

Meta-Data for a lot of LOD

Editor(s): Aidan Hogan, Universidad de Chile, Chile

Solicited review(s): Jürgen Umbrich, Vienna University of Economics and Business (WU Wien), Austria; Sebastian Hellmann, University of Leipzig, Germany; Christopher Krauss, Fraunhofer FOKUS, Germany; Aidan Hogan, Universidad de Chile, Chile

Laurens Rietveld^a, Wouter Beek^a, Rinke Hoekstra^{a,b} and Stefan Schlobach^a

^a *Department of Computer Science, VU University Amsterdam, The Netherlands*

E-mail: {laurens.rietveld,w.g.j.beek,stefan.schlobach,rinke.hoekstra}@vu.nl

^b *Leibniz Center for Law, Faculty of Law, University of Amsterdam, The Netherlands*

E-mail: hoekstra@uva.nl

Abstract.

This paper introduces the LOD Laundromat meta-dataset, a continuously updated RDF meta-dataset that describes the documents crawled, cleaned and (re)published by the LOD Laundromat. This meta-dataset of over 110 million triples contains structural information for more than 650.000 documents (and growing). Dataset meta-data is often not provided alongside published data, it is incomplete or it is incomparable given the way they were generated. The LOD Laundromat meta-dataset provides a wide range of structural dataset properties, such as the number of triples in LOD Laundromat documents, the average degree in documents, and the distinct number of Blank Nodes, Literals and IRIs. This makes it a particularly useful dataset for data comparison and analytics, as well as for the global study of the Web of Data. This paper presents the dataset, its requirements, and its impact.

Keywords: Dataset Meta-data, Linked Data, Dataset Descriptions

1. Introduction

In this paper we present the LOD Laundromat meta-dataset, a uniform collection of dataset meta-data that describes the structural properties of very many (over 650.000) Linked Data documents that together contain over 38 billion triples. These Linked Data documents range from large compressed data dumps and RDFa embedded in web pages, to dereferenceable URIs and SPARQL CONSTRUCT query references. This meta-dataset is unique in its scale (both in terms of the number of datasets it describes, and the number of meta-data properties included), the consistent way in which meta-data properties are calculated, the explicit description of the computational processes used to calculate these properties, and the use cases it supports. The meta-dataset uniquely facilitates the analysis and comparison of many datasets, and supports research scenarios in which algorithms make innovative use of meta-data values to improve performance [19].

Analyzing, comparing, and using multiple Linked Open Datasets currently is a hassle: finding download locations, hoping the downloaded data dumps are valid, and parsing the data in order to analyze or compare it based on some criteria. It is even more difficult to search for datasets based on characteristics that are relevant for machine-processing, such as syntactic conformance and structural properties such as the average outdegree of nodes. What is needed is a uniform representation of the *dataset* and a uniform representation of *dataset descriptions*.

The LOD Laundromat [5] realizes the first: it (re)publishes the largest (collection of) dataset(s) on the Web of Data (over 38 billion triples and counting). Every dataset is published in the same format that is fully conformant with Linked Open Data (LOD) publication standards for machine-processability. The purpose of the LOD Laundromat is to drastically simplify the task of data preprocessing for the data consumer.

However, the creation of meta-data that uniformly describes the datasets is still left to the original data publisher. We see that many data publishers do not publish a dataset description that can be found by automated means, and that dataset descriptions that *can* be found do not always contain all (formally or de-facto) standardized meta-data. More importantly, the meta-data values are generally not comparable between datasets since different data publishers may interpret and calculate the same meta-data property differently. For instance, it is not generally the case that a dataset with a higher value for the `void:triples` property contains more triples: this value might be outdated with respect to the original dataset, or it might have been incorrectly calculated. Because of such incompatibilities between existing dataset descriptions, it is difficult to reliably analyze and compare datasets on a large scale.

Therefore, in addition to the uniform *dataset* representations published by the LOD Laundromat, we need the same uniform representation for publishing *dataset meta-data*. Secondly, even straightforward meta-data should come with provenance annotations that describe how meta-data was generated. The LOD Laundromat meta-dataset presented here, brings exactly this: a collection of dataset descriptions, linked to the same canonical dataset representation, all modeled, created, and published in the same manner, and with provenance annotations that reflect how the meta-data was generated.

Section 2 gives an overview of comparable datasets. In section 3 we identify shortcomings of existing meta-data standards and collections, and formulate a set of requirements for a dataset that allows large collections of datasets to be analyzed, compared, and used. Section 4 presents the meta-data we publish, the model that underlies it, dependencies on external (standard) vocabularies, a discussion in the context of the five stars of Linked Data Vocabulary use [12], and clarification of how the LOD Laundromat meta-dataset is generated and maintained. Section 5 shows the applications and use cases that the LOD Laundromat meta-dataset supports. We conclude with section 6.

2. Similar Datasets

SPARQL Endpoint Status¹ [6] presents an overview of dataset descriptions that can be found by automated

means. It shows that even the uptake of the core meta-data properties, such as the ones from the VoID [1] specification, is still quite low: only 12,9% of analyzed SPARQL endpoints is described using VoID. Because of this apparent lack of LOD meta-data, several initiatives tried to fill this gap by creating uniform meta-data descriptions for multiple datasets. We discuss two of these: LODStats and Sindice.

LODStats [2] provides statistical information for all Linked Open Datasets that are published in the CKAN-powered² Datahub³ catalog. It offers a wide range of statistics, e.g., including the number of blank nodes in a dataset and the average outdegree of subject terms. Unfortunately, it only publishes a small subset of these statistics as Linked Data. Secondly, Sindice [17] provides statistical information that is similar to LODStats, but mostly analyzes smaller datasets that are crawled from Web pages. The meta-data provided by Sindice are similar to those in the VoID specification but they are not published in a machine-readable format such as RDF.

Although Sindice and LODStats provide a step in the right direction by uniformly creating meta-data descriptions for many Linked Datasets, they 1) only support a subset of existing meta-data properties, they 2) do not publish exhaustive meta-data descriptions as Linked Data, and 3) they do not publish structural information on the meta-data generation procedure. Finally, they are limited to Linked Datasets that are published in only certain locations, web pages and Datahub, respectively.

3. Meta-Data Requirements

In this section we present a requirements analysis for a dataset that satisfies our goal of supporting the meaningful analysis, comparison, and use, of very many datasets.

We explain problems with respect to meta-data specifications (section 3.1), dataset descriptions (section 3.2) and collections of dataset descriptions (section 3.3). Based on these considerations, the requirements are presented in section 3.4.

¹<http://sparql.es.ai.wu.ac.at/>

²<http://ckan.org/>

³<http://datahub.io/>

3.1. Meta-data specifications

Existing dataset vocabularies include VoID [1], VoID-ext [14], DCAT⁴, and Bio2RDF [7]. VoID is a vocabulary for expressing meta-data about Linked Datasets. It supports generic meta-data (e.g., the homepage of a dataset), access meta-data (e.g., which protocols are available), links to other datasets, exemplary resources, as well as dataset statistics (e.g., the number of triples). Only some of the VoID meta-data properties can be automatically generated. Others can only be given by human authors, –such as exemplary resources– since they depend on interpretation. Bio2RDF presents a collection of dataset meta-data properties that extends the set of VoID properties and provides more detail. For example, Bio2RDF includes properties that describe how often particular types are used in the subject position and in the object position for a given property; e.g. property `ex:livesIn` links 10 subjects of type `ex:Person` to 6 objects of type `ex:City`. The use of such descriptive properties can increase the size of a meta-dataset significantly when the described dataset has a large number of classes and properties.

