

LEAPS: A Semantic Web and Linked data framework for the Algal Biomass Domain

Monika Solanki¹ and Johannes Skarka²

¹ Aston University, UK
m.solanki@aston.ac.uk

² Karlsruhe Institute of Technology, ITAS, Germany
johannes.skarka@kit.edu

Abstract. In this paper we present, *LEAPS*, a Semantic Web and Linked data framework for searching and visualising datasets from the domain of Algal biomass. *LEAPS* provides tailored interfaces to explore algal biomass datasets via REST services and a SPARQL endpoint for stakeholders in the domain of algal biomass. The rich suite of datasets include data about potential algal biomass cultivation sites, sources of CO₂, the pipelines connecting the cultivation sites to the CO₂ sources and a subset of the biological taxonomy of algae derived from the world's largest online information source on algae.

1 Motivation

Algal biomass holds huge promises. The use of microalgae as a food source for humans has been considered for overpopulated countries and for space travel since as early as 1961 [3]. If algae is grown under proper environmental conditions, the protein yield from it may be quite high. Algae have been collected for more than 4000 years in China and Japan for use as human food ³,

Recently the idea that algae biomass based biofuels could serve as an alternative to fossil fuels has been embraced by councils across the globe. Major companies [1,2], government bodies [4] and dedicated non-profit organisations such as ABO (Algal Biomass Organisation) ⁴ and EABA(European Algal Biomass Association)⁵ have been pushing the case for research into clean energy sources including algae biomass based biofuels.

It is quickly evident that because of extensive research being carried out, the domain itself is a very rich source of information. Most of the knowledge is however largely buried in various formats of images, spreadsheets, proprietary data sources and grey literature that are not readily machine accessible/interpretable. A critical limitation that has been identified is the lack of a knowledge level infrastructure that is equipped with the capabilities to provide semantic

³ <http://www.botgard.ucla.edu/html/botanytextbooks/economicbotany/Algae/index.html>

⁴ <http://www.algalbiomass.org/>

⁵ <http://www.eaba-association.eu/>

grounding to the datasets for algal biomass so that they can be interlinked, shared and reused within the biomass community.

Integrating algal biomass datasets to enable knowledge representation and reasoning requires a technology infrastructure based on formalised and shared vocabularies. Stakeholders in the domain who would benefit from such a structured, unambiguous and machine interpretable representation of data include researchers, algae producers and users, biofuels producers, oil companies, airline, cars and aerospace industry, national public authorities, international organisation and NGOs amongst others.

In this paper, we present *LEAPS*⁶, a Semantic Web/Linked data framework for the representation and visualisation of knowledge in the domain of algal biomass. One of the main goals of *LEAPS* is to enable the stakeholders of the algal biomass domain to interactively explore, via linked data, potential algal sites and sources of their consumables across NUTS (Nomenclature of Units for Territorial Statistics)⁷ regions in North-Western Europe.

Some of the objectives of *LEAPS* are,

- motivate the use of Semantic Web technologies and LOD for the algal biomass domain.
- laying out a set of ontological requirements for knowledge representation that support the publication of algal biomass data.
- elaborating on how algal biomass datasets are transformed to their corresponding RDF model representation.
- interlinking the generated RDF datasets along spatial dimensions with other datasets on the Web of data.
- visualising the linked datasets via an end user LOD REST Web service.
- visualising the scientific classification of the algae species as large network graphs.

The paper is structured as follows: Section 2 presents a brief overview of the dataset transformation process. Section 3 presents a description of the system architecture. Section 4 presents an overview of the querying mechanism underlying the *LEAPS* interface.

2 *LEAPS* Datasets

The transformation of the raw datasets to linked data takes place in two steps. The first part of the data processing and the potential calculation are performed in a GIS-based model which was developed for this purpose using ArcGIS⁸ 9.3.1.

The second step of lifting the data from XML to RDF is carried out using a bespoke parser that exploits XPath⁹ to selectively query the XML datasets and generate linked data using the ontologies. While in most cases, transforming

⁶ <http://www.semanticwebservices.org/enalgae>

⁷ <http://bit.ly/I7y5st>

⁸ <http://www.esri.com/software/arcgis/index.html>

⁹ <http://www.w3.org/TR/xpath/>

XML datasets to their linked data counterparts is done assuming a simplistic one-to-one mapping between the XML elements and RDF entities, in our scenario, the original data sources had several limitations and a one-to-one transformation was not possible. In order to produce a linked data representation of the datasets, that directly interlinked the resources of sites, sources, pipelines and region potential to each other and their NUTS regions of location, a bespoke parser that utilised a complex underlying data structure to facilitate the transformation was implemented.

The transformation process yielded four datasets which were stored in distributed triple store repositories: Biomass production sites, CO₂ sources, pipelines and region potential. We stored the datasets in separate repositories to simulate the realistic scenario of these datasets being made available by distinct and dedicated dataset providers in the future. While a linked data representation of the NUTS regions data¹⁰, was already available there was no SPARQL endpoint or service to query the dataset for region names. We retrieved the dataset dump and curated it in our local triple store as a separate repository. The NUTS dataset was required to link the biomass production sites and the CO₂ sources to regions where they would be located and to the dataset about the region potential of biomass yields. The transformed datasets interlinked resources defining sites, CO₂ sources, pipelines, regions and NUTS data using link predicates defined in the ontology network.

