



Provided by the author(s) and University of Galway in accordance with publisher policies. Please cite the published version when available.

Title	SPORTAL: Profiling the Content of Public SPARQL Endpoints
Author(s)	Hasnain, Ali; Mehmood, Qaiser; Sana e Zainab, Syeda; Hogan, Aidan
Publication Date	2016
Publication Information	Hasnain, A., Mehmood, Q., Sana e Zainab, S., & Hogan, A. (2016). SPORTAL: Profiling the Content of Public SPARQL Endpoints. <i>International Journal on Semantic Web and Information Systems (IJSWIS)</i> , 12(3), 134-163. doi:10.4018/IJSWIS.2016070105
Publisher	IGI Global
Link to publisher's version	<a href="http://dx.doi.org/10.4018/IJSWIS.2016070105">http://dx.doi.org/10.4018/IJSWIS.2016070105</a>
Item record	<a href="http://hdl.handle.net/10379/6476">http://hdl.handle.net/10379/6476</a>
DOI	<a href="http://dx.doi.org/10.4018/IJSWIS.2016070105">http://dx.doi.org/10.4018/IJSWIS.2016070105</a>

Downloaded 2024-04-23T20:27:45Z

Some rights reserved. For more information, please see the item record link above.



# SPORTAL: Profiling the Content of Public SPARQL Endpoints

Ali Hasnain, INSIGHT Centre for Data Analytics, National University of Ireland, Galway, Ireland

Qaiser Mehmood, INSIGHT Centre for Data Analytics, National University of Ireland, Galway, Ireland

Syeda Sana e Zainab, INSIGHT Centre for Data Analytics, National University of Ireland, Galway, Ireland

Aidan Hogan, Center for Semantic Web Research, Department of Computer Science, University of Chile, Santiago, Chile

## ABSTRACT

Access to hundreds of knowledge bases has been made available on the Web through public SPARQL endpoints. Unfortunately, few endpoints publish descriptions of their content (e.g., using VoID). It is thus unclear how agents can learn about the content of a given SPARQL endpoint or, relatedly, find SPARQL endpoints with content relevant to their needs. In this paper, the authors investigate the feasibility of a system that gathers information about public SPARQL endpoints by querying them directly about their own content. With the advent of SPARQL 1.1 and features such as aggregates, it is now possible to specify queries whose results would form a detailed profile of the content of the endpoint, comparable with a large subset of VoID. In theory it would thus be feasible to build a rich centralised catalogue describing the content indexed by individual endpoints by issuing them SPARQL (1.1) queries; this catalogue could then be searched and queried by agents looking for endpoints with content they are interested in. In practice, however, the coverage of the catalogue is bounded by the limitations of public endpoints themselves: some may not support SPARQL 1.1, some may return partial responses, some may throw exceptions for expensive aggregate queries, etc. The authors' goal in this paper is thus twofold: (i) using VoID as a bar, to empirically investigate the extent to which public endpoints can describe their own content, and (ii) to build and analyse the capabilities of a best-effort online catalogue of current endpoints based on the (partial) results collected.

## KEYWORDS

Catalogue, Self-Descriptive Queries, SPARQL, SPORTAL (SPARQL Portal)

## 1. INTRODUCTION

Linked Data aims at making data available on the Web in an interoperable format so that agents can discover, access, combine and consume content from different sources with higher levels of automation than would otherwise be possible (Heath & Bizer, 2011). The envisaged result is a “Web of Data”: a Web of structured data with rich semantic links where agents can query in a unified manner -across sources- using standard languages and protocols. Over the past few years, hundreds of knowledge bases with billions of facts have been published according to the Semantic Web standards (using RDF as a data model and RDFS and OWL to provide explicit semantics) following the Linked Data principles.

As a convenience for consumer agents, Linked Data publishers often provide a SPARQL endpoint for querying their local content (Jentzsch, Cyganiak, & Bizer, 2011). SPARQL is a declarative query language for RDF in which graph pattern matching, disjunctive unions, optional clauses, dataset construction, solution modifiers, etc., can be used to query RDF knowledge bases; the recent

SPARQL 1.1 release adds features such as aggregates, property paths, sub-queries, federation, and so on (Harris, Seaborne, & Prud'hommeaux, 2013). Hundreds of public endpoints have been published in the past few years for knowledge bases of various sizes and topics (Buil-Aranda, Hogan, Umbrich, & Vandenbussche, 2013; Jentzsch et al., 2011). Using these endpoints, clients can receive direct answers to complex queries using a single request to the server.

However, it is still unclear how clients (be they human users or software agents) should find endpoints relevant for their needs in the first place (Buil-Aranda et al., 2013; Paulheim & Hertling, 2013). A client may have a variety of needs when looking for an endpoint, where they may, for example, seek endpoints with data:

1. About a given resource, e.g., *MICHAEL JACKSON*;
2. About instances of a particular type of class, e.g., *PROTEINS*;
3. About a certain type of relationship between resources, e.g., *DIRECTS-MOVIE*;
4. About certain types of values associated with resources, e.g., *RATING*;
5. About resources within a given context or with specific values, for example, *CRIMES WITH LOCATION U.K. IN YEAR 1967* or *RAT GENES AND DISEASE STRAINS*;
6. A combination of one or more of the above.

Likewise, a client may vary in how they are best able to specify these needs: some clients may only have keywords; others may know the specific IRI(s) of the resource, class or property they are interested in; some may be able to specify concrete queries or sub-queries that they wish to answer.

We argue that a service offering agents the ability to find relevant public endpoints on the Web would serve as an important part of the SPARQL querying infrastructure, enabling ad-hoc discovery of datasets over the Web. However, realising such a service over the current SPARQL infrastructure on the Web is challenging. Looking at the literature (in particular, works on the related problem of federated querying (Acosta, Vidal, Lampo, Castillo, & Ruckhaus, 2011; Harth et al., 2010; Hasnain et al., 2014, 2016; Quilitz & Leser, 2008; Schwarte, Haase, Hose, Schenkel, & Schmidt, 2011)), we can find two high-level approaches that have been investigated thus far:

- **Runtime Queries:** The first option is to take an agent's request and query the endpoints directly at runtime to determine if they have relevant metadata or not (Schwarte et al., 2011). For example, if the agent were interested in instances of `mo:MusicalWork`,<sup>1</sup> one could issue a list of endpoints the following query:

```
Ask WHERE { ?s a mo:MusicalWork }
```

Any endpoint returning true for this query would contain information relevant to the original agent. Likewise, more complex queries could be used depending on the user's need. For example, if a user were interested in endpoints with more than 100 such instances, the service could issue:

```
SELECT (COUNT(DISTINCT ?s) AS ?c)  
WHERE { ?s a mo:MusicalWork }
```

Any endpoint returning a result greater than 100 would be relevant.

- **Published Content Descriptions:** The second option is to rely on a static description of the content of each endpoint (Acosta et al., 2011; Harth et al., 2010; Quilitz & Leser, 2008; Schwarte et al., 2011). These works either assume that a description is available in a popular format, such as the Vocabulary of Interlinked Datasets (Alexander, Cyganiak, Hausenblas, & Zhao, 2009), or a custom format (Acosta et al., 2011; Harth et al., 2010; Quilitz & Leser, 2008). For example, the VoID vocabulary allows for defining class partitions that not only state which classes are in a dataset, but how many instances it has, which properties appear, and so forth (Alexander et al., 2009). These descriptions can then be used directly to find endpoints with relevant content.

However, these approaches are themselves problematic. With respect to the first approach, each user request would require a query to be sent to several hundred public endpoints, which would incur very slow response times and could flood public services with too many requests. Likewise, users would need to know the IRIs of the resources, classes and/or properties they are interested in where supporting keyword search would be cumbersome to support: (i) although SPARQL does support functions such as REGEX that could be used to find relevant terms in literals, these functions are often executed as post-filtering operations, incurring unpredictable performance; (ii) although many SPARQL engines do build and maintain inverted indexes for efficient full-text search with keywords, this support is non-standard, different engines support full-text search in different manners, and determining the engine powering a SPARQL endpoint is non-trivial (Buil-Aranda et al., 2013).

With respect to the second approach, Buil-Aranda et al. (Buil-Aranda et al., 2013) previously observed that only one third of public SPARQL endpoints give static descriptions of their content in a standard location using suitable vocabularies such as VoID, and even where they are provided, it is unclear what level of detail these descriptions contain or indeed how accurate or up-to-date these descriptions are. Over the past several years, at least 159 distinct websites have begun hosting SPARQL endpoints (Buil-Aranda et al., 2013). Putting the burden on publishers to provide static descriptions of their endpoints' content or to otherwise change how they host data would incur a prohibitive technical and social cost.

For these reasons, in this paper we propose and explore the feasibility of a third approach:

- **Computing Content Descriptions:** Rather than relying on publishers to compute and keep content descriptions up to date, we propose to compute such descriptions directly from the endpoints themselves. In particular, we propose to design a set of queries that can be issued to endpoints to learn about their content, where the results of these queries can then be used to build a catalogue that enables clients to find endpoints with relevant content.

This approach offers a number of useful trade-offs when compared with the previous two approaches discussed earlier.

Comparing the use of computed content descriptions with runtime queries, in the former case, the client will query a centrally indexed catalogue, which incurs a lower cost, both for the client in terms of response time, and also for the remote SPARQL infrastructure in terms of the number of requests generated. However, the client will be restricted to finding endpoints using the metadata collected in the catalogue. For this reason, it is important for the catalogue to capture general descriptions of content.

When compared with using published content descriptions, we do not need to assume that such descriptions are provided by the publishers of SPARQL endpoints independently of the endpoint itself. Also, by computing the content descriptions, we ensure that the endpoint is still available (since it needs to answer the queries we send it), that the description is at least as recent as the last time the descriptions were computed, and that the statistics have a simple SPARQL query that acts as provenance (rather than using descriptions provided by the publishers themselves that could be produced by tools with, e.g., different interpretations of statistics or that may include manual approximations). However, it is not clear if public endpoints would be able to support the type of

complex SPARQL query required to compute detailed content descriptions, and indeed certain types of descriptors (e.g., licence) may not be automatically computed from the endpoint but rather require the perspective of the publisher.

In this paper, we explore the feasibility of computing content descriptions directly from SPARQL endpoints. More concretely, we propose SPOTAL (SPARQL PORTAL): a centralised catalogue indexing content descriptions computed from individual SPARQL endpoints. The goal of SPOTAL is to help both human and software agents find public SPARQL endpoints relevant for their needs. The system makes minimal assumptions about how data are hosted: SPOTAL relies only on SPARQL queries to gather information about the content of each endpoint and hence only assumes a working SPARQL interface rather than requiring the publishers hosting endpoints to provide additional descriptions of the datasets. Rather than send a query to each public endpoint at runtime, we issue each endpoint queries offline to gather metadata about its content, which are later used to find relevant endpoints. Taking a simple example, instead of querying each endpoint every time an agent is looking for a given class, we can occasionally query each endpoint (on a fortnightly basis) for an up-to-date list of their classes and use that list to find relevant endpoints for the agent at runtime.

