



HAL
open science

AgroLD API. Une architecture orientée services pour l'extraction de connaissances dans la base de données liées AgroLD

Gildas Tagny Ngompe, Aravind Venkatesan, Nordine El Hassouni, Manuel Ruiz, Pierre Larmande

► To cite this version:

Gildas Tagny Ngompe, Aravind Venkatesan, Nordine El Hassouni, Manuel Ruiz, Pierre Larmande. AgroLD API. Une architecture orientée services pour l'extraction de connaissances dans la base de données liées AgroLD. *Revue des Sciences et Technologies de l'Information - Série ISI: Ingénierie des Systèmes d'Information*, 2016, 21 (5-6), pp.133-158. 10.3166/isi.21.5-6.133-157 . hal-01411532

HAL Id: hal-01411532

<https://hal.science/hal-01411532>

Submitted on 7 Dec 2016

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.

AgroLD API : Une architecture orientée services pour l'extraction de connaissances dans la base de données liées AgroLD

Gildas Tagny Ngompe¹, Aravind Venkatesan¹,
Nordine El Hassouni^{1,2}, Manuel Ruiz^{1,2}, Pierre Larmande^{1,3}

1. Institut de Biologie Computationnelle (IBC)
Université de Montpellier, 860 rue St Priest, 34095 Montpellier Cedex 5, France
{aravindvenkatesan}{tagnyngompe}@gmail.com
2. AGAP, Plateforme SouthGreen, CIRAD, Avenue Agropolis, 34398 Montpellier,
France
{manuel.ruiz}{nordine.el_hassouni}@cirad.fr
3. UMR DIADE, Plateforme SouthGreen, IRD, 911 Avenue Agropolis, 34394
Montpellier, France
pierre.larmande@ird.fr

RÉSUMÉ. L'agronomie est un domaine important constitué de divers domaines de recherche tels que la génétique, la biologie moléculaire des plantes, l'écologie et les sciences de la terre. Les dernières décennies ont vu le succès du développement de technologies à haut débit qui ont révolutionnés et transformés la recherche agronomique. L'application de ces technologies a généré de grandes quantités de données et de ressources sur le Web. Dans la plupart des cas, ces sources restent autonomes et déconnectés. Le projet Agronomic Linked Data (AgroLD) est une base de connaissances du web sémantique conçu pour intégrer des données provenant de diverses sources de données centrées sur des plantes disponibles publiquement. L'objectif du projet AgroLD est de fournir un portail web pour les bioinformaticiens et les experts du domaine afin d'exploiter les données homogénéisées et permettre de combler les connaissances dans ce domaine.

ABSTRACT. Agronomy is an overarching field constituting various research areas such as genetics, plant molecular biology, ecology and earth science. The last several decades has seen the successful development of high-throughput technologies that have revolutionised and transformed agronomic research. The application of these technologies have generated large quantities of data and resources over the web. In most cases these sources remain autonomous and disconnected. The Agronomic Linked Data project (AgroLD) is a Semantic Web knowledge base designed to integrate data from various publicly available plant centric data sources. The aim of AgroLD project is to provide a portal for bioinformaticians and domain experts to exploit the

homogenized data towards enabling to bridge the knowledge.

MOTS-CLÉS : Biologie Moléculaire, Agronomie, Web Sémantique, Données Liées, RDF, SPARQL, Services Web RESTFul

KEYWORDS: Molecular Biology, Agronomy, Semantic Web, Linked Data, RDF, SPARQL, REST-Ful Web Services

DOI:10.3166/ISI.28.2-3.1-25 © 2014 Lavoisier

1. Introduction

L'agronomie est un champ multidisciplinaire qui comprend des domaines de recherche tels que la biologie moléculaire végétale, la physiologie et l'agroécologie. La recherche agronomique vise à répondre entre autres, à des questions sur l'amélioration de la production des cultures, leur résistance face aux maladies et l'étude de l'impact environnemental sur les cultures. Pour cela, les chercheurs ont besoin de comprendre en détail les implications des différents processus biologiques, par conséquent, de lier les données à différentes échelles (e.g. la génomique, la protéomique et le phéno). Nous assistons actuellement à des progrès rapides dans les technologies de production de données biologiques conduisant à un flot d'information dans les domaines mentionnés ci-dessus. Cependant, une grande partie de cette information est dispersée dans différentes bases de données spécifiques à un domaine et en outre, ces données sont mises à disposition dans des formats de données variées. Par conséquent, l'utilisation efficace de ces ressources et l'adoption d'une approche intégrée reste un défi majeur.

Parmi les nombreux axes de recherche que compte le domaine bioinformatique, la gestion des connaissances est devenue un domaine de recherche important (Goble, Stevens, 2008), axé sur l'interconnexion de l'information et la représentation des connaissances (Antezana *et al.*, 2009). Dans ce domaine, les ontologies sont devenues une pierre angulaire dans la représentation des connaissances biologiques. Elles fournissent la structure nécessaire pour représenter des concepts biologiques et leurs relations. Aujourd'hui, de nombreuses applications exploitent les avantages offerts par les ontologies biologiques telles que Gene Ontology (GO) (Ashburner *et al.*, 2000), Sequence Ontology (SO) (Mungall *et al.*, 2011) et Plant Ontology (PO) (Cooper *et al.*, 2013) pour n'en citer que quelques-unes.

De plus, la gestion efficace des connaissances exige l'adoption de méthodes robustes d'intégration de données. Cela implique une intégration efficace des sources de données hétérogènes distribuées, représentées dans un format interprétable par les machines. La technologie du Web sémantique (SW) proposée par Tim Berners-Lee (?), offre une solution pour faciliter cette intégration et permettre l'interopérabilité entre les machines.

Ces dernières années ont vu l'utilisation croissante du SW dans la communauté biomédicale, dont un certain nombre d'initiatives ont démontré le potentiel de ce der-

nier. Certaines initiatives notables incluent Bio2RDF (Belleau *et al.*, 2008) , OpenPHACTS (Williams *et al.*, 2012) , Life Data Linked (Momtchev *et al.*, 2009) , KUPKB (Jupp *et al.*, 2011) et EBI Plate-forme RDF (Jupp *et al.*, 2014). Pris dans son ensemble une application SW mis en œuvre avec succès permet aux scientifiques de poser des questions très complexes par le biais d'une requête ou un ensemble de requêtes qui retourneront des réponses très pertinentes (Luciano *et al.*, 2011), ce qui facilite la formulation d'hypothèses de recherche (Venkatesan *et al.*, 2014).

Aujourd'hui, il y a une prise de conscience croissante similaire dans le domaine agronomique. Bien que les ontologies agronomiques sont majoritairement utilisées pour annoter les données contenues dans diverses bases de données, à la différence du domaine bio-médical, l'utilisation du SW en agronomie doit encore être mieux exploité. Afin de contribuer à cet effort, nous avons mis au point une base de connaissances sémantique, nommée Agronomic Linked Data (AgroLD) (www.agrold.org). Le but de notre démarche est de fournir un portail d'information agronomique intégrée pour aider les experts du domaine à répondre aux questions biologiques pertinentes.

AgroLD est conçue pour intégrer des informations disponibles sur diverses espèces végétales. Le cadre conceptuel de la connaissance est basé sur des ontologies bien établies dans le domaine. En outre, compte tenu de la portée de l'effort, nous avons décidé de construire AgroLD en plusieurs phases. La phase actuelle (phase un) couvre les informations sur les gènes, les protéines, les prédictions de genes homologues, les voies métaboliques, des phenotypes de plantes et le matériel génétique. Des informations sur les bases de données intégrées peut être trouvées sur la page de documentation d'AgroLD ¹.

De manière standard, pour interroger une base de triplets RDF, il faut faire exécuter des requêtes SPARQL par des points d'accès SPARQL (« sparql endpoint »). Les points d'accès SPARQL ont l'avantage de fournir une API standard. Même si le langage SPARQL est assez efficace pour construire les requêtes (maximise l'expressivité des requêtes (Ferré, 2014)), il reste difficile à prendre en main vue sa large variété de fonctionnalités.

Dans cet article, nous proposons un modèle d'architecture comprenant des systèmes implémentant des paradigmes des systèmes de recherche sémantique (Haag *et al.*, 2014). Chacun de ces paradigmes a des atouts, mais aussi des limites que leur intégration peut combler. Le reste du papier est organisé comme suit: dans la section 2, nous passons en revue les travaux connexes par rapport à notre domaine d'intérêt; section 3 présente le modèle d'architecture et les détails de son implementation; section 4 examine nos cas d'utilisation et leurs résultats et la section 5 conclut sur la mise en place de l'architecture.