VoID-ext extends the set of meta-data properties that are found in VoID as well. It includes the in- and out-degree of entities, the number of blank nodes, the average string length of literals, and a partitioning of the literals and URIs based on string length. The Data Catalog Vocabulary (DCAT) is a vocabulary for describing datasets on a higher level; i.e., it includes properties such as the dataset title, description and publishing/modification date. Such information is difficult to reliably extract from the dataset in an automated fashion.

We observe the following problems with these existing meta-data specifications:

Firstly, some existing meta-data properties are subjective. For example, `void:entities` is intended to denote a subset of the IRIs of a dataset based on “arbitrary additional requirements” imposed by the authors of the dataset description. Since different authors may impose different requirements, the number of entities of a dataset may vary between zero and the number of resources.

Secondly, some existing meta-data properties are defined in terms of undefined concepts. For example, LODStats specifies the set of vocabularies that are

reused by a given dataset. The notion of a ‘reused vocabulary’ is itself not formally defined but depends on heuristics about whether or not an IRI belongs to another dataset. LODStats calculates this set by using relatively simple string operations according to which IRIs of the form

`http://<authority>/<string>/<value>` are assumed to belong to the vocabulary denoted by

`http://<authority>/<string>`. Although this is a fair attempt at identifying reused vocabularies, there is not always a bijective map between datasets and URI substrings that occur in datasets. The number of links to other datasets suffers from the same lack of a formal definition.

3.2. Dataset descriptions

We observe the following problems with existing dataset descriptions: Firstly, uptake of dataset descriptions that can be found by automated means is still quite low (section 2). Secondly, for reasons discussed above, the values of meta-data properties that do not have a well-founded definition cannot be meaningfully compared across datasets. For instance, if two dataset descriptions contain different values for the `void:entities` property it is not clear whether this denotes an interesting difference between the two datasets or whether this is due to the authors having different criteria for identifying the set of entities. Thirdly, even the values of well-defined meta-data may have been calculated in different ways by different computational procedures. We observe that there are significant discrepancies between meta-data which occurs *in* the original dataset description and those from the LOD Laundromat. For example, a dataset about a Greek fire brigade contains 3.302.302 triples according to its original VoID description⁵, but 4.134.725 triples according to the LOD Laundromat meta-dataset⁶.

Similar discrepancies exist between meta-data values that occur in different dataset description *collections*, e.g. between LODStats and the LOD Laundromat meta-dataset.⁷

Since it is difficult to assess whether a computational procedure that generates meta-data is correct, we

⁵See <http://greek-lod.auth.gr/Fire/void.ttl>

⁶See <http://lodlaundromat.org/resource/0ca7054f382b29319c82796a7f9c3899>

⁷E.g., according to LODStats the dataset located at <http://www.open-biomed.org.uk/open-biomed-data/bdgp-images-all-20110211.tar.gz> contains 1.080.060 triples while the LOD Laundromat meta-dataset states 1.070.072.

⁴See <http://www.w3.org/TR/vocab-dcat/>

believe it is necessary that all generated meta-data is annotated with provenance information that describes the used computational procedure. Although relatively verbose, this approach circumvents the arduous discussion of which version of what tool is correct/incorrect for calculating a given meta-data value. We assume that there will always be multiple values for the same meta-data property. The fact that there are different values, and that these have been derived by different means, is something that has to be made transparent to the consumer of this meta-data. The onus is on the data consumer to trust one computational procedure for calculating a specific meta-data value more than another. This requires provenance that details the mechanism behind the calculated meta-data.

3.3. Dataset description collections

We observe two problems with existing collections of dataset descriptions: Firstly, even though the meta-data may be calculated consistently within a collection, the computational procedure that is used is not described in a machine-processable format (if at all). This means that values can only be compared within the collection, but not with dataset descriptions external to the collection (e.g. occurring in other collections). Secondly, meta-data that is calculated within existing collections is not always published in a machine-interpretable format (e.g. LODStats).

3.4. Requirements

Based on the above considerations, we formulate the following requirements which allow multiple datasets to be meaningfully compared based on their meta-data:

1. The LOD Laundromat meta-dataset must cover very many datasets in order to improve data comparability.
2. The meta-dataset should reuse official and de-facto meta-data standards as much as possible, in order to be compatible with other dataset descriptions and to promote reuse.
3. The meta-dataset must be generated algorithmically in order to assure that values are calculated in the same way for every described dataset.
4. Only those meta-data properties must be used that can be calculated efficiently, because datasets can have peculiar properties that may not have been anticipated when the meta-data properties were first defined.
5. The LOD Laundromat meta-dataset must contain provenance annotations that explain how and when the meta-data was calculated.
6. The meta-data must be disseminated as LOD and must be accessible via a SPARQL endpoint.
7. The LOD Laundromat meta-dataset must be able to support a wide range of real-world use cases that involve analyzing and/or comparing datasets such as Big Data algorithms that process LOD.

4. The LOD Laundromat meta-dataset

In this section we present the meta-data we publish, the model we use, and how we generate this dataset.

4.1. Published Meta-Data

The LOD Laundromat meta-dataset is generated in adherence to the requirements formulated in section 3. Since there are multiple ways in which these requirements can be prioritized and made concrete, we will now discuss the considerations that have guided the generation of the meta-data.

Firstly, there is a trade-off between requirements 2 and 3: since the meta-dataset has to be constructed algorithmically, only well-defined meta-data properties can be included.

Secondly, there is a conflict between requirements 1 and 4 on the one hand, and requirement 2 on the other: since the LOD Laundromat meta-dataset must describe many datasets, some of which are relatively large, and we want calculations to be efficient, we chose to narrow down the set of meta-data properties to those that can be calculated by *streaming* the described datasets. This excludes properties that require loading (large parts of) a dataset into memory, e.g. in order to perform joins on triples.

Thirdly, because of the scale at which the LOD Laundromat meta-dataset describes datasets, it is inevitable that some datasets will have atypical properties. This includes datasets with extremely long literals, datasets where the number of unique predicate terms is close to the total number of predicate terms, or datasets where the number of unique literal datatype equals the total number of literals. It is only when meta-data is systematically generated on a large scale, that one finds such corner cases. These corner cases can make dataset descriptions impractically large. This is especially true for meta-data properties that consist of enumerations. For instance, for

some datasets the partition of all properties – as defined by VoID-ext and Bio2RDF – is only (roughly) a factor 3 smaller than the described dataset itself (and this is only one meta-data property). Or, take as example the `void-ext:subjectPartition`, that refers to a partition that contains triples for a certain subject. Using such partitions for all the subjects in a dataset would generate a meta-dataset that equals the size of the original dataset. Therefore, in order to keep data descriptions relatively small w.r.t. the dataset described, the meta-dataset does not include properties whose values are dataset partitions.

