

LinkLion: A Link Repository for the Web of Data

Markus Nentwig¹, Tommaso Soru², Axel-Cyrille Ngonga Ngomo², and Erhard Rahm¹

¹ Database Group, Department of Computer Science, University of Leipzig

² AKSW, Department of Computer Science, University of Leipzig

Abstract. Links between knowledge bases build the backbone of the Web of Data. Consequently, numerous applications have been developed to compute, evaluate and infer links. Still, the results of many of these applications remain inaccessible to the tools and frameworks that rely upon it. We address this problem by presenting LinkLion, a repository for links between knowledge bases. Our repository is designed as an open-access and open-source portal for the management and distribution of link discovery results. Users are empowered to upload links and specify how these were created. Moreover, users and applications can select and download sets of links via dumps or SPARQL queries. Currently, our portal contains 12.6 million links of 10 different types distributed across 3184 mappings that link 449 datasets. In this demo, we will present the repository as well as different means to access and extend the data it contains. The repository can be found at <http://www.linklion.org>.

1 Introduction

In addition to being central for question answering across several datasets, links also play a key role in various other domains such as data fusion and federated SPARQL queries. It is a well-known problem that links make up less than 3% of the RDF triples on the Web of Data [4]. This problem is being addressed by link discovery and ontology matching tools and frameworks [3,2]. However, due to the architectural choices behind the Web of Data, the results of a link discovery (LD) framework cannot be added directly to the datasets involved in the link discovery process. Further, the direct addition of links to a knowledge base fails to provide means to track the source of these links for later reference. Moreover, the availability of some endpoints still remains a major issue,³ making the direct addition of linking results to some endpoints unattractive.

We address these drawbacks by presenting the open-source link repository LINKLION. The main goal of LINKLION is to facilitate the publication, retrieval and use of links between knowledge bases. Our repository thus provides dedicated functionality for the upload, storage, querying and download of large sets of links. Currently, it contains 63 million triples which describe 12.6 million links of 10 different types (e.g., `owl:sameAs`, `dbo:spokenIn`, `foaf:made`, `spatial:P`)⁴ distributed on 3184 mappings that link 449 datasets. These links were retrieved from the Web as well as computed by tools such as LIMES [3] and Silk [6]. Our repository provides a SPARQL

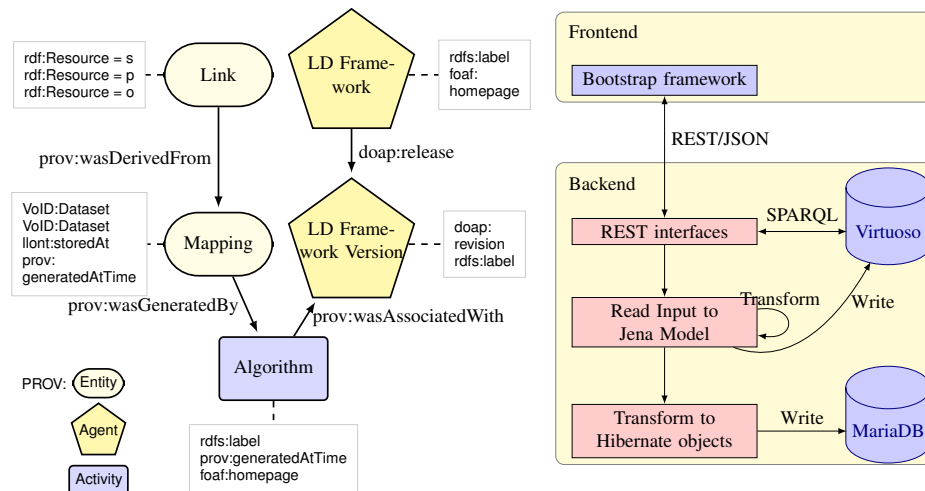
³ <http://labs.mondeca.com/sparqlEndpointsStatus.html>

⁴ We used the prefixes available at <http://prefix.cc>.

query interface as well as commodity interfaces to access the mappings. In contrast to other portals such as BioPortal⁵, LINKLION focuses exclusively on links and provides dedicated functionality for manipulating them. Moreover, we do not limit ourselves to a single domain such as the life sciences. In the following, we give a brief overview of the repository and show the use cases that will be presented during the demo. The repository can be accessed at <http://www.linklion.org>. The code of the repository is available at <http://github.com/AKSW/LinkingLodPortal>. The SPARQL endpoint can be found at <http://www.linklion.org:8890/sparql>.

