



GEOLIFT - Spatial mapping framework for enriching RDF datasets with Geo-spatial information.

Abstract:

This manual presents the spatial mapping component dubbed GEOLIFT. The goal of GEOLIFT is to enrich RDF datasets with geo-spatial information. To achieve this goal, GeoLit relies on three atomic modules based on dereferencing, linking and NLP. GEOLIFT was implemented in Java, is open-source and can be accessed at <https://github.com/GeoKnow/GeoLift/>.

Contents

1	Introduction	4
2	Assumptions	4
3	Technical Approach	5
3.1	Architecture	5
3.2	Using Dereferencing	6
3.3	Using Linking	7
3.4	Using Named Entity Recognition	9
4	Developers' Manual	10
4.1	Dereferencing module	11
4.1.1	Input	11
4.1.2	Output	12
4.1.3	Process	12
4.1.4	Sample code to run the module	12
4.2	Linking module	14
4.2.1	Input	14
4.2.2	Output	14
4.2.3	Process	14
4.2.4	Sample code to run the module	15
4.3	NLP module	17
5	Conclusions	19

1 Introduction

Manifold RDF data contain implicit references to geographic data. For example, music datasets such as *Jamendo* include references to locations of record labels, places where artists were born or have been, etc. The aim of the spatial mapping component, dubbed GEOLIFT, is to retrieve this information and make it explicit. In the following, we begin by presenting the basic assumptions that influence the development of the first component of GEOLIFT. Then, we present the technical approach behind GEOLIFT. Finally, we present the detailed developers' manual of GEOLIFT.

2 Assumptions

Geographical information can be mentioned in three different ways within Linked Data:

1. *Through dereferencing*: Several datasets contain links to datasets with explicit geographical information such as DBpedia or LinkedGeoData. For example, in a music dataset, one might find information such as `http://example.org/Leipzig`
`owl:sameAs`
`http://dbpedia.org/resource/Leipzig`.

We call this type of reference *explicit*. We can now use the semantics of RDF to fetch geographical information from DBpedia and attach it to the resource in the other ontology as `http://example.org/Leipzig` and `http://dbpedia.org/resource/Leipzig` refer to the same real-world object.

2. *Through linking*: It is known that the Web of Data contains an insufficient number of links. The latest approximations suggest that the Linked Open Data Cloud alone consists of 31+ billion triples but only contains approximately 0.5 billion links (i.e., less than 2% of the triples are links between knowledge bases). The second intuition behind our approach is thus to use link discovery to map resources in an input knowledge base to resources in a knowledge that contains explicit geographical information. For example, given a resource `http://example.org/Athen`, GEOLIFT should aim to find a resource such as `http://dbpedia.org/resource/Athen` to map it with. Once having established the link between the two resources, GEOLIFT can then resolve to the approach defined above.

3. *Through Natural Language Processing*: In some cases, the geographic information is hidden in the objects of data type properties. For example, some datasets contain biographies, textual abstracts describing resources, comments from users, etc. The idea here is to use this information by extracting Named Entities and keywords using automated Information Extraction techniques. Semantic Web Frameworks such as FOX¹ have the main advantage of providing URIs for the keywords and entities that they detect. These URIs can finally be linked with the resources to which the datatype properties were attached. Finally, the geographical information can be dereferenced and attached to the resources whose datatype properties were analyzed.

The idea behind GEOLIFT is to provide a generic architecture that contains means to exploit these three characteristics of Linked Data. In the following, we present the technical approach underlying GEOLIFT.