One of the main design questions for SPOTAL then is: what content descriptions should such a system try to compute from endpoints? Ideally the content descriptions should be as general as possible, supporting a variety of different types of clients and searches. With respect to the information collected, SPARQL is a powerful query language that can be used to learn about the underlying knowledge base of the endpoint. With the advent of novel features in SPARQL 1.1 like aggregates, it is now possible to formulate queries that ask, e.g., how many triples the knowledge base contains, which classes or properties are used, how many unique instances of each class appears, which properties are used most frequently with instances of which classes, and so on. In this sense, we argue that – at least in theory – SPARQL endpoints can be considered *self-descriptive*: they can describe their own content.

On the other hand, SPOTAL is limited in what it can collect by practical thresholds on the amount of data that a SPARQL endpoint will return. Buil-Aranda et al. (Buil-Aranda et al., 2013) found that many endpoints return a maximum of 10,000 results: given that many endpoints contain millions of resources and text literals, this rules out, for example, building a complete inverted index over the content of an individual endpoint, or indexing all resources that an endpoint mentions. In any case, the goal of SPOTAL is to compute concise content descriptions rather than mirroring remote endpoint content (which would be prohibitively costly for both SPOTAL and the remote endpoints, particularly to keep up-to-date). Thus, we focus on computing concise, schema-level descriptions of endpoints. Using such descriptions, we can directly find relevant endpoints given queries of type 2, 3, 4 mentioned earlier, and can indirectly help with other forms of queries (e.g., to find endpoints that contain instances of GENE, though they may not necessary be from a rat). In particular, we focus on computing extended Vocabulary of Interlinked Datasets (VoID) descriptions from endpoints: VoID has become the de-facto standard for describing datasets in RDF (Alexander et al., 2009), and is also used in federated scenarios to find relevant endpoints (Acosta et al., 2011; Akar, Halaç, Ekinci, & Dikenelli, 2012; Basca & Bernstein, 2014; Hasnain et al., 2014, 2016; Quilitz & Leser, 2008; Schwarte et al., 2011).

SPOTAL is further limited by the inability of some endpoints to return answers to complex queries. Buil-Aranda et al. (Buil-Aranda et al., 2013) previously reported that endpoints may exhibit performance and reliability issues, may return partial results, etc. Some endpoints may not support SPARQL 1.1, some may be hosted on underpowered machines, others may index very large and/or diverse datasets over which complex aggregates cannot be successfully executed, and so forth. This again creates a practical limit with respect to how detailed a content description SPOTAL can generate for certain endpoints. For example, in later results we will show that while 93.8% of operational public endpoints respond successfully when asked for a list of all classes in their dataset, only 40.2% respond successfully when additionally asked how many instances those classes have. Thus, the SPOTAL catalogue would include metadata about the classes that appear in 93.8% of

the catalogued endpoints, but only in 40.2% cases would the catalogue have information about how many instances appear in those classes.

Rather than limiting ourselves to building uniform descriptions of each endpoint based on information that can be computed from, say, >90% of endpoints, SPORTAL also considers more complex queries in its scope: while most endpoints cannot return responses to such queries, as we will show, a non-trivial percentage of endpoints do respond. In the interest of collecting as much data as possible from these latter endpoints, we include these more complex queries. Likewise we would hope that as SPARQL implementations mature, the percentage of endpoints responding to more complex queries may grow over time. As a result, the descriptive metadata available for an individual endpoint may differ from others depending on its ability to answer increasingly complex queries over its dataset. A core contribution of this paper is thus to evaluate the ability of public SPARQL endpoints to answer increasingly complex self-descriptive queries, which reflects the coverage of high-level metadata available to the SPORTAL catalogue (and similar agents) using only the SPARQL interface.

More specifically, our working hypothesis in this paper is that – despite problems with endpoint reliability and performance – by computing content descriptions using self-descriptive queries issued directly to endpoints, we can create a catalogue with (i) broader coverage and (ii) more up-to-date information than existing catalogues of SPARQL endpoints that rely on currently-available content descriptions provided by the publishers themselves. Towards investigating the validity of this hypothesis, this paper is structured as follows:

- We begin in *Section 2* with some background on related areas: Linked Data access methods, proposals for describing RDF datasets, proposals for schemes to help find relevant SPARQL endpoints, as well as discussion on how the problem could be viewed from the perspective of Linked Data as a Distributed System.
- In *Section 3*, we look at how the content of endpoints can be described in a general-purpose, automated manner. In order to extract a *description* of the content of each endpoint, we propose to use a set of 29 self-descriptive SPARQL (1.1) queries that capture a large “computable” subset of a VoID description (Alexander et al., 2009) as well as some additional features.<sup>2</sup>
- In *Section 4*, we first investigate, in a controlled environment, how well current SPARQL 1.1 engines are able to process these self-describing queries, some of which involve aggregation across an entire dataset and thus may require a prohibitive amount of processing, especially for large datasets. We run experiments over four datasets of increasing size and complexity using four SPARQL engines – 4store (Harris, Lamb, & Shadbolt, 2009), Jena/Fuseki<sup>3</sup>, Sesame (Broekstra, Kampman, & van Harmelen, 2002) and Virtuoso (Erling & Mikhailov, 2009) – to see how well our self-descriptive queries perform. These engines are mostly commonly used to power public SPARQL endpoints (Buil-Aranda et al., 2013).
- With an idea of how the queries perform in a local environment for a variety of datasets and engines, in *Section 5*, we investigate how effectively public SPARQL endpoints process these queries. We take a list of 526 public endpoints and investigate the ratio that can answer each of the self-descriptive queries and characterise the typical performance we can expect in a realistic, uncontrolled environment. Our results show that, depending on the query, the ratio of operational endpoints<sup>4</sup> returning non-empty (but possibly partial) responses vary from 25 – 94%.
- In *Section 6*, we introduce the SPORTAL catalogue based on the results collected from the remote endpoints. We describe the manner in which it can help both human and software agents to find public endpoints on the Web that may be relevant for their needs. Based on the results of the previous questions, we discussed the (in)completeness of the catalogue and both the capabilities and limitations of the system. We also provide a high-level comparison of the SPORTAL catalogue with two catalogues based on publisher-provided content descriptions: VoID Store and DataHub.
- We conclude in *Section 7* by recapitulating the main results of the paper and lessons learnt with respect to the goal of building a central catalogue of public SPARQL endpoints.

## 2. BACKGROUND

Before we continue to the core of the paper, we provide some brief background on (1) methods for accessing Linked Data, (2) the problem of peer discovery in the area of Distributed Systems, (3) works on finding relevant SPARQL endpoints, and (4) techniques for describing/ summarising RDF datasets.

### 2.1. Linked Data Access Methods

Traditionally there have been three methods provided for consumer agents to access content from knowledge bases published as Linked Data: *dereferencing*, where IRIs of interest are looked up via HTTP; *dumps*, where the entire content of a dataset is made available for download; and *SPARQL endpoints*, where a query interface is provided over the local content. A more recent proposal – *Linked Data Fragments* (Verborgh et al., 2014) – has recently begun to gain attention.

Both dereferencing and dumps are lightweight methods in-tune with current practices on the Web; however, they can be inefficient for agents to use. Consider an agent wishing to retrieve the populations of Asian capitals from DBpedia. An agent has no direct way of finding the correct IRIs to dereference; even if they did, DBpedia specifies a Crawl-delay of 10 seconds: assuming that the *DBpedia* IRIs of 49 Asian capitals needed dereferencing, a polite agent would require 8 minutes to retrieve the respective documents and would ultimately use one triple out of potentially hundreds of thousands in each document. Using a dump would entail downloading an entire dataset to get at 49 triples; hosting a local dump mirror would require constant refreshing.

Hence publishers provide SPARQL endpoints as a convenient alternative to dereferencing or dumps. To get the populations of Asian capitals, an agent could run the following query against the DBpedia SPARQL endpoint<sup>5</sup>:

```
SELECT ?pop ?city WHERE
{ ?city dct:subject dbc:Capitals_in_Asia ;
  dbo:populationTotal ?pop . }
```

All going well, the query will return populations in less than a second. Likewise, only the data that the client is interested in will be transferred. However, SPARQL endpoints push the burden from data consumers to producers: hosting such a public query service is expensive and as a result, endpoints may not be able to answer all queries for all consumer agents (Buil-Aranda et al., 2013). Despite problems with reliability, SPARQL endpoints still offer an appealing method for consumer agents to interact with remote Linked Data knowledge bases where endpoints such as DBpedia serve millions of queries for clients (Gallego, Fernández, Martínez-Prieto, & de la Fuente, 2011).

As an alternative to SPARQL endpoints, Verborgh et al. (Verborgh et al., 2014), propose methods for providing and organising multiple access methods to a Linked Dataset, including a lightweight “triple pattern fragment”, which allows clients to request all triples matching a single pattern, the goal of which is to allow publishers to host highly reliable but greatly simplified query services, thus trying to strike a better balance between the costs on the client and server side. Although their Linked Data Fragments (LDF) proposal offers a valuable compromise between client and server costs, being a recent proposal, SPARQL endpoints still greatly outnumber the number of LDF servers on the Web.

### 2.2. SPARQL Endpoints as a Distributed System

Viewed from the perspective of Distributed Computing, each SPARQL endpoint on the Web involves a client- server architecture, where numerous clients use the SPARQL protocol to interface with a single external server.<sup>6</sup> However, when hundreds of public SPARQL endpoints are viewed collectively, they can be seen as forming a decentralised peer-to-peer (P2P) system. In particular, with the advent of SPARQL 1.1 Federation (Prud’hommeaux & Buil-Aranda, 2013), endpoints can query each other and thus may perform computation on behalf of other peers.

In this light, the goal of finding relevant SPARQL endpoints relates to the core problem of peer discovery in the P2P area, wherein a peer wishes to find another peer with a particular piece of data. To make this task more efficient, *structured P2P systems* impose an overall organisation on the network overlay to ensure rapid peer discovery; the most common structure is a Distributed Hash Table (DHT), which is effectively a distributed map where keys are hashed to determine on which peer(s) a given set of key – value pairs should be stored (Ratnasamy, Francis, Handley, Karp, & Shenker, 2001; Rowstron & Druschel, 2001; Stoica et al., 2003; Zhao et al., 2004). However, all such structured schemes assume that peers in the network can be assigned data, which is not true of SPARQL endpoints where peers themselves decide which datasets they wish to index.

As such, public SPARQL endpoints collectively form an *unstructured P2P system*, where, since there is no correlation imposed between a peer and the data it indexes, peer discovery would necessarily involve one of two options: a separate search index that records the content at each peer (e.g., trackers in BitTorrent (Qiu & Srikant, 2004), or blindly flooding the network with queries looking for the desired data from peers in a “brute force manner” (e.g., Gnutella (Ripeanu, Iamnitchi, & Foster, 2002)).

