

# LDIF - Linked Data Integration Framework

Andreas Schultz<sup>1</sup>, Andrea Matteini<sup>2</sup>, Robert Isele<sup>1</sup>, Christian Bizer<sup>1</sup>, and  
Christian Becker<sup>2</sup>

1. Web-based Systems Group, Freie Universität Berlin, Germany  
a.schultz@fu-berlin.de, mail@robertisele.com, chris@bizer.de

2. MediaEvent Services GmbH & Co. KG, Germany  
a.matteini@mes-info.de, c.becker@mes-info.de

**Abstract.** The LDIF - Linked Data Integration Framework can be used within Linked Data applications to translate heterogeneous data from the Web of Linked Data into a clean local target representation while keeping track of data provenance. LDIF provides an expressive mapping language for translating data from the various vocabularies that are used on the Web into a consistent, local target vocabulary. LDIF includes an identity resolution component which discovers URI aliases in the input data and replaces them with a single target URI based on user-provided matching heuristics. For provenance tracking, the LDIF framework employs the Named Graphs data model. This paper describes the architecture of the LDIF framework and presents a performance evaluation of a life science use case.

**Keywords:** Linked Data, Data Integration, Data Translation, Identity Resolution

## 1 Motivation

The Web of Linked Data grows rapidly but the development of Linked Data applications is still cumbersome due to the heterogeneity of the Web of Linked Data [1]. Two major roadblocks for Linked Data applications are vocabulary heterogeneity and URI aliases. A fair portion of the Linked Data sources reuse terms from widely-deployed vocabularies to describe common types of entities such as people, organizations, publications and products. For more specialized, domain-specific entities, such as genes, pathways, descriptions of subway lines, statistical and scientific data, no wide-spread vocabulary agreement has evolved yet. Data sources in these domains thus use proprietary terms [2]. A second problem are identity links. Some data sources set *owl:sameAs* links pointing at data about the same entity in other data sources. Many other data sources do not [2].

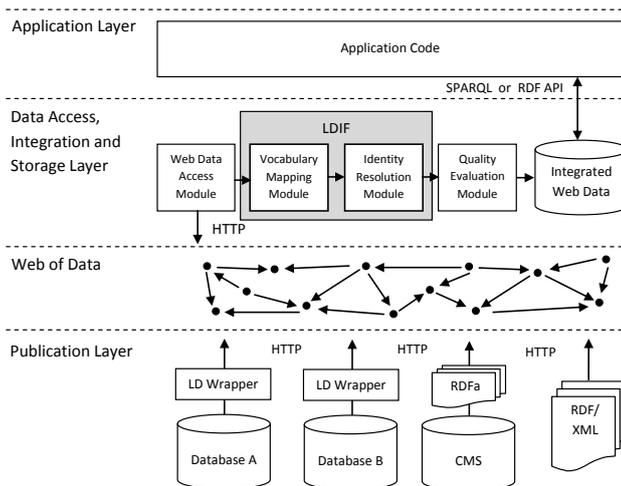
In contrast to the heterogeneity of the Web, it is beneficial in the application context to have all data describing one class of entities being represented using

the same vocabulary. Instead of being confronted with URI aliases which refer to data that might or might not describe the same entity, Linked Data applications would prefer all triples describing the same entity to have the same subject URI as this eases many application tasks including querying, aggregation and visualization.

In order to ease using Web data in the application context, it is thus advisable to translate data to a single target vocabulary (vocabulary mapping) and to replace URI aliases with a single target URI on the client side (identity resolution), before doing any more sophisticated processing.

There are various open source tools available that help application developers with either data translation or identity resolution. But to date, there are hardly any integrated frameworks available that cover both tasks. With LDIF, we try to fill this gap and provide an open-source Linked Data integration framework that provides for data translation and identity resolution while keeping track of data provenance.

Figure 1 shows the schematic architecture of a Linked Data application that implements the crawling/data warehousing pattern [1]. The figure highlights the steps of the data integration process that are currently supported by LDIF.



**Fig. 1.** Role of LDIF within the architecture of a Linked Data application.

The LDIF framework is implemented in Scala and can be downloaded from the project website<sup>1</sup> under the terms of the Apache Software License. In the following, we explain the architecture of the LDIF framework and present a performance evaluation along the example of a life science use case.