3 Technical Approach

3.1 Architecture

GEOLIFT was designed to be a modular tool which can be easily extended and re-purposed. In its first version, it provides two main types of artifacts:

1. *Modules*: These artifacts are in charge of generating geographical data based on RDF data. To this aim, they implement the three intuitions presented above. The input for such a module is an RDF dataset (in Java, a *Jena Model*). The output is also an RDF dataset enriched with geographical information (in Java, an enriched *Jena Model*). Formally, a module can thus be regarded as a function $\mu : \mathcal{R} \rightarrow \mathcal{R}$, where \mathcal{R} is the set of all RDF datasets.
2. *Operators*: The idea behind operators is to enable users to define a workflow for processing their input dataset. Thus, in case a user knows the type of enrichment that is to be carried out (using linking and then links for example), he can define the sequence of modules that must be used to process his dataset. Note that the format of the input and output of modules is identical. Thus, the user is empowered to create workflows of arbitrary complexity by simply connecting modules. Formally, an operator can be regarded as a function $\varphi : \mathcal{R} \cup \mathcal{R}^2 \rightarrow \mathcal{R} \cup \mathcal{R}^2$.

¹<http://fox.aksw.org>

The corresponding architecture is shown in Figure 1. The input layer allows reading RDF in different serializations. The enrichment modules are in the second layer and allow adding geographical information to RDF datasets by different means. The operators (which will be implemented in the future version of GEOLIFT) will combine the enrichment modules and allow defining a workflow for processing information. The output layer serializes the results in different format. The enrichment procedure will be monitored by implementing a controller, which will be added in the future version of GEOLIFT.

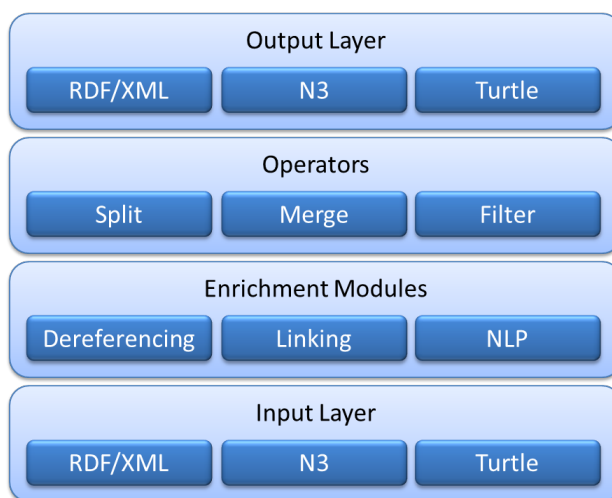


Figure 1: Architecture of GEOLIFT

In the following, we present the implementation of the three intuitions presented above in GEOLIFT.

3.2 Using Dereferencing

For datasets which contain `owl:sameAs` links, we deference all links from the dataset to other datasets by using a content negotiation on HTTP as shown in Figure 2. This returns a set of triples that needs to be filtered for relevant geographical information. Here, we use a predefined list of attributes that links to geographical information. Amongst others, we look for `geo:lat`, `geo:long`, `geo:lat_long`, `geo:line` and `geo:polygon`. The list of retrieved property values can be configured.



Figure 2: Content Negotiation as used by GEOLIFT (courtesy of W3C)

3.3 Using Linking

As pointed out before, links to geographical resources do not occur in several knowledge bases. Here, we rely on the metrics implemented in the LIMES framework² [5, 4, 6] to link the resources in the input dataset with geographical datasets. LIMES, the **Link** Discovery Framework for **Metric** Spaces, is a framework for discovering links between entities contained in Linked Data sources. LIMES is a hybrid framework [4] that combines the mathematical characteristics of metric spaces as well prefix-, suffix- and position filtering to compute pessimistic approximations of the similarity of instances. These approximations are then used to filter out a large amount of those instance pairs that do not suffice the mapping conditions. By these means, LIMES can reduce the number of comparisons needed during the mapping process by several orders of magnitude and complexity without losing a single link. The architecture of LIMES is shown in Figure 3

Linking using LIMES [4, 3] can be achieved in three ways:

1. *Manually*, by the means of a link specification [4], which is an XML-description of (1) the resource in the input and target datasets that are to be linked and (2) of the similarity measure that is to be employed to link these datasets.
2. *Semi-automatically* based on active learning [7, 8, 9]. Here, the idea is that if the user is not an expert and thus unable to create a link specification, he can simply provide the framework with positive and negative examples iteratively. Based on these examples, LIMES can compute links for mapping resources with high accuracy.
3. *Automatically* based on unsupervised machine learning. Here, the user can simply specify the sets of resources that are to be linked with each other. LIMES implements both a deterministic and non-deterministic