Rather than requiring a complete global structure or accepting zero structure, other proposals aim to strike a balance by imposing a limited form of structure over nodes. For example, routing indices (Crespo & Garcia-Molina, 2002) allow nodes to index whatever data they wish, but require that each peer must additionally store pointers to a neighbouring peer that is closer to the desired data; this avoids blind flooding of queries during peer discovery, instead allowing peers to be *routed* to relevant peer(s). Likewise, routing indexes avoid the need for a central index of peer content.

However, our goal in this paper is to enable peer discovery of SPARQL endpoints without changing the current infrastructure; we feel that it is important to explore options over the current infrastructure first before proposing that hundreds of stakeholders change how they host their data. For example, we do not presume that publishers will agree to add and maintain routing indexes towards the endpoints of external publishers. Hence we assume that no structure is imposed on the peers, but rather that each SPARQL endpoint indexes its own data. Thus the scenario is effectively unstructured: we have no guarantees about which data may appear at which endpoint/peer. Our hypothesis instead is that we can use the SPARQL query interface to learn about the content at each peer.

### 2.3. Describing/Summarising RDF Datasets

With respect to building a central search service for endpoints based on their content, it would seem infeasible to index all of the data from the endpoint; hence some form of summary or schema overview must be indexed. A variety of works have proposed methods to describe and/or summarise RDF datasets.

In terms of describing metadata about RDF datasets, Cyganiak et al. (Cyganiak, Stenzhorn, Delbru, Decker, & Tummarello, 2008) propose Semantic Sitemaps to mark the locations of different Linked Data access points; however, information captured is limited to broad concepts such as change frequency. Alexander et al. (Alexander et al., 2009) later proposed VoID for describing RDF datasets and the links between them. As we will see, the vocabulary provides terms for describing high-level statistics about a dataset, as well as about the instances of specific classes and the usage of specific properties. A number of works have proposed extensions to the VoID vocabulary. Mountantonakis et al. (Mountantonakis et al., 2014) propose to extend VoID with metrics about the connectivity of pairs of data sources to capture, for example, the number of common RDF terms used in both sources, the increase in average node degree with both sources are combined, etc. Omitola et al. (Omitola et al., 2011) propose to extend VoID to allow publishers to describe in more depth the provenance of their dataset.

With respect to computing dataset descriptions, or profiling datasets, Bohm et al. (Böhm, Lorey, & Naumann, 2011) demonstrated that computing a VoID description for large datasets is feasible using MapReduce techniques. As part of the LOD Laundromat service – which aims to clean up and republish existing datasets in a more uniform manner – Beek et al. (Beek, Rietveld,

Bazoobandi, Wielemaker, & Schlobach, 2014) compute a VoID description for each dataset indexed. More recently, Fetahu et al. (Fetahu et al., 2014) propose extracting topics from a dataset based on a combination of information retrieval techniques such as PageRank and HITS and Named Entity Recognition applied offline over the dataset. Abejan et al. (Abedjan, Grütze, Jentzsch, & Naumann, 2014) propose ProLOD++: a system to profile Linked Datasets that applies clustering techniques, statistical analysis, and association rules to find semantically related groups of entities, statistical distributions, properties that together uniquely identify resources, as well as suggested changes to the dataset/ontology. Mihindikulasooriya et al. (Mihindikulasooriya, Poveda-Villalón, García-Castro, & Gómez-Pérez, 2015) propose Loupe: a system that extracts a schema-level summary of a dataset similar to that captured by VoID (e.g., number of triples, number of instances per class, etc.), with additional information about namespaces, ontological definitions, etc.

Closer to our own contribution, various works have proposed using SPARQL to extract high-level information about an RDF dataset. Auer et al. (Auer, Demter, Martin, & Lehmann, 2012) propose LODstats, which applies analytics over a stream of RDF data but which uses SPARQL filters to (reject)/select (ir)relevant triples; use of SPARQL is limited to filters. Langegger & Wöß propose RDFStats (Langegger & Wöß, 2009), which uses a pipeline of SPARQL (1.0) queries to generate a histogram on a per-class basis, representing the predicates and types of values associated with its instances. Holst & Höfig (Holst & Höfig, 2013) propose the use of SPARQL 1.1 queries to discover specific aspects of an RDF dataset, but the authors do not consider VoID and only run local experiments over three datasets. Mountantonakis et al. (Mountantonakis et al., 2014) propose a set of SPARQL 1.1 queries that can compute the connectivity metrics with which they extend VoID. Mäkelä (Mäkelä, 2014) propose Aether: a system for extracting extended VoID descriptions from a SPARQL 1.1 compliant endpoint; this work is perhaps most similar in spirit to ours, however, the focus is more on getting an overview of a known endpoint rather than building a catalogue that can be used by clients to find endpoints of interest.

The SPARQL 1.1 Service Description (SD) (Williams, 2013), vocabulary was recently recommended by the W3C; however, unlike previously discussed works, which focus on describing the content of datasets, SD describes technical aspects of an endpoint, such as features supported, dataset configurations, etc.

Other works have focused on summarising the content of RDF datasets (rather than describing them using a high-level RDF description). Umbrich et al. (Umbrich, Hose, Karnstedt, Harth, & Polleres, 2011) propose to use an approximate, hash-based indexing structure, called a QTree, to aid in source selection; the QTree allows for determining which sources are likely to contain matches for a given RDF triple pattern but at a fraction of the size of the original dataset. Khatchadourian & Consens (Khatchadourian & Consens, 2010) propose creating bisimulation labels that capture connectivity in an RDF graph on the level of the namespaces of the instance URIs and the schema used. Campinas et al. (Campinas, Delbru, & Tummarello, 2013) propose using existing graph summary algorithms to summarise RDF graphs, where nodes that are equivalent per some relation – e.g., having the same types, or the same attributes – are collapsed into a single node to create a smaller summary graph.

## 2.4. Discovering SPARQL Endpoints

As previously discussed in the introduction, there are two high-level options for discovering SPARQL endpoints with relevant data: (1) flood the endpoints with queries, or (2) build a central search index. For example, federated SPARQL engines employ one or both of these strategies (Acosta et al., 2011; Akar et al., 2012; Basca & Bernstein, 2014; Quilitz & Leser, 2008; Schwarte et al., 2011). Our goal is to build a central catalogue based on data collected from endpoints through their SPARQL interfaces.

Paulheim & Hertling (Paulheim & Hertling, 2013) looked at how to find a SPARQL endpoint containing content about a given Linked Data URI: using VoID descriptions and the DataHub catalogue, the authors could find suitable endpoints for about 15% of the sample of ten thousand URIs considered. Mehdi et al. (Mehdi et al., 2014) looked at the problem of discovering endpoints

that may be relevant to a set of domain-specific keywords: their approach involved generating a list of RDF literals from the keywords and flooding queries against endpoints to see if they contained, e.g., case or language-tag variations of the literals.

Buil-Aranda et al. (Buil-Aranda et al., 2013) propose SPARQLES as a catalogue of SPARQL endpoints, but focus on performance and stability metrics rather than cataloguing content; they do however remark that they could only find static descriptions for the content of about one third of the public endpoints surveyed, making endpoint discovery difficult. Likewise, the analysis by Lorey (Lorey, 2014) of public endpoints focused on characterising the performance offered by these services rather than on the problem of discovery.

There are a variety of locations online where lists of public endpoints can be found and searched over. For example, DataHub<sup>7</sup> provides a list of hundreds of Linked Datasets, which can be filtered to find those that offer SPARQL endpoint locations. One can, for example, search for datasets relating to uk crime and filter to only show those with SPARQL endpoints. However, the search functionality provided is limited in most cases to keyword search over the dataset title, or to browsing datasets with a given tag, etc. Still, the service often provides links to VoID files that could be used to catalogue the content of endpoints. Unlike SPORAL however, these VoID files are provided by publishers rather than being computed from the endpoints. Hence we will later compare our catalogue with that formed by collecting the VoID files that DataHub links to for each dataset.

As part of the RKBexplorer infrastructure (Glaser, Millard, & Jaffri, 2008), the VoID Store allows for performing searches over VoID files submitted to the system.<sup>8</sup> A service is also provided to find endpoints that index content about a given resource (using the REGEX patterns sometimes provided in VoID). This catalogue could thus be used by clients to find relevant SPARQL endpoints. Currently the store contains information related to 118 endpoints. Like DataHub – but unlike SPORAL – the VoID files indexed by VoID Store are again computed and uploaded by publishers.

## 2.5. Novelty

We focus on the problem of helping clients find relevant SPARQL endpoints. To the best of our knowledge, there are two online services that clients could use to try to find SPARQL endpoints based on their content: DataHub and/or VoID Store. However, both of these services rely on static content descriptions provided by publishers themselves. As noted by Buil-Aranda et al. (Buil-Aranda et al., 2013), many of the endpoints listed in the DataHub have been offline for years; also, of the endpoints surveyed, VoID descriptions are only available in the DataHub for 33.3% and in the VoID Store for 22.4%. We instead propose to compute extended VoID descriptions for public endpoints directly through their SPARQL interface. We provide a high level comparison between SPORAL, DataHub and VoID Store in *Section 6*.

## 3. SELF-DESCRIPTIVE QUERIES

With respect to describing the content of an endpoint, in this section, we list the set of SPARQL 1.1 queries that we use to compute a VoID-like description from the content indexed by an endpoint.

### 3.1. Functionality

First we wish to filter unavailable endpoints and to determine those that (partially) support SPARQL 1.1. We consider an endpoint available if it is accessible through the HTTP SPARQL protocol, it responds to a SPARQL-compliant query, and it returns a response in an appropriate SPARQL format; for this, we use query  $Q_{A1}$  (see *Table 1*), which should be trivial for an endpoint to compute, returning a single binding for any triple. We consider an endpoint SPARQL 1.1 aware if it likewise responds to a query valid only in SPARQL 1.1; for this, we use query  $Q_{A2}$  (see *Table 1*), which tests two features unique to SPARQL 1.1: sub-queries and the count aggregate function.<sup>9</sup>

### 3.2. Dataset-level Statistics

Second we list a set of queries to capture high-level “dataset-level” statistics that form a core part of VoID. We issue five queries, as listed in *Table 2*, to ascertain the number of triples ( $Q_{B1}$ ), and the number of distinct classes ( $Q_{B2}$ ), properties ( $Q_{B3}$ ), subjects ( $Q_{B4}$ ), and objects ( $Q_{B5}$ ). These queries require support for SPARQL 1.1 COUNT and sub-query features (as tested in  $Q_{A2}$ ). The  $\langle D \rangle$  term refers to an IRI constructed from the SPARQL endpoint’s URL to indicate the dataset it indexes.

Once these statistics are catalogued for public SPARQL endpoints, agents can use them to find endpoints indexing datasets that fall within a given range of triples in terms of overall size, or, for example, to find the endpoints with the largest datasets. Counts may be particularly useful – in combination with later categories – to order the endpoints; for example, to find the endpoints with a given class (using data from the next category) and order them by the total number of triples they index.