Under these restrictions, the meta-dataset is able to include a large number of datasets while still being relatively efficient to construct. Implementation-wise, the generation of the meta-dataset takes into account the many advantages that come from the way in which LOD Laundromat (re)publishes datasets. LOD Laundromat allows datasets to be opened as gzip-compressed streams of lexicographically sorted N-Triples and N-Quads. Since these streams are guaranteed to contain no syntax error nor any duplicate occurrences of triples, they can be processed on a line-by-line / triple-by-triple basis, making it convenient to generate meta-data for inclusion in the LOD Laundromat meta-dataset.

Table 1 gives an overview of the meta-data properties included in the LOD Laundromat meta-dataset, together with those that are included in existing dataset description standards. As can be seen from the table, the only meta-data properties that are excluded from our dataset (because of computational issues) are the distinct number of classes that occur in either the subject, predicate, or object position, as specified in VoID-ext. These three meta-data properties cannot be calculated by streaming the data a single time. In addition, all meta-data properties whose values must be represented as partitions are excluded in order to preserve brevity for all dataset descriptions, and to maintain scalability. Considering these limitations, the meta-data properties presented in Bio2RDF are similar to those in VoID and VoID-ext. Therefore, Bio2RDF is not referenced in our vocabulary. The generation of several statistics (e.g. the distinct number of URIs) requires in-memory lists. To reduce this memory consumption, we use an efficient in-memory dictionary (RDF Vault [3]).

Since we want the LOD Laundromat meta-dataset to be maximally useful for a wide range of use cases (requirement 7), we have added several meta-data properties that do not occur in existing specifications:

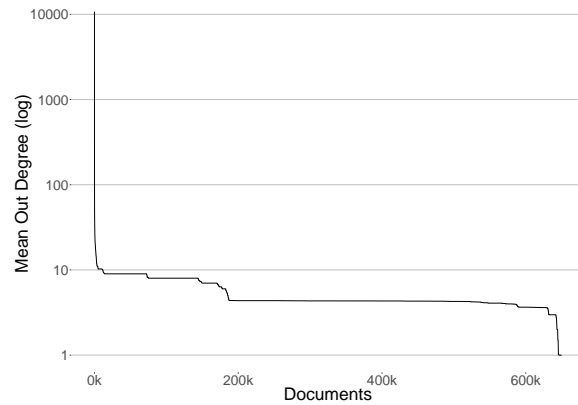


Fig. 1. Average out degree distribution of LOD Laundromat documents

1. Next to the number of distinct IRIs, blank nodes and literals (i.e., *types*), we also include the number of (possibly non-distinct) occurrences (i.e., *tokens*).
2. Existing vocabularies specify the number of properties and classes (although they do so incorrectly, see section 3). The meta-dataset also includes the number of classes and properties that are *defined* in a dataset, such as `<prop> rdf:type rdf:Property`
3. Existing dataset description vocabularies such as VoID-ext use arithmetic means to describe number series such as the literal lengths in given document. The LOD Laundromat meta-dataset uses more detailed descriptive statistics, that include the median, minimum, maximum and standard deviation values as well.
4. Similar statistics are provided for network characteristics such as Degree, In Degree and Out Degree.

Considering that only 0,5% of the datasets publish a corresponding dataset license via RDF, we exclude this information for now. We expect these dataset licenses to increase in use and popularity though, and will include this meta-data in a future crawl.

Figure 1 illustrates one of the published meta-data properties: the average out degree of datasets. The figure illustrates our previous remark that analyzing many datasets will inevitably include datasets with a-typical properties or ‘corner cases’. For instance, the dataset with the highest average out degree, contains 10.004 triples, and only one subject, thereby strongly skewing the dataset distribution. Such a-typical properties of datasets are potentially important as e.g. a means of explaining deviating evaluation results between datasets.

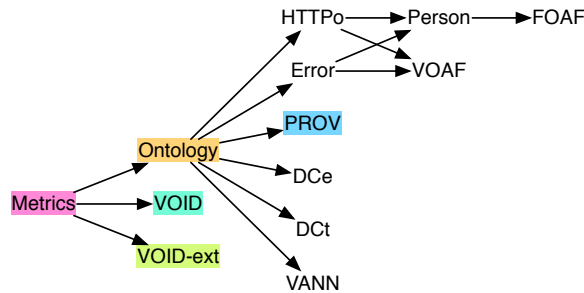


Fig. 2. Dependencies of LOD Laundromat meta-dataset vocabulary

Note that generating the data behind this figure requires the following SPARQL query, illustrating the ease of use:

```
SELECT * {[] l1m:outDegree/l1m:mean ?mean}
```

Besides publishing the meta-data, and in line with requirement 5, the meta-dataset contains a provenance trail of how the meta-data was generated. The provenance trail includes a reference to the code that was used to generate the meta-data. For this we use a Git commit identifier in order to uniquely identify the exact version that was used. The provenance trail also includes all the steps that preceded the calculation of the meta-data:

1. Where the file was downloaded (either the original URL or the archive that contained the file).
2. When the file was downloaded (date and time).
3. Meta-data on the download process, such as the status code and headers from the original HTTP reply. For archived data the applied compression techniques (possibly multiple ones) are enumerated as well.
4. Detailed meta-data on the data preparation tasks performed by the LOD Laundromat in order to clean the data. This includes the number of bytes that were read (not necessarily the same as the value for `Content-Length` HTTP header) and syntax errors that were encountered (e.g., malformed syntax, unrecognized encoding, undefined prefixes).
5. The number of duplicate triples in the original dataset.
6. A reference to the online location where the cleaned file is stored, and from which the meta-data is derived.

4.2. Model

The meta-data is specified in the LOD Laundromat meta-dataset (See Table 2). Of the 26 meta-data prop-

erties that are included, 22 are linked to one or more other dataset description vocabularies. Figure 2 shows the dependencies between our meta-dataset vocabulary and other vocabularies. The referenced dataset description vocabularies are VoID and VoID-ext. Figure 3⁸ shows an example dataset description that illustrates the structure of this meta-dataset⁹. The meta-dataset also includes information about the vocabulary *itself*, such as its license (Creative Commons¹⁰), last modification date, creators, and homepage. As such, it implements the first 4 of the 5 stars for vocabulary re-use [12]. The fifth star (re-use *by* other vocabularies) is not reached yet because the vocabulary is quite recent. However, the LOD Laundromat meta-dataset has been submitted to the Linked Open Vocabulary catalog¹¹, thereby hopefully supporting its re-use and findability.

The provenance information of datasets is described using the PROV-O vocabulary [15], a W3C recommendation. Figure 4 presents an overview on how PROV-O is used by the LOD Laundromat meta-dataset. Similar vocabularies exist, such as the VoiDp [16] vocabulary which matches the provenance of Linked Datasets with the VoID vocabulary. However, because VoiDp uses a predecessor of the PROV-O standard, we model our provenance in PROV-O directly. The Provenance Vocabulary [9] aims to describe the provenance of Linked Datasets as well, but is too specific for our use considering the wide range of provenance (see below) we describe.