<sup>1</sup> <http://www4.wiwiw.fu-berlin.de/bizer/ldif/>

## 2 Architecture

The LDIF framework consists of a runtime environment and a set of pluggable modules. The runtime environment manages the data flows between the modules. The pluggable modules are organized as data access components, data transformation components and data output components. So far, we have implemented the following modules:

### 2.1 Data Access: N-Quads Loader

The current version of LDIF expects input data to be represented as Named Graphs and be stored in N-Quads format. The graph URI is used for provenance tracking. Provenance meta-information describing the graphs can be provided as part of the input data within a specific provenance graph. The name of this provenance graph can be set in the LDIF configuration file. LDIF does not make any assumptions about the provenance vocabulary that is used to describe the graphs, meaning that you can use your provenance vocabulary of choice.

### 2.2 Transformation: R2R Data Translation

LDIF employs the R2R Framework [3] to translate Web data that is represented using terms from different vocabularies into a single target vocabulary. Vocabulary mappings are expressed using the R2R Mapping Language. The language provides for simple transformations as well as for more complex structural transformations (1-to-n and n-to-1) and property value transformations such as normalizing different units of measurement or complex string manipulations. So-called modifiers make it possible to change the language tag or data type of a literal or the RDF node type (URI  $\leftrightarrow$  literal). The syntax of the R2R Mapping Language is very similar to the SPARQL query language, which eases the learning curve. The expressivity of the language enabled us to deal with all requirements that we have encountered so far when translating Linked Data from the Web into a target representation [3].

### 2.3 Transformation: Silk Identity Resolution

LDIF employs the Silk Link Discovery Framework [4] to find different URIs which identify the same real-world entity. Silk is a flexible identity resolution framework that allows the user to specify identity resolution heuristics using the declarative *Silk - Link Specification Language* (Silk-LSL). In order to specify the condition which must hold true for two entities to be considered a duplicate, the user may apply different similarity metrics, such as string, date or URI comparison methods, to multiple property values of an entity or related entities. The Link Specification Language provides a variety of data transformations to normalize the data prior to comparing it. The resulting similarity scores can be combined and weighted using various similarity aggregation functions.

Silk uses a novel blocking approach which removes definite non-duplicates early in the matching process, thereby significantly increasing its efficiency. For each set of duplicates which have been identified by Silk, LDIF replaces all URI aliases with a single target URI within the output data. In addition, it adds *owl:sameAs* links pointing at the original URIs, which makes it possible for applications to refer back to the original data sources on the Web. If the LDIF input data already contains *owl:sameAs* links, the referenced URIs are normalized accordingly.

## 2.4 Data Output: N-Quads Writer

The N-Quads writer dumps the final output of the integration work flow into a single N-Quads file. This file contains the translated versions of all graphs from the input graph set as well as the contents of the provenance graph.

## 2.5 Runtime Environment

The runtime environment manages the data flow between the modules and the caching of the intermediate results. In order to parallelize processing, data is partitioned into entities prior to supplying it to a transformation module. An entity represents a Web resource together with all data that is required by a transformation module to process this resource. Entities consist of one or more graph paths and include a provenance URI for each node. Each transformation module specifies which paths should be included into the entities it processes. By splitting the data set into fine-grained entities, LDIF is able to parallelize the workload on machines with multiple cores. In the next release, it will allow the workload to be parallelized on multiple machines using Hadoop.

## 3 Performance Evaluation

We evaluated the performance of LDIF using two life science data sets: KEGG GENES<sup>2</sup>, a collection of gene catalogs generated from publicly available resources, and UniProt<sup>3</sup>, a data set containing protein sequence, genes and functions.

We defined R2R mappings for translating genes, diseases and pathways from KEGG GENES and genes from UniProt into a proprietary target vocabulary<sup>4</sup>. The mappings employ complex structural transformations. The prevalent value transformations rely on regular expressions, e.g. for extracting an integer value from a URI, and modify the target data types. We defined Silk linkage rules for identifying equivalent genes in both datasets. For the benchmark, we generated subsets of both data sources together amounting to 25 million, 50 million, and

<sup>2</sup> <http://www.genome.jp/kegg/genes.html>

<sup>3</sup> <http://www.uniprot.org/>

<sup>4</sup> <http://www4.wiwiw.fu-berlin.de/bizer/ldif/resources/Wiki.owl>

100 million quads. The R2R mappings, Silk linkage rules as well as the evaluation data sets can be downloaded from the LDIF website.

