

# A More Decentralized Vision for Linked Data

Axel Polleres<sup>a,b,\*</sup>, Maulik Rajendra Kamdar<sup>c</sup> and Javier David Fernández<sup>a,b</sup> Tania Tudorache<sup>c</sup>  
Mark Musen<sup>c</sup>

<sup>a</sup> *Institute for Information Business, Vienna University of Economics and Business, Vienna, Austria*

<sup>b</sup> *Complexity Science Hub Vienna, Vienna, Austria*

*E-mails: axel.polleres@wu.ac.at, javier.fernandez@wu.ac.at*

<sup>c</sup> *Department first, University or Company name, Abbreviate US states, Country*

*E-mails: tudorache@stanford.edu, musen@stanford.edu*

**Abstract.** In this *deliberately provocative* position paper, we claim that more than ten years into Linked Data there are still (too?) many unresolved challenges towards arriving at a truly machine-readable *and* decentralized Web of data. We take a deeper look at key challenges in usage and adoption of Linked Data from the ever-present “LOD cloud” diagram.<sup>1</sup> Herein, we try to highlight and exemplify both key technical and non-technical challenges to the success of LOD, and we outline potential solution strategies. We hope that this paper will serve as a discussion basis for a fresh start towards more actionable, truly decentralized Linked Data, and as a call to the community to join forces.

**Keywords:** Linked Data, Decentralization, Semantic Web

## 1. Decentralization Myths on the Semantic Web

Let us start with a rant, arguing that the Semantic Web may well be considered a story of failed promises with regards to decentralization:

- We had hopes (as a community) to revolutionize Social Networks in a way that every data subject owns and controls their social network data in **decentralized FOAF** [2] files published in their personal Web space – we got siloed, centralized social networks (Facebook, LinkedIn). Attempts to re-decentralize the Social Web, for instance, through the work of the W3C Social Web WG<sup>2</sup> appear not to have found major adoption at a level comparable with these siloed sites.<sup>3</sup>

- We envisioned a **decentralized network of ontologies on the Web** that would enable smart agents to seamlessly talk to each other, and that would enable easy integration of data by following the guiding principles of ontology engineering and Gruber’s often cited vision of ontologies as shared conceptualizations [3].<sup>4</sup> While there are indeed certain areas in which ontologies are used to share conceptualizations of a domain, mostly these are insular efforts that do the job well for a particular community. However, on Web scale, ontology and vocabulary reuse is still extremely limited. Instead, we got a main *central* schema (schema.org), and fast-growing community projects like Wikidata [4] refusing to buy

\*Corresponding author. E-mail: axel.polleres@wu.ac.at.

<sup>1</sup>The Linking Open Data cloud diagram, available at <http://lod-cloud.net/>, which has been regularly updated since 2007 by Andreas Abele, John P. McCrae, Paul Buitelaar, Anja Jentzsch and Richard Cyganiak, with its latest version having been created in March 2019 [1].

<sup>2</sup><https://www.w3.org/wiki/Socialwg>

<sup>3</sup>While there is some hope left, in ActivePub being picked up by several implementations, cf. <https://en.wikipedia.org/w/index.php?title=ActivityPub&oldid=841568831#Implementations>, reversing

the network effects that have drawn a critical mass of users to these siloed sites seems still far away.

<sup>4</sup>or, as Dan Brickley, one of the inventors of FOAF stated slightly sarcastically in personal communication: “we took one useful feature of RDF/RDFS (fine grained vocabulary composition) and elevated it to a cult-like holy law, to the extent that anyone who created a useful RDF vocabulary and wanted to keep improving it, found themselves pushed instead into combining it with dozens of other half-finished, poorly documented efforts that weren’t really designed to fit together nicely.”

into the need for re-using terms from other ontologies.<sup>5</sup>

- We put a lot of effort into **formal semantics and clean axiomatization** of those ontologies – we got logical inconsistency.<sup>6</sup> Whereas, serious attempts to apply such reasoning about Web Data in the wild have either had to restrict themselves to lightweight ontologies or have not been further developed in the past five years, with (a) the semantics of OWL [9] and even parts of RDF(S) [10] turning out to be too hard to grasp for normal Web users and developers to survive in the World Wild Web [11, 12]; and (b) the DL community mostly having turned their back to seriously taking the challenge of decentralized applications at Web scale.
- Berners-Lee et al. in their original Semantic Web article [13] promised **Web-scale automation**: automated calendar synchronisation, personalised health care assistance, home automation – some of these applications are a reality now (Amazon Alexa’s home control, or Google’s and Apple’s widely used services), but rather than relying on a decentralized Semantic Web, use single companies’ curated knowledge bases – also now called “Knowledge Graphs” – that enhance these companies’ services’ backend systems.
- More specifically, we see **knowledge graphs** evolve and embrace them as a success story of the Semantic Web. Yet a good definition of what a Knowledge Graph is and what differentiates it from an “ontology” is still to be provided – apart from the single distinguishing feature that all known examples of knowledge graphs (Google’s, Bing’s, and Yahoo’s knowledge graphs as well as their open pendants DBpedia and Wikidata) are NOT decentralized.

So, here we are... however, there is one lighthouse project that clearly has implemented the vision of a decentralized Semantic Web, this single project that we,

<sup>5</sup>The main reason for Wikidata not to prescribe existing vocabularies was to leave the community freedom to link and use what they deem useful within one consistent scheme/namespace: one of the reasons was to avoid the needed buy-in to existing ontologies, the popularity of which or agreement about could shift over time. Therefore, they “left it to the community to choose a stronger semantics - like OWL - or a weaker semantic - like SKOS[5] or not” (personal communication Denny Vrandečić).

<sup>6</sup>Even within DBpedia [6, 7], the central crystallization point of the LOD cloud [8].

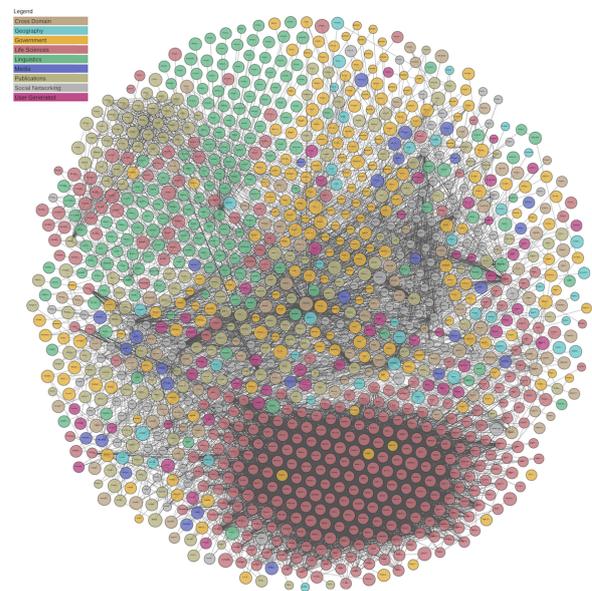


Figure 1. The “LOD cloud” diagram [1], from April 2018, counting 1,184 datasets.

as a community, hinge upon and tend to accept as a clear success to wipe away all the failed promises mentioned above is: Linked Data [14]. The promise to be able to publish structured data in a truly decentralized fashion, with a couple of simple principles to enable the automatic retrieval and integration by just “following your nose”, i.e., dereferencing HTTP links. This principle is the most powerful promise that filled the community with new enthusiasm through the so-called “LOD cloud”, cf. Fig. 1. If we measure the number of datasets published according to the *four linked data principles* [15] and that link to each other, we find evidence of growth and prosperity (cf. Fig. 2), and hope to finally make the vision of a decentralized Web of data come true. Meanwhile, indeed this “cloud” contains over 1,184 datasets, which should be considered good news.