Datasets about algae cultivation can become more meaningful and useful to the biomass community, if they are integrated with datasets about algal strains. This can help the plant operators in taking judicious decisions about which strain to cultivate at a specific geospatial location. Algaebase¹¹ provides the largest online database of algae information. While Algaebase does not make RDF versions of the datasets directly available through its website, they can be programmatically retrieved via their LSIDs (Life Science Identifiers) from the LSID Web resolver¹² made available by Biodiversity Information Standards (TDWG)¹³ working group.

We retrieved RDF metadata for 113061 species of algae¹⁴ and curated in our triple store. We then used the Semantic import plugin with Gephi to visualise the biological taxonomy of the algae species.

3 System Description

LEAPS provides an integrated view over multiple heterogeneous datasets of potential algal sites and sources of their consumables across NUTS regions in North-Western Europe. Figure 1 illustrates the conceptual architecture of *LEAPS*. The main components of the application are

¹⁰ <http://nuts.geovocab.org/>

¹¹ <http://www.algaebase.org/about/>

¹² <http://lsid.tdwg.org/>

¹³ <http://www.tdwg.org/>

¹⁴ The retrieval algorithm ran on an Ubuntu server for three days

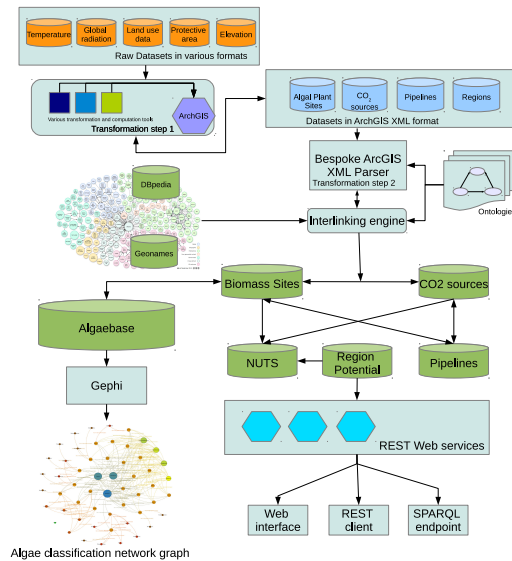


Fig. 1. Architecture of *LEAPS*

- **Parsing modules:** As shown in Figure 1, the parsing modules are responsible for lifting the data from their original formats to RDF. The lifting process takes place in two stages to ensure uniformity in transformation.
- **Linking engine:** The linking engine along with the bespoke XML parser is responsible for producing the linked data representation of the datasets. The linking engine uses ontologies, dataset specific rules and heuristics to generate interlinking between the five datasets. From the LOD cloud, we currently provide outgoing links to DBpedia¹⁵ and Geonames¹⁶.
- **Triple store:** The linked datasets are stored in a triple store. We use OWLIM SE 5.0¹⁷.
- **Web services:** Several REST Web services have been implemented to provide access to the linked datasets.
- **SPARQL endpoints:** SPARQL endpoints that provide access to individual dataset repositories are available. Snorql has been customised as the front end for the endpoint. An endpoint for federated queries is planned to be implemented as part of future work.
- **Ontologies:** A suite of OWL ontologies for the algal biomass domain have been designed and made available.
- **Interfaces:** The Web interface provides an interactive way to explore various facets of sites, sources, pipelines, regions, ontologies and SPARQL endpoints.

¹⁵ <http://dbpedia.org/About>

¹⁶ <http://sws.geonames.org/>

¹⁷ <http://www.ontotext.com/owlim/editions>

The map visualisation has been rendered using Google maps. Besides the SPARQL endpoint and the interactive Web interface, a REST client has been implemented for access to the datasets. Query results are available in RDF/XML, JSON, Turtle and XML formats.

- **Biological taxonomy visualisation:** A subset of the Algaebase database which is the largest information source of algae on the Web, has been retrieved and curated in our triple store. This dataset when integrated with the dataset for algal cultivation site, can inform stakeholders about the strains of algae that can be harvested on that site. Further, the Semantic Import plugin¹⁸ of Gephi¹⁹ has been exploited to visualise the biological taxonomy of algae. This visualisation is also made available via the *LEAPS* interface.

4 Application access

*LEAPS*²⁰ is available on the Web. The interface currently provides visualisation and navigation of the algae cultivation datasets in a way most intuitive for the phycologists. The application has been demonstrated to several stakeholders of the community at various algae-related workshops and congresses. They have found the navigation very useful and made suggestions for future dataset aggregation. At the time of this writing, data retrieval is relatively slow for some queries because of their federated nature, however optimisation work on the retrieval mechanism is in progress to enable faster retrieval of information.

Acknowledgments

The research described in this paper is partly supported by the Energetic Algae project (EnAlgae), a 4 year Strategic Initiative of the INTERREG IVB North West Europe Programme.

References

1. A. H. Claire Smith. Research needs in ecosystem services to support algal biofuels, bioenergy and commodity chemicals production in the uk. Technical report, NNFCC, 2011.
2. Oilgae. Oilgae comprehensive report, energy from algae: Products, market, processes and strategies. Technical report, Oilgae, 2011.
3. R. C. Powell and E. M. Nevels. Algae feeding in humans. *Journal of Nutrition*, 1961.
4. U.S. Department of Energy. National Algal Biofuels Technology Roadmap. Technical report, accessed June 2012.

¹⁸ <http://wiki.gephi.org/index.php/SemanticWebImport>

¹⁹ <https://gephi.org/>

²⁰ <http://www.semanticwebservices.org/enalgae>