We ran the performance tests on a machine with an Intel i7 950, 3.07GHz (4 cores) processor and 24GB of memory out of which we assigned 20GB to LDIF.

Table 1 summarizes the LDIF runtimes for the different data set sizes. The overall runtime is split according to the different processing steps of the integration process.

**Table 1.** Runtimes of the integration process for different input data set sizes.

	25M	50M	100M
Load and build entites for R2R	128.1 <i>sec</i>	297.2 <i>sec</i>	1059.7 <i>sec</i>
R2R data translation	169.9 <i>sec</i>	515.0 <i>sec</i>	1109.2 <i>sec</i>
Build entities for Silk	15.3 <i>sec</i>	36.8 <i>sec</i>	107.4 <i>sec</i>
Silk Identity Resolution	103.0 <i>sec</i>	568.5 <i>sec</i>	2954.9 <i>sec</i>
Final URI rewriting	8.1 <i>sec</i>	27.0 <i>sec</i>	65.0 <i>sec</i>
Overall execution time	7.0 <i>min</i>	24.0 <i>min</i>	88.3 <i>min</i>

Table 2 provides statistics about the data integration process. The original number of input quads decreases in the process as LDIF was configured to discards input quads which are irrelevant for the defined mappings, and therefore can not be translated into the target vocabulary. The number decreases again after the actual translation, as the input data uses more verbose vocabularies and as multiple quads from the input data are thus combined into single quads in the target vocabulary.

**Table 2.** Data integration statistics for different input data set sizes.

	25M	50M	100M
Number of input quads	25,000,000	50,000,000	100,000,000
Number of quads after irrelevance filter	13,576,394	25,397,310	44,249,757
Number of quads after mapping	4,419,410	11,398,236	24,972,112
Number of pairs of equivalent entities resolved	24,782	113,245	213,062

## 4 Related Work

We are aware of two other Linked Data integration frameworks that also provide for data translation and identity resolution: The *ALOE - Assisted Linked Data Consumption* framework<sup>5</sup> developed at the Universität Leipzig and the

<sup>5</sup> <http://aksw.org/projects/aloe>

*Information Workbench*<sup>6</sup> [5] developed by fluid Operations. Compared to both frameworks, LDIF provides more expressive languages for data translation and identity resolution as well as an end-to-end concept for provenance tracking. Additionally, both other frameworks have not published performance numbers for use cases involving larger amounts of triples yet.

## 5 Outlook

Over the next months, we will extend LDIF along the following lines:

1. Implementing a Hadoop version of the runtime environment in order to be able to distribute processes and data over a cluster of machines and thereby allowing to scale to really large amounts of input data.
2. Adding Web data access modules (Linked Data crawler, SPARQL endpoint reader) as well as a scheduling component which provides for regularly updating the local input data cache.
3. Adding a data quality evaluation and data fusion module which allows Web data to be filtered according to different data quality assessment policies and provides for fusing Web data according to different conflict resolution methods.

## 6 Acknowledgments

This work was supported in part by Vulcan Inc. as part of its Project Halo ([www.projecthalo.com](http://www.projecthalo.com)) and by the EU FP7 project LOD2 - Creating Knowledge out of Interlinked Data (<http://lod2.eu/>, Ref. No. 257943).

## References

1. Heath, T., Bizer C.: *Linked Data: Evolving the Web into a Global Data Space*. Synthesis Lectures on the Semantic Web: Theory and Technology, Morgan & Claypool Publishers, ISBN 978160845431, 2011.
2. Bizer, C., Jentzsch, A., Cyganiak, R.: *State of the LOD Cloud*. <http://www4.wiwi.fu-berlin.de/lodcloud/state/>, August 2011.
3. Bizer, C., Schultz, A.: *The R2R Framework: Publishing and Discovering Mappings on the Web*. 1st International Workshop on Consuming Linked Data (COLD 2010), Shanghai, November 2010.
4. Isele, R., Jentzsch, A., Bizer, B.: *Silk Server - Adding missing Links while consuming Linked Data*. 1st International Workshop on Consuming Linked Data (COLD 2010), Shanghai, November 2010.
5. Haase, P., Schmidt, M., Schwarte, A.: *The Information Workbench as a Self-Service Platform for Linked Data Applications*. 2nd International Workshop on Consuming Linked Data (COLD 2011), Bonn, Oktober 2011.

---

<sup>6</sup> <http://www.fluidops.com/information-workbench/>