However, as we will discuss in the present paper, there are still serious barriers to consume and use this data. Thus, we would like to take a step back and assess the situation. We will identify some serious challenges in consuming and using Linked Data from the “cloud”, wherein we have to question the usefulness of the current LOD cloud, and, finally, we call for a more principled and, in our opinion, more useful restart and for more collaboration and decentralization in the community itself.

Along these lines, in the remainder of this paper, we start with some background on the genesis of the cur-

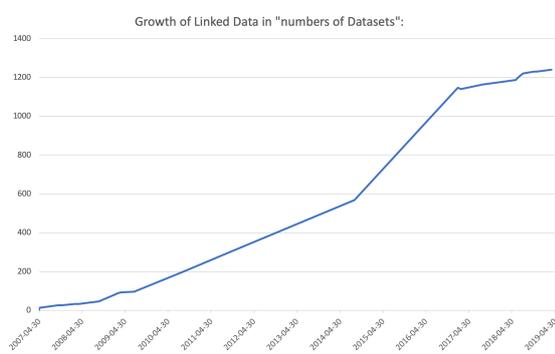


Figure 2. The growth of the “LOD cloud” in number of datasets seems to indicate steady, while not rapid or even overwhelming adoption; we still have to view this as opposed to the probably much more rapid growth of other parts of the Web in the same time period [16]

rent LOD cloud in Section 2. We will then highlight five perceived main challenges we deem important to be addressed to make Linked Data more usable and, therefore, useful. These challenges will be presented – by examples and discussing their implications – in the course of Section 3. Finally, we conclude with a call to collaboratively and openly address these challenges as a community in order to (re-)decentralize the Semantic Web again in Section 4.

## 2. Background: The genesis of the LOD Cloud

The creation of a complete Web index is a never-ending story. Since the early days of the Linked Data Web, several attempts have been created and failed to sustain exhaustive Linked Data Search engines, such as Sindice [17], SWSE [18], Watson [19], Swoogle [20], just to name a few. Typically based on bespoke, crawler-based architectures, these search engines relied on either (i) collecting data published under the Linked Data principles and particularly applying the “follow-your-nose” approach enabled through these principles (i.e., find more Linked Data by dereferencing links appearing in Linked Data), and sometimes (ii) relying on “registry” or “pingback” services to collect and advertise linked data assets, such as Semantic Pingback [21]. In the meantime, unfortunately all of these search engines have been discontinued, and we are not aware of any active, public Semantic Pingback services. As more recent efforts, the LOD-Laundromat project [22] offers an URL lookup

service<sup>7</sup> generated from/for the (accessible parts of) the LOD cloud, and also the LOD Cache by Open-LinkSW<sup>8</sup> remain available for LOD entity lookups and SPARQL queries, although it does not provide a detailed specification of which datasets it indexes.

Both of these more recent efforts though, claim to refer to datasets in the LOD cloud: the LOD cloud diagram [1] took a different approach, that is, it has been generated from metadata provided by the community at a (CKAN-driven) Open Data portal, namely <http://datahub.io>. Interestingly, this is only confirmed for its prior version in August 2017, as references to datahub.io have been removed from current, later LOD diagram versions; also note that datahub.io has moved in the meantime and the “old” LOD-cloud dataset metadata descriptions are only available via the suggestively “deprecated” URL <https://old.datahub.io/>. While the LOD-cloud initiative itself seems to have been suffering from starvation as well, the current noble effort is depending on a few individuals, such as John McCrae et al. (the creators and maintainers of the LOD cloud diagram), which seems to put the initiative at risk; at least, there is recent active development, with regular at least bimonthly updates on the lod-cloud between April 2018 and April 2019.

Still, the LOD-cloud at lod-cloud.net and the metadata at datahub.io seem to remain the single most popular entry points to Semantic Web data (with the exception of domain specific portals such as Bioportal [23]), and therefore a bottleneck.

The metadata the LOD cloud relies on, comprises metadata fields such as:<sup>9</sup>

- **tags**, where as a pre-filter, only those datasets are included in the cloud that have the tag “lod”,
- **link descriptions**, i.e. declarations of numbers of links to other datasets,
- **resources**, that is, URLs to access the dataset in the form of e.g. dumps, as SPARQL endpoints, or semantic descriptions (e.g. in the form of a Void [24] descriptions) or an XML sitemap.

Apart from the LOD cloud, a similar effort exists to collect and register Linked Data *vocabularies* and document their interconnections in the Linked Open

<sup>7</sup><http://lotus.lodlaundromat.org/>

<sup>8</sup>[lod.openlinksw.com/](http://lod.openlinksw.com/)

<sup>9</sup>Disclaimer: Note that our observations base on metadata from datahub.io collected in April 2018; since then, lod-cloud.net has discontinued on datahub.io and now provides an own form-based submission system for metadata on its webpage.

Vocabularies (LOV) project by Vandenbussche et al. [25]. As opposed to the purely meta-data based approach of the LOD cloud collection, LOV relies on curation and quality checks, verification of parsable vocabulary descriptions, etc. We note that the distinction between Linked “vocabularies” and “data” is not always straightforward, with for instance the entries of the BioPortal [23], a registry of ontologies (which could by definition be considered as well as vocabularies), being (an in fact significant) part of the LOD cloud, but not being present in LOV.

So, where does this leave us? We have seen a lot of resources being put into publishing and using Linked Data, but yet a publicly widely bisible “killer app” is still missing. The reason for this, in the opinion and experiences of the authoes lies all to ofen in the frustrating experiences when trying to actually use Linked Data for building actual applications. Many attempts and projects end up in the end to still use a centralized warehousing approach, integrating a handful of data sets directly from their raw data sources, rather than being able to leverage their “lifted” Linked Data versions. Linked Data still too often proofs to be largely insufficient to sustain any real application, plus if a central warehouse is used, the use and benefits of RDF and Linked Data over conventional databases and warehouses technologies, where more trained people are available, remain questionable. In the following, we will elaborate on the main reasons for this current state, as we perceive them, however, with a hopeful perspective, that is, sketching solution paths to overcome these challengenges.

### 3. Key Challenges in usage and adoption of Linked Data

Reasons for LOD not yet having reached its full potential are manifold and not simple, and we do not claim to be exhaustive herein; yet, we would like to provide a list from the experiences of the authors to help explain some major challenges in the current state of affairs around LOD. We have chosen to divide reasons into technical and non-technical underlying challenges.