1. Documentation d'AgroLD: <http://volvestre.cirad.fr:8080/agrold/documentation.jsp>

2. Les approches existantes d'interrogation des Bases de Données de triplets RDF

Dans cette section, nous proposons un état de l'art des approches nous semblant contribuer à l'interrogation des bases de données de triplets. Pour aider à la recherche sémantique, certains outils sont souvent proposés soit pour assister les utilisateurs dans la construction de leurs requêtes, soit pour cacher le langage SPARQL et fournir une interface plus facile à utiliser pour l'interrogation des bases cibles.

2.1. Approches d'aide à la construction de requêtes

2.1.1. Les éditeurs syntaxiques de requêtes SPARQL

Les fournisseurs de données RDF proposent une interface web pour l'édition et l'exécution de requêtes SPARQL sur leurs données. Les systèmes de gestion de graphes RDF tels que OpenLink Virtuoso propose déjà un formulaire HTML simple pour indiquer les paramètres d'exécution au point d'accès SPARQL.

Certains publieurs de données améliorent l'utilisabilité de leurs clients SPARQL en intégrant des bibliothèques d'édition de code (texte de la requête) et de gestion des résultats. Par exemple, UniProt² et « linked life data »³ joignent une grammaire de SPARQL à la librairie CodeMirror⁴ pour proposer la coloration et vérification syntaxique, et l'auto-complétion. Parmi les bibliothèques les plus avancées, on retrouve YASQE, de la famille YASGUI (Rietveld, Hoekstra, 2015), qui est basée aussi CodeMirror mais propose un éditeur adapté à la syntaxe SPARQL. YASR (Rietveld, Hoekstra, 2015) associé à YASQE permet de gérer les messages d'erreur et de mieux afficher les résultats des requêtes.

Des interfaces fournissant plus de fonctionnalités, comme YASGUI⁵ et Flint SPARQL Editor⁶, sont disponibles en ligne pour interroger n'importe quel serveur SPARQL. YASGUI propose par exemple, l'interrogation simultanée de plusieurs serveurs SPARQL, la recherche de SPARQL endpoint, la conservation d'une requête entre deux sessions web.

2.1.2. Les langages visuels de requêtes

Il s'agit des systèmes permettant aux utilisateurs d'effectuer les requêtes de manière graphique. Leur but commun est d'augmenter l'intuitivité et l'utilisabilité du langage SPARQL. Leur principe est de construire les requêtes SPARQL en passant par des représentations graphiques.

2. Point d'accès SPARQL d'UniProt: <http://sparql.uniprot.org/>

3. Point d'accès SPARQL de « linked life data »: <http://linkedlifedata.com/sparql>

4. CodeMirror: <https://codemirror.net/>

5. YASGUI: <http://yasgui.org/>

6. Flint SPARQL Editor: <http://openuplabs.tso.co.uk/demos/sparqleditor>

Parmi les systèmes les plus récents et avancés, nous pouvons citer ReVeaLD (Kamdar *et al.*, 2014), QueryVOWL (Haag *et al.*, 2015) et SPARQLGraph (Schweiger *et al.*, 2014).

ReVeaLD⁷ est un système d'interrogation fédérée de données de recherches sur le cancer. Son principe est de partir d'un langage graphique, extensible, avec des concepts et liens spécifiques au domaine, pour permettre à l'utilisateur de représenter des questions complexes de son domaine. Un serveur SPARQL central reçoit donc la requête SPARQL construite, puis la traduit en plusieurs requêtes fédérées pour les différents serveurs de données.

QueryVOWL⁸ est défini par une représentation graphique des différents concept du langage SPARQL (Concepts, individus, disjonction, filtre, etc.). Le graphe construit par l'utilisateur, dans ce cas, ne représente pas la requête mais un patron des données. Chaque nœud correspond à une requête SPARQL et donc renvoie une certaine information par rapport à cet unique nœud. Le serveur SPARQL est interrogé régulièrement pour assister l'utilisateur avec la liste des valeurs possibles pour un nœud, ou la description d'un individu par exemple. Ces systèmes ne couvrent pas toutes les fonctionnalités du langage SPARQL mais peut servir à la construction rapide du début d'une requête complexe. L'utilisateur pourrait affiner cette dernière par la suite dans un éditeur textuel. L'avantage de ReVeaLD par rapport à QueryVOWL, c'est qu'il propose plus de flexibilité dans la sélection des informations à retourner (QueryVOWL par exemple ne retourne pas des lignes de résultats où on peut voir la relation entre les différentes valeurs des variables de requête, mais uniquement le label du nœud sélectionné).

SPARQLGraph⁹ (Schweiger *et al.*, 2014), comme ReVeaLD (Kamdar *et al.*, 2014), a été développé dans le contexte des sciences du vivant en prenant en compte leur spécificité. Sa particularité est de laisser le soin à l'utilisateur de choisir la source pour une certaine information, d'intégrer lui-même les différents services dans sa requête. A la différence de ReVeaLD qui passe par un médiateur de requête SPARQL, SPARQLGraph permet d'interroger plusieurs sources simultanément à partir d'une seule requête fédérée SPARQL. En plus, il permet un partage de requêtes visuelles et une collaboration entre les experts biologistes.

2.1.3. Interrogation basée sur les phrases en langage naturel

Les systèmes basés sur le langage naturel assistent l'utilisateur dans l'interrogation des bases de données en acceptant des questions en langage naturel. Ainsi, l'utilisateur n'a besoin de connaître que la terminologie des données, et non plus d'apprendre le langage SPARQL.

7. ReVeaLD: <http://srvgal78.deri.ie/reveald/>

8. QueryVOWL: <http://vowl.visualdataweb.org/queryvowl/index.html>

9. SPARQLGraph: <http://sparqlgraph.i-med.ac.at>

Dans le prototype actuel de LODQA¹⁰ (Kim, Cohen, 2013), l'utilisateur fournit une question sur un domaine de connaissance préalablement sélectionnée. La question passe ensuite par des étapes de traitement de langage naturel et de représentation intermédiaire. Les identifiants (URI) correspondant aux mots clés de la question sont recherchés dans des sources ontologiques ou de données spécifiques. Si ces URI sont trouvés, le système génère enfin plusieurs requêtes SPARQL, pouvant correspondre à la question posée.

Sparklis¹¹ (Ferré, 2014) assiste de manière interactive l'utilisateur dans la construction de sa requête en langage naturelle. Le point d'accès SPARQL est interrogé à chaque étape pour orienter l'utilisateur avec des informations possibles pour compléter la requête. La requête SPARQL est construite en arrière-plan. La difficulté pour l'utilisateur est de retrouver à chaque étape l'information à sélectionner car la quantité proposée est souvent très grande. Son avantage est l'exploration à facettes personnalisée.

2.2. *Approches de recherche d'informations*

2.2.1. *Exploration interactive*

Pour rechercher de l'information, on a généralement recours à la navigation dans un ensemble de données. La recherche d'information est généralement basée sur des mots clés et peut se faire suivant deux modes: les filtres et les facettes. Ces deux modes, par principe, analysent un grand ensemble de contenu et excluent les éléments qui ne respectent pas un certain critère (Whitenton, s. d.). La navigation par facettes utilise plusieurs filtres pour donner différentes vues afin que l'utilisateur ait une meilleure compréhension des données. Toutefois, la réalisation des facettes demande plus d'effort que celle des filtres.

Le service web Virtuoso Facets facilite le développement d'application web sémantique de navigation à facettes. La requête et la réponse étant sous forme d'arbre XML, la partie de la requête SPARQL nécessaire à la navigation à facettes est mieux précisée (Erling, Mikhailov, 2009) que dans le cas d'une requête en simple texte (SPARQL).

Le système SPARQLFilterFlow¹² (Haag *et al.*, 2014) étend le modèle filtre/flux pour l'adapter au SPARQL. Il consiste à faire passer toutes les données de la source à travers un réseau de flux d'un nœud d'entrée vers un nœud de sortie. Chaque nœud représente un filtre ne laissant passer que les données respectant le critère correspondant. Dans ce système, chaque nœud peut avoir plusieurs points de connexions en entrée et en sortie .

10. Linked Open Data Question Answering: <http://lodqa.org/>

11. Sparklis: <http://www.irisa.fr/LIS/ferre/sparklis/osparklis.html>

12. SPARQLFilterFlow: <http://sparql.visualdataweb.org/>

2.2.2. API de type service web

Pour l'accès aux données, les API sous formes de services web sont généralement utilisées. Dans ce cas, les services basées sur les principes RESTful présentent plus d'avantages (FIGURE 1 (Pautasso, 2014)) que les services dits WS-* (basés sur le protocole SOAP). REST (REpresentational State Transfer) (Fielding, Taylor, 2002) est un style d'architecture décrivant un ensemble de contraintes que doit respecter un service Web. Le plus intéressant dans ce style est sa simplicité, la possibilité d'utiliser des technologies standards du web (HTTP, URL, ...), et l'existence de plusieurs formats pour retourner les résultats des requêtes.