As the LOD Laundromat cleaning process is part of the provenance trail, we model this part of the dataset using separate vocabularies: Firstly, the LOD Laundromat vocabulary (See Table 2) describes the crawling and cleaning process of LOD Laundromat. This description includes the download time and date of the original document, and therefore specifies which version of the original document is described by the meta-dataset. Secondly, the HTTP vocabulary (See Table 2) describes HTTP status codes. Thirdly, the error ontology (See Table 2) models all exceptions and warnings, and is used by the LOD Laundromat vocabulary to rep-

⁸For brevity, figures in tables in this paper do not contain the following LOD Laundromat namespaces:

```
PREFIX l1e <http://lodlaundromat.org/errors/ontology/>
PREFIX l1m <http://lodlaundromat.org/metrics/ontology/>
PREFIX l1o <http://lodlaundromat.org/ontology/>
```

⁹For brevity, only a subset of the available meta-data properties are included in this figure

¹⁰See <http://creativecommons.org/licenses/by/3.0/>

¹¹<http://lov.okfn.org/>

Table 1

An overview of dataset meta-data properties, grouped by the vocabularies that define them and dataset description collections that include them. For brevity's sake, properties whose values are dataset partitions and properties that require manual intervention are excluded

Meta-data Property	VoID	Bio2RDF	VoID-ext	LOD Laundromat	DataType / Range
Triples	v	v	v	v	xsd:integer
Entities	v	v	v	v	xsd:integer
Distinct Classes	v	v	v	v	xsd:integer
Distinct Properties	v	v	v	v	xsd:integer
Distinct Subject	v	v	v	v	xsd:integer
Distinct Objects	v	v	v	v	xsd:integer
Distinct RDF Nodes			v	v	xsd:integer
Distinct IRIs			v	v	xsd:integer
IRIs				v	xsd:integer
Distinct Blank Nodes			v	v	xsd:integer
Blank Nodes				v	xsd:integer
Distinct Literals	v		v	v	xsd:integer
Literals				v	xsd:integer
Distinct URIs in subject position			v	v	xsd:integer
Distinct Blank Nodes in subject position			v	v	xsd:integer
Distinct URIs in object position			v	v	xsd:integer
Distinct Blank Nodes in object position			v	v	xsd:integer
Distinct Literal Data-Types			v	v	xsd:integer
Distinct Literal Languages			v	v	xsd:integer
Length statistics of IRIs			v	v	xsd:integer
Length statistics of IRIs in subject position			v	v	llm:DescriptiveStatistics
Length statistics of IRIs in predicate position			v	v	llm:DescriptiveStatistics
Length statistics of IRIs in object position			v	v	llm:DescriptiveStatistics
Length statistics of Literals			v	v	llm:DescriptiveStatistics
Degree Statistics				v	llm:DescriptiveStatistics
Indegree Statistics				v	llm:DescriptiveStatistics
Outdegree Statistics				v	llm:DescriptiveStatistics
Defined Classes				v	xsd:integer
Defined Properties				v	xsd:integer
Distinct Classes occurring in the subject position			v		
Distinct Classes occurring in the predicate position			v		
Distinct Classes occurring in the object position			v		

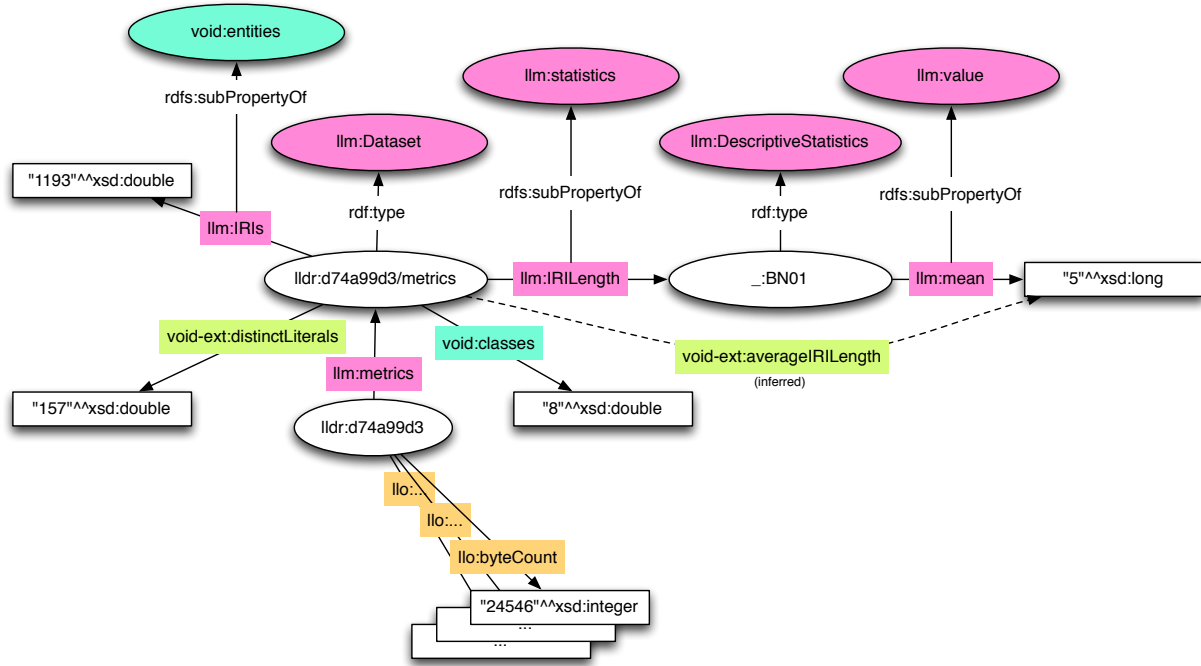


Fig. 3. Example (partial) dataset meta-data description, color-coded using vocabularies from Figure 2

resent errors that occur during the crawling and cleaning process. Each of these vocabularies are linked to other vocabularies. For instance, the HTTP vocabulary is an extension of the W3C HTTP in RDF vocabulary¹².

4.3. Naming Scheme

The LOD Laundromat meta-dataset uses the following naming scheme. As a running example, we take a Semantic Web Dog Food file that is crawled by LOD Laundromat¹³.

- The LOD Laundromat document identifier for this dataset is generated by appending an MD5 hash of the data source IRI to `http://lodlaundromat.org/resource/`¹⁴.
- The calculated structural properties of this dataset are accessible by appending `/metrics` to the LOD Laundromat document identifier¹⁵.

- Provenance that describes the procedure behind the metrics calculation is accessible by appending `metricCalculation` to the LOD Laundromat identifier¹⁶.

The prefixes we use (further in this paper as well), are:

```
ll <http://lodlaundromat.org/ontology/>
llm <http://lodlaundromat.org/metrics/ontology/>
lle <http://lodlaundromat.org/errors/ontology/>
```

4.4. Dissemination

All the online resources related to the LOD Laundromat are shown in Table 2. The LOD Laundromat [5] continuously crawls and analyses Linked Data dumps. In order to get a maximum coverage of the LOD Cloud, it searches both linked data catalogs and the LOD Laundromat datasets themselves for references to datadumps. Because it does not claim to have a complete seed list that links to all LOD in the world, users have the option to manually or algorithmically add seed-points to the LOD Laundry Basket (See Table 2).