#### 3.1. Technical challenges

The current model of collection of LOD by meta-data published once-off by the creators of datasets has lead to mainly a nice drawing, rather than mak-

ing Linked Data accessible and usable. In fact, we see the following major challenges when attempting to use Linked Data, parts of which we underpin by some preliminary analyses on the metadata from old.datahub.io; we are obviously not the first ones to recognize these as such, wherefore we will accompany them with similar analyses and references where available. Yet, we focus on challenges which we believe to need a solution first, before we can dream about federated queries or optimizing query answering over linked data (which is what we do mostly in our research papers now — without practical applications over *several datasets in real existing Linked Data*).

##### 3.1.1. Availability and resource limits.

As a result of a recent analysis we did over the metadata on datahub.io, we unfortunately confirmed a very low level of availability of resources, which was already identified as one of the main challenges in the biomedical domain: among the mentioned 5435 resources in the 1281 “LOD”-tagged datasets on datahub.io, there are only 1917 resources URLs that could be dereferenced. Among all the datasets only 646 dataset descriptions contain such dereferenceable (not counting links to PDF, CSV, TSV files) resource URLs; i.e., almost half, 635 dataset descriptions contain no dereferenceable resource URLs that would point to data at all. We applied a best effort here, that is dereferencing both HTTP and FTP URLs with a timeout of 10 seconds awaiting a potential response, counting all 2xx return codes for a HEAD request for HTTP (and following redirects), or, resp. LIST requests for the containing directory for FTP as positives. This confirms the similar experiments by Debattista in his thesis [26, Section 9] and in a more recent article [27]; many LOD cloud datasets are indeed not even being mentioned in his quality assessment framework<sup>10</sup>, which only covers 130 accessible datasets.

We note that even a best effort of availability could be viewed as optimistic, if we look in a finer grained analysis of the different different formats in these URLs, cf. Figure 3, e.g. concerning SPARQL endpoints: indeed our small experiment reconfirms that, among the mentioned 444 potential SPARQL endpoint URLs in metadata, only 252 responded at all, and only 195 responded “true” to a simple ASK { ?S ?P ?O } query.<sup>11</sup> Table 1 shows the numbers for re-

<sup>10</sup><http://jerdeb.github.io/lodqa>

<sup>11</sup>Also, some endpoint implementations returned non-SPARQL-protocol-conformant results such as <http://identifiers.org/services/>

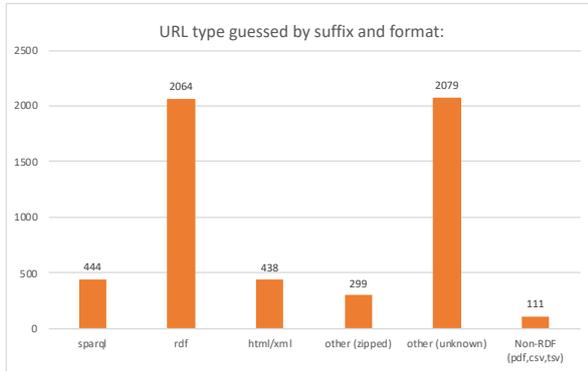


Figure 3. types of URLs in the “LOD cloud” guessed by declared metadata format and suffixes.

query	conformant responses
ASK {?S ?P ?O }	195 (true) + 7 (false)
ASK { }	150 (true) + 7 (false)
ASK {GRAPH ?G { ?S ?P ?O } }	192 (true) + 9 (false)
ASK {GRAPH ?G { } }	146 (true) + 11 (false)
SELECT (count(*) AS ?C) WHERE {?S ?P ?O }	143 (137 non-zero)
SELECT (count(*) AS ?C) WHERE { GRAPH ?G { ?S ?P ?O } }	134 (132 non-zero)

Table 1

SPARQL protocol conformant responses out of the 251 of overall 440 endpoints that responded at all.

sponding endpoint (without timeouts) to a set of test queries, which seem to indicate a considerable number of non-responding and also non-SPARQL-protocol-conformant endpoints.

*Towards a solution path:* As a part of a solution path, we view regular monitoring frameworks like SPARQLES [28],<sup>12</sup> or the Dynamic Linked Data Observatory[29],<sup>13</sup> as essential, which both (i) assess which parts of the LOD cloud are still “alive” and also (ii) could notify the providers and publishers about potential problems. Similarly, Debattista’s fine-grained quality framework mentioned above, aimed originally at re-assessing and testing LOD data regularly could be a valid starting point, but seems to not have been updated since 2016. The same applies for the LOD Laundromat crawl, which is not updated on a regularly basis.

sparql which returns “false” on above ask query, although clearly its default graph is not empty.

<sup>12</sup><http://sparqles.ai.wu.ac.at/>

<sup>13</sup><http://km.aifb.kit.edu/projects/dyldo/>

Outdated, as well as non-available data is worthless and the frustrating experiences of not finding half the resources when trying to retrieve Linked Data, rather jeopardizes the LOD initiative than inviting externals to our own close community to buy in to the ideas of Linked Data. That is, the LOD cloud itself needs to be “live” and providers that do not comply with minimal availability over a certain duration should be notified and removed. Also, notoriously outdated, stale data should not be listed.

### 3.1.2. Size and Scalability.

The situation in terms of dataset sizes have changed dramatically since the early days of semantic search engines, where relatively small amounts of triples could be feasibly managed in a single triple store: few datasets generated from big databases reach dramatic sizes. For instance, the latest edition of DBpedia (2016-10), consists of more than 13 billion triples, Wikidata comprises +5B triples and the whole LOD-Laundromat project, which attempts to process and cleanse the accessible part of the LOD cloud, reports at the moment 38.8b indexed triples.

We also note that, to the best of our knowledge, current triple stores on commodity servers do not scale up to more than 50b triples, apart from lab experiments on hardware probably not yet available to most research labs in our community: AllegroGraph and Oracle triple stores have reported dealing with up to 1 trillion triples.<sup>14</sup>

We already see sizes of triples reported on the LOD cloud diverging from what a simple SELECT (COUNT (\*) AS ?C) WHERE {?S ?P ?O} to their respective endpoints reports in various examples, just to name some: the Pubmed-Bio2RDF endpoint<sup>15</sup>, reports 1.37b triples on the query above,<sup>16</sup> whereas the dump<sup>17</sup> reports 1.8b triples. Yet again, on a side note, different to both of that, the metadata at datahub.io reports 5b triples for the same dataset,<sup>18</sup> where however it cannot be easily determined in how far these numbers refer to different versions or subsets of the

<sup>14</sup>cf. <https://www.w3.org/wiki/LargeTripleStores>, last retrieved 2018-05-16, where we note that these experiments have been conducted on synthetic LUBM data, which does not necessarily reflect the characteristics of Linked Data “in the wild”.

<sup>15</sup><http://pubmed.bio2rdf.org/sparql>

<sup>16</sup>The same number is returned on a query for quads, i.e. SELECT (COUNT (\*) AS ?C) WHERE {GRAPH ?G {?S ?P ?O}}, which is of course not necessarily the case for all SPARQL endpoints.

<sup>17</sup><http://download.bio2rdf.org/#/release/4/pubmed/>

<sup>18</sup><https://old.datahub.io/dataset/bio2rdf-pubmed>

dataset. Likewise, Wikidata’s query service responds to the same query a number of 5.2b triples, which is significantly lower than the 5.7b triples we retrieved from the dump mentioned above.