2 Implementation

An overview of LINKLION's architecture is given in Figure 1b. The back end consists of a triple store in which we save data according to the vocabulary shown in Figure 1a.



(a) Overview of the LINKLION ontology. New classes such as Link, Mapping, Algorithm and LD Framework to store the mappings in the Virtuoso and are specified as subclasses of the PROV vocabulary. (b) Visualization of front and back end to store the mappings in the Virtuoso and MariaDB.

Fig. 1: Overview LINKLION ontology and architecture.

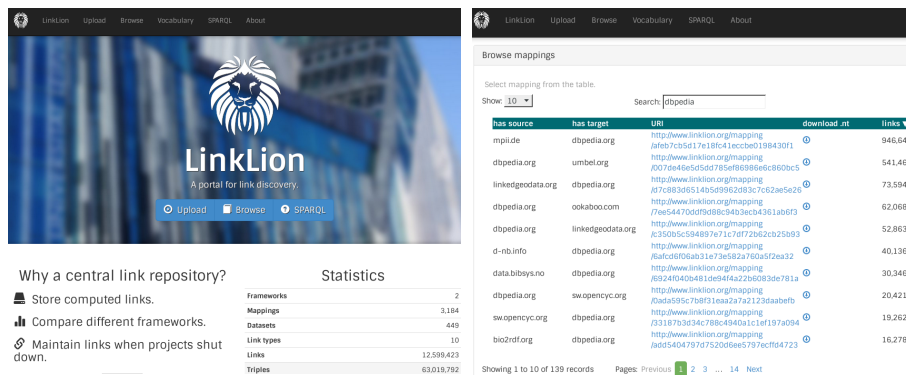
The ontology⁶ was designed with usability and reuse in mind. Especially, we wanted to allow end users of the portal to select dedicated portions of certain mappings at will. This meant designing an ontology that allowed amongst others (1) retrieving all links that pertain to a particular resource or set of resources, (2) gathering all mappings between datasets of interest as well as (3) getting aggregated information on how particu-

⁵ <http://www.bioontology.org/BioPortal>

⁶ Available at <http://www.linklion.org/ontology>.

lar links came about. We implemented this vision by storing the output of a link discovery tool under an instance of the mapping class. Individual mappings can be described by metadata including the datasets that they link, the tool (incl. a version number) used to generate the links and the creation date of the mapping. We refrained from using blank nodes for links. Instead, we gave each link a unique ID. Note that we reused existing vocabularies (especially PROV⁷, VoID⁸ and DOAP⁹) as much as we could. The use of a triple store pays off as end users can choose to provide more metadata such as the link specification used or parameters of the algorithm they used to discover the link without us having to alter our schema. For the sake of scalability, we yet also provide the core of the data in the triple store as SQL dump. The functionality of the back end is exposed by RESTful interfaces, which allow a programmatic access to LINKLION from code written in virtually any modern programming language.

The front end of our repository provides an easy way to use some of the functionality of LINKLION (see Figure 2a). First, it allows users to upload new mappings. Users are asked to provide a source file in the N-Triples format¹⁰, well as the algorithm used within this framework have to be provided (note that we consider humans to also be linking frameworks). The data (and especially the mapping) given by the user is then checked for consistency and uploaded into the underlying triple store. The content of the triple store can be browsed directly from the web page (see Figure 2b). Especially, the front end includes search functionality and pagination which allow end users to search for mappings that link to or from a dataset of interest. The upload and browsing functionality will be presented during the demo.



(a) LINKLION Homepage with statistics regarding the content. (b) Mapping Browser. The search for DBpedia returns 139 mappings.

Fig. 2: Front-end views.

⁷ <http://www.w3.org/TR/prov-o/>

⁸ <http://www.w3.org/TR/void/>

⁹ <https://github.com/edumbill/doap/>

¹⁰ <http://www.w3.org/2001/sw/RDFCore/ntriples/>

3 Use Cases

In this section, we present and motivate a selection of use cases that will be presented during the demo session.