¹²See <http://www.w3.org/2011/http>
¹³See <http://data.semanticweb.org/dumps/conferences/iswc-2013-complete.rdf>
¹⁴See <http://lodlaundromat.org/resource/05c4972cf9b5ccc346017126641c2913>
¹⁵See <http://lodlaundromat.org/resource/05c4972cf9b5ccc346017126641c2913/metrics>

¹⁶See <http://lodlaundromat.org/resource/05c4972cf9b5ccc346017126641c2913/metricCalculation>

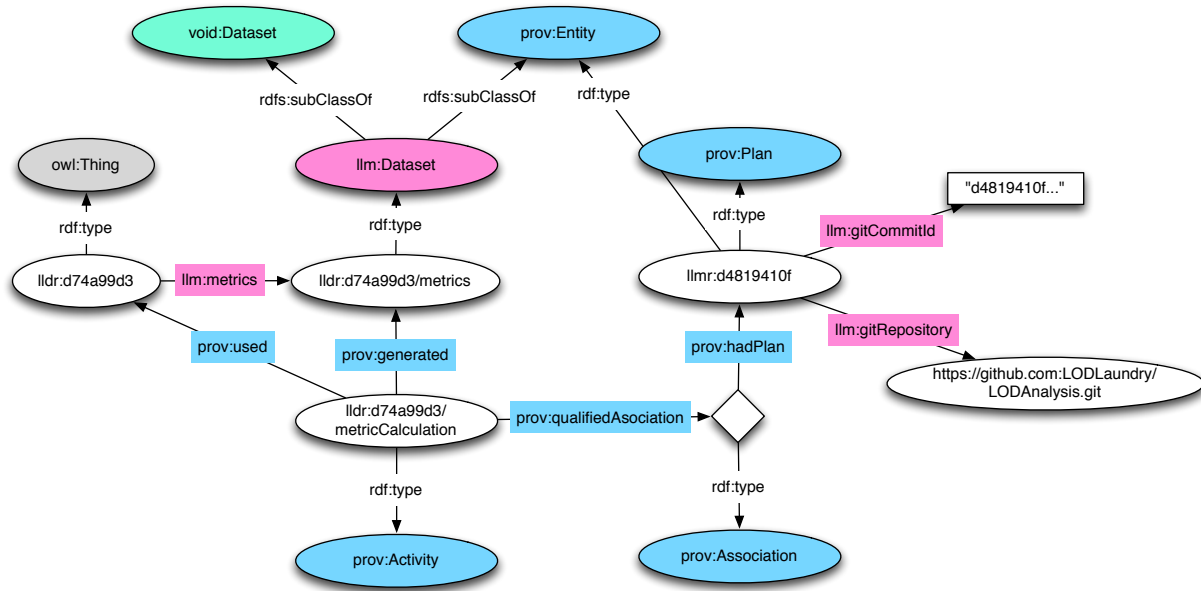


Fig. 4. Provenance model illustration

Table 2

Online LOD Laundromat Resources

Website

Main Website	http://lodlaundromat.org
Documents Overview	http://lodlaundromat.org/wardrobe
Seed-List	http://lodlaundromat.org/basket
Visualizations	http://lodlaundromat.org/visualizations

APIs

SPARQL Endpoint	http://lodlaundromat.org/sparql
Linked Data Fragments	http://ldf.lodlaundromat.org
Resource Index API	http://index.lodlaundromat.org
Text Index (LOTUS)	http://lotus.lodlaundromat.org

Vocabularies

LOD Laundromat Vocabulary	http://lodlaundromat.org/ontology/
meta-data Vocabulary	http://lodlaundromat.org/metrics/ontology/
HTTP Vocabulary	http://lodlaundromat.org/http/ontology/
Error Vocabulary	http://lodlaundromat.org/error/ontology/

Open Source Code

meta-data Generation	https://github.com/LOD-Laundromat/LODAnalysis
Linked Data Crawler & Cleaner	https://github.com/LOD-Laundromat/LOD-Laundromat

The code used to generate the LOD Laundromat meta-dataset runs immediately after a document is crawled and cleaned by the LOD Laundromat, and is published via a public SPARQL endpoint (See Table 2). The time it takes between a LOD Laundromat crawl and a published document with corresponding meta-data depends on the original serialization format, the size of the dataset, the number of errors encoun-

tered during the cleaning phase, and other idiosyncrasies a dataset might have such as only one subject for thousands of triples. The crawling and cleaning phase processes on average 460.000 triples per second, and the meta-data application is able to stream and analyze 400.000 triples per second (both on a server with 5TB SSD disk space, 32-core CPU, and 256GB of memory). Considering the costs for downloading the

original document files, and a latency (with a maximum of 10 minutes) between the cleaning and meta-data generation, most datasets and the corresponding meta-data are public 15 minutes after the crawl began.

SPARQL is preferred over HDT as publishing method for the meta-data, because HDT files are static and do not support updates. In line with requirement 6, a nightly version of the meta-dataset is copied from the meta-data SPARQL endpoint and published as data dump in the same standardized N-Quad serialization format of the LOD Laundromat.

Considering some meta-data is too verbose and expensive to host as RDF, we publish non-RDF data as well. Specifically, we publish the mapping between all the IRIs and namespaces to the corresponding LOD Laundromat documents these occur in. We provide access to this meta-data via a custom RESTful API (See Table 2).

4.5. Statistics and Use

This section briefly describes some of the characteristics of the dataset, and gives an indication as to how it is being used.

4.5.1. Dataset Characteristics

As mentioned before, the LOD Laundromat crawled and re-publishes over 650.000 documents containing over 38.000.000 triples. The meta-data of these crawled documents are published in the LOD Laundromat meta-dataset, that now contains over 110 million triples, accessible via a data dump and SPARQL endpoint.

Figure 5 shows the top 10 most frequently instantiated classes. It shows that the descriptive statistics class is used frequently, which is not surprising considering that each dataset has several descriptive statistics blank-nodes to describe e.g. in degree, out degree, and IRI/literal lengths. Other frequent classes are related to provenance. This is not a surprise either, considering the extensive provenance model that the LOD Laundromat uses (See figure 4).

Figure 6 shows the top 10 most frequently occurring predicates in the dataset. These statistics show a similar pattern as the class instance frequency: properties related to descriptive statistics and provenance appear frequently.