In addition to that, it is mostly impossible to indeed retrieve all triples from a SPARQL endpoint, due to result size restrictions that many endpoints apply, either in the form of timeouts or only returning a certain maximum number of results/triples. For details, see also [30], which discusses some of these restrictions, and also explains, why in general they cannot be trivially circumvented, e.g. by “paging” results with LIMIT and OFFSET. As another example of related problems, Uniprot, reported to have +39b triples served on its public endpoint, cf. Footnote 14, times out on the simple query to count its triples mentioned above.

Another potential challenge in terms of size and scalability is the amount of duplicates in current dumps: as an example, the PubMed RDF dump from Bio2RDF we mentioned above, cf. Footnote 17, consists of +7.22b nquads spread over 1151 dump files. A lot of triples are actually duplicated across these dump files from the same dataset; downloading all of these and de-duplicating them locally both wastes bandwidth and makes processing such dumps unnecessarily cumbersome.

*Towards a solution path:* It seems that in order to avoid both such discrepancies and bottlenecks for downloads and query processing, a combination of (i) dumps provided in HDT [31], a compressed and queryable RDF format, as well as (ii) Triple Pattern Fragments (TPF) endpoints [32] as the standard access method for Linked Datasets could alleviate some of these problems: the triple-patterns fragment interface – essentially limits queries to an endpoint to simple triple matching queries which offloads processing of complex joins and other operations to the client-side, while still not having to download complete dumps. HDT,<sup>19</sup> on the other hand is an already compressed dump-format that allows such triple pattern queries without decompression and also guarantees duplicate-freeness. Notably, there are already several TPF endpoints available,<sup>20</sup> most of them powered by HDT in the backend, thus creating a small server-footprint and -load, for either answering triple pattern queries or downloading the whole dump. HDT has also been recently extended

to handle also quads besides RDF triple dumps, thus also being usable for datasets consisting of different (sub)graphs [33]; an analogous extension of the TPF interface to quads would be straightforward. Lastly, we note that e.g. the number of triples it encoded and stored during dump generation in the metadata header of HDT files, thus providing a single, reliable entry to the dataset size.

### 3.1.3. Findability and (Meta-)Data Formats.

The current metadata available on the LOD cloud does not tell us a lot about how to access the single datasets.

Over time, various dataset description formats and mechanisms have been proposed, typically (i) VoID descriptions, (ii) (Semantic) Sitemaps, and (iii) SPARQL service endpoint descriptions. In the following, we analyze the current state of affairs in the LOD cloud.

*The Vocabulary of Interlinked Datasets (VoID)* had been designed as a minimalistic entry point for describing datasets and how to access them, containing properties for locating dumps (`void:dataDump`), finding SPAQL endpoints (`void:sparqlEndpoint`) or describing the size of the dataset in terms of numbers of triples (`void:triples`) and other structural statistics. In order to find the VoID description, it is suggested to place the dataset description under `/.well-known/void` in the root directory of a Web-server.

There are various problems with this approach: firstly, different datasets hosted under one common domain/server cannot provide different dataset descriptions; as illustration, obviously for Github hosted data, `https://github.com/.well-known/void` would not return a valid VoID description, although github is gaining popularity for hosting Linked Data sets. Secondly, even the “epicenter” of the LOD-cloud, `dbpedia.org` does not follow the rules and provides a VoID description at the non-obviously findable URL `http://dbpedia.org/void/page/Dataset` instead. Lastly, indeed, among all 881 hostnames mentioned in URLs in `datahub.io`’s metadata, 159 respond to an HTTP Get with this recipe, at least 75 of which though seem to be HTML responses, and only 56 valid RDF;<sup>21</sup> without going into further detail, even if the HTML contained RDFa (which in the cases we inspected it did not), it seems that easy to parse RDF re-

<sup>19</sup><http://rdfhdt.org>

<sup>20</sup><http://linkeddatafragments.org/>

<sup>21</sup>We tested all hosts from the URLs that provided non-error results

sults with valid VoID descriptions seem to be the exception.

(*Semantic*) *Sitemaps* XML Sitemaps<sup>22</sup> seem to be a more commonly implemented pattern to discover data and pages accessible via an HTTP server, not least because of their recommendation by search engines. It is a simple XML format that should guide crawlers across sites, where Tummarello et al. had even proposed an extension of the Sitemaps protocol to link to RDF datasets specifically [34], that has been implemented in Sindice [17]. Sitemaps are expected to be found under the root of a dataset’s directory on a host in a file called ‘sitemap.xml’, that is, not necessarily directly underneath root directory of the host address. datahub.io’s metadata contains hints (by filename) to such sitemaps for 57 datasets, 56 indeed returning valid sitemaps, and 55 of which indeed use the semantic sitemap extension [34] (52 containing a `sc:dataDump` attribute and 53 containing a `sc:sparqlEndpoint` field). So, overall, while semantic sitemaps are only used for a marginal 5% of the datasets in datahub.io, they seem to be fairly consistent.

*SPARQL service endpoint descriptions* according to the SPARQL1.1 specification, “*SPARQL services made available via the SPARQL Protocol SHOULD return a service description document at the service endpoint when dereferenced using the HTTP GET operation without any query parameter strings provided. This service description MUST be made available in an RDF serialization, MAY be embedded in (X)HTML by way of RDFa [RDFa], and SHOULD use content negotiation [CONNEG] if available in other RDF representations.*” Yet, out of the 251 potential respondent endpoint addresses mentioned above only 136 respond to this recipe, out of which in fact 63 return HTML (mostly query forms), even if attempting CONNEG.<sup>23</sup>

We note that while some of these mentioned HTML responses *might* contain RDFa, it is still an extra step to extract and parse and each such extra step will bloat a potential consuming client unnecessarily. Similarly, when attempting to find data dumps, without a semantic sitemap or a VoID file in place, our best guess would be to guess and try parsers from “format” descriptors in the metadata or from filename suffixes.

<sup>22</sup><https://www.sitemaps.org/protocol.html>

<sup>23</sup>with sending an ‘Accept: text/turtle, application/n-triples, application/trig, application/n-quads, application/rdf+xml, \*’ header.

An additional complication here are compressed formats, where attempting different decompression formats (gzip, bzip, tar, zip, just to name a few), sometimes even used in combination, further complicate accessibility. Some of the the guessed formats we found in all URLs are listed again in Fig. 3 above.

We also note that by manual inspection, some endpoint addresses or accessibility of datasets could be recovered, but since we herein would like to emphasize on machine accessibility, manual “recovery” seems an undesirable option.

*Towards a solution path:* We feel that as for automatic findability, Semantic Sitemaps with pointers to a VoID description, with concrete pointers to primarily a dump, preferably in HDT as well as (optionally) a pointer to a SPARQL endpoint (or TPF endpoint) should be the commonly to be agreed upon practice. We note here, that the use of HDT makes this task even simpler, as indeed the Header part of an HDT dump file holds a place for metadata descriptions about the dataset readily.<sup>24</sup> Also, SPARQL endpoints should provide service descriptions in easily accessible RDF (not RDFa) available via CONNEG, where again these SPARQL service descriptions should describe service limitations (such as e.g. result size limits or connection limits and timeouts). Also, the service description should declare potential differences between the data in the dump and in the endpoint, if any. We emphasize here, that to the best of our knowledge there is no agreed upon vocabulary for SPARQL endpoint restrictions and capabilities.