²<http://limes.sf.net>

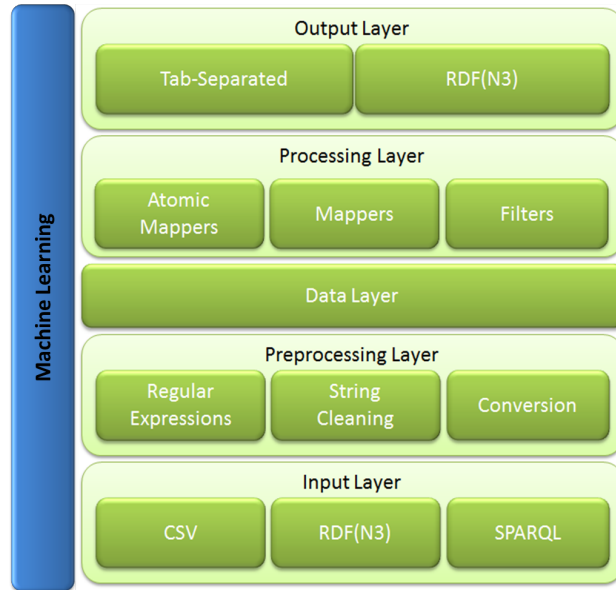


Figure 3: Architecture of LIMES

machine-learning approaches that optimize a pseudo-F-measure to create a one-to-one mapping.

The techniques implemented by LIMES can be accessed via the SAIM user interface³, of which a screenshot is shown in Figure 4.

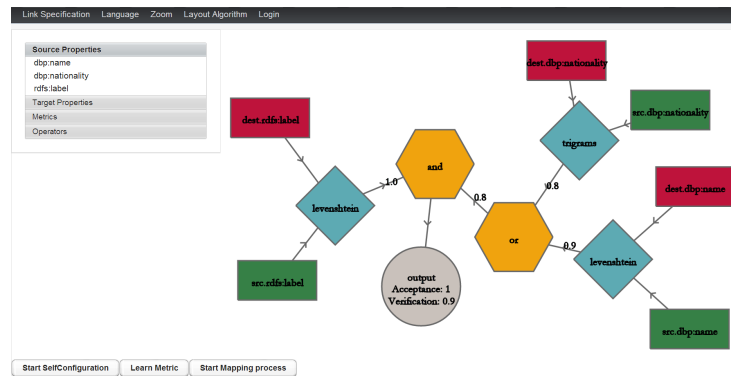


Figure 4: Screenshot of SAIM

³<http://saim.aksw.org>

3.4 Using Named Entity Recognition

The geographical information hidden in datatype properties is retrieved by using Named Entity Recognition. In the first version of GEOLIFT, we rely on the FOX framework. The FOX framework is a stateless and extensible framework that encompasses keyword extraction and named entity recognition. Its architecture consists of three layers as shown in Figure 5.

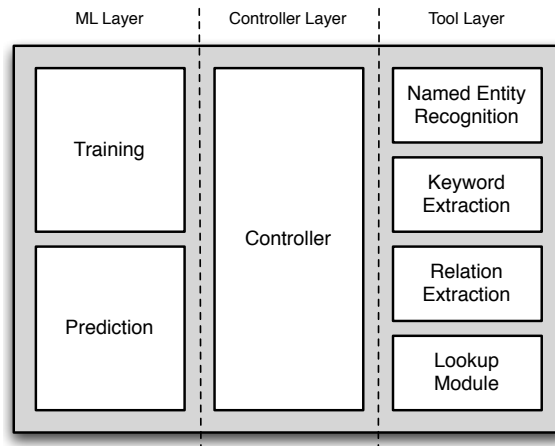


Figure 5: Architecture of the FOX framework.