4.5.2. Usage of the Dataset

Since the release of LOD Laundromat in September 2014 and the release of the meta-dataset in January 2015, we registered a total of 8.158 unique visitors to

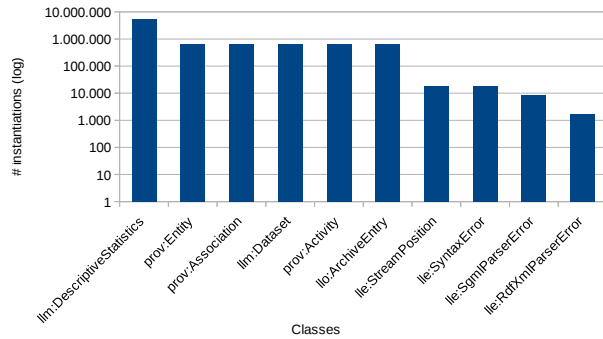


Fig. 5. Top 10 class instances

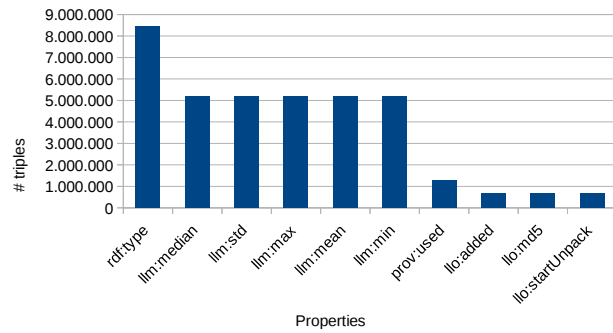


Fig. 6. Top 10 most frequently occurring properties

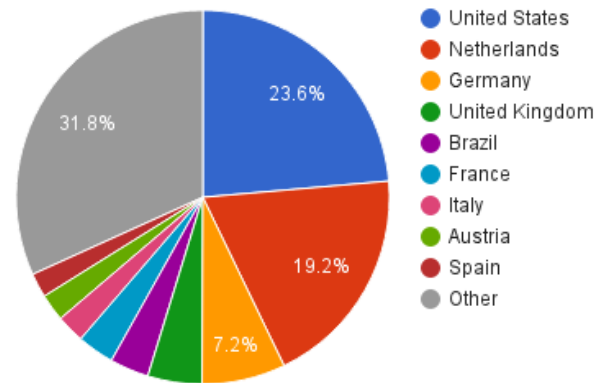


Fig. 7. Geo-Location visitors <http://lodlaundromat.org>

the LOD Laundromat website, of which 23,6% of the requests originated in the United States, and 19,2% of the requests originated from The Netherlands (see Figure 7).

The SPARQL endpoint receives the following four types of SPARQL requests.

1. Queries executed while browsing the LOD Laundromat website. These queries can vary in length, depending on the user interaction with the website.
2. Queries executed by the `Frank` tool, described in the next section. These queries can vary in length as well, and may contain user-provided custom triple-patterns.
3. Queries executed by others via custom scripts (e.g. python, bash, java)
4. Queries that are manually written and executed via the public LOD Laundromat SPARQL graphical user interface, or via interfaces such as <http://yasgui.org>.

Distinguishing these types of interactions with the SPARQL endpoint is not trivial: request headers are insufficient to separate these categories, and query features (e.g. number of triple patterns or projection variables) are difficult to translate to different types of SPARQL requests [18,20,21]. We do expect a large part of the logged SPARQL queries to come from the `Frank` tool. This server log goes back to June 2015, and shows a total of 12.376.978 SPARQL queries, coming from 2.887 unique IPs¹⁷. Figure 8 shows the number of queries where a particular IRI occurs in the predicate position. This figure shows that the provenance information such as the original download location (`llo:url`, 6.660.932 queries) and the time a dataset was added to the laundromat (`llo:added`, 2.826.283 queries) are most often used in queries. A typical dataset metric such as the number of triples (`llo:triples`, 993.550 queries) is popular as well. These statistics only paint part of the picture (queries with a bound predicate), but they illustrate the variety of accessed information.

The index API that exposes the mappings between IRIs/Namespaces and documents, received 56.727 requests, coming from 124 unique IPs. Of these requests, 445 used the namespace-to-document index, where the remaining 56.282 used the resource-to-document index.

The server registered 5.640.331 downloads, coming from 881 IPs. The unique number of downloads is 649.908, which equals the number of available LOD Laundromat documents.

These server logs show that a relatively small number of IPs (less than 3.000) are responsible for millions of requests. In the next section we present three use cases and applications that can (at least partly) account for this large amount of use.

5. Use Cases

The LOD Laundromat meta-dataset is intended to support a wide array of non-trivial use cases. The first use case we present is the evaluation of Semantic Web (SW) algorithms. In contemporary SW research novel algorithms are usually evaluated against only a handful of – often the same – datasets (i.e., mainly DBpedia, Freebase, and Billion Triple Challenge). The risk of this practice is that – over time – SW algorithms will be optimized for datasets with specific distributions, but not for others. In [19], we re-evaluate parts of three SW research papers using `Frank` [4], a bash interface interfacing with the LOD Laundromat. We showed how the LOD Laundromat meta-dataset can be used to relate datasets to their overall structural properties, and how SW evaluations can be performed on a much wider scale, leading to results that are more indicative of the *entire* LOD Cloud. For example, the re-evaluation of RDF HDT [8] (a binary compressed representation for RDF) showed a –previously unknown– relation between the degree of datasets and the RDF HDT compression ratio. This use case combines the strength of both the LOD Laundromat collection of documents and the LOD Laundromat meta-dataset. The following SPARQL query was used in the RDF HDT re-evaluation to find documents with a low average out degree:

```
SELECT * WHERE {
  ?datadoc llm:metrics
           /llm:outDegree
           /llm:mean ?outDegree .
  FILTER(?outDegree < 5)
}
```

Similarly to *evaluating* SW algorithms, the LOD Laundromat meta-dataset can also be used to *tune* Linked Data applications or prune datasets with the desired property at an early stage, i.e., without having to load and interpret them. An example of this is `PrefLabel`¹⁸, an online service that returns a human-readable label for a given resource-denoting IRI. The

¹⁷The number of unique IPs may not correspond to the actual number of unique users, as a single IP can be shared by employees of the same institute or company, or the same user may access the endpoint from different machines.