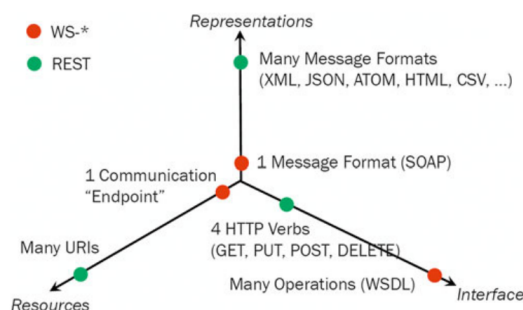


Figure 1. Avantages des services RESTful sur les services basés sur SOAP (WS-*)

Dans le domaine de l'intégration de données liées, certains projets ont démontré les avantages de fournir l'accès aux données via des API de type RESTful.

Le projet Open PHACTS Discovery Platform (Groth *et al.*, 2014), par exemple, exploite les données liées pour fournir un accès intégré à des bases de connaissance pharmaceutiques à partir d'API de type RESTful. Les données de plus de 12 sources sont présentes dans le système Open PHACTS (Uniprot, ENZYME, DrugBank, ...). Son architecture (Groth *et al.*, 2014) étend l'architecture standard des applications de données liées et ouvertes. Les modules d'accès aux données Web, de mise en correspondance de vocabulaire, de résolution d'identité, et d'évaluation de qualité servent à la création de la base d'intégration de données web. Ces modules sont décrits dans la section ci-dessous. Les autres modules relevant du fonctionnement de l'API sont :

- le médiateur qui est l'implémentation de l'API et donc effectue les requêtes SPARQL nécessaires vers la base de données et les services externes pour répondre aux requêtes venant des applications tierces;
- les définitions déclaratives qui décrivent précisément les liens entre les services et les données;
- les documentations des services de l'API qui permettent aux développeurs d'applications tierces de mieux comprendre comment utiliser les services proposés par l'API, et même de les tester ;

- les bibliothèques, spécifique à un langage de programmation, réduisent l'effort de développement d'applications tierces.

- les services externes sont les diverses sources de données interfacées par l'API; ils sont interrogés pour enrichir l'information retournée.

Cette architecture facilite l'implémentation des services de l'API de manière itérative et incrémentale car chaque module a une fonction bien précise. De plus, Open PHACTS Discovery Platform dispose d'un contrôle d'accès à travers le gestionnaire d'API, 3scale¹³, et d'un pilote pour l'intégration dans le système de Workflow KNIME¹⁴. Le contrôle d'accès peut être important lorsqu'on veut gérer des aspects comme (Groth *et al.*, 2014) des services payants, le suivi de l'utilisation de l'API, ou encore le respect de la politique d'exploitation.

SADI (Semantic Automated Discovery and Integration) (González *et al.*, 2014) est un ensemble de principes d'implémentation et d'exposition des services web REST. A la différence d'Open Phacts Discovery Platform qui utilise les standards du protocole HTTP, SADI requiert un fichier RDF/XML, respectant une sémantique particulière, pour passer les arguments ("input") aux services et, les résultats ("output") sont aussi retournés uniquement en RDF/XML. Il expose les services web sur internet à partir du registre SHARE (Semantic Health And Research Environment). Les services ainsi exposés peuvent être découverts par des applications comme le plugin SADI-Galaxy (Aranguren *et al.*, 2014) qui intègre les services SADI dans l'environnement de workflow Galaxy.

3. Développement de l'infrastructure d'accès aux données liées

3.1. Architecture proposée

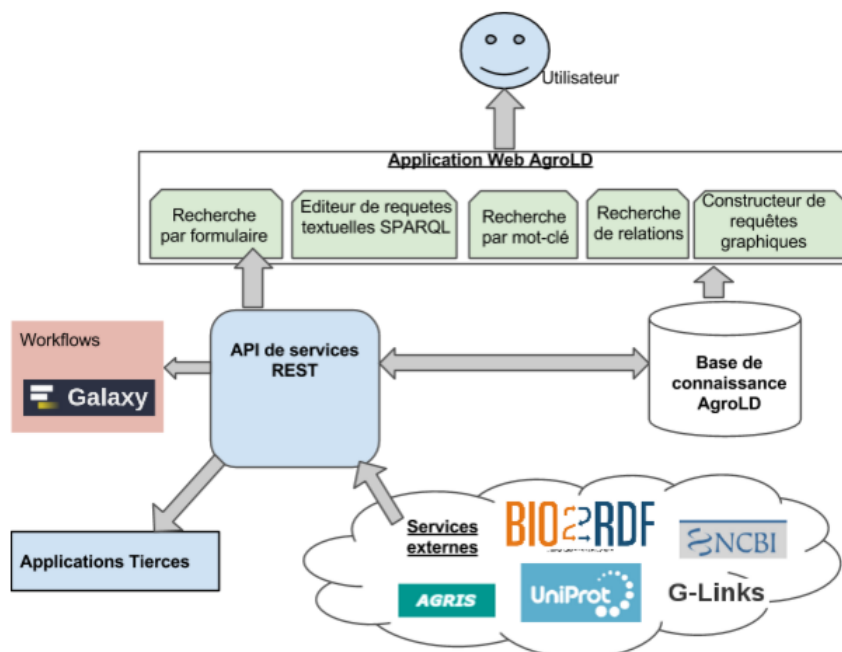
Les composants de notre architecture et leurs articulations sont présentés à la FIGURE 2. L'architecture reprend le modèle classique 3-tiers avec l'application web, l'API de Services Web et la base de données liées.

Au niveau de l'application Web, le module de recherche par mot-clé est destiné à explorer rapidement la base pour découvrir les données et ontologies stockées dans AgroLD à partir de mot-clés. Le constructeur de requêtes graphiques doit aider à la construction rapide de requêtes SPARQL à partir d'une représentation des termes des experts correspondant aux classes de la base. La recherche de relations permet de retrouver rapidement les relations existantes entre deux ou plusieurs entités.

L'API comprend les requêtes correspondant à des questions biologiques. Elle est utilisée lors de la recherche par formulaire et dans les workflows Galaxy. L'API peut interroger des services externes pour compléter les résultats trouvés dans AgroLD.

13. 3scale: www.3scale.net

14. KNIME: www.knime.org/



(Les flèches indiquent le sens des données)

Figure 2. Architecture proposée pour l'application web d'AgroLD

La base de connaissance AgroLD peut ainsi être enrichie pendant l'exécution des requêtes.

Des applications tierces peuvent être développées à base des services de l'API. Notre cas d'étude ici est la recherche par formulaire qui utilise l'API. Nous proposons une recherche par formulaire basée sur le modèle itératif et exploratoire défendu par V. Uren et Al (Uren *et al.*, 2007) pour une meilleure utilisabilité. Il s'agit en effet d'une recherche qui propose à l'utilisateur de nouvelles voies de recherche en fonction des résultats des recherches précédentes. L'utilisateur peut ainsi découvrir plus d'information relative à sa recherche.

3.2. Prototype Implémenté

Parmi les systèmes existants décrits dans l'état de l'art, certains peuvent être réutilisés pour accéder à AgroLD. Cependant, d'autres doivent être développés pour des besoins bien spécifiques. Afin de proposer une plateforme accessible aussi bien aux biologistes qu'aux bioinformaticiens, nous proposons d'intégrer plusieurs systèmes implémentant les différents paradigmes cités en section 2. Chaque utilisateur pourra ainsi utiliser les systèmes lui fournissant un bon compromis entre utilisabilité, expressivité, et convivialité des résultats.

3.2.1. *Intégration et adaptation de systèmes existants*

3.2.1.1. Système de recherche rapide

Pour ce système, nous pouvons intégrer le système de recherche à facettes, **Open-Link Faceted**, déjà disponible dans Openlink Virtuoso. Bien que les facettes présentées ici ne soient pas spécifiques aux concepts biologiques, Faceted permet de rapidement retrouver toutes les entités ayant un attribut de type littéral (textuel) contenant l'expression ou le mot-clé recherché.