### 3.3. Class-based Statistics

Third we ascertain similar statistics about the instances of each class following the notion of *class partitions* in VoID: a subset of the data considering only triples where instances of that class are in the subject position. *Table 3* lists the six queries we use. The first query ( $Q_{C1}$ ) merely lists all class partitions. The other five queries ( $Q_{C2-6}$ ) count the triples and distinct classes, predicates, subjects and objects for each class partition; e.g.,  $Q_{C2}$  retrieves the number of triples where instances of that class are in the subject position. Queries  $Q_{C2-6}$  introduce COUNT, sub-queries and also GROUPBY features from SPARQL 1.1.

Once catalogued, agents can use statistics describing class partitions of the datasets to find endpoints mentioning a given class, where they can additionally (for example) sort results in descending order according to the number of unique instances of that class, or triples used to define such instances, and so forth. Hence the counts computed by ( $Q_{C2-6}$ ) help agents to distinguish endpoints that may only have one or two instances of a class to those with thousands or millions. Likewise criteria can be combined arbitrarily for multiple classes, or with the overall statistics computed previously.

Table 1. Queries for basic functionalities

No	Query
$Q_{A1}$	SELECT * WHERE { ?s ?p ?o } LIMIT 1
$Q_{A2}$	SELECT (COUNT(*) as ?c) WHERE { SELECT * WHERE { ?s ?p ?o } LIMIT 1 }

Table 2. Queries for dataset-level VoID statistics

No	Query
$Q_{B1}$	CONSTRUCT { <D> v:triples ?x } WHERE { SELECT (COUNT(*) AS ?x) WHERE { ?s ?p ?o } }
$Q_{B2}$	CONSTRUCT { <D> v:classes ?x } WHERE { SELECT (COUNT(DISTINCT ?o) AS ?x) WHERE { ?s a ?o } }
$Q_{B3}$	CONSTRUCT { <D> v:properties ?x } WHERE { SELECT (COUNT(DISTINCT ?p) AS ?x) WHERE { ?s ?p ?o } }
$Q_{B4}$	CONSTRUCT { <D> v:distinctSubjects ?x } WHERE { SELECT (COUNT(DISTINCT ?s) AS ?x) WHERE { ?s ?p ?o } }
$Q_{B5}$	CONSTRUCT { <D> v:distinctObjects ?x } WHERE { SELECT (COUNT(DISTINCT ?o) AS ?x) WHERE { ?s ?p ?o } }

### 3.4. Property-based Statistics

Fourth we look at property partitions in the dataset, where a property partition refers to the set of triples with that property term in the predicate position. Queries are listed in *Table 4*. As before,  $Q_{D1}$  lists the property partitions. ( $Q_{D2-4}$ ), count the number of triples, distinct subjects and distinct objects. We do not count classes (which would be ‘0’ for all properties except `rdf:type`) or properties (which would always be ‘1’).

Using these statistics about property partitions in the catalogue, agents can, for example, retrieve a list of public endpoints using a given property ordered by the number of triples using that specific property. Likewise criteria can be combined arbitrarily for multiple properties, or with the dataset- or class-level metadata previously collected; for example, an agent may wish to order endpoints by the ratio of triples using a given property (where the count from  $Q_{D2}$  for the property in question can be divided by the total triple count from  $Q_{B1}$ ), or to find endpoints where all subjects have an `rdfs:label` value (where the count computed from  $Q_{D3}$  for that property should match the count for  $Q_{B4}$ ).

### 3.5. Nested Class–property Statistics

Fifth we look at how classes and properties are used together in a dataset, gathering statistics on property partitions nested within class partitions: these statistics detail how properties are used for instances of specific classes. *Table 5* lists the four queries used. [ $Q_{E1}$ ] lists the property partitions nested inside the class partitions, and [ $Q_{E2-4}$ ] count the number of triples using a given predicate for instances of that class, as well as the number of distinct subjects and objects those triples have. In

**Table 3. Queries for statistics about classes**

No	Query
$Q_{C1}$	CONSTRUCT { <D> v:classPartition [ v:class ?c ] WHERE { ?s a ?c }
$Q_{C2}$	CONSTRUCT { v:classPartition [ v:class ?c ; v:triples ?x ] WHERE { SELECT (COUNT(?p) AS ?x) ?c WHERE { ?s a ?c ; ?p ?o } GROUP BY ?c }
$Q_{C3}$	CONSTRUCT { v:classPartition [ v:class ?c ; v:classes ?x ] WHERE { SELECT (COUNT(DISTINCT ?d) AS ?x) ?c WHERE { ?s a ?c, ?d } GROUP BY ?c }
$Q_{C4}$	CONSTRUCT { v:classPartition [ v:class ?c ; v:properties ?x ] WHERE { SELECT (COUNT(DISTINCT ?p) AS ?x) ?c WHERE { ?s a ?c ; ?p ?o } GROUP BY ?c }
$Q_{C5}$	CONSTRUCT { v:classPartition [ v:class ?c ; v:distinctSubjects ?x ] WHERE { SELECT (COUNT(DISTINCT ?s) AS ?x) ?c WHERE { ?s a ?c } GROUP BY ?c }
$Q_{C6}$	CONSTRUCT { v:classPartition [ v:class ?c ; v:distinctObjects ?x ] WHERE { SELECT (COUNT(DISTINCT ?o) AS ?x) ?c WHERE { ?s a ?c ; ?p ?o } GROUP BY ?c }

**Table 4. Queries for statistics about properties**

No	Query
$Q_{D1}$	CONSTRUCT { v:propertyPartition [ v:property ?p ] WHERE { ?s ?p ?o }
$Q_{D2}$	CONSTRUCT { v:propertyPartition [ v:property ?p ; v:triples ?x ] WHERE { SELECT (COUNT(?o) AS ?x) ?p WHERE { ?s ?p ?o } GROUP BY ?p }
$Q_{D3}$	CONSTRUCT { v:propertyPartition [ v:property ?p ; v:distinctSubjects ?x ] WHERE { SELECT (COUNT(DISTINCT ?s) AS ?x) ?p WHERE { ?s ?p ?o } GROUP BY ?p }
$Q_{D4}$	CONSTRUCT { v:propertyPartition [ v:property ?p ; v:distinctObjects ?x ] WHERE { SELECT (COUNT(DISTINCT ?o) AS ?x) ?p WHERE { ?s ?p ?o } GROUP BY ?p }

Table 5. Queries for nested property/class statistics

No	Query
Q <sub>E1</sub>	CONSTRUCT { v:classPartition [ v:class ?c ; v:propertyPartition [ v:property ?p ] ] } WHERE { ?s a ?c ; ?p ?o }
Q <sub>E2</sub>	CONSTRUCT { v:classPartition [ v:class ?c v:propertyPartition [ v:property ?p ; v:triples ?x ] ] } WHERE { SELECT (COUNT(?o) AS ?x) ?p WHERE { ?s a ?c ; ?p ?o } GROUP BY ?c ?p }
Q <sub>E3</sub>	CONSTRUCT { v:classPartition [ v:class ?c ; v:propertyPartition [ v:distinctSubjects ?x ] ] } WHERE { SELECT (COUNT(DISTINCT ?s) AS ?x) ?c ?p WHERE { ?s a ?c ; ?p ?o } GROUP BY ?c ?p }
Q <sub>E4</sub>	CONSTRUCT { v:classPartition [ v:class ?c ; v:propertyPartition [ v:distinctObjects ?x ; v:property ?p ] ] } WHERE { SELECT (COUNT(DISTINCT ?o) AS ?x) ?c ?p WHERE { ?s a ?c ; ?p ?o } GROUP BY ?c ?p }

Table 6. Queries for miscellaneous statistics

No	Query
Q <sub>F1</sub>	CONSTRUCT { e:distinctIRIReferenceSubjects ?x } WHERE { SELECT (COUNT(DISTINCT ?s) AS ?x) WHERE { ?s ?p ?o FILTER(isIri(?s)) } }
Q <sub>F2</sub>	CONSTRUCT { e:distinctBlankNodeSubjects ?x } WHERE { SELECT (COUNT(DISTINCT ?s) AS ?x) WHERE { ?s ?p ?o FILTER(isBlank(?s)) } }
Q <sub>F3</sub>	CONSTRUCT { e:distinctIRIReferenceObjects ?x } WHERE { SELECT (COUNT(DISTINCT ?o) AS ?x) WHERE { ?s ?p ?o FILTER(isIri(?o)) } }
Q <sub>F4</sub>	CONSTRUCT { e:distinctLiterals ?x } WHERE { SELECT (COUNT(DISTINCT ?o) AS ?x) WHERE { ?s ?p ?o FILTER(isLiteral(?o)) } }
Q <sub>F5</sub>	CONSTRUCT { e:distinctBlankNodeObjects ?x } WHERE { SELECT (COUNT(DISTINCT ?o) AS ?x) WHERE { ?s ?p ?o FILTER(isBlank(?o)) } }
Q <sub>F6</sub>	CONSTRUCT { e:distinctBlankNodes ?x } WHERE { SELECT (COUNT(DISTINCT ?b) AS ?x) WHERE { { ?s ?p ?b } UNION { ?b ?p ?o } FILTER(isBlank(?b)) } }
Q <sub>F7</sub>	CONSTRUCT { e:distinctIRIReferences ?x } WHERE { SELECT (COUNT(DISTINCT ?u) AS ?x) WHERE { { ?u ?p ?o } UNION { ?s ?u ?o } UNION { ?s ?p ?u } FILTER(isIri(?u)) } }
Q <sub>F8</sub>	CONSTRUCT { e:distinctRDFNodes ?x } WHERE { SELECT (COUNT(DISTINCT ?n) AS ?x) WHERE { { ?n ?p ?o } UNION { ?s ?n ?o } UNION { ?s ?p ?n } } }
Q <sub>F9</sub>	CONSTRUCT { v:propertyPartition [ v:property ?p ; s:subjectTypes [ s:subjectClass ?sType ; s:distinctMembers ?x ] ] } WHERE { SELECT (COUNT(?s) AS ?x) ?p ?sType WHERE { ?s ?p ?o ; a ?sType . } GROUP BY ?p ?sType }
Q <sub>F10</sub>	CONSTRUCT { v:propertyPartition [ v:property ?p ; s:objectTypes [ s:objectClass ?oType ; s:distinctMembers ?x ] ] } WHERE { SELECT (COUNT(?o) AS ?x) ?p ?oType WHERE { ?s ?p ?o . ?o a ?oType . } GROUP BY ?p ?oType }

terms of technical features, these queries involve GROUP BY over multiple terms. In general, the queries listed in this section are quite complex where we would expect that many endpoints would struggle to return metadata about their content at this detailed level of granularity.

