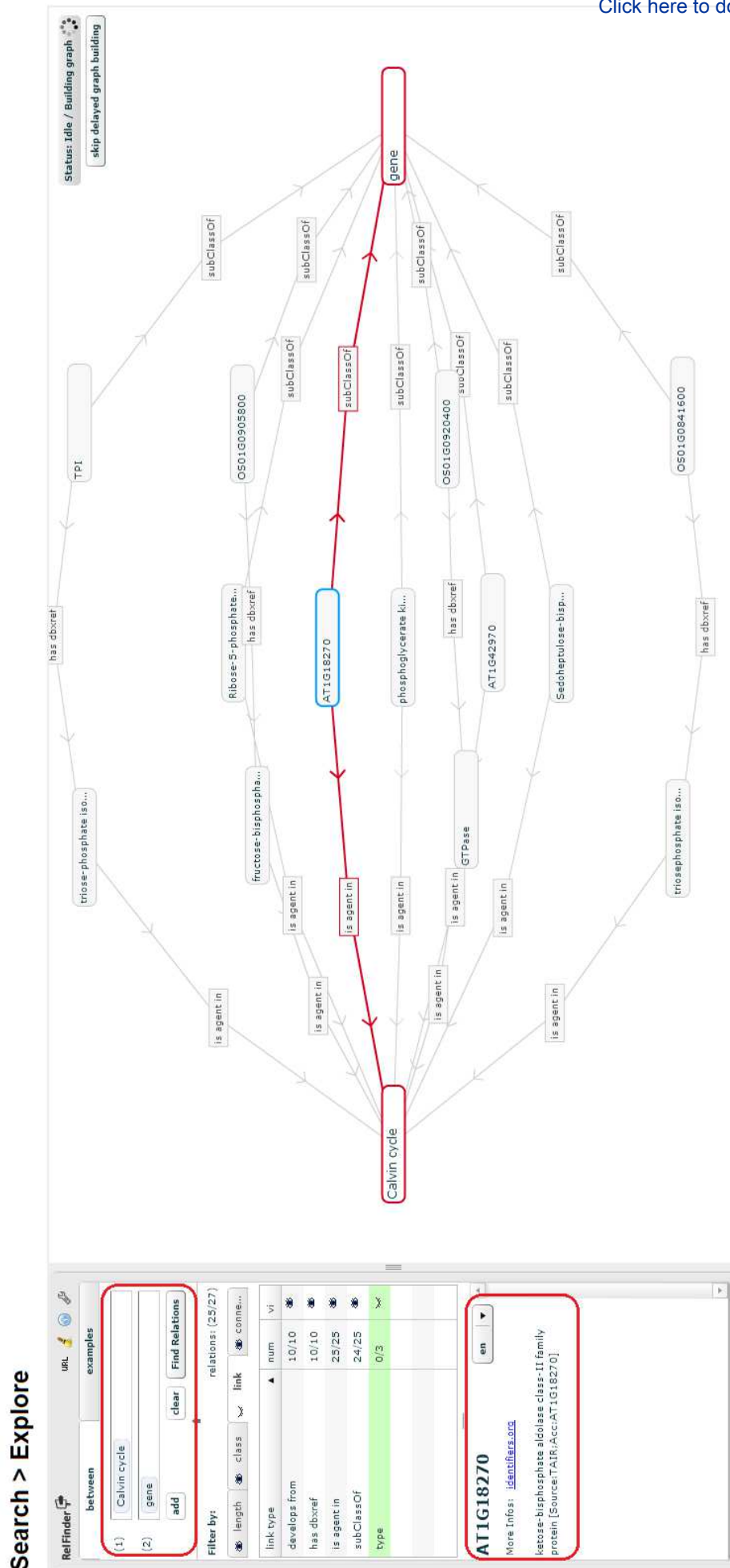


Fig 6

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'

a)

Search pathway with keyword "Calvin cycle"

Pathway

Search: Show entries

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

Showing 1 to 1 of 1 entries



PATHWAY : CALVIN-PWY / Calvin cycle

URI: <http://www.southgreen.fr/agroid/biocyc/pathway/CALVIN-PWY>

b)

Participating genes

geneid	gene_name	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)	



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](#) ±

Click here to download Figure Fig6.eps

c)