### 3.1.4. “RDF Data Quality” of Datasets and the “Semantics of Links”.

The linked data principles define rough guidelines on dereferenceability and linkage of datasets, yet in order for RDF datasets, once downloaded, to be truly machine-processable and being able to traverse and interpret those links fruitfully, more detailed guidelines seem to be indispensable: in an early approach, Hogan et al. proposed the “Pedantic Web”[35] alongside with an in the meanwhile discontinued tool, RDFAlerts, to check and assess the quality, dereferenceability, and finally syntactical (e.g. use of ill-defined literals) log-

<sup>24</sup>In fact, some automatically computable VoID properties are already computed and included in HDT’s header per default, and it is well possible to add additional properties such as pointers to (SPARQL or Linked Data fragments) endpoints, or used namespaces within this header, as a single point of access through an HDT dump file.

ical consistency (in terms of RDFS/OWL inferences, use of literals in place of object properties, availability of definitions for used properties and classes, etc.) of RDF datasets. A lot of these checks though, were not necessarily designed to scale to datasets of billions of triples, or, resp. should be reassessed in terms of feasibility. Again, HDT could serve as a basis for scalable, out-of-the-box implementations of such checks on a dataset level.

Besides the aforementioned “semantic heterogeneity” issue in the biomedical domain, as a particular additional example of checks that should be automatically performed on a dataset level, we mention the links in the LOD cloud diagram, shall indicate in how far one dataset links to another dataset; to the best of our knowledge, these links and their strength, have been created so far from datahub.io’s metadata field `links:<Dataset-acronym>`, i.e. been typically manually specified by the contributors of said metadata: the definition for how such links should be declared on lod-cloud.net provides the following inclusion/exclusion criterion for datasets in the LOD cloud: “The dataset must be connected via RDF links to a dataset that is already in the diagram. This means, either your dataset must use URIs from the other dataset, or vice versa. We arbitrarily require at least 50 links.” An older version of the page also provided a slightly more concrete definition of what is meant by a link here: “A link, for our purposes, is an RDF triple where subject and object URIs are in the namespaces of different datasets.” We however find this definition hard to assess. Since so concrete guideline with regards to “ownership” of name spaces is provided here, any attempt to compute such links automatically is doomed to fail. As from our observation when investigating different datasets, it is by no means always clear

1. to which namespace a URI belongs, or
2. to which dataset a namespace belongs

As for 1, we note that in many cases it is not even clear entirely purely from the RDF data which part of the URIs in a dataset denote namespaces: namespaces and qnames in RDF have no special status as in XML, they simply denote prefixes; while certain “recipes” for such prefixes exist, such as most commonly used ‘/’ and ‘#’ prefixes, some ontologies use completely different recipes to separate identifiers from prefixes. In fact, various datasets “mint” URIs with differing recipes, for instance, we find the prefix scheme `http://bioonto.de/sbml.owl#Uniprot:` within the BIOMODELS ontology from Bioportal, with 562

identifiers using this scheme, e.g.

`http://bioonto.de/sbml.owl#Uniprot:Q9UJX6.`

In this case, what is the namespace prefix? It seems intuitive that this URI minting scheme is referring to UNIPROT which indeed means the dereferenceable URL

`https://www.uniprot.org/uniprot/Q9UJX6.`

Now, at a closer look this example<sup>25</sup> illustrates several problems at once:

- it is unclear which prefix denotes the “namespace”: `http://bioonto.de/sbml.owl#` or rather `http://bioonto.de/sbml.owl#Uniprot:?`
- the same entities exist in the LOD cloud under different, disconnected namespace prefixes, such as the Uniprot identifier Q9UJX6, the “official” prefix (as per the authoritative pay-level-domain `uniprot.org`) of which is `http://purl.uniprot.org/uniprot/Q9UJX6`.
- likewise, the overall “#namespace” `http://bioonto.de/sbml.owl#` does not refer to a dereferenceable URI; the data itself comes in fact from a dataset dump in an old version of bioportal, that has been fixed in the meantime, but nonetheless it serves for illustration; a detailed analysis of present such quality issues in the LOD-cloud is still on our agenda, but we have reason to believe that many such issues still persist also in the current LOD cloud. In fact, the example BIOMODELS ontology dataset now exists on different places in the LOD cloud, within BIO2RDF, within BIOPORTAL, but also as an RDF dataset directly published by EBI<sup>26</sup> in three different “RDF exports” of the same database.

While – depending on the serialisation – namespaces could be filtered out based on being explicitly represented (e.g. marked with XML namespaces in RDF/XML or by @prefix declaration in Turtle, respectively, this seems not to be a reliable way of recognizing all used namespaces within an RDF datadump in a declarative machine-readable manner. Plus, as the example illustrates, even if we had all namespaces occurring within a dataset, various URL schemes used refer to either non-dereferenceable or non-RDF publishing third-party namespaces, that cannot be simple assigned to “belonging” to a single dataset. More issues about

<sup>25</sup>which is one of many, we emphasize it is not our intention to point fingers to anyone!

<sup>26</sup> at `ftp://ftp.ebi.ac.uk/pub/databases/RDF/biomodels/`

1 URI schemes and namespaces and term (non-)re-use  
2 have been described in [36] and [37].

3 Last, but not least, as an open problem, links in  
4 one dataset always refer to a particular *version* of the  
5 linked dataset, which if not archived cannot be guaran-  
6 teed to persist or being dereferenceable in the future.  
7 For a more sustainable version of Linked Open Data,  
8 we therefore deem versioned Linked Data as well as  
9 archives a necessity.

10 *Towards a solution path:* We feel that in order to  
11 avoid such issues, to be established best practices for  
12 Linked Data publishing would need to provide more  
13 guidelines for URL minting and reuse. Namespace and  
14 ID minting should probably be restricted to machine-  
15 recognizable patterns (such as strict adherence to ‘/’  
16 and ‘#’-namespaces), with dereferenceable names-  
17 pace URLs). Ownership of a namespace could – for in-  
18 stance – be restricted to pay-level-domain, that is, def-  
19 inition of namespaces being restricted to the own pay-  
20 level domain, and URL and namespace schemas given  
21 a clear machine-readable ownership relation. We leave  
22 a concrete definition of such a machine-readable and  
23 assessable ownership open for now, but refer to simi-  
24 lar concepts and thoughts about URI “authority” hav-  
25 ing been discussed before in the context of ontolog-  
26 ical inference by Hogan in his thesis [38, Section 5]  
27 as a potential starting point. Hogan’s thesis also con-  
28 tains some details on scalable implementations of the  
29 above-mentioned checks that have been described in  
30 RDFAlerts [35] earlier, which we believe could be im-  
31 plemented directly and efficiently on top of indexed  
32 compressed formats HDT, which we leave to future  
33 work on our agenda for now.