FOX takes text or HTML as input. Here we use the objects of datatype properties, i.e., plain text. This data is sent to the *controller layer*, which implements the functionality necessary to clean the data, i.e., remove HTML and XML tags as well as further noise. Once the data has been cleaned, the controller layer begins with the orchestration of the tools in the *tool layer*. Each of the tools is assigned a thread from a thread pool, so as to maximize usage of multi-core CPUs. Every thread runs its tool and generates an event once it has completed its computation. In the event that a tool does not complete after a set time, the corresponding thread is terminated. So far, FOX integrates tools for KE, NER and RE. The KE is realized by tools such as KEA⁴ and the Yahoo Term Extraction service⁵. In addition, FOX integrates the Stanford Named Entity Recognizer⁶ [2], the Illinois Named Entity Tagger⁷ [10] and Alchemy⁸ for NER.

The results from the tool layer are forwarded to the *prediction module* of

⁴<http://www.nzdl.org/Kea/>

⁵<http://developer.yahoo.com/search/content/V1/termExtraction.html>

⁶<http://nlp.stanford.edu/software/CRF-NER.shtml>

⁷http://cogcomp.cs.illinois.edu/page/software_view/4

⁸<http://www.alchemyapi.com>

the *machine-learning layer*. The role of the prediction module is to generate FOX's output based on the output the tools in FOX's backend. For this purpose, it implements several ensemble learning techniques [1] with which it can combine the output of several tools. Currently, the prediction module carries out this combination by using a feed-forward neural network. The neural network inserted in FOX was trained by using 117 news articles. It reached 89.21% F-Score in an evaluation based on a ten-fold-cross-validation on NER, therewith outperforming even commercial systems such as Alchemy.

Once the neural network has combined the output of the tool and generated a better prediction of the named entities, the output of FOX is generated by using the vocabularies shown in Figure 6. These vocabularies extend the two broadly used vocabularies Annotea⁹ and Autotag¹⁰. In particular, we added the constructs explicated in the following:

- `scms:beginIndex` denotes the index in a literal value string at which a particular annotation or keyphrase begins;
- `scms:endIndex` stands for the index in a literal value string at which a particular annotation or keyphrase ends;
- `scms:means` marks the URI assigned to a named entity identified for an annotation;
- `scms:source` denotes the provenance of the annotation, i.e., the URI of the tool which computed the annotation or even the system ID of the person who curated or created the annotation and
- `scmsann` is the namespace for the annotation classes, i.e, location, person, organization and miscellaneous.

4 Developers' Manual

GEOLIFT contains three basic *Java* packages:

- `IO package` which deals with input/output operations using the *Reader* and *Writer* classes.
- `Operators package` will contains implementation of different operators like *merge*, *split*, *filter*, ...

NOTE: Operators package will be implemented in future version of GEOLIFT .

⁹<http://www.w3.org/2000/10/annotation-ns#>

¹⁰<http://commontag.org/ns#>

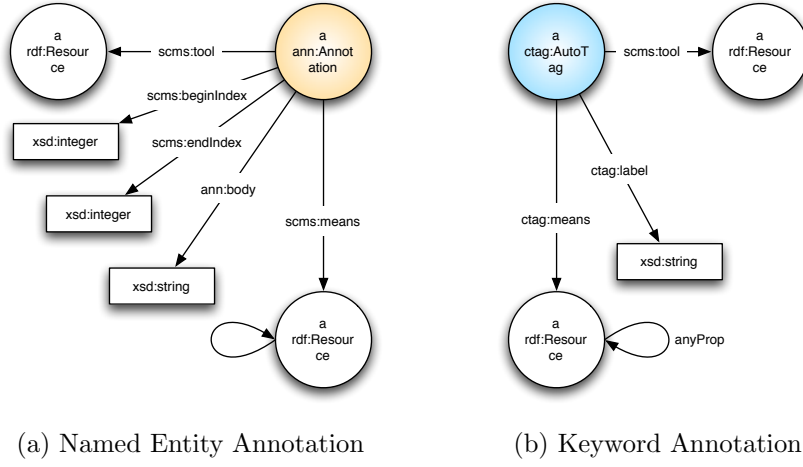


Figure 6: Vocabularies used by FOX for representing named entities (a) and keywords (b)

- **Modules package** contains the **GeoLiftModule** interface which is implemented by the basic classes: **Dereference** class which handles dereferencing geographical information extending process. **Linking** class which handles linking geographical information extending processes. and **NLP** class which handles named entity extraction process.