3.2.1.2. Recherche de relations entre entités

La recherche des relations entre ressources peut être utile pour découvrir les liens existants entre des concepts ou entités biologiques. RelFinder (Heim *et al.*, 2009) présente l'avantage de donner une représentation graphique et de permettre une recherche rapide des entités par mot-clé lorsqu'on n'ignore leur URI. Cependant, la version originale de RelFinder a été développée (en ActionScript¹⁵) et configurée pour DBpedia. Nous avons proposé une configuration et une modification du système RelFinder qui sont adaptées à AgroLD. La configuration concerne principalement le point d'accès SPARQL, les propriétés à considérer pour la recherche d'entités et pour la description des ressources.

Dans notre implémentation nous avons ajouté quelques exemples pour guider les utilisateurs. L'exécution du premier exemple permet de voir graphiquement les relations entre le gène "adenosylmethionine decarboxylase" et les deux pathways "spermine biosynthesis" et "spermidine biosynthesis" auxquels il participe (FIGURE 14).

3.2.1.3. Système de requêtes visuelles

Nous proposons d'intégrer un système basé sur ReVeaLD (Kamdar *et al.*, 2014) mais avec un langage spécifique aux données d'AgroLD. Le langage est spécifié dans un fichier OWL en indiquant les différents concepts et propriétés des ressources d'AgroLD. ReVeaLD a ainsi besoin en entrée: des types de propriétés, des classes des sujets (« domain ») et des types possibles de valeurs (« range ») de ces propriétés. Le principal objectif est de fournir un moyen pour construire rapidement une requête SPARQL. ReVeaLD ne permet pas de faire toutes les interrogations possibles sur la base cible (par exemple, il est limité par la portée du langage qui le définit). La requête peut être visualisée et complétée dans l'éditeur textuel de requête SPARQL.

Ce système n'a pas pu être complètement développé pendant notre travail car nous ne disposions pas d'assez de temps pour définir le langage orienté utilisateur de concepts biologiques.

15. <http://www.adobe.com/devnet/actionsript.html>

3.2.1.4. Dé-référencement d'URI

L'application Pubby¹⁶, utilisé par DBpedia, nous a paru être une bonne solution pour répondre rapidement au problème de dé-référencement d'URI. Il s'agit en fait d'une application web Java, qui s'utilise dans le cas de point d'accès SPARQL supportant les requêtes DESCRIBE (de SPARQL). Nous l'avons aussi choisi parce que le serveur Virtuoso d'AgroLD supporte les requêtes DESCRIBE.

Toutefois à l'avenir, nous souhaiterions utiliser des systèmes plus récents comme TheDataTank (Vander Sande *et al.*, 2012).

3.2.2. Développement de nouvelles fonctionnalités

Globalement, notre prototype est une application web développée en Java, HTML et JavaScript. Les systèmes que nous avons implémenté sont présentés dans cette section.

3.2.2.1. Editeur de requêtes SPARQL

Il est beaucoup plus destiné à l'usage des bioinformaticiens (développeurs d'application biologiques). Ce système comprend plusieurs fonctionnalités (FIGURE 3) dont:

1. une zone d'édition des requêtes: En utilisant YASQE (Rietveld, Hoekstra, 2015), l'édition est facilitée avec les fonctionnalités d'auto-complétion, de vérification et coloration syntaxique, des raccourcis clavier, etc.
2. la sauvegarde des requêtes et résultats : Les requêtes peuvent être sauvegardées dans des fichiers textes pour être utilisés plus tard. Les résultats quant à eux, peuvent être téléchargés dans les formats RDF, CSV, JSON, etc.
3. un composant de remplacement des paramètres de patrons de requêtes, pour utiliser un exemple avec différentes valeurs de paramètres;
4. une liste de patrons de requêtes correspondants à des questions biologiques; ces patrons servent de bases à l'édition de requête
5. Une zone d'affichage des résultats, utilisant YASR (Rietveld, Hoekstra, 2015) pour la bonne présentation des résultats et des messages d'erreurs.

3.2.2.2. API de services web

ARCHITECTURE L'architecture globale de l'API de services Web a été inspirée de celle d'Open PHACTS Discovery Platform (Groth *et al.*, 2014) car elle est modulaire et basée sur REST. Elle présente ainsi plus d'avantages que l'architecture SADI (González *et al.*, 2014) au niveau des formats des résultats, de l'utilisation des standards, etc.

16. <http://wifo5-03.informatik.uni-mannheim.de/pubby/>

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. [Watch how!](#)

Set values of parameters: Replace "plant height" by: **3**

KEYBOARD COMMANDS

Query Text **1**

```

1 BASE <http://www.southgreen.fr/agro1d/>
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 vocab <vocabulary/>
5
6 SELECT DISTINCT ?term_id ?term_name ?graph
7 WHERE {
8   GRAPH ?graph {
9     ?term_id rdfs:label ?term_name .
10    FILTER regex(str(?term_name), 'plant height')
11  }
12 }

```

Execution timeout: 20000 milliseconds (values less than 1000 are ignored) **2** Results Format: **RDF/XML** **Download Results**

Filename to Save As: **Save File** **No file chosen**

Query: sparql **Save Query** **Load Selected Query File**

Query Patterns **4**

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

Results **5**

Raw Response Table Pivot Table Search: Show 50 entries

term_id	term_name	graph
http://purl.obolibrary.org/obo/TO_0000207	"plant height" ^{hexd-string}	http://www.southgreen.fr/agro1dto
http://purl.obolibrary.org/obo/TO_0001003	"plant height uniformity" ^{hexd-string}	http://www.southgreen.fr/agro1dto
http://purl.obolibrary.org/obo/TO_0001034	"relative plant height" ^{hexd-string}	http://www.southgreen.fr/agro1dto

Showing Top 3 of 3 entries.

Figure 3. Editeur de requêtes textuelles SPARQL

Certains développeurs partagent leurs expériences de la conception d'API accessible en ligne, en publiant des principes et astuces à prendre en compte. Nous pouvons citer entre autres ¹⁷:

- gérer les versions de l'API
- Utiliser correctement les codes d'états de HTTP
- Utiliser les méthodes HTTP: GET, POST, PUT, DELETE
- Renvoyer des messages d'erreurs compréhensibles et dans le même format que les résultats
- Retourner les résultats des requêtes en page pour éviter de lourds transfert de données sur le réseau
- Prévenir la surcharge du serveur en limitant le nombre d'accès simultanés
- fournir une bonne documentation

17. <https://localize-software.phraseapp.com/posts/best-practice-10-design-tips-for-apis/>

L'API que nous avons proposée comprend donc:

- **le module serveur**: correspondant au médiateur d'Open PHACTS, il implémente les services web. Il interroge le point d'accès SPARQL d'AgroLD et les services externes pour répondre aux requêtes. En plus de quelques services généraux utiles, ce module est constitué en majeure partie des services correspondants à des questions biologiques (diagramme de classe FIGURE 4). Il a été implémenté à l'aide de la librairie Jersey¹⁸ (version 2.17) parce qu'elle est open source et est destinée au développement de services REST en Java.

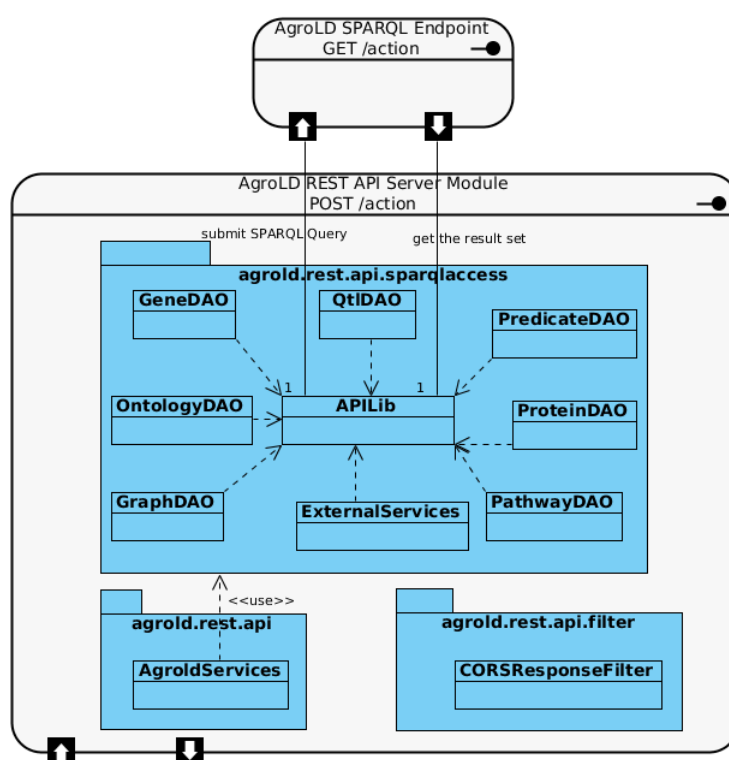


Figure 4. Module serveur de l'API d'AgroLD

- **la documentation interactive des services web**: elle donne une description complète pour comprendre et utiliser les services (rôle, paramètres, requête SPARQL exécutée, format de retour des résultats, signification des code HTTP, ...). Son interactivité permet aussi de tester les services et d'observer le comportement des services dans diverses situations. Les services sont décrits dans un fichier JSON¹⁹. Ce dernier

18. <https://jersey.java.net/>

19. Description de l'API d'AgroLD: <http://volvestre.cirad.fr:8080/aldp/swagger/agrold.json>

est utilisé par le framework SWAGGER²⁰(2.0) pour générer la page web de documentation (FIGURE 16).