34 As for archiving and versions, we refer to [39]  
35 and references therein in terms of starting points; al-  
36 though no single agreed proposal exists at this point  
37 for how to publish versioned RDF archives we again  
38 refer to possible HDT-based solutions, particularly en-  
39 abled through the recent extension of HDT to handle  
40 quads [33].

### 41 3.2. Non-Technical Challenges

42  
43  
44  
45 Even if we will be able to solve all the above tech-  
46 nical challenges, there are several pertinent issues that  
47 are in the critical path to the success of LOD. That is,  
48 we also see many non-technical challenges that should  
49 be fixed in order to stimulate adoption of linked data, a  
50 non-exhaustive list of which we briefly describe here-  
51 after.

#### 1 3.2.1. Completeness/Consistency.

2 Several well-known and important RDF datasets  
3 are missing in the LOD cloud, e.g. EBI RDF is not  
4 there (plus various other well-known data bases from  
5 the biomedical and life sciences domain), which have  
6 gone through the effort of publishing RDF, but not  
7 taken the additional hurdle of manually adding and up-  
8 dating their metadata in yet another centralized catalog  
9 such as datahub.io. For similar reasons, e.g., Wikidata  
10 is not a dataset in the LOD cloud, although it is clearly  
11 linked well with several datasets present.

12 Overall, the burden of manually and pro-actively  
13 needing to provide and maintain LOD cloud metadata  
14 on the publisher-side has proven unsustainable.

#### 15 3.2.2. Trust.

16 Besides the pervasive issues of availability and reli-  
17 ability, developers are rightfully worried that the pub-  
18 lished data in the cloud is not kept up to date, and  
19 as such the technical issues mentioned above might  
20 overall give rise to (or have already given rise to,  
21 possibly) doubts on the technology and principles of  
22 Linked Data. Stale datasets, while still available, but  
23 with outdated, once-off RDF exports of in the mean-  
24 time evolved databases, likewise raise trustworthiness  
25 issues in Linked Data.

26 While it seems to have been a sufficient incentive  
27 to “appear” in the LOD cloud to publish datasets ad-  
28 hering to Linked Data principles, a similarly strong in-  
29 centive to sustain and maintain quality of published  
30 datasets seems to be missing.

31 It is therefore important for us as a community to  
32 keep this project up and alive, by creating sustainable  
33 publishing and monitoring processes.<sup>27</sup>

#### 34 3.2.3. Governance.

35 We note that not only trust in the LOD cloud itself,  
36 but also mutual trust between LOD providers may be a  
37 problem that is difficult to circumvent. For instance the  
38 presence of various different unlinked “RDF dumps”  
39 or LOD datasets that actually arise from exports of  
40 the same legacy database (BIOMODELS given as *one*  
41 illustrative example of many above) could be poten-  
42 tially related to many of our exports and datasets hav-  
43 ing been created in isolation, by closed groups, without  
44 inviting collaboration or being based on infrastructures  
45 to share and evolve those exports jointly. We feel that

46  
47  
48  
49  
50  
51  
<sup>27</sup>Of course, with the alternative to eventually re-brand it under  
a different name after survival of an “LOD winter” from unfulfilled  
expectations)

1 this issue can only be solved by a more collaborative,  
2 and truly open governance.

### 3.2.4. Documentation and Usability.

3 Besides the technical accessibility discussed above,  
4 usability issues and documentation standards have  
5 been long overlooked in many Linked Data projects.  
6 Industry-strength tools to consume and use Linked  
7 Data with sufficient documentation are still under-  
8 developed.

9 We believe this issue can be ameliorated by: (1) bet-  
10 ter metadata for describing the datasets; (2) better doc-  
11 umentation for using the datasets, including sample  
12 queries; (3) better tool support for enabling reuse of ex-  
13 isting vocabularies; and (4) Supporting and promoting  
14 the use of developer-friendly formats, such as JSON-  
15 LD.

16 In addition, in terms of positive examples, we would  
17 again like to name the aforementioned HDT and TPF  
18 projects, as well as useful SPARQL query editing tools  
19 such as YASGUI [40] or Wikidata's query interface,  
20 which have appeared in the last two years; we need  
21 more tools like those.

### 3.2.5. Funding & Competition.

22 Last, but not least, while the EU and other funding  
23 agencies have supported our endeavor to create a Web  
24 of data greatly, we also feel that there are problem-  
25 atic side effects which need discussion and counter-  
26 strategies:

- 27 – cross-continental research initiatives are not be-  
28 ing funded
- 29 – EU project consortia are typically being judged  
30 by complementary partner expertise

31 Both these factors, which prevent research groups  
32 working on overlapping topics from collaboration, and  
33 rather stimulate an environment of isolated closed re-  
34 search than open collaboration to jointly address the  
35 issues mentioned so far.

36 Lack of collaboration may in other cases also just  
37 be caused by the disconnect of research communities:  
38 this is for instance exemplified by the Semantic Web  
39 in Life Sciences community, for instance seemingly  
40 having recently started efforts very similar to SPAR-  
41 QLES [28] in building up a completely independent  
42 SPARQL endpoint monitoring framework [41],<sup>28</sup> not  
43 even citing SPARQLES (sic!), which seems unneces-  
44 sarily duplicating efforts instead of collaboratively de-  
45 veloping and maintaining such services.

<sup>28</sup>available at <http://yummydata.org/endpoints>

## 4. Conclusions and Next Steps

1 So, is Linked Data doomed to fail? In this paper we  
2 did not present a lot of new insights, but our delib-  
3 eratively provocative articulation of rethinking Linked  
4 Open Data and its principles. It is not too late to coun-  
5 teract and join forces. We hope that our summary of  
6 problems and challenges, reminders of valuable past  
7 attempts to address them, and outline of potential so-  
8 lution strategies can serve as a discussion basis for  
9 a fresh starts ahead towards more actionable Linked  
10 Data.

11 On the bright side, specific communities, such as  
12 the biomedical community have been very successful  
13 in using OWL and Semantic Web technologies for the  
14 management of large biomedical vocabularies and on-  
15 tologies, for a detailed overview of successes in this  
16 area we refer to [42]. Main factors for success projects  
17 are: (1) Having a dedicated and very active develop-  
18 ment team behind it with continuous funding over sev-  
19 eral years; (2) Actively building a strong community  
20 of domain users from different areas, and using their  
21 needs as the driver for the ontology development; (3)  
22 Having an exemplary documentation, about both on-  
23 tology, but also about how to use Linked Data in appli-  
24 cations targeted to domain users, as well as documen-  
25 tation about the processes for building and maintain-  
26 ing collaboratively generated Linked Data sources; (4)  
27 Using a principled approach for developing the under-  
28 lying ontology and maintaining the vocabularies used;  
29 (5) Using automated pipelines to check and ensure data  
30 and vocabulary quality.

31 Our hope is that the Linked Data community can  
32 learn from such specific projects, and that it will try to  
33 apply some of the same approaches that proved to be so  
34 successful. We believe the community needs to work  
35 on those by joining forces, rather than by competition.  
36 We also argued that HDT, a compressed and queryable  
37 dump format for Linked Datasets, could play a central  
38 role as a starting point to address some (but not all) of  
39 the technical challenges we have outlined, i.e., implic-  
40 itly suggesting a "fifth Linked Data principle" [15]:

- 41 5. Publish your dataset as an **HDT dump**, in-  
42 cluding **VOID metadata** as part of its header and  
43 declaring (i) the (authoritatively) **owned names-**  
44 **spaces**, (ii) links to previous and most current **ver-**  
45 **sions** of the dataset, (iii) and – whenever you use  
46 namespaces owned by other datasets or ontolo-  
47 gies – the **links to specific versions of these other**  
48 **datasets**.