¹⁸<http://preflabel.org>

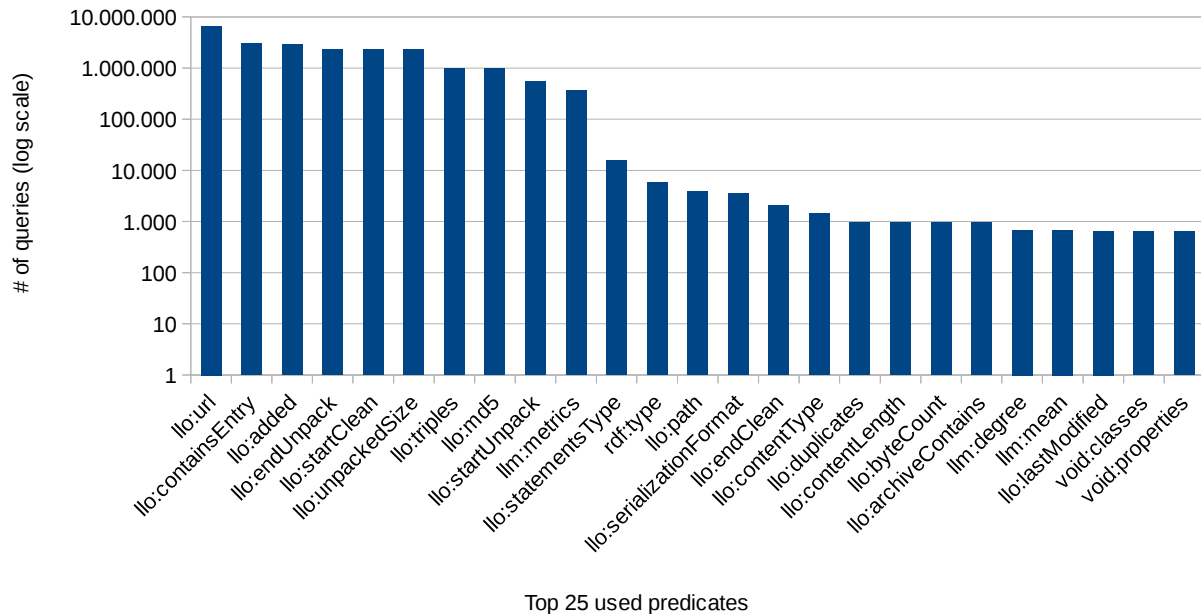


Fig. 8. Frequency of predicates in SPARQL query logs

index behind the PrefLabel Web service is populated by streaming and analyzing LOD Laundromat datasets for RDFS label statements in datasets. PrefLabel uses the LOD Laundromat meta-dataset by pruning for datasets that do not contain RDF literals at all. This crude way of using the meta-dataset already excludes 20% of all the triples that are in the LOD Laundromat today, thereby significantly optimizing the algorithm. The following SPARQL query is used by PrefLabel to prune the list of documents:

```
SELECT ?doc WHERE {
  ?doc llo:metrics
    /llo:literals ?lit;
  FILTER(?lit = 0)
}
```

Another use case involves using the LOD Laundromat meta-dataset to analyze and compare datasets, e.g., in order to create an overview of the state of the LOD Cloud at a given moment in time. A common approach (see e.g. [10,13,22]) is to crawl Linked Data via dereferenceable URIs using tools such as LD-spider [11], and/or to use catalogs such as datahub to discover the Linked Datasets. Both dereferenceable URIs and dataset catalogs come with limitations: most Linked Data URIs are not dereferenceable, and the dataset catalogs only cover a subset of the LOD Cloud. The LOD Laundromat on the other hand pro-

vides access to more than dereferenceable URIs only, and aims to provide a complete as possible dataset collection. The corresponding meta-dataset provides a starting point for e.g. finding datasets by Top Level Domain, serialization format, or structural properties such as number of triples. In [19] we re-evaluate (part of) exactly such a Linked Data Observatory paper [22], where we use the meta-dataset and LOD Laundromat to find the documents and extract namespace statistics.

Next to the *structural* meta-data properties, the *provenance* meta-data provides an interesting data source as well. Such provenance enables e.g. an analysis of common RDF serialization formats, as shown in figure 9. The following SPARQL query fetches the serialization information used by this figure:

```
SELECT ?format (SUM(?t) AS ?count)
WHERE {
  [] llo:serializationFormat ?format;
    llo:triples ?t .
}
GROUP BY ?format
```

The provenance information can be used to measure publishing best practices as well, such as whether the content length specified by the HTTP response matches the actual content length of the file. This is visualized in figure 9, which uses the following SPARQL query to fetch the data:

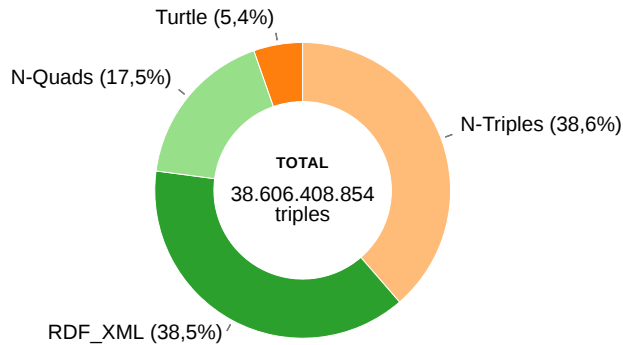


Fig. 9. Serialization formats

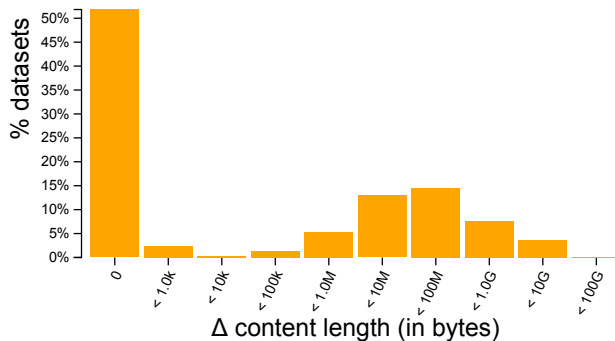


Fig. 10. Invalid HTTP Content Lengths

```

SELECT ?clength ?bcount ?t WHERE {
  [] llo:contentLength ?clength ;
     llo:byteCount ?bcount ;
     llo:triples ?t .
}

```

6. Conclusion

The dataset presented in this paper offers access to a large set of uniformly represented datasets descriptions, acting as an enabler for large scale Linked Data research: finding or comparing linked datasets with certain structural properties is now as easy as executing a SPARQL query. And even better: because the dataset descriptions are linked to their uniform dataset representations, the access to the underlying data is extremely easy as well.

We are exploring the possibilities of storing snapshots of both the meta-dataset and the corresponding cleaned datasets, effectively creating snapshots of the state of the LOD Cloud. At this point, we consider this future work though.

Another future improvement we consider is to publish partitions of the datasets via more scalable and efficient ways than SPARQL. As explained in section 4, corner-cases in the LOD cloud can drastically increase the corresponding meta-data. Therefore, an efficient and scalable method is required for hosting such partitions. We consider publishing a selection of such partitions using non-SPARQL APIs with a stronger focus on scalability and efficiency.

Acknowledgements

This work was supported by the Dutch national program COMMIT.