– **Une librairie cliente pour application JavaScript:** elle est fournie par la librairie JavaScript `swagger-client.js` (de SWAGGER) à partir de la description JSON des services web.

SERVICES WEB DISPONIBLES Une trentaine de services sont actuellement accessibles à partir de l'API. Ils sont organisés en six catégories: *gene*, *graphs*, *ontologies*, *pathway*, *protein*, *QTL*. Nous présentons quelques exemples de services dans la TABLE 1.

Tableau 1. Exemples de services web disponibles dans l'API d'AgroLD

L'URI de base est <http://volvestre.cirad.fr:8080/agrold/api/1.0>

URI	Information recherchée
/genes/publications/byId	des publications à partir de l'identifiant d'un gène
/genes/participatingInPathway	les gènes qui participent à un pathway dont l'identifiant est donné
/graphs/predicates	liste des URIs des prédicats d'AgroLD
/ontologies/terms/children/byId	les descendant directs d'un concept ontologique étant donné son Id
/ontologies/terms/associatedWithQtl	les concepts ontologiques associés à un QTL donné et la nature de leur association
/proteins/EncodedByGene	les protéines encodées par un gène donné
/qtls/associatedWithProteinId	les QTLs associés à une protéine donnée

3.2.2.3. Formulaire dynamique de recherche

Il implémente la recherche par formulaire de notre architecture. L'idée était de développer un formulaire dynamique basée sur l'API pour faire des recherches d'information dans AgroLD. Le système proposé permet en plus d'explorer les informations fournies par les services. A partir de la librairie cliente de l'API, nous l'avons implémenté en JavaScript et JSP pour cinq types d'entités: gène, protéine, QTL, pathway, classe d'ontologies. L'exploration est représentée par le diagramme d'activité de la FIGURE 5. Par exemple, après avoir trouvé un QTL, on peut découvrir les informations relatives aux protéines ou concepts ontologiques auxquels il est associé.

3.2.2.4. Accès aux données à partir de Galaxy

Pour l'accès aux services web, Galaxy dispose du plugin **webservice_toolsuite**²¹. Ce dernier lit un fichier **WADL** (XML) décrivant les ressources RESTful disponibles et les rend accessibles dans Galaxy sous forme de sources de données. `webservice_toolsuite`

20. <http://swagger.io/>

21. https://toolshed.g2.bx.psu.edu/view/ganjoo/webservice_toolsuite/d5cd409b8a18

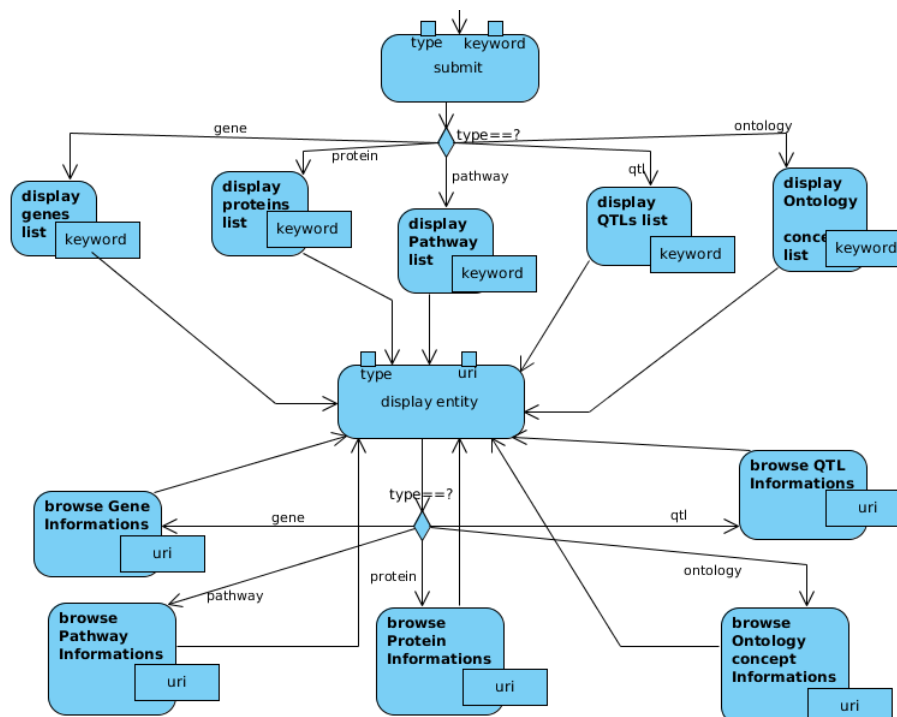


Figure 5. Activités de navigation avec le formulaire dynamique

n'utilise que la méthode POST de HTTP. C'est la raison principale pour laquelle nous avons défini la méthode POST pour les services de l'API au lieu de GET qui est plus indiqué lorsque les services font de la recherche d'information.

Nous avons décrit les services dans un fichier **WADL**²². Ensuite, tout utilisateur ayant installé l'outil `webservice_toolsuite` dans Galaxy peut passer par ses étapes d'intégration pour installer les ressources de l'API qu'il souhaite utiliser.

3.2.2.5. Enrichissement de l'information d'AgroLD

Dans notre travail, nous avons mis l'accent sur les résultats retournés à l'utilisateur. Nous proposons une intégration de services Web au niveau du formulaire dynamique. L'avantage par rapport à l'intégration des données directement dans AgroLD est de pouvoir utiliser aussi bien des données en RDF ou que dans un autre format. Pour répondre aux questions nécessitant des données externes, nous avons enveloppé l'accès aux services externes dans des services de l'API (classe `ExternalServices` de la FIGURE 4). Par exemple, pour la recherche de publications scientifiques relatives à une protéine, l'API passe par deux services:

22. http://volvestre.cirad.fr:8080/agrold/agrold_api.wadl

1. **G-Links**²³: pour obtenir leur identifiant (résultats en JSON);
2. puis **Europe PubMed Central**²⁴ pour obtenir leurs informations (en XML) : titre, auteurs, journal, année de publication, etc.;

4. Utilisation de l'application web AgroLD

4.1. Utilisation d'AgroLD par des utilisateurs humains

Pour l'évaluation des systèmes de recherche sémantique, Khadija Elbedweihy et al. (Elbedweihy *et al.*, 2012) recommandent une enquête auprès des utilisateurs à partir de quelques requêtes à soumettre au système. Nous nous sommes basés sur cette approche pour intégrer à notre application un formulaire de recueil des avis, remarques et suggestions des utilisateurs. Le formulaire comprend quelques questions du modèle de questionnaires "System Usability Scale"(SUS)(Brooke, 1996) et d'autres questions que nous avons jugées importantes.

Les trois principaux critères analysés sont

1. l'utilisabilité qui concerne la facilité d'entrer la requête (nombre de tentatives, temps nécessaire) et, la présentation des résultats;
2. l'expressivité qui définit pour un système de recherche quelles requêtes un utilisateur est capable de poser;
3. la performance en rapport avec la capacité à retourner rapidement des résultats corrects et complets.

En attendant les réponses, des utilisateurs, qui guideront l'évolution de l'application, nous analysons ici quelques scénarios d'utilisation des fonctionnalités de recherche du prototype implémenté:

1. « recherche des entités liées au terme "plant height" »: est une requête générale où l'utilisateur n'a pas une idée précise du type d'entité qu'il recherche;
2. « recherche des QTLs associés à l'identifiant ontologique "EO:0007403" »: est une requête un plus précise mais plus complexe que la précédente;
3. « recherche des publications relatives à la protéine "TBP1" »: est une requête qui nécessite l'interrogation de services externes car la base AgroLD ne contient pas de publications.
4. « recherche des relations existantes entre le gène AT3G25570 et les deux pathways "spermine biosynthesis" et "spermidine biosynthesis" »

Nous analysons surtout les moyens les plus pratiques d'obtenir des résultats. Nous avons fait nos analyses avec l'instance en développement de notre prototype (<http://volvestre.cirad.fr:8080/aldp/>).