An agent could use the resulting metadata to find endpoints describing instances of specific classes with specific properties, with filtering or sorting criteria based on, e.g., the number of triples. For example, an agent might be specifically interested in images of people, where they would be looking for the class-partition foaf:Person with the nested property-partition foaf:depicts. Using the previous statistics, it would have been possible to find endpoints that have data for the class foaf:Person and triples with the property foaf:depicts, but not that the images were defined for people. The counts from [Q<sub>E2-E4</sub>] again allow an agent to filter or order endpoints by the amount of relevant

### 3.6. Miscellaneous Statistics

In our final set of experiments, we look at queries that yield statistics not supported by VoID as listed in *Table 6*. In particular, we experiment to see if endpoints can return a subset of statistics from the VoID Extension Vocabulary<sup>10</sup>, which include counts of different types of unique RDF terms in different positions: subjects IRIs ( $Q_{F1}$ ), subject blank nodes ( $Q_{F2}$ ), objects IRIs ( $Q_{F3}$ ), literals ( $Q_{F4}$ ), object blank nodes ( $Q_{F5}$ ), all blank nodes ( $Q_{F6}$ ), all IRIs ( $Q_{F7}$ ), and all terms ( $Q_{F8}$ ). Inspired by the notion of “schema maps” as proposed by Kinsella et al. (Kinsella, Bojars, Harth, Breslin, & Decker, 2008), we also count the classes that the subjects and objects of specific properties are instances of ( $Q_{F9-10}$ ); these are “inverses” of queries ( $Q_{E3-4}$ ).<sup>11</sup>

The resulting data could serve a number of purposes for agents looking for public endpoints. For example, the agent in question could look for datasets without any blank nodes or for datasets where a given number of the objects of a given property are of a certain type. Likewise, the user can combine these criteria with earlier criteria; for example, to find endpoints with more than ten million triples where at least 30% of the unique object terms are literals.

## 4. LOCAL EXPERIMENTS

In our first set of experiments, we test whether or not SPARQL implementations can locally answer the queries we specified in the previous section. These implementations are used to power individual endpoints and hence we would like to see if running these queries is feasible in a locally controlled environment before running remote experiments.

Along these lines, given the 29 self-descriptive queries mentioned previously, we test four popular SPARQL query engines: *Virtuoso* (07.10.3207), *Fuseki* (1.0.2), *Astore* (4s-httpd/v1.1.4) and *Sesame* (2.7.12). Given that the cost of the self-descriptive queries listed previously depends directly on the size and nature of the dataset indexed, for each engine, we perform experiments with respect to the four real-world datasets listed in *Table 7*, representing a mix of datasets at a variety of scales and with a variety of diversity in predicates and classes used. The experiments are run on a server with *Ubuntu 14*, a *4x Intel Core i5 CPU (M540@2.53 GHz)* processor and 8 GB of *RAM*. A timeout of 10 minutes was set for the first result to return. Result-size thresholds were switched off where applicable.

All of the engines passed the functionality tests. In *Table 8*, we list the runtimes for all other queries (spanning *Table 2-6*), for the four datasets and the four engines. We manually inspected the results so as to only include runtimes where the correct response was returned. With respect to the (partially) failed queries, in the table, we differentiate between:

- **Empty results** (—) where zero results are returned, most commonly caused by a 10-minute timeout;
- **Partial results** (~) where the stream of results returned is correct but ends prematurely;
- **Incorrect results** (X) where the results returned are false, most commonly caused by counts not considering all results or by query processor bugs.

**Table 7. High-level statistics for datasets used in local experiments**

	Triples	Subjects	Predicates	Objects	Classes
DRUGBANK	517,023	19,693	119	276,142	8
JAMENDO	1,049,647	335,925	26	440,686	11
KEGG	1,090,830	34,260	21	939,258	4
DBPEDIA	114,456,676	11,194,893	53,200	27,518,753	447

In the following, we draw high-level conclusions from these results.

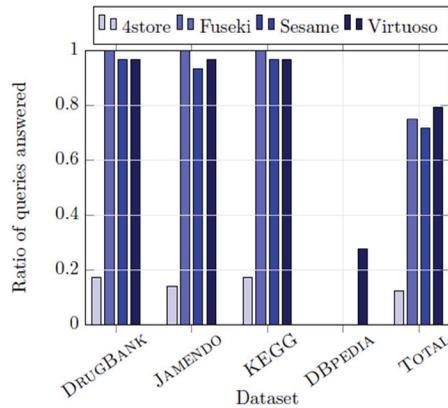
#### 4.1. Datasets

We see in Table 8 that while *Fuseki*, *Sesame* and *Virtuoso* successfully run almost all queries for *DRUGBANK*, *JAMENDO* and *KEGG* – datasets around or below a million triples – all engines struggle for the *DBpedia* dataset, which is two orders of magnitude larger and contains two orders of magnitude more classes, predicates, objects and subjects. This is better illustrated by Figure 1, where the difference between *DBPEDIA* and the other datasets is evident in terms of success rate. Only *Virtuoso* managed to return correct results for some queries over *DBPEDIA*, including counts for triples, classes, properties, triples per property partition, blank node subjects, blank node objects and blank nodes in any position<sup>12</sup> and object IRIs. We posit that with the available memory, queries

**Table 8. Local query runtimes for self-descriptive queries (times in millisecond; engines are keyed as 4store, Fuseki, Sesame, Virtuoso; ‘.’ indicates empty results, ‘~’ partial results, ‘X’ incorrect results)**

No	4	DrugBank			Jamendo			KEGG			DBpedia					
		F	S	V 4	F	S	V 4	F	S	V 4	F	S	V			
Q <sub>B1</sub>	X	5,128	6,962	95	X	4,637	6,001	127	X	2,018	3,059	130	-	-	-	6,175
Q <sub>B2</sub>	27	3,140	658	43	X	8,278	2,258	158	42	2,773	325	55	X	-	-	6,260
Q <sub>B3</sub>	29	1,664	7,029	155	34	1,840	6,478	278	35	1,016	3,140	284	-	-	-	25,130
Q <sub>B4</sub>	X	6,029	7,324	252	X	7,987	6,503	1,493	X	2,455	3,263	498	-	-	-	-
Q <sub>B5</sub>	X	17,141	8,764	1,269	X	4,714	6,824	1,736	X	7,174	3,715	2,523	-	-	-	-
Q <sub>C1</sub>	X	4,201	21,845	3,560	X	16,868	166,727	33,559	X	3,125	19,282	4,453	~	-	-	-
Q <sub>C2</sub>	X	5,342	9,305	374	X	5,336	10,086	291	X	3,097	5,609	360	-	-	-	-
Q <sub>C3</sub>	X	949	758	99	X	2,946	5,153	390	X	801	526	119	-	-	-	-
Q <sub>C4</sub>	X	2,844	9,423	949	X	3,940	10,028	1,135	X	1,990	5,933	1,110	-	-	-	-
Q <sub>C5</sub>	X	227	479	178	X	991	2,952	2,255	X	249	258	226	-	-	-	-
Q <sub>C6</sub>	X	17,612	10,393	2,979	X	5,530	10,339	3,472	X	3,483	6,802	4,301	-	-	-	-
Q <sub>D1</sub>	X	80,119	754,030	-	X	106,486	-	-	X	35,902	323,709	-	-	-	-	-
Q <sub>D2</sub>	X	13,468	8,787	109	X	7,675	740	90	X	5,022	3,917	111	-	-	-	41,188
Q <sub>D3</sub>	X	2,718	90,936	1,425	X	3,870	8,322	3,105	X	1,580	4,120	1,764	-	-	-	-
Q <sub>D4</sub>	X	15,504	10,252	1,710	X	4,983	8,355	2,298	X	2,258	4,356	3,298	-	-	-	-
Q <sub>E1</sub>	X	60,310	1,296,659	17,899	X	39,836	2,663,655	34,926	X	17,825	683,787	37,204	-	-	-	-
Q <sub>E2</sub>	X	13,644	9,410	445	X	5,086	12,886	279	X	2,947	6,469	240	-	-	-	-
Q <sub>E3</sub>	X	4,113	10,029	2,815	X	6,246	10,522	3,464	X	3,505	7,021	2,763	-	-	-	-
Q <sub>E4</sub>	X	20,489	10,603	3,415	X	6,449	10,649	3,448	X	3,359	7,756	4,596	-	-	-	-
Q <sub>F1</sub>	X	2,587	9,990	428	X	4,181	9,131	1,943	X	1,375	4,622	822	-	-	-	-
Q <sub>F2</sub>	30	3,066	8,992	52	36	2,407	8,178	42	45	1,410	4,151	54	-	-	-	378
Q <sub>F3</sub>	X	17,559	11,069	486	X	3,663	9,416	1,467	X	1,711	4,470	1,815	-	-	-	65,398
Q <sub>F4</sub>	X	15,862	9,580	1,100	X	3,061	9,775	644	X	1,761	4,592	1,160	-	-	-	-
Q <sub>F5</sub>	29	14,686	8,877	122	40	3,678	7,755	162	38	1,652	3,880	168	-	-	-	16,192
Q <sub>F6</sub>	52	17,208	17,820	130	58	4,914	15,968	173	60	2,049	7,691	204	-	-	-	16,323
Q <sub>F7</sub>	X	23,360	31,675	1,333	X	11,497	27,428	4,183	X	4,456	13,401	3,131	-	-	-	-
Q <sub>F8</sub>	X	26,517	27,651	1,788	X	10,565	23,549	6,728	X	5,122	11,642	3,615	-	-	-	-
Q <sub>F9</sub>	X	4,725	9,546	414	X	6,746	9,875	315	X	3,210	6,822	271	-	-	-	-
Q <sub>F10</sub>	X	1,272	-	127	X	5,290	-	619	X	1,228	-	149	-	-	-	-

Figure 1. Ratio of successful queries per dataset/ engine



over the smaller datasets could be processed largely in-memory whereas queries over *DBPEDIA* may have led to a lot of on-disk processing.

## 4.2. Engines

With respect to the success rate of the four implementations, from Table 8 and Figure 1, we can see that *4store* struggled the most with the self-descriptive queries specified, returning correct results only for counts of classes, properties and blank nodes. *Fuseki* was the most reliable engine for the smaller datasets, successfully answering all queries over *DRUGBANK*, *JAMENDO* and *KEGG*, whereas *Sesame* and *Virtuoso* struggled on queries  $[Q_{D1}]$  and  $[Q_{F10}]$ . From further investigation, we discovered that for  $[Q_{D1}]$  – list all property partitions – the engines were returning two output triples for every triple indexed, producing a large volume of non-lean data, as opposed to returning two output triples for every unique property. To get around this, a sub-query specifying *DISTINCT* on *?p* could be used at the cost of requiring SPARQL 1.1 support. On the other hand,  $[Q_{F10}]$  would seem on face value to be the most expensive query in our collection, requiring an open join and aggregation step that may naturally fail even for small-to-medium-sized datasets.