3.1 Gather all links and mappings to a given resource

Gathering and fusing all information on a resource of interest is of central importance to applications such as Question Answering systems, Linked Data Browsers and Quality Assessment tools. LINKLION allows to gather all links pertaining to a particular resource (`dbpedia:Thailand` in our example) by means of the following SPARQL query. By using this information, novel repair-based algorithms for link discovery such as COLIBRI can find errors or inconsistencies in the data [5].

```
SELECT ?link WHERE { { ?link rdf:subject dbpedia:Thailand }
                    UNION { ?link rdf:object dbpedia:Thailand } }
```

The portal also allows gather all mappings that contain links pertaining to a particular resource, e.g.. `dbpedia:Thailand`, as shown below.

```
SELECT DISTINCT ?mapping WHERE { ?link prov:wasDerivedFrom ?mapping .
                                  { ?link rdf:subject dbpedia:Thailand }
                                  UNION { ?link rdf:object dbpedia:Thailand } }
```

3.2 Get support for a link

Ensemble learning techniques have been shown to improve the results of manifold machine-learning applications such as Named Entity Recognition frameworks. Our repository facilitates the use of ensemble learning for combining the results of different link discovery tools. Especially, LINKLION allows us to retrieve (if any) the list of mappings that contain a given link, as well as the algorithms and the frameworks that generated it. In the following example, the support for `dbpedia:Thailand owl:sameAs <http://sws.geonames.org/1605651/>` is queried.

```
SELECT ?mapping ?algorithm ?framework WHERE {
  ?mapping prov:wasGeneratedBy ?algorithm .
  ?algorithm prov:wasAssociatedWith ?framework .
  ?link prov:wasDerivedFrom ?mapping ;
        rdf:predicate owl:sameAs .
  { ?link rdf:subject dbpedia:Thailand;
    rdf:object <http://sws.geonames.org/1605651/> }
  UNION { ?link rdf:object dbpedia:Thailand;
          rdf:subject <http://sws.geonames.org/1605651/> } }
```

3.3 Link Composition

With the growth of the Linked Data Web, it becomes ever more important to regard link discovery as a holistic process that goes beyond linking a pair of knowledge bases. Algorithms based on composition can exploit sequences of links to enrich their mapping composition graphs [1]. Moreover, algorithms which link several knowledge bases at

the same time [5] can achieve higher accuracies. By using LINKLION, composition and concurrent linking algorithms are now enabled to gather the data they require without having to manage all the links by themselves. In the query below, all resources related to `dbpedia:Thailand` over two links are retrieved from the repository.

```
SELECT DISTINCT ?resource WHERE { {  
  ?link rdf:subject dbpedia:Thailand ; rdf:object ?x .  
  ?link2 rdf:subject ?x ; rdf:object ?resource }  
  UNION { ?link rdf:object dbpedia:Thailand ; rdf:subject ?x .  
          ?link2 rdf:object ?x ; rdf:subject ?resource } }
```

4 Summary

This paper presents LINKLION, a repository for links between knowledge bases of the Web of Data. The repository enables users to upload results of a link discovery process and furthermore allows them to add information on how the results were created. LINKLION therefore provides management and distribution capabilities through a both open-access and open-source web interface. Resulting sets of links can be reviewed in the portal and via SPARQL queries, additionally the results can be downloaded via dumps. In future work, LINKLION will be extended to support a closely collaboration with external link discovery frameworks and applications.

References

1. Michael Hartung, Anika Groß, and Erhard Rahm. Composition Methods for Link Discovery. In *BTW*, 2013.
2. Toralf Kirsten, Anika Gross, Michael Hartung, and Erhard Rahm. GOMMA: a component-based infrastructure for managing and analyzing life science ontologies and their evolution. *J. Biomedical Semantics*, 2:6, 2011.
3. Axel-Cyrille Ngonga Ngomo. On Link Discovery using a Hybrid Approach. *J. Data Semantics*, 1(4):203–217, 2012.
4. Axel-Cyrille Ngonga Ngomo and Sören Auer. LIMES: a time-efficient approach for large-scale link discovery on the web of data. In *Proc. of the 22nd international joint conference on Artificial Intelligence - Volume Three, IJCAI'11*, pages 2312–2317. AAAI Press, 2011.
5. Axel-Cyrille Ngonga Ngomo, Mohamed Ahmed Sherif, and Klaus Lyko. Unsupervised Link Discovery Through Knowledge Base Repair. In *Extended Semantic Web Conference (ESWC 2014)*, 2014.
6. Julius Volz, Christian Bizer, Martin Gaedke, and Georgi Kobilarov. Silk - A Link Discovery Framework for the Web of Data. In Christian Bizer, Tom Heath, Tim Berners-Lee, and Kingsley Idehen, editors, *LDOW*, volume 538 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2009.