References

- [1] Keith Alexander and Michael Hausenblas. Describing Linked Datasets - On the Design and Usage of VoiD, the Vocabulary of Interlinked Datasets. In *In Linked Data on the Web Workshop (LDOW 09), in conjunction with 18th International World Wide Web Conference (WWW), Madrid, Spain, April 2009*, 2009.
- [2] Sören Auer, Jan Demter, Michael Martin, and Jens Lehmann. LODStats – an extensible framework for high-performance dataset analytics. In Annette ten Teije, Johanna Völker, Siegfried Handschuh, Heiner Stuckenschmidt, Mathieu d’Aquin, Andriy Nikolov, Nathalie Aussenac-Gilles, and Nathalie Hernandez, editors, *Knowledge Engineering and Knowledge Management - 18th International Conference, EKAW 2012, Galway City, Ireland, October 8-12, 2012. Proceedings*, volume 7603 of *Lecture Notes in Computer Science*, pages 353–362. Springer, 2012.
- [3] Hamid R. Bazoobandi, Steven de Rooij, Jacopo Urbani, Annette ten Teije, Frank van Harmelen, and Henri E. Bal. A Compact In-Memory Dictionary for RDF Data. In Fabien Gandon, Marta Sabou, Harald Sack, Claudia d’Amato, Philippe Cudré-Mauroux, and Antoine Zimmermann, editors, *ESWC, volume 9088 of Lecture Notes in Computer Science*, pages 205–220. Springer, 2015.
- [4] Wouter Beek and Laurens Rietveld. Frank: The LOD Cloud at your Fingertips. In Ruben Verborgh and Miel Vander Sande, editors, *Proceedings of the ESWC Developers Workshop 2015 co-located with the 12th Extended Semantic Web Conference (ESWC 2015), Portoroz, Slovenia, May 31, 2015*, volume 1361 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2015.
- [5] Wouter Beek, Laurens Rietveld, Hamid R. Bazoobandi, Jan Wielemaker, and Stefan Schlobach. LOD Laundromat: A Uniform Way of Publishing Other People’s Dirty Data. In Peter Mika, Tania Tudorache, Abraham Bernstein, Chris Welty, Craig A. Knoblock, Denny Vrandečić, Paul T. Groth, Natasha F. Noy, Krzysztof Janowicz, and Carole A. Goble, editors, *Semantic Web Conference (1)*, volume 8796 of *Lecture Notes in Computer Science*, pages 213–228. Springer, 2014.
- [6] Carlos Buil-Aranda, Aidan Hogan, Jürgen Umbrich, and Pierre-Yves Vandenbusche. SPARQL Web-Querying Infras-

- structure: Ready for Action? In *Proceedings of the 12th International Semantic Web Conference, Sydney, Australia*, pages 277–293. Springer, 2013.
- [7] Alison Callahan, José Cruz-Toledo, Peter Ansell, and Michel Dumontier. Bio2RDF Release 2: Improved Coverage, Interoperability and Provenance of Life Science Linked Data. In Philipp Cimiano, Oscar Corcho, Valentina Presutti, Laura Hollink, and Sebastian Rudolph, editors, *ESWC*, volume 7882 of *Lecture Notes in Computer Science*, pages 200–212. Springer, 2013.
- [8] Javier D. Fernández, Miguel A. Martínez-Prieto, Claudio Gutiérrez, Axel Polleres, and Mario Arias. Binary RDF representation for publication and exchange (HDT). *Journal of Web Semantics*, 19:22–41, 2013.
- [9] Olaf Hartig and Jun Zhao. Publishing and Consuming Provenance Metadata on the Web of Linked Data. In Deborah L. McGuinness, James Michaelis, and Luc Moreau, editors, *IPAW*, volume 6378 of *Lecture Notes in Computer Science*, pages 78–90. Springer, 2010.
- [10] Tom Heath and Christian Bizer. *Linked Data: Evolving the Web into a Global Data Space*, volume 1 of *Synthesis Lectures on the Semantic Web: Theory and Technology*. Morgan & Claypool, February 2011.
- [11] Robert Isele, Jürgen Umbrich, Christian Bizer, and Andreas Harth. LDspider: An Open-source Crawling Framework for the Web of Linked Data. In Axel Polleres and Huajun Chen, editors, *ISWC Posters & Demos*, volume 658 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2010.
- [12] Krzysztof Janowicz, Pascal Hitzler, Benjamin Adams, Dave Kolas, and Charles Vardeman. Five stars of Linked Data vocabulary use. *Semantic Web*, 5(3):173–176, 2014.
- [13] Tobias Käfer, Jürgen Umbrich, Aidan Hogan, and Axel Polleres. DyLDO: Towards a Dynamic Linked Data Observatory. In Christian Bizer, Tom Heath, Tim Berners-Lee, and Michael Hausenblas, editors, *LDOW*, volume 937 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2012.
- [14] Eetu Mäkelä. *The Semantic Web: ESWC 2014 Satellite Events: ESWC 2014 Satellite Events, Anissaras, Crete, Greece, May 25-29, 2014, Revised Selected Papers*, chapter Aether – Generating and Viewing Extended VoID Statistical Descriptions of RDF Datasets, pages 429–433. Springer International Publishing, Cham, 2014.
- [15] Deborah McGuinness, Timothy Lebo, and Satya Sahoo. PROV-O: The PROV ontology. W3C recommendation, W3C, April 2013.
- [16] Tope Omitola, Landong Zuo, Christopher Gutteridge, Ian Millard, Hugh Glaser, Nicholas Gibbins, and Nigel Shadbolt. Tracing the provenance of linked data using void. In Rajendra Akerkar, editor, *WIMS*, page 17. ACM, 2011.
- [17] Eyal Oren, Renaud Delbru, Michele Catasta, Richard Cyganiak, Holger Stenzhorn, and Giovanni Tummarello. Sindice.com: a Document-Oriented Lookup Index for Open Linked Data. *International Journal of Metadata, Semantics and Ontologies*, 3(1):37–52, 2008.
- [18] Aravindan Raghuv eer. Characterizing machine agent behavior through SPARQL query mining. In *Proceedings of 2nd International Workshop on Usage Analysis and the Web of Data (USEWOD 2012), co-located with the 21st International World Wide Web Conference 2012 (WWW 2012), Lyon, France, april 2012, 2012*.
- [19] Laurens Rietveld, Wouter Beek, and Stefan Schlobach. LOD Lab: Experiments at LOD Scale. In Marcelo Arenas, Oscar Corcho, Elena Simperl, Markus Strohmaier, Mathieu d’Aquin, Kavitha Srinivas, Paul T. Groth, Michel Dumontier, Jeff Heflin, Krishnaprasad Thirunarayan, and Steffen Staab, editors, *International Semantic Web Conference (2)*, volume 9367 of *Lecture Notes in Computer Science*, pages 339–355. Springer, 2015.
- [20] Laurens Rietveld and Rinke Hoekstra. Man vs. machine: Differences in sparql queries. In Bettina Berendt, Laura Dragan, Laura Hollink, Markus Luczak-Rösch, Elena Demidova, Stefan Dietze, Julian Szymanski, and John G. Breslin, editors, *Proceedings of the ESWC2014 workshop on Usage Analysis and the Web of Data (USEWOD), Anissaras, Crete, Greece, May 2014, 2014*.
- [21] Laurens Rietveld and Rinke Hoekstra. YASGUI: Feeling the Pulse of Linked Data. In Krzysztof Janowicz, Stefan Schlobach, Patrick Lambrix, and Eero Hyvönen, editors, *EKAW*, volume 8876 of *Lecture Notes in Computer Science*, pages 441–452. Springer, 2014.
- [22] Max Schmachtenberg, Christian Bizer, and Heiko Paulheim. Adoption of the Linked Data Best Practices in Different Topical Domains. In Peter Mika, Tania Tudorache, Abraham Bernstein, Chris Welty, Craig A. Knoblock, Denny Vrandečić, Paul T. Groth, Natasha F. Noy, Krzysztof Janowicz, and Carole A. Goble, editors, *Semantic Web Conference (1)*, volume 8796 of *Lecture Notes in Computer Science*, pages 245–260. Springer, 2014.