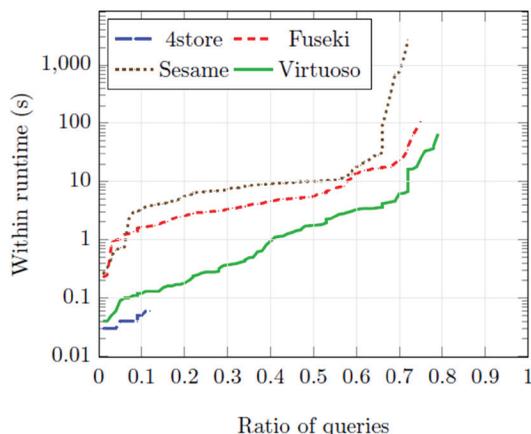
## 4.3. Runtimes

We see quite a large variance in runtimes between the different engines, varying in orders of magnitude. In order to get a better insight into the differences in performance, in Figure 2, we plot the ratio of all 116 queries (29 queries  $\times$  4 datasets) that ran below a certain runtime, where, for example, we can see that *Virtuoso* successfully ran 40% of the queries in less than one second and 72% of the queries in less than ten seconds. The plots end where queries began to fail. Interestingly, although *4store* was the most unreliable engine, it offered the fastest runtimes for the simpler queries it did answer, suggesting some index may have been used for optimisation purposes. Although *Sesame* and *Fuseki* were faster for certain queries, the trend in Figure 2 suggests that overall; *Virtuoso* was fastest for most queries. We also see that many queries continued to stream results well in excess of the one-minute connection timeout, with *Sesame* having some of the slowest successful query executions (the slowest being 44 minutes).

## 4.4. Errors

Although all engines returned empty results, *4store* was the only engine that was found to return partial or incorrect answers where, for count queries, the engine seemed to return a partial count of what it had found up to a certain point.<sup>13</sup> Otherwise, the other engines tended to have fail-stop

Figure 2. Ratio of queries executed within given runtime



errors, meaning that they either returned full correct results or no results at all. With respect to public SPARQL endpoints, although we would expect partial results due to result-size limits (Buil-Aranda et al., 2013), the pattern of errors in *Table 8* suggests that *if* one of the latter three engines successfully returns a count result, then that value is likely to be correct.

#### 4.5. Summary

In general, we see that Fuseki, Sesame and Virtuoso are capable of describing – in considerable detail – the content of small-to-medium-sized datasets under controlled conditions, returning correct results for almost all queries over DRUGBANK, JAMENDO and KEGG. These results suggest that when deployed as public SPARQL endpoints, these implementations could provide a rich catalogue of the content of such datasets. However, we would not expect to derive as rich or as trustworthy a description from 4store-powered endpoints, nor from larger datasets or datasets with more diverse schema terms.

### 5. REMOTE EXPERIMENTS

We now look at how public SPARQL endpoints themselves perform for the list of self-descriptive queries we have previously enumerated. Along these lines, we collected a list of 540 SPARQL endpoints registered in the DataHub in April 2015<sup>14</sup>. We likewise collected a list of 137 endpoints from Bio2RDF releases 1-3. In total, we considered 618 unique endpoints (59 endpoints were present in both lists). The results are based on experiments we performed in April 2015.

#### 5.1. Implementations Used

With respect to the previous local experiments, we are first interested to see if we can determine which implementations are used by the in-scope endpoints. As per the observations of Buil-Aranda et al. (Buil-Aranda et al., 2013), although there is no generic exact method of determining the engine powering a SPARQL endpoint, the HTTP header may contain some clues in the Server field. Hence our first step was to perform a lookup on the endpoint URLs. In *Table 9*, we present the response codes of this step, where we see that quite a large number of endpoints return error codes 4xx, 5xx, or some other exception. This indicates that a non-trivial fraction of the endpoints from our list are offline; we will return to this issue later.

Table 9. HTTP response

Response	No
200 (successful)	307
200(unsuccesful)	43
400	56
404	66
500	4
502	0
503	32
Unknown host	51
Time out	23
Connection refused	18
Not responding	7

With respect to the server names returned by those URLs that returned a HTTP response, *Table 10* enumerates the main prefixes that we discovered. Although some of the server names denote generic HTTP servers – more specifically *Apache*, *nginx*, *Jetty*, *GlassFish*, *Restlet* and *lighttpd* – we also see some names that indicate SPARQL implementations – namely *Virtuoso*, *Fuseki* and *4s-httpd* (4store). Interestingly, we see that two of the engines that performed quite well in our local experiments – *Virtuoso* and *Fuseki* – are quite prevalent amongst SPARQL endpoints<sup>15</sup>.

## 5.2. Availability and Version

Based on the previous experiment, we suspect some of the endpoints in our list may be offline. Hence we next look at how many endpoints respond to the basic availability query [Q<sub>A1</sub>].

Given that we run queries in an uncontrolled environment, we perform multiple runs to help mitigate temporary errors and remote server loads: the core idea is that if an endpoint fails at a given moment of time, a catalogue could simply reuse the most recent successful result. Along these lines, we ran three weekly experiments in the month of April 2015. In total, 306 endpoints (49.5%) responded to [Q<sub>A1</sub>] at least once in the three weeks; we deem these endpoints to be operational and others to

Table 10. Server Names

Server-field	No
Apache	203
Virtuoso	174
Nginx	38
Jetty	25
Fuseki	15
Glassfish	3
4s-httpd	2
Lighttpd	1
Empty	130

be *offline*. Of the operational endpoints, 7 (1.1%) responded successfully exactly once to  $[Q_{A1}]$ , 28 (4.5%) responded successfully exactly twice, and 272 (44.1%) responded successfully thrice. In the most recent run, 298 endpoints responded to  $[Q_{A1}]$ . Of these, 168 (56.4%) also responded with a single result for  $[Q_{A2}]$ , indicating some support for SPARQL 1.1 in about half of the operational endpoints.

Moving forward, to mitigate the issue of temporary errors, for each endpoint, we consider the most recent non-empty results returned for each endpoint and each query over the three runs.

### 5.3. Success Rates

We first focus on the overall success rates for each query, looking at the ratio of the 307 endpoints that return non-empty results. The results are illustrated in Figure 3, where we see success rates varying from 25% for  $[Q_{E3}]$  on the lower end, to 94% for  $[Q_{C1}]$  on the higher end. The three queries with the highest success rates require only SPARQL 1.0 features to run: list all class partitions ( $Q_{C1}$ ), all property partitions ( $Q_{D1}$ ), and all nested partitions ( $Q_{E1}$ ). Hence we see that – as expected given that only 49% could respond to the SPARQL 1.1 test query  $[Q_{A1}]$  – more endpoints can answer queries not requiring novel SPARQL 1.1 features such as counts or sub-queries. The query with the highest success rate that involved SPARQL 1.1 features was  $Q_{B1}$ , where 51% of endpoints responded with a count of triples. In general, queries deriving counts within partitions had the lowest success rates.

### 5.4. Result Sizes

Next we focus on the size of results returned for each query. To illustrate this, in *Figure 4* we show result sizes in log scale for individual queries at various percentiles considering all endpoints that returned a non-empty result. As expected, queries that return a single count triple return one result across all percentiles. For other queries, the result sizes extended into the tens of thousands. One may note that the higher percentiles are quite compressed for certain queries, indicating the presence of result thresholds. For example, for  $[Q_{C1}]$ , a common result-size was precisely 40,000, which would appear to be the effect of a result-size threshold. Hence, unlike the local experiments where result thresholds could be switched off, we see that for public endpoints, partial results are sometimes returned.

Figure 3. Ratio of endpoints returning non-empty results per query

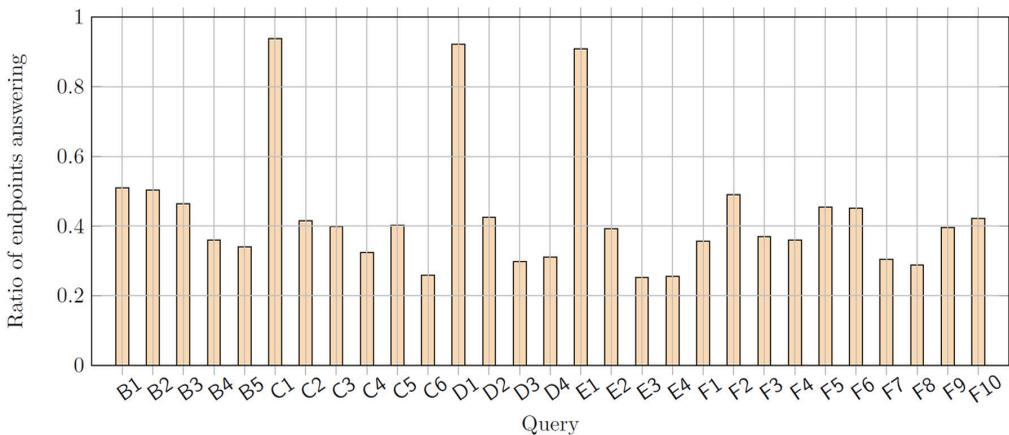
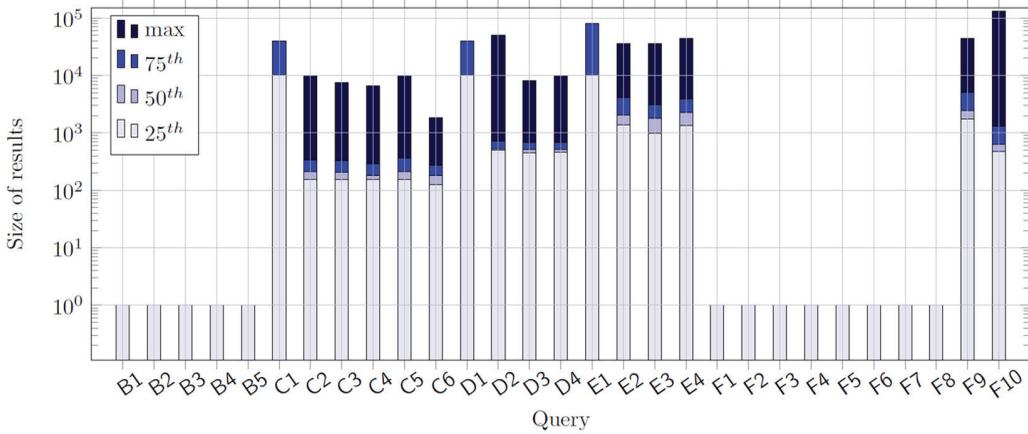


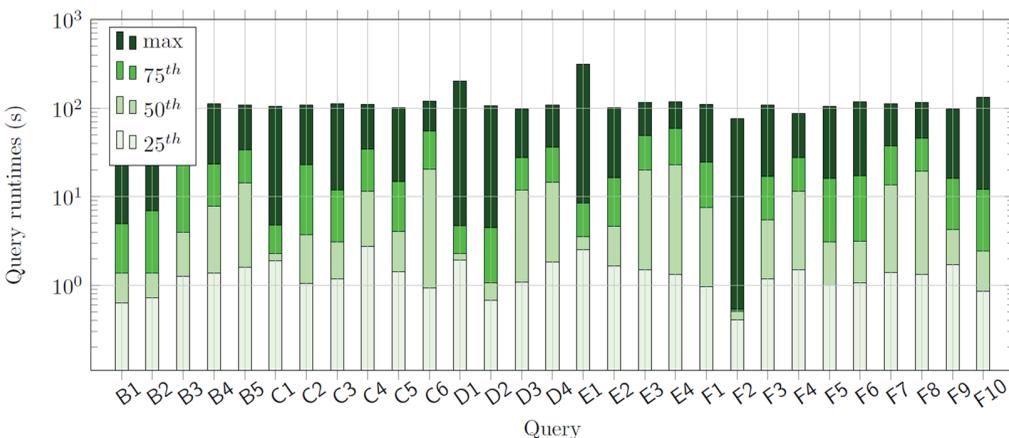
Figure 4. Sizes of results for different queries taking 25th, 50th (median), 75th and 100th (max) percentiles, inclusive, across all endpoints returning non-empty results



### 5.5. Runtimes

Finally, we focus on runtimes for successfully executed queries, incorporating the total response time for issuing the query and streaming all results. In *Figure 5*, we again present the runtimes for each query considering different percentiles across all endpoints returning non-empty results in log scale. We see quite a large variance in runtimes, which is to be expected given that different endpoints host datasets of a variety of sizes and schemata on servers with a variety of computational capacity. In general, we see that the 25<sup>th</sup> percentile roughly corresponds with the one second line, but that slower endpoints may take tens or hundreds of seconds. The at max trend seems to be the effect of remote timeout policies, where query runtimes often maxed out at between 100 – 120 seconds, likely returning partial results.