All modules implement the **GeoLiftModule** interface's two methods **getParameters()** and **process()**. Method **getParameters()** returns set of input parameters was given to the implementing module. The **process()** method takes input model to be enriched in addition to parameters used while processing in the form of Map structure then use them to start the implementing module process.

4.1 Dereferencing module

4.1.1 Input

- **Data model** contains the triples of the dataset to be enriched (This data model can be an output from previous stage or can be loaded from file directly before using the module).
- **Predicates list** list of interested predicates to be added as enrichment to the data model.

Table 1 provides details about the **Dereferencing** module's parameters

Table 1: Dereferencing parameters description

Parameter Name	Default value	Description
<code>model</code>	<code>null</code>	Original Model to be enriched with geographical information.
<code>predicates list</code>	<code>null</code>	List of interesting predicates to enrich the model with them and their Objects' values, e.g. <code>http://www.w3.org/2003/01/geo/wgs84_pos#lat</code> . The list of predicates is given in form of Map structure where the key and value of each entry are the predicate itself.

4.1.2 Output

- **Data model** data model enriched with additional geographic information using the interested predicates given. The information are represented as triples include interested predicate and its value.

4.1.3 Process

In this module, a dataset model is given and a list of interested predicates as inputs. The purpose is to extend the model's geographical information by set of information through the given predicates. This is done by iterating over the model's resources (dubbed as original resources) and for each original resource an extraction of the predicates' values (objects) that are in the form of URI is performed. These URIs (dubbed as dereferenced resources) are more filtered to be the resources used in dbpedia. The dereferenced resources are supposed to a dereference operation in order to find the interested predicates list for them. Such predicates and their objects' values are fetched and added to the the original resource to extend its information.

4.1.4 Sample code to run the module

This is sample code shows how to use the URIDereference module:

```
void main()
{
    /*
```

```

Load list of needed parameters usinf function
    getConfigured(). It is written to read from a file list
    of parameters. The path to the file is given in args[0].
The parameters include:
File with to load the original dataset from it.
List of interested predicates list
*/
List<String> configurations= getConfigured(args[0]);
//Load model with required dataset
Model
    model=org.aksw.geolift.io.Reader.readModel(configurations.get(0));
    //model is loaded with dataset from specified file
List<String> predicates=configurations.subList(1,
    configurations.size()); //load targeted predicates to be
    added to enrich information in dataset
//Pre-processing step : Collect list of targeted predicates
    into Map as process method get them as Map structure
Map<String, String> parameters= new HashMap<String,
    String>();
for (String predicate : predicates)
{
    parameters.put(predicate, predicate);
}
//Create the 'Dereference' object
URIDereferencing u = new URIDereferencing();
//Run the dereferencing process it requires model contains
    the dataset and list of targeted predicates to enrich
    the model
Model resultedModel = u.process(model, parameters);
    try
    {
        //Write the enriched model into .ttl file
        org.aksw.geolift.io.Writer.writeModel(resultedModel,
            "TTL",
            "src/main/resources/dereferencing/DereferencingEnriched.ttl");
    } catch (IOException e)
    {
        // TODO Auto-generated catch block
        e.printStackTrace();
    }
}

```

4.2 Linking module

4.2.1 Input

- **Data model** contains the triples of the dataset to be enriched (This data model can be an output from previous stage or can be loaded from file directly before using the module).
- **Parameters list** is a list of parameters that will be used during the process. These parameters include:

Specification file path, the path to the spec.xml file contains the linking specifications

Links file path, the path to the file contains the resulted links

URI position, represents the original model's URI position as source/left or target/right in the linking specifications.