23. <http://www.g-language.org/wiki/glinks>

24. <https://europepmc.org/RestfulWebService>

4.1.1. Entrée des requêtes et expressivité

Pour le scénario 1, il est plus facile de soumettre la requête par la fonctionnalité "Quick Search" que par "Advanced Search" et "SPARQL Query". Il suffit de saisir le mot-clé et de cliquer sur le bouton "Search" (Figure 6).

Search Text

Figure 6. Scénario 1: entrée de la requête

Par contre, pour les requêtes plus complexes, la simple recherche par mot-clé ("Quick Search") est moins facile à manipuler que les deux autres. Pour le scénario 2, l'absence de résultat par la recherche de "EO:0007403" dans "Quick Search", limite déjà l'utilisateur qui ne peut pas continuer la construction de sa requête. Avec "Advanced Search", la construction de la requête se fera seulement en 2 étapes (Figure 7) : « trouver le concept ontologique "EO:0007403" » puis « accéder aux QTLs associés ».

Ontology

Search ontology with keyword "EO:0007403"

1- trouver le concept ontologique "EO:0007403"

Id	Name	Description	URI
1 EO:0007403 (display)	unknown environment		http://purl.obolibrary.org/obo/EO_0007403 (in Sparql)

ONTOLOGY : EO:0007403 / unknown environment

URI: http://purl.obolibrary.org/obo/EO_0007403

Parents ±

Children ±

Proteins associated ±

QTL associated ±

2- accéder aux QTLs

Figure 7. Scénario 2: entrée de la requête dans le formulaire dynamique

Pour préciser les types d'informations qu'on souhaite avoir sur ces QTLs, il serait préférable de passer par l'éditeur de requêtes SPARQL (Figure 8). En plus cet éditeur peut être agrandi et il contient une assistance syntaxique (YASQUE (Rietveld, Hoekstra, 2015)) pour un meilleur confort dans la saisie.

Les informations sur les publications n'étant pas accessibles depuis un point d'accès SPARQL, nous ne pouvons pas utiliser de requêtes SPARQL fédérées dans le scénario 3. Par contre, avec "Advanced Search", après avoir trouver la protéine avec le mot clé "TBP1", l'utilisateur peut demander la recherche des publications concernant cette protéine (Figure 9).



Figure 8. Scénario 2: entrée de la requête dans l'éditeur de requête SPARQL

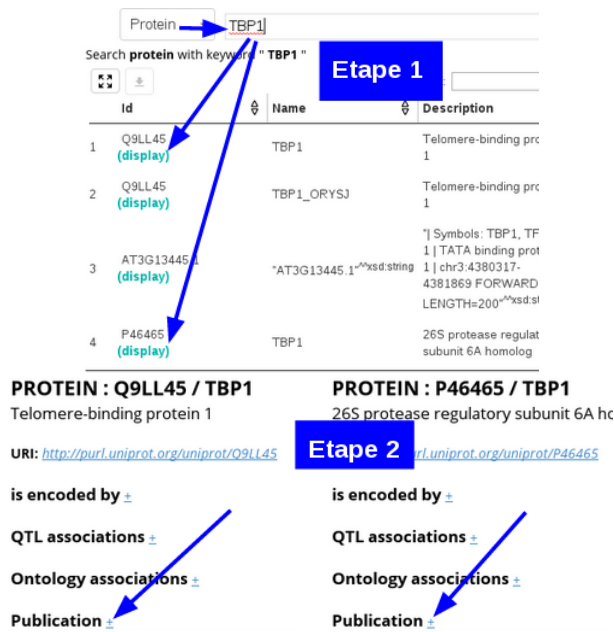


Figure 9. Scénario 3: entrée de la requête dans le formulaire dynamique

Dans le dernier scénario, il est possible d'utiliser une requête SPARQL, mais elle sera très difficile à construire. En utilisant la fonctionnalité de recherche de relations, l'utilisateur retrouve rapidement les ressources à partir des noms du gène et des pathways, et soumet sa requête (Figure 10)

4.1.2. Exécution des requêtes et temps de réponse

L'outil OpenLink Faceted utilisé par la fonctionnalité "Quick Search" retourne le temps de recherche effectué. Nous observons que la requête du scénario 1 s'exécute en environ 452 millisecondes; un temps relativement proche de celui offert par les mo-

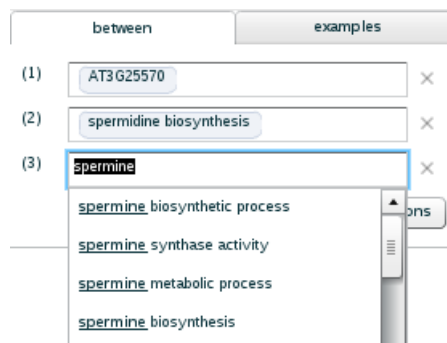


Figure 10. Scénario 4: entrée de la requête

teurs de recherche classiques. Par contre Faceted ne fournit pas de moyen d'effectuer la recherche sur l'aspect sémantique (par exemple rechercher uniquement les gènes).

La durée de l'exécution d'une requête par l'éditeur de requête SPARQL dépend des clauses de la requête et des performances du point d'accès SPARQL. Plus la requête sera complexe (nombre de clauses, ou type de filtres, ...), plus ce temps sera long.

Nous avons noté que le temps d'exécution lors de la recherche par formulaire est de l'ordre de 18 secondes (13 puis 5 pour les 2 étapes) dans le scénario 2 (recherche dans la base AgroLD). Il est de 24 secondes (20 puis 4 pour les deux étapes) dans le scénario 3 (interrogation de services web externes). Ce qui indique que la recherche par mot-clé prend beaucoup plus de temps par rapport aux autres actions que peut effectuer l'utilisateur pendant son exploration.

Nous n'avons pas recueilli le temps d'exécution de la recherche de relation, mais elle semble aussi rapide que le "Quick Search". Cependant, l'utilisateur doit patienter quelques secondes pendant l'affichage des relations qui se fait de manière séquentielle, l'une après l'autre.

4.1.3. Présentation des résultats

Dans le scénario 1 avec "Quick Search", nous observons une bonne présentation (Figure 11) sous forme de tableau avec les mot-clés mis en gras pour indiquer où ils ont été retrouvés. Les liens hypertexte sur les URIs permettent de continuer l'exploration des données (Figure 11).

Par contre l'utilisateur rencontrera beaucoup de difficultés à parcourir tous les résultats retournés s'ils sont très nombreux (5537 pour le scénario 1).

Par le scénario 2, nous observons que ce problème peut être résolu en permettant à l'utilisateur de préciser le type d'entité qu'il recherche, mais en faisant un choix seulement parmi ceux qui sont proposés par la fonctionnalité "Advanced Search". Les résultats obtenus ici sont plus précis (Figure 12).

Displaying Ranked Entity Names and Text summaries where:
 ?sl has any Attribute with Value "plant height" Drop.

View query as SPARQL Facet permalink

Go to: Show 150 651 - 800 of 5537 total

Entity	Title	Named Graph	plant height.
http://www.identifie...gramene.qtl/AQBB053	BNL6.32	http://www.southgree...old/qtl.annotations	plant height.
http://www.identifie...gramene.qtl/AQEA376	BNL6.32	http://www.southgree...old/qtl.annotations	plant height.
http://www.identifie...gramene.qtl/COZ18	BNL6.32	http://www.southgree...old/qtl.annotations	plant height.
http://www.identifie...gramene.qtl/AQEA346	BNL6.32	http://www.southgree...old/qtl.annotations	plant height.
http://www.identifie...gramene.qtl/AQBB070	BNL6.32	http://www.southgree...old/qtl.annotations	plant height.
http://www.identifie...gramene.qtl/AQEA320	BNL6.32	http://www.southgree...old/qtl.annotations	plant height.

Entity Relationship Filters
 Type
 Attributes
 Values
 Distinct (Count)
 Places Any location

Figure 11. Scénario 1: présentation des résultats avec la recherche rapide par mot-clé

QTL associated Next page>>

Search: Show 30 entries

qtlId	Association	URI
1 AQAA003 (display)	observed_in	http://www.identifiers.org/gramene.qtl/AQAA003 (in Sparql)
2 AQAA015 (display)	observed_in	http://www.identifiers.org/gramene.qtl/AQAA015 (in Sparql)

Figure 12. Scénario 2 : présentation des résultats

La présentation des résultats dans "Advanced Search" et l'éditeur de requêtes SPARQL est aussi faite sous forme de tableaux pour une bonne lisibilité. Des liens y sont disponibles pour obtenir des informations complémentaires (liens "display" et URLs dans la Figure 12). La recherche par formulaire dynamique ("Advanced Search") propose aussi un lien pour une description de ressource dans l'éditeur SPARQL (lien "in SPARQL" dans la Figure 12).