Figure 5. Runtimes for different queries taking 25th, 50th (median), 75th and 100th (max) percentiles inclusive, across all endpoints returning non-empty results



## 5.6. Summary

Although we see a high success rate in asking for class and property partitions where we would expect to have such data for over 90% of the endpoints, the success rate for queries using novel SPARQL 1.1 features drops to 25 – 50%. We also noted that for queries generating larger result sizes, thresholds and timeouts would likely lead to only partial results being returned. But based on local experiments and the implementations most prominently used by endpoints, we posit that the partial data returned by these endpoints is likely to be accurate even if incomplete.

## 6. SPARQL PORTAL

Our primary motivation in this paper is to investigate a method for cataloguing the content of public SPARQL endpoints without requiring them to publish separate, static descriptions of their content—or indeed, for publishers to offer any additional infrastructure other than the query interface itself. In the previous sections, we performed a variety of experiments that characterised the feasibilities and limitations of collecting metadata about the content of endpoints by directly querying them. In this section, we describe the **SPORTAL** catalogue itself, including its interfaces, capabilities and limitations. A prototype of **SPORTAL** is available online at <http://www.SPORTALproject.org>.

### 6.1. Building the Catalogue

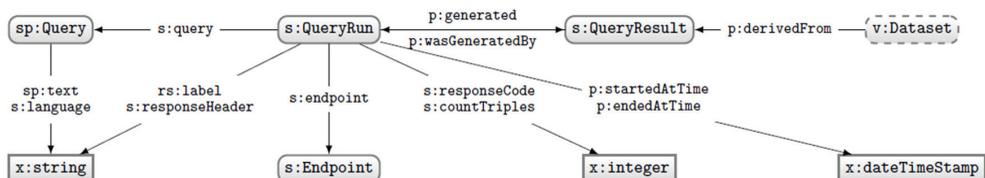
The results of the self-descriptive queries are used to form a content description for each endpoint, which collectively form the **SPORTAL** catalogue. This catalogue is indexed in a local SPARQL endpoint that agents can access. The result for each self-descriptive query over each endpoint is loaded into a dedicated Named Graph and annotated with provenance information using the model illustrated in Figure 6, which follows the recommendations of the *W3C PROV-O* ontology (Lebo, Sahoo, & McGuinness, 2013) (based on the notion of activities and entities). Each query run is (implicitly)<sup>16</sup> considered to be an activity, with an associated start time and end time. This activity uses a query entity and an endpoint entity to generate a query-result entity (a Named Graph with the results). Each VoID dataset is *derived from* potentially multiple query-results. We also keep track of other information, such as HTTP response codes, the number of triples generated by the query, whether the query is SPARQL 1.0 or SPARQL 1.1, the text of the query, etc.

Box 1 provides a real-world example output from executing  $Q_{BI}$  over an endpoint, with provenance information following the model previously described.

### 6.2. SPARQL Interface

**SPORTAL** itself provides a public SPARQL endpoint, where the RDF triples produced by the CONSTRUCT clauses of the self-descriptive queries issued against public endpoints can themselves be queried. This allows users with specific requirements in mind to interrogate our catalogue in a flexible manner.

Figure 6. SPORTAL provenance data model



**Box 1. An Example RDF output of QB1 with provenance metadata**

```

@prefix ep: <http://www.linklion.org:8890/sparql#> .
@prefix p: <http://www.w3.org/ns/prov#> .
@prefix s: <http://vocab.deri.ie/sad#> .
@prefix sp: <http://spinrdf.org/spin#> .
@prefix rs: <http://www.w3.org/2000/01/rdf-schema#> .
@prefix v: <http://rdfs.org/ns/void#> .
@prefix x: <http://www.w3.org/2001/XMLSchema#> .

## PROVENANCE metadata
# each query run is implicitly a p:Activity
ep:totalNumberOfTriplesQueryRun a s:QueryRun ;
  p:generated ep:totalNumberOfTriplesResult ;
  s:responseCode 200 ; s:countTriples 1 ;
  p:startedAtTime "2015-05-11T21:20:54.065Z"^^x:dateTimeStamp ;
  p:endedAtTime "2015-05-11T21:20:55.511Z"^^x:dateTimeStamp ;
  rs:label "Extracting triple count from 'http://www.linklion.org:8890/sparql' on 2015-05-11T21:20:54.065Z"@en ;
  s:resultDataset ep:dataset ;
  s:responseHeader "StatusCode=[HTTP/1.1 200 OK] & Server=[Virtuoso/07.00.3203 (Linux) x86_64-suse-linux-gnu]" ;
  s:endpoint <http://www.linklion.org:8890/sparql> ; # s:endpoint sub-property of p:used
  s:query ep:totalNumberOfTriplesQuery . # s:query sub-property of p:used

# each query is implicitly a p:Entity
ep:totalNumberOfTriplesQuery a sp:Query ;
  ep:text ""PREFIX void: <http://rdfs.org/ns/void#>
  CONSTRUCT { <http://www.linklion.org:8890/sparql#dataset> void:triples ?count }
  WHERE { SELECT (COUNT(*) AS ?count) WHERE { ?s ?p ?o } }"" ;
  s:language "SPARQL1.1" .

# represents the RDF graph returned by the query, implicitly a p:Entity
ep:totalNumberOfTriplesResult a s:QueryResult ;
  p:wasGeneratedBy ep:totalNumberOfTriplesQueryRun .

# connects the dataset mentioned in the results and the graph storing the results
ep:dataset p:wasDerivedFrom s:totalNumberOfTriplesResult .

## RDF GENERATED BY THE QUERY
# loaded into the Named Graph ep:totalNumberOfTriplesResult
ep:dataset v:triples 77455301 .
    
```

To take a first example, a client could pose the following query asking for the SPARQL endpoints for which the catalogue has the top 5 largest triple counts:

```

SELECT DISTINCT ?endpoint ?triples
WHERE { ? dataset v:triples ?triples ; v:sparqlEndpoint ?endpoint . }
ORDER BY DESC(?triples) LIMIT 5
    
```

This will return the answer presented in Table 11.<sup>17</sup>

As another example, referring back to the second client scenario mentioned in the introduction, take a user who is interested in data about proteins and asks for endpoints with at least 50,000 instances of bp:Protein, with results in descending order of number of instances. This user could ask:

**Table 11. SPARQL endpoints with top 5 largest triple counts**

?endpoint	?triple
http://commons.dbpedia.org/sparql_	1,229,690,546
http://lod.b3kat.de/sparql#dataset_	981,672,146
http://www.linklion.org:8890/sparql#dataset_	727,421,750
http://live.dbpedia.org/sparql#dataset_	560,701,025
http://linked.opendata.cz/sparql#dataset_	555,666,667

Table 12. List of SPARQL endpoints with at least 50,000 instances of “bp:Protein”

?endpoint	?instances
https://www.ebi.ac.uk/rdf/services/reactome/sparql	260,546

Table 13. List of SPARQL endpoints with at least 50 unique images of people

?endpoint	?imgs
http://eu.dbpedia.org/sparql	4,517
http://eudbpedia.deusto.es/sparql	4,517
http://data.open.ac.uk/query	311
http://apps.morelab.deusto.es/labman/sparql	78

```
SELECT DISTINCT ?endpoint ?instances
WHERE { ?dataset v:classPartition
[ v:class bp:Protein ; v:distinctSubjects ?instances ] ;
v:sparqlEndpoint ?endpoint . FILTER(?instances > 50000) }
ORDER BY DESC (?instances)
```

This query will return the answer presented in Table 12.

As a final example combining scenarios 2 and 3 in the introduction, consider an agent looking for SPARQL endpoints with at least 50 unique images of people, where this agent may ask:

```
SELECT DISTINCT ?endpoint ?imgs
WHERE {
?dataset v:classPartition [ v:class f:Person ; v:propertyPartition [
v:property f:depiction; v:distinctObjects ?imgs ] ] ;
v:sparqlEndpoint ?endpoint . FILTER(?imgs > 50) }
ORDER BY DESC (?imgs)
```

This returns the result presented in Table 13.

Of course, this is just to briefly highlight three examples of the capabilities of SPORAL and the kinds of results it can return. One could imagine various other types of queries that a user could be interested in posing over the SPORAL catalogue, which supports a variety of types of queries referring to high-level dataset statistics as well as schema-level information. However, the catalogue does not support finding endpoints mentioning a specific resource or value, nor does it currently support keyword search on the topic of the dataset.

### 6.3. User Interface

In order to use the SPARQL interface, the agent must first be familiar with SPARQL, and second must know the IRI of the particular classes and/or properties that they are interested in. To help non-expert users, SPORAL also provides an online user interface with a number of functionalities.

First, users can search for specific endpoints by their URL, by the classes in their datasets, and/or by the properties in their datasets. These features are offered by means of auto-completion on keywords, meaning that the agent need not know the specific IRIs they are searching for. Taking a simple example, if a user wishes to find endpoints with instances of drugs, they may type “*drug*”

into the search bar and then select one of the presented classes matching that search; once a class is selected, the user is presented with a list of public endpoints mentioning that class, ordered by the distinct subjects for that class partition (as available).

If a user clicks on or searches for an endpoint, they can retrieve all the information available about that endpoint as extracted by the queries previously described, providing an overview of how many triples it contains, how many subjects, how many classes, etc. (as available).

The SPORAL user interface also includes some graphical visualisations of some of the high-level features of the catalogue, such as the most popular classes and properties based on the number of endpoints in which they are found, the most common server headers, and so forth. While this may not be of use to a user with a specific search in mind, it offers a useful overview of the content available across all endpoints on the Web, and the schema-level terms that are most often instantiated.