In fact, we would argue that more principled Linked Data publishing could allow to auto-generate LOD clouds from a set of such HDT dumps, which to demonstrate is on our agenda for future work.

Apart from technical challenges, other issues arose, that seem equally important, such as the establishment of collaborative and shared research infrastructures to guarantee sustainable funding and persistence of Linked Data assets, as we have seen many promising efforts and initiatives mentioned in this paper having discontinued unfortunately. In the meanwhile, we also emphasize that initiatives like the recently US-founded “Open Knowledge Network”<sup>29</sup> initiative or Dagstuhl seminar on “New Directions for Knowledge Representation on the Semantic Web”<sup>30</sup> have provided platforms to openly discuss such a fresh start, in the context of new trends and efforts around Knowledge Graphs [43] and the FAIR principles [44], that parallel and complement the Linked Data movement.

#### Acknowledgements

A preliminary version of this position paper[45] was presented at the DESEMWEB2018 workshop, where we gained some valuable feedback from the participants. We thank Dan Brickley, Sarven Capadisli, and Denny Vrandečić for comments on the first revision of this paper [46]. Axel Polleres’ work was supported under Stanford University’s Distinguished Visiting Austrian Chair program. Javier Fernández’ work was supported by the EC under the H2020 project SPECIAL and by the Austrian Research Promotion Agency (FFG) under the project “CitySpin”.

#### References

- [1] A. Abele, J.P. McCrae, P. Buitelaar, A. Jentzsch and R. Cyganiak, Linking Open Data cloud diagram (2018-04-30), 2018, From <http://lod-cloud.net/>; retr. 2018/06/01.
- [2] D. Brickley and L. Miller, FOAF Vocabulary Specification 0.99, 2014, <http://xmlns.com/foaf/0.1/>.
- [3] T.R. Gruber, Toward principles for the design of ontologies used for knowledge sharing?, *International journal of human-computer studies* **43**(5–6) (1995), 907–928.
- [4] D. Vrandečić and M. Krötzsch, Wikidata: A Free Collaborative Knowledgebase, *Commun. ACM* **57**(10) (2014), 78–85. <http://doi.acm.org/10.1145/2629489>.
- [5] A. Miles and S. Bechhofer, Simple Knowledge Organization System Reference, Recommendation, W3C, 2009.
- [6] H. Paulheim and A. Gangemi, Serving DBpedia with DOLCE - More than Just Adding a Cherry on Top, in: *The Semantic Web - ISWC 2015 - 14th International Semantic Web Conference, Bethlehem, PA, USA, October 11-15, 2015, Proceedings, Part I*, 2015, pp. 180–196.
- [7] S. Bischof, M. Krötzsch, A. Polleres and S. Rudolph, Schema-Agnostic Query Rewriting in SPARQL 1.1, in: *Proceedings of the 13th International Semantic Web Conference (ISWC 2014)*, LNCS, Springer, 2014. <http://www.polleres.net/publications/bisc-etal-2014iswc.pdf>.
- [8] C. Bizer, J. Lehmann, G. Kobilarov, S. Auer, C. Becker, R. Cyganiak and S. Hellmann, DBpedia - A crystallization point for the Web of Data, *J. Web Sem.* **7**(3) (2009), 154–165.
- [9] M.K. Smith, C. Welty and D.L. McGuinness, OWL Web Ontology Language Guide, 2004, <http://www.w3.org/TR/owl-guide/>.
- [10] D. Brickley and R.V. Guha, RDF Schema 1.1, 2014, <http://www.w3.org/TR/rdf-schema/>.
- [11] B. Glimm, A. Hogan, M. Krötzsch and A. Polleres, OWL: Yet to arrive on the Web of Data?, in: *WWW2012 Workshop on Linked Data on the Web (LDOW)*, 2012. <http://www.polleres.net/publications/glim-etal-2012LDOW.pdf>.
- [12] A. Mallea, M. Arenas, A. Hogan and A. Polleres, On Blank Nodes, in: *Proceedings of the International Semantic Web Conference (ISWC)*, LNCS, Vol. 7031, Springer, 2011. <http://www.polleres.net/publications/mall-etal-2011ISWC.pdf>.
- [13] T. Berners-Lee, J. Hendler and O. Lassila, The Semantic Web, *Scientific American* (2001), 29–37.
- [14] C. Bizer, T. Heath and T. Berners-Lee, Linked Data - The Story So Far, *Int. J. Semantic Web Inf. Syst.* **5**(3) (2009), 1–22.
- [15] T. Berners-Lee, Linked Data, 2006, From <http://www.w3.org/DesignIssues/LinkedData.html>; retr. 2018/06/01.
- [16] A. Polleres, A. Hogan, A. Harth and S. Decker, Can we ever catch up with the Web?, *Semantic Web* **1**(1–2) (2010), 45–52. doi:10.3233/SW-2010-0016.
- [17] E. Oren, R. Delbru, M. Catasta, R. Cyganiak, H. Stenzhorn and G. Tummarello, Sindice.com: a document-oriented lookup index for open linked data., *IJMSO* **3**(1) (2008), 37–52.
- [18] A. Hogan, A. Harth, J. Umbrich, S. Kinsella, A. Polleres and S. Decker, Searching and browsing Linked Data with SWSE: The Semantic Web Search Engine, *J. Web Sem.* **9**(4) (2011), 365–401. <http://www.deri.ie/fileadmin/documents/DERI-TR-2010-07-23.pdf>.
- [19] M. d’Aquin and E. Motta, Watson, More Than a Semantic Web Search Engine, *Semant. web* **2**(1) (2011), 55–63. <http://dl.acm.org/citation.cfm?id=2019470.2019476>.
- [20] L. Ding, T. Finin, A. Joshi, R. Pan, R.S. Cost, Y. Peng, P. Reddivari, V. Doshi and J. Sachs, Swoogle: A Search and Metadata Engine for the Semantic Web, in: *Proceedings of the ACM International Conference on Information and Knowledge Management (CIKM)*, ACM, 2004, pp. 652–659. ISBN 1-58113-874-1. doi:10.1145/1031171.1031289.
- [21] S. Tramp, P. Frischmuth, T. Ermilov and S. Auer, Weaving a Social Data Web with Semantic Pingback, in: *Proceedings of the Knowledge Engineering and Knowledge Management by the Masses (EKAW)*, Vol. 6317 of LNAI, Springer, 2010, pp. 135–149.
- [22] W. Beek, L. Rietveld, H.R. Bazoobandi, J. Wielemaker and S. Schlobach, LOD laundromat: a uniform way of publishing other people’s dirty data, in: *Proceedings of the Inter-*