Les résultats retournés dans le scénario 3 sont présentés dans un style proche des citations de documents (Figure 13). Ainsi la reconnaissance des différentes informations relatives à la publication est intuitive. Un lien associé permet de se rendre à la source de ces informations pour plus de détails.

Publication

1. Yu EY, Kim SE, Kim JH, Ko JH, Cho MH, Chung IK., " **Sequence-specific DNA recognition by the Myb-like domain of plant telomeric protein RTBP1.** ", *J Biol Chem*, 2000
 More at: <http://www.ncbi.nlm.nih.gov/pubmed/10811811>
2. International Rice Genome Sequencing Project., " **The map-based sequence of the rice genome.** ", *Nature*, 2005
 More at: <http://www.ncbi.nlm.nih.gov/pubmed/16100779>
3. Rice Annotation Project, Tanaka T, Antonio BA, Kikuchi S, Matsumoto T, Nagamura Y, Numa H, Sakai H, ..., " **The Rice**

Figure 13. Scénario 3 : présentation des résultats

Dans le dernier scénario, les relations sont observées sous forme d'arête entre les nœuds d'un graphe pour une plus rapide lisibilité et compréhension. L'utilisateur peut filtrer les types de liens qu'il souhaite observer (les liens "type" sont cachés dans la Figure 14).

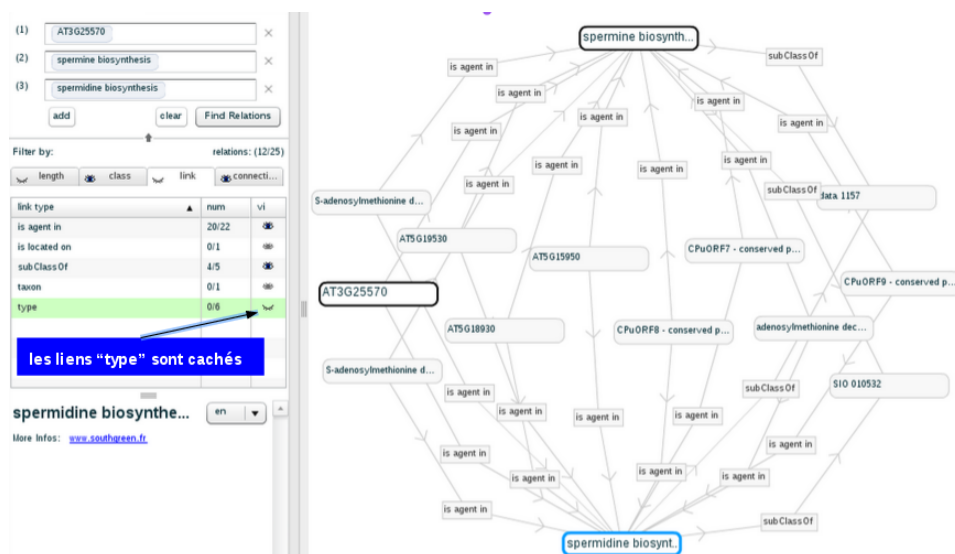


Figure 14. Scénario 4 : Relations découvertes entre le gène "adenosylmethionine decarboxylase" (AT3G25570) et les deux pathways "spermine biosynthesis" et "spermidine biosynthesis"

4.2. Utilisation d'AgroLD dans des applications

Le développement du formulaire dynamique était aussi l'occasion de démontrer l'utilisation en programmation de l'API d'AgroLD. Comme nous l'avons mentionné dans le chapitre, la description JSON des services de l'API peut être utilisée avec une bibliothèque cliente de SWAGGER (par exemple **swagger-client.js** pour le JavaScript).

Avec **swagger-client.js**, il suffit de préciser l'URL du fichier JSON, puis les services de l'API peuvent être appelés comme les méthodes d'un objet (**swagger.apis**) sans se préoccuper des URLs des services ou de la méthode HTTP qu'ils supportent (Figure 15).

En plus de cette facilité, la documentation interactive des services (Figure 16) permet aux développeurs d'applications d'avoir toutes les informations sur chaque service pour savoir l'utiliser.

```

window.swagger
url: url,
success:
conso
displayHoldMessage("result");
switch (type) {
case "gene":
    swagger.apis.gene.getGenesByKeyword({format: "sparql-json", keyword: keyword, _page:
n/json}), function (data) {
    }
}
    
```

url du fichier JSON de définition déclarative de l'API

paramètres du service

appel au service sous forme d'appel de méthode

Figure 15. Utilisation du service de recherche de gène par mot-clé dans un programme JavaScript

Méthode HTTP acceptée

Request URL

Requête SPARQL exécutée

```

SELECT DISTINCT ?proteinId
WHERE
{
  SELECT ?ontoEtr
  WHERE
  {
    ?ontoEtr rdfs:subClassOf ?ontoEtrClass.
    FILTER REGEX(STR(?ontoEtr), CONCAT(".*", ".*", ".*"))
    limit 1
  }
  protein ?predicate ?ontoEtr.
  ?protein rdfs:subClassOf+ http://purl.obolibrary.org/obo/GO_0001104-.
  BIND(REPLACE(?protein, ".*", "AS ?proteinId").
  ORDER BY ?proteinId
  LIMIT ?page_size
  OFFSET ?page
  }
    
```

Paramètres obligatoires et optionnels

Parameter	Value	Description	Parameter	Data Type
keyword	CC0003204	The ID of the ontological identifier (e.g. 000000000)	query	string
_page		Number (1-10) of the page of the result set to display (0 = the first page by default)	query	integer
_page_size		Size of a page of the results	query	integer

URL complète de la ressource retournée

Résultats

```

{
  "proteinId": "Association",
  "uri": "http://www.southgreen.fr/agrold/greenphy/sequence/At1g19930.1",
  "name": "Name",
  "name_function": "http://purl.uniprot.org/uniprot/211K1A?",
  "211K1A?": "http://purl.uniprot.org/uniprot/211K1A?",
  "211K1A?": "http://purl.uniprot.org/uniprot/211K1A?",
  "211K1A?": "http://purl.uniprot.org/uniprot/211K1A?",
  "211K1A?": "http://purl.uniprot.org/uniprot/211K1A?",
  "211K1A?": "http://purl.uniprot.org/uniprot/211K1A?",
  "211K1A?": "http://purl.uniprot.org/uniprot/211K1A?"
}
    
```

Statut de la réponse HTTP

Response Headers

```

{
  "date": "Thu, 24 Sep 2015 07:08:43 GMT",
  "server": "Apache/2.4.18.1",
  "Content-Type": "application/json"
}
    
```

Format des résultats

Figure 16. Documentation du service de recherche des protéines associées à un identifiant ontologique

5. Conclusion

Parmi les nombreux axes de recherche que compte le domaine bioinformatique, la gestion des connaissances est devenue un domaine de recherche important, axé sur l'interconnexion de l'information et la représentation des connaissances. AgroLD est l'une des premières initiatives connues pour appliquer les pratiques du SW au domaine génétique et moléculaire en agromomie. Aujourd'hui, elle joue un rôle complémentaire des approches intégratives adoptées par cette communauté. Dans cet article, nous proposons un modèle d'architecture implémentant les paradigmes de recherche d'information sémantique (Haag *et al.*, 2014). Certains sont proposés soit pour assister les utilisateurs dans la construction de leurs requêtes, soit pour cacher le langage SPARQL et fournir une interface plus facile à utiliser pour l'interrogation des bases cibles. Chacun de ces paradigmes a des atouts, mais aussi des limites que leur intégration dans une seule infrastructure peut combler.

A l'avenir, nous allons améliorer les performances et étendre les fonctionnalités de l'API de services afin de couvrir d'autres entités représentées dans AgroLD (par exemple: annotation génomique et l'information d'homologie). Nous proposons éga-

lement d'intégrer un système basé sur ReVeLD (Kamdar *et al.*, 2014) mais avec un langage spécifique aux données d'AgroLD. Le principal objectif est de fournir aux utilisateurs un moyen pour construire rapidement une requête SPARQL. De plus, nous prévoyons d'explorer les fonctionnalités offertes par les outils d'indexation et recherche dans le texte telles qu'ElasticSearch²⁵ et SolR²⁶. Ceci permettra d'améliorer la fonctionnalité de recherche rapide en récupérant plus d'informations textuelles et ainsi de cacher les détails techniques.