The parameters, other than the data model parameter, are collected in Map structure form where each entry's value in the Map represents the parameter itself. Table 2 provides details about the **Linking** module's parameters .

4.2.2 Output

- **Data model** data model enriched with additional geographic information URIs represented in owl:sameAs predicates.

4.2.3 Process

In this module an input model is given and list of parameters for used files during the process. The process starts by generating links between the dataset model and another dataset as second partner. This is done using LIMES interlinking tool by specifying the linking specification file given as parameter. The links are generated in accept.nt file that is used after to combine such links with their original resources in our dataset model as owl:sameAs predicate's objects. The result is a dataset model enriched with geographical information links.

Another forward step is to feed the Dereference module with the resulted enriched model from Linking module as input. The previously generated links in the model in addition to other Objects in the URIs form will be dereferenced adding more and detailed geographical information.

Table 2: Linking parameters description

Parameter Name	Default value	Description
model	null	Original Model to be enriched with geographical information.
Specification file	N/A	The path to specification file used for linking process, the original dataset to be enriched must be on the source dataset , e.g. <code>linkingModuleData/linking/spec.xml</code> . The parameter's entry in th Map structure has key 'specFilePath'.
Links file	N/A	The path to links file resulted from the linking process. This file's path is the same as the one specified in LIME's specifications file as output file, e.g. <code>linkingModuleData/linking/links.nt</code> . The parameter's entry in th Map structure has key 'linksFilePath'.
Original URI position	N/A	represents the original model's URI position as source/left or target/right in the linking specifications. Its value is either 'source' or 'target'. The parameter's entry in th Map structure has key 'linksPart'.

4.2.4 Sample code to run the module

This is sample code shows how to use the Linking module:

```

public static void main(String[] args)
{
Map<String, String> parameters=new HashMap<String, String>();

// The path to the dataset file to be loaded
parameters.put("datasetFilePath",args[0]);

//The path to the spec.xml file contains the linking specifications
parameters.put("specFilePath",args[1]);

// The path to the file contains the resulted links

```

```

parameters.put("linksFilePath",args[2]);// The path to the file
contains the resulted links

/*The position of the Original URI to be enriched in the generated
links
(right side/source or left side/target), so the otherside is the
enriching
owl:sameAs link to be added to it */
parameters.put("linksPart",args[3]);

//Load model to be enriched from the specified file path
Model
    model=org.aksw.geolift.io.Reader.readModel(parameters.get("datasetFilePath"));

//Create linking object
Linking l= new Linking();

//Run the the Linking Model to enrich the original model with
links generated
model=l.process(model, parameters);
try
{
    org.aksw.geolift.io.Writer.writeModel(model, "TTL",
        "src/main/resources/linking/datasetUpdated.nt");
} catch (IOException e) {
    // TODO Auto-generated catch block
    e.printStackTrace();
}
//Further step is combining the Linking module with Dereferncing
module
URIDereferencing d= new URIDereferencing();

//prepare Map of interested predicates to enrich by
Map<String, String> PredicatesParameters= new HashMap<String,
String>();
PredicatesParameters.put("http://www.w3.org/2003/01/geo/wgs84_pos#lat",
"http://www.w3.org/2003/01/geo/wgs84_pos#lat");
PredicatesParameters.put("http://www.w3.org/2003/01/geo/wgs84_pos#long",
"http://www.w3.org/2003/01/geo/wgs84_pos#long");
PredicatesParameters.put("http://www.w3.org/2003/01/geo/wgs84_pos#geometry",
"http://www.w3.org/2003/01/geo/wgs84_pos#geometry");
PredicatesParameters.put("http://www.w3.org/2003/01/geo/wgs84_pos#lat_long",
"http://www.w3.org/2003/01/geo/wgs84_pos#lat_long");