<sup>29</sup><http://ichs.ucsf.edu/open-knowledge-network/>

<sup>30</sup><https://www.dagstuhl.de/en/program/calendar/semhp/?semnr=18371>

- national Semantic Web Conference (ISWC), Springer, 2014, pp. 213–228.
- [23] M. Salvadores, P.R. Alexander, M.A. Musen and N.F. Noy, BioPortal as a dataset of linked biomedical ontologies and terminologies in RDF, *Semantic Web* 4(3) (2013), 277–284. doi:10.3233/SW-2012-0086.
- [24] K. Alexander, R. Cyganiak, M. Hausenblas and J. Zhao, Describing Linked Datasets with the VoID Vocabulary, 2011, From <https://www.w3.org/TR/void/>; retr. 2018/06/01.
- [25] P. Vandenbussche, G. Atemezing, M. Poveda-Villalón and B. Vatant, Linked Open Vocabularies (LOV): A gateway to reusable semantic vocabularies on the Web, *Semantic Web* 8(3) (2017), 437–452. doi:10.3233/SW-160213.
- [26] J. Debattista, Scalable Quality Assessment of Linked Data, PhD thesis, Rheinische Friedrich-Wilhelms-Universität Bonn, 2016. <http://hss.ulb.uni-bonn.de/2017/4720/4720.htm>.
- [27] J. Debattista, C. Lange, S. Auer and D. Cortis, Evaluating the Quality of the LOD Cloud: An Empirical Investigation, *Semantic Web* (2017), 1–42.
- [28] P.-Y. Vandenbussche, J. Umbrich, L. Matteis, A. Hogan and C. Buil-Aranda, SPARQLES: Monitoring public SPARQL endpoints, *Semantic Web* 8(6) (2017), 1049–1065.
- [29] T. Käfer, A. Abdelrahman, J. Umbrich, P. O’Byrne and A. Hogan, Observing linked data dynamics, in: *Proceedings of Extended Semantic Web Conference (ESWC)*, Springer, 2013, pp. 213–227.
- [30] C. Buil-Aranda, A. Polleres and J. Umbrich, Strategies for Executing Federated Queries in SPARQL1.1, in: *Proceedings of the International Semantic Web Conference (ISWC)*, Springer, 2014.
- [31] J.D. Fernández, M.A. Martínez-Prieto, C. Gutiérrez, A. Polleres and M. Arias, Binary RDF Representation for Publication and Exchange (HDT), *J. Web Sem.* 19(2) (2013). <http://www.websemanticsjournal.org/index.php/ps/article/view/328>.
- [32] R. Verborgh, M.V. Sande, O. Hartig, J.V. Herwegen, L.D. Vocht, B.D. Meester, G. Haesendonck and P. Colpaert, Triple Pattern Fragments: A low-cost knowledge graph interface for the Web, *J. Web Sem.* 37-38 (2016), 184–206. doi:10.1016/j.websem.2016.03.003.
- [33] J.D. Fernandez, M.A. Martínez-Prieto, A. Polleres and J. Reindorf, HDTQ: Managing RDF Datasets in Compressed Space, in: *Proceedings of the European Semantic Web Conference (ESWC)*, Springer, 2018. <http://polleres.net/publications/fern-et-al-ESWC2017.pdf>.
- [34] R. Cyganiak, H. Stenzhorn, R. Delbru, S. Decker and G. Tumarello, Semantic Sitemaps: Efficient and Flexible Access to Datasets on the Semantic Web, in: *Proceedings of the European Semantic Web Conference (ESWC)*, 2008, pp. 690–704.
- [35] A. Hogan, A. Harth, A. Passant, S. Decker and A. Polleres, Weaving the Pedantic Web, in: *International Workshop on Linked Data on the Web (LDOW) at WWW*, 2010.
- [36] M.R.e.a. Kamdar, An Empirical Meta-analysis of the Life Sciences (Linked?) Open Data Cloud, 2018, Unpublished Manuscript, available at <http://onto-apps.stanford.edu/lslodminer>. <http://onto-apps.stanford.edu/lslodminer>.
- [37] M.R. Kamdar, T. Tudorache and M.A. Musen, A systematic analysis of term reuse and term overlap across biomedical ontologies, *Semantic web* 8(6) (2017), 853–871.
- [38] A. Hogan, Exploiting RDFS and OWL for Integrating Heterogeneous, Large-Scale, Linked Data Corpora, PhD thesis, Digital Enterprise Research Institute, National University of Ireland, Galway, 2011, From <http://aidanhogan.com/docs/thesis/>; retr. 2010/10/27.
- [39] J.D. Fernandez, J. Umbrich, A. Polleres and M. Knuth, Evaluating Query and Storage Strategies for RDF Archives, *Semantic Web* (2018), to appear (accepted for publication). <http://www.semantic-web-journal.net/content/evaluating-query-and-storage-strategies-rdf-archives-0>.
- [40] L. Rietveld and R. Hoekstra, The YASGUI family of SPARQL clients, *Semantic Web* 8(3) (2017), 373–383. doi:10.3233/SW-150197.
- [41] A.S. Yasunori Yamamoto Atsuko Yamaguchi, Umaka-Yummy Data: A Place to Facilitate Communication between Data Providers and Consumers, in: *Proceedings of the International Conference Semantic Web Applications and Tools for Life Sciences (SWAT4LS)*, CEUR, Vol. 1795, 2016.
- [42] M.R. Kamdar, A WEB-BASED INTEGRATION FRAMEWORK OVER HETEROGENEOUS BIOMEDICAL DATA AND KNOWLEDGE SOURCES, PhD thesis, 2019, defended.
- [43] P.A. Bonatti, S. Decker, A. Polleres and V. Presutti, Knowledge Graphs: New Directions for Knowledge Representation on the Semantic Web (Dagstuhl Seminar 18371), P.A. Bonatti, S. Decker, A. Polleres and V. Presutti, eds, Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, Dagstuhl, Germany, 2019, pp. 29–111. ISSN 2192-5283. doi:10.4230/DagRep.8.9.29. <http://drops.dagstuhl.de/opus/volltexte/2019/10328>.
- [44] M.D. Wilkinson, M. Dumontier, I.J. Aalbersberg, G. Appleton, M. Axton, A. Baak, N. Blomberg, J.-W. Boiten, L.B. da Silva Santos, P.E. Bourne et al., The FAIR Guiding Principles for scientific data management and stewardship, *Scientific data* 3 (2016).
- [45] A. Polleres, M.R. Kamdar, J.D. Fernández, T. Tudorache and M.A. Musen, A More Decentralized Vision for Linked Data, in: *Decentralizing the Semantic Web (Workshop of ISWC2018)*, CEUR Workshop Proceedings, Vol. 2165, CEUR-WS.org, 2018, An extended technical report of this paper is available at <http://epub.wu.ac.at/6371/>. <http://ceur-ws.org/Vol-2165/paper1.pdf>.
- [46] A. Polleres, M.R. Kamdar, J.D. Fernandez, T. Tudorache and M.A. Musen, A More Decentralized Vision for Linked Data, Technical Report, 02/2018, Working Papers on Information Systems, Information Business and Operations, 2018, Available at <http://epub.wu.ac.at/6371/>. [epub.wu.ac.at/6371/](http://epub.wu.ac.at/6371/).