Bibliographie

- Antezana E., Kuiper M., Mironov V. (2009). Biological knowledge management: the emerging role of the Semantic Web technologies. *Briefings in Bioinformatics*, vol. 10, n° 4, p. 392–407.
- Aranguren M. E., González A. R., Wilkinson M. D. (2014). Executing SADI services in Galaxy. *Journal of biomedical semantics*, vol. 5, n° 1, p. 42.
- Ashburner M., Ball C. A., Blake J. A., Botstein D., Butler H., Cherry J. M. *et al.* (2000). Gene Ontology: tool for the unification of biology. *Nature genetics*, vol. 25, n° 1, p. 25–29.
- Barrell D., Dimmer E., Huntley R. P., Binns D., O'Donovan C., Apweiler R. (2009). The GOA database in 2009 - An integrated Gene Ontology Annotation resource. *Nucleic Acids Research*, vol. 37, n° SUPPL. 1.
- Belleau F., Nolin M.-A., Tourigny N., Rigault P., Morissette J. (2008). Bio2RDF: towards a mashup to build bioinformatics knowledge systems. *Journal of biomedical informatics*, vol. 41, n° 5, p. 706–16.
- Brooke J. (1996). SUS-A quick and dirty usability scale. *Usability evaluation in industry*, vol. 189, n° 194, p. 4–7.
- Conte M. G., Gaillard S., Lanau N., Rouard M., Périn C. (2008). GreenPhylDB: a database for plant comparative genomics. *Nucleic Acids Research*, vol. 36, n° Database issue, p. D991–998.
- Cooper L., Walls R. L., Elser J., Gandolfo M. A., Stevenson D. W., Smith B. *et al.* (2013). The plant ontology as a tool for comparative plant anatomy and genomic analyses. *Plant & cell physiology*, vol. 54, n° 2, p. e1.
- Droc G., Ruiz M., Larmande P., Pereira A., Piffanelli P., Morel J. B. *et al.* (2006). OryGenesDB: a database for rice reverse genetics. *Nucleic Acids Research*, vol. 34, n° Database issue, p. D736–40.
- Elbedweihy K., Wrigley S. N., Ciravegna F., Reinhard D., Bernstein A. (2012). Evaluating semantic search systems to identify future directions of research. In *The semantic web: Eswc 2012 satellite events*, p. 148–162. Springer.
- Erling O., Mikhailov I. (2009). Faceted Views over Large-Scale Linked Data. In *Ldow*.

25. ElasticSearch: <https://www.elastic.co>

26. SolR: <https://lucene.apache.org/solr/>

- Ferré S. (2014). Expressive and scalable query-based faceted search over SPARQL endpoints. In *The semantic web—iswc 2014*, p. 438–453. Springer.
- Fielding R. T., Taylor R. N. (2002). Principled design of the modern Web architecture. *ACM Transactions on Internet Technology*, vol. 2, n° 2, p. 115–150.
- Goble C., Stevens R. (2008). State of the nation in data integration for bioinformatics. *Journal of biomedical informatics*, vol. 41, n° 5, p. 687–93.
- González A., Callahan A., Cruz-Toledo J., Garcia A., Egaña Aranguren M., Dumontier M. *et al.* (2014). Automatically exposing OpenLifeData via SADI semantic Web Services. *J Biomed Semantics*, vol. 5, n° 1, p. 46.
- Groth P., Loizou A., Gray A. J., Goble C., Harland L., Pettifer S. (2014). API-centric Linked Data integration: The Open PHACTS Discovery Platform case study. *Web Semantics: Science, Services and Agents on the World Wide Web*, p. 1–7.
- Haag F., Lohmann S., Bold S., Ertl T. (2014). Visual SPARQL Querying based on Extended Filter/Flow Graphs. In *Proceedings of the 12th international working conference on advanced visual interfaces (avi '14)*, p. 305–312. New York, NY, USA: ACM.
- Haag F., Lohmann S., Siek S., Ertl T. (2015). Visual Querying of Linked Data with {Query-VOWL}. In *Joint proceedings of sumpre 2015 and hswi 2014-15*. CEUR-WS.
- Hamelin C., Sempere G., Jouffe V., Ruiz M. (2012). TropGeneDB, the multi-tropical crop information system updated and extended. *Nucleic acids research*, p. gks1105.
- Heim P., Hellmann S., Lehmann J., Lohmann S., Stegemann T. (2009). RelFinder: Revealing relationships in RDF knowledge bases. In *Semantic multimedia*, p. 182–187. Springer.
- Jupp S., Klein J., Schanstra J., Stevens R. (2011). Developing a kidney and urinary pathway knowledge base. *Journal of biomedical semantics*, vol. 2, n° 2, p. 1.
- Jupp S., Malone J., Bolleman J., Brandizi M., Davies M., Garcia L. *et al.* (2014). The EBI RDF platform: linked open data for the life sciences. *Bioinformatics*, vol. 30, n° 9, p. 1338–1339.
- Kamdar M. R., Zeginis D., Hasnain A., Decker S., Deus H. F. (2014). ReVealD: A user-driven domain-specific interactive search platform for biomedical research. *Journal of Biomedical Informatics*, vol. 47, p. 112–130.
- Kim J.-D., Cohen K. B. (2013). Natural language query processing for SPARQL generation: A prototype system for SNOMED CT. In *Proceedings of biolink*, p. 32–38.
- Larmande P., Gay C. C., Lorieux M., Périn C., Bouniol M., Droc G. G. *et al.* (2008). Oryza Tag Line, a phenotypic mutant database for the Genoplante rice insertion line library. *Nucleic Acids Research*, vol. 36, n° Database issue, p. D1022–7.
- Luciano J. S., Andersson B., Batchelor C., Bodenreider O., Clark T., Denney C. K. *et al.* (2011). The Translational Medicine Ontology and Knowledge Base: driving personalized medicine by bridging the gap between bench and bedside. *Journal of biomedical semantics*, vol. 2 Suppl 2, n° Suppl 2, p. S1.
- Magrane M., Consortium U. P. (2011). UniProt Knowledgebase: A hub of integrated protein data. *Database*, vol. 2011.
- Momtchev V., Peychev D., Primov T., Georgiev G. (2009). Expanding the Pathway and Interaction Knowledge in Linked Life Data. In *International semantic web challenge*.

- Monaco M. K., Stein J., Naithani S., Wei S., Dharmawardhana P., Kumari S. *et al.* (2014). Gramene 2013: Comparative plant genomics resources. *Nucleic Acids Research*, vol. 42, n° D1.
- Mungall C. J., Batchelor C., Eilbeck K. (2011). Evolution of the Sequence Ontology terms and relationships. *Journal of biomedical informatics*, vol. 44, n° 1, p. 87–93.
- Networks S.L 3scale. (2011). *What is an API? Your guide to the Internet Business (R)evolution*. Consulté sur <http://www.3scale.net/wp-content/uploads/2012/06/What-is-an-API-1.0.pdf> (Disponible à <http://www.3scale.net/wp-content/uploads/2012/06/What-is-an-API-1.0.pdf>)
- Pautasso C. (2014). RESTful Web Services: Principles, Patterns, Emerging Technologies. *Web Services Foundations, Springer Science+Business Media New York*, p. 31–51.
- Rietveld L., Hoekstra R. (2015). The YASGUI Family of SPARQL Clients. *Semantic Web Journal*.
- Schweiger D., Trajanoski Z., Pabinger S. (2014). SPARQLGraph: a web-based platform for graphically querying biological Semantic Web databases. *BMC bioinformatics*, vol. 15, n° 1, p. 279.
- Uren V., Lei Y., Lopez V., Liu H., Motta E., Giordanino M. (2007). The usability of semantic search tools: a review. *The Knowledge Engineering Review*, vol. 22, n° 04, p. 361–377.
- Vander Sande M., Colpaert P., Van Deursen D., Mannens E., Walle R. de. (2012). The datatank: an open data adapter with semantic output. In *21st international conference on world wide web, proceedings*, p. 4.
- Venkatesan A., Tripathi S., Sanz de Galdeano A., Blondé W., Lægreid A., Mironov V. *et al.* (2014). Finding gene regulatory network candidates using the gene expression knowledge base. *BMC bioinformatics*, vol. 15, n° 1, p. 386.
- Whitenton K. (s. d.). *Filters vs. Facets: Definitions*.
- Williams A. J., Harland L., Groth P., Pettifer S., Chichester C., Willighagen E. L. *et al.* (2012). *Open PHACTS: Semantic interoperability for drug discovery*.