```



```

PredicatesParameters.put("http://www.w3.org/2003/01/geo/wgs84_pos#line",
    "http://www.w3.org/2003/01/geo/wgs84_pos#line");
PredicatesParameters.put("http://www.w3.org/2003/01/geo/wgs84_pos#polygon",
    "http://www.w3.org/2003/01/geo/wgs84_pos#polygon");
//Run the Dereferencing operation
model=d.process(model, PredicatesParameters);
try
{
    //Write the resulted double enriched model in .ttl file
    org.aksw.geolift.io.Writer.writeModel(model, "TTL",
        "src/main/resources/linking/datasetLinkingDereferenced.nt");
} catch (IOException e) {
    // TODO Auto-generated catch block
    e.printStackTrace();
}

System.out.println("Finished");
}

```

4.3 NLP module

The `getParameters()` method of the NLP module returns a list of parameters, which can be set by the user to provide custom control of the *Named entity extraction* provided by the implemented *FOX framework*.

The `process()` method of the NLP module takes as input a *Jena* model and a `Map` of different parameters, and it generates also a *Jena* model as output. Table 3 provides details about the NLP module's parameters .

Table 3: NLP parameters description

Parameter Name	Default value	Description
literalProperty	Top Ranked	Literal property used by FOX for NER, if not set the top ranked property is pecked by LiteralPropertyRanker , which ranks the lateral properties of a model according to the average size of each literal property divided by the number of instances of such property.
useFoxLight	false	Use the light version of FOX, setting it generates faster execution time but less accurate results)
askEndPoint	false	Ask the <i>DBpedia</i> endpoint for each location returned by FOX (setting it generates slower execution time but more accurate results)
foxType	TEXT	FOX input type : { text url }
foxTask	NER	FOX task : {NER} for Named Entity Recognition
foxInput	" "	FOX actual input as text or an URL
foxOutput	TURTLE	FOX output format: { JSONLD N3 N-TRIPLE RDF/{ JSON XML XML-ABBREV } TURTLE }
foxUseNif	false	FOX generates NIF: { true false }
foxReturnHtml	false	FOX returns HTML: { true false }

5 Conclusions

In this manual, we presented the GEOLIFT component for enriching RDF datasets with geo-spatial data. In future work, we aim to implement a graphical user interface on top of GEOLIFT to enable users to specify their workflows graphically. Moreover, we aim to implement workflow checking functionality.

References

- [1] Thomas G. Dietterich. Ensemble methods in machine learning. In *MCS*, pages 1–15, London, UK, 2000. Springer-Verlag.
- [2] J. Finkel, T. Grenager, and C. Manning. Incorporating non-local information into information extraction systems by gibbs sampling. In *ACL*, pages 363–370, 2005.
- [3] Axel-Cyrille Ngonga Ngomo. Link discovery with guaranteed reduction ratio in affine spaces with minkowski measures. In *Proceedings of ISWC*, 2012.
- [4] Axel-Cyrille Ngonga Ngomo. On link discovery using a hybrid approach. *Journal on Data Semantics*, 1:203 – 217, December 2012.
- [5] Axel-Cyrille Ngonga Ngomo and Sren Auer. A time-efficient approach for large-scale link discovery on the web of data. In *IJCAI*, page 2011, 2011.
- [6] Axel-Cyrille Ngonga Ngomo, Lars Kolb, Norman Heino, Michael Hartung, Sören Auer, and Erhard Rahm. When to reach for the cloud: Using parallel hardware for link discovery. In *Proceedings of ESCW*, 2013.
- [7] Axel-Cyrille Ngonga Ngomo, Jens Lehmann, Sören Auer, and Konrad Höffner. Raven – active learning of link specifications. In *Proceedings of OM@ISWC*, 2011.
- [8] Axel-Cyrille Ngonga Ngomo and Klaus Lyko. Eagle: Efficient active learning of link specifications using genetic programming. In *Proceedings of ESWC*, 2012.
- [9] Axel-Cyrille Ngonga Ngomo, Klaus Lyko, and Victor Christen. Coala – correlation-aware active learning of link specifications. In *Proceedings of ESWC*, 2013.

- [10] Lev Ratinov and Dan Roth. Design challenges and misconceptions in named entity recognition. In *CONLL*, pages 147–155, 2009.