## 6.4. Updates

An important aspect of the SPORAL service is to keep up-to-date information about current SPARQL endpoints. Along these lines, we currently recompute the content descriptions every 15 days: we perform a backup of the old catalogue and simply recompute everything from scratch. One shortcoming of this approach is that the catalogue may miss endpoints that were temporarily unavailable during the computation. Currently we do not implement any special workaround for this issue, but we could in future consider importing data from the previous catalogue for endpoints, with a fixed limit for how long into the past we are willing to still consider content descriptions as valid.

## 6.5. Comparison

In Table 14, we compare the SPORAL catalogue with two other publicly available services that could be used to find relevant SPARQL endpoints using VoID descriptions: DataHub and VoID Store. Unlike SPORAL, both of these services rely on publisher-contributed VoID descriptions.

In the comparison, we include all endpoints that had an associated VoID description in the given service. For DataHub and VoID Store, it is possible to have multiple VoID descriptions associated with an endpoint, and multiple endpoints associated with a VoID file. We count a SPARQL endpoint as available if it could respond with a valid SPARQL response to the query (as used, for example, by the SPARQLES system (Buil-Aranda et al., 2013)):

```
SELECT ?s WHERE { ?s ?p ?o } LIMIT 1
```

For DataHub, VoID les are not hosted locally, where links are provided instead. We used the LDspider v1.3 (Isele, Umbrich, Bizer, & Harth, 2010) crawler to download the VoID files from these URLs,<sup>18</sup> from which we extract the availability (number of VoID files successfully downloaded) and the content for later statistics. To give a brief comparison of the coverage of the catalogues, we also display the number of unique classes and unique properties that are associated in each catalogue with at least one endpoint; in more detail, we count the unique classes and unique properties that would be returned for the following queries over the catalogues, respectively:

Table 14. A comparison of the availability and coverage of SPORAL, DataHub and VoID Store

Service	Endpoints		Descriptions		Classes	Properties
	Total	Available	Total	Available		
<b>SPORAL</b>	307	231 (75%)	298	298 (100%)	19,216	46,313
<b>DataHub</b>	200	115 (58%)	260	162(62%)	1,636	829
<b>VoID Store</b>	118	69 (58%)	148	148 (100%)	30	217

```
SELECT DISTINCT ?c WHERE { ?s v:class ?c }  
SELECT DISTINCT ?p WHERE { ?s v:property ?p }
```

Although this only partially captures the full wealth of information available in VoID, it gives an overview of the diversity of domain terms indexed from endpoints.

From the results, with respect to endpoints, we see that SPORTAL has the broadest coverage: unlike DataHub and VoID Store, it does not require publishers to compute and submit VoID descriptions but rather computes them automatically. For this reason, we see that SPORTAL indexes twice as many available endpoints as DataHub and more than three times that of VoID Store. We also see that the endpoints that SPORTAL indexes have the highest availability ratio: for DataHub and VoID Store, many of the indexed descriptions refer to endpoints that are long dead.

For both SPORTAL and VoID Store, descriptions are hosted locally, meaning that they are always available when the respective catalogue is available; however, for DataHub, 38% of the VoID links provided could not be resolved to RDF content by LDspider.

With respect to the class and property terms mentioned, we see that the SPORTAL catalogue contains orders of magnitude more unique classes and properties than either DataHub or VoID Store.

From these results, we conclude that when compared to DataHub and VoID Store, clients using SPORTAL can expect to find a broader range of relevant endpoints for (e.g.) a broader range of classes and properties, and that the endpoints returned are more likely to be available and to still contain the content in question. Thus we see the benefits of a catalogue based on computing content descriptions rather than relying on those provided by publishers.

## 6.6. Limitations

SPORTAL naturally inherits many of the limitations raised during earlier experiments. For instance, the previous example queries would probably miss endpoints that could not return results for the relevant self-descriptive queries. In general, the catalogue should be considered a best-effort initiative to collect as much metadata about the content of endpoints as possible, rather than a 100% complete catalogue.

Another limitation is that SPORTAL can only help to find endpoints based on the metadata collected from self-describing queries, which mainly centres on the schema terms used. For example, the system cannot help to find endpoints that mention a given literal, or a given subject IRI (which is partially supported by VoID Store using REGEX patterns), or to find endpoints based on the text of the description or the tags associated with the relevant dataset (which is supported by DataHub), etc.

We must also note that by focusing on the problem of finding relevant SPARQL endpoints, SPORTAL may miss relevant Linked Datasets that do not offer a SPARQL endpoint. According to statistics by Jentzsch et al. (Jentzsch et al., 2011), only 68% of the Linked Datasets surveyed provided a SPARQL endpoint. Hence, in addition to missing out on endpoints that cannot answer the self-descriptive queries that SPORTAL issues, we also do not cover Linked Datasets without SPARQL endpoints. However, our focus is specifically on the problem of relevant SPARQL endpoints, which we argue is a sufficiently noteworthy problem in and of itself: a problem that merits specialised methods such as those proposed in this paper.

## 7. CONCLUSION

In this paper, we proposed a novel cataloguing scheme for helping agents to find public SPARQL endpoints relevant to their needs. Given that the endpoints in question are made available by hundreds of different parties, we chose to investigate a cataloguing system that works with the existing SPARQL infrastructure and, for each endpoint indexed, only requires a working SPARQL interface. We ruled out the option of flooding runtime requests to public SPARQL endpoints looking for the desired content since this would lead to long runtimes and could generate a lot of traffic to public

endpoints. Instead, we proposed to use self-descriptive queries to incrementally generate high-level descriptions of the content of public endpoints. We experimented with the performance of running these queries for four datasets and four engines, showing that although Fuseki, Sesame and Virtuoso could successfully answer the queries over small-to-medium-sized datasets, only Virtuoso managed to return results to some queries over larger datasets. We then looked at what sort of success rate public endpoints had in answering these queries, where out of 306 operational endpoints, the ratio of non-empty responses ranged from 25–94% depending on the query in question. Finally we presented details of the SPORAL prototype that uses the catalogue we have extracted from public endpoints to help users find interesting datasets on the Web.

Although SPORAL has its limitations, we have shown that it compares favourably with existing services to help clients find SPARQL endpoints: when compared with DataHub and VoID Store, the SPORAL catalogue has better coverage of available endpoints and, for example, indexes a much broader range of the class and property terms used in the data of remote endpoints. However, it lacks some of the features of these other services: for example, exploring Linked Datasets (and not just SPARQL endpoints) using tags, keyword search over dataset abstracts, searching by resource IRIs, etc.

Our goal in the immediate future is to build upon the existing prototype by seeking feedback from the Linked Data community on what features they feel might be useful, and to gather feedback on the usability of the system. We would also like to investigate fall-back methods of extracting metadata directly from endpoints, such as incremental methods that query, e.g., for statistics about one class/property partition at a time.<sup>19</sup>

The SPORAL prototype is available online at <http://www.SPORALproject.org/>. Project Code is available at <https://github.com/SALiHasnain/sparqlautodescription>.

## ACKNOWLEDGMENT

This publication was supported in part by a research grant from Science Foundation Ireland (SFI) under Grant Number SFI/12/RC/2289, by the Millennium Nucleus Center for Semantic Web Research under Grant NC120004, and by Fondecyt Grant No. 11140900

## REFERENCES

- Abedjan, Z., Grütze, T., Jentzsch, A., & Naumann, F. (2014). Profiling and mining RDF data with ProLOD++. *Proceedings of the International Conference on Data Engineering ICDE* (pp. 1198–1201). <http://doi.org/doi:10.1109/ICDE.2014.6816740> doi:10.1109/ICDE.2014.6816740
- Acosta, M., Vidal, M.-E., Lampo, T., Castillo, J., & Ruckhaus, E. (2011). ANAPSID: An Adaptive Query Processing Engine for SPARQL Endpoints. *Proceedings of the International Semantic Web Conference (ISWC)* (pp. 18–34). Springer. [http://doi.org/doi:10.1007/978-3-642-25073-6\\_2](http://doi.org/doi:10.1007/978-3-642-25073-6_2) doi:10.1007/978-3-642-25073-6\_2
- Akar, Z., Halaç, T. G., Ekinci, E. E., & Dikenelli, O. (2012). Querying the Web of Interlinked Datasets using VOID Descriptions. In *Linked Data On the Web (LDOW)*. CEUR.
- Alexander, K., Cyganiak, R., Hausenblas, M., & Zhao, J. (2009). Describing Linked Datasets. In *Linked Data On the Web (LDOW)*. CEUR.
- Auer, S., Demter, J., Martin, M., & Lehmann, J. (2012). LODStats – An Extensible Framework for High-Performance Dataset Analytics. In *Knowledge Engineering and Knowledge Management EKAW* (pp. 353–362). Springer.
- Basca, C., & Bernstein, A. (2014). Querying a messy web of data with Avalanche. *Journal of Web Semantics*, 26, 1–28. doi:10.1016/j.websem.2014.04.002

- Beek, W., Rietveld, L., Bazoobandi, H. R., Wielemaker, J., & Schlobach, S. (2014). LOD Laundromat: A Uniform Way of Publishing Other People's Dirty Data. *Proceedings of the International Semantic Web Conference (ISWC)* (pp. 213–228). Springer. doi:10.1007/978-3-319-11964-9\_14
- Böhm, C., Lorey, J., & Naumann, F. (2011). Creating VoID descriptions for Web-scale data. *Journal of Web Semantics*, 9(3), 339–345. doi:10.1016/j.websem.2011.06.001
- Broekstra, J., Kampman, A., & van Harmelen, F. (2002). Sesame: A Generic Architecture for Storing and Querying RDF and RDF Schema. *Proceedings of the International Semantic Web Conference (ISWC)* (pp. 54–68). Springer.
- Buil-Aranda, C., Hogan, A., Umbrich, J., & Vandenbussche, P.-Y. (2013). SPARQL Web-Querying Infrastructure: Ready for Action? *Proceedings of the International Semantic Web Conference (ISWC)* (pp. 277–293). Springer.
- Campinas, S., Delbru, R., & Tummarello, G. (2013). Efficiency and precision trade-offs in graph summary algorithms. *Proceedings of the International Database Engineering & Applications Symposium (IDEAS)* (pp. 38–47). doi:doi:10.1145/2513591.2513654 doi:10.1145/2513591.2513654
- Crespo, A., & Garcia-Molina, H. (2002). Routing Indices for Peer-to-Peer Systems. *Proceedings of the International Conference on Distributed Computing Systems (ICDCS)* (pp. 23–32). doi:
- Cygniak, R., Stenzhorn, H., Delbru, R., Decker, S., & Tummarello, G. (2008). Semantic Sitemaps: Efficient and Flexible Access to Datasets on the Semantic Web. *Proceedings of the European Semantic Web Conference ESWC* (pp. 690–704). Springer. doi:10.1007/978-3-540-68234-9\_50
- Erling, O., & Mikhailov, I. (2009). RDF Support in the Virtuoso DBMS. In *Networked Knowledge – Networked Media*. Springer. doi:10.1007/978-3-642-02184-8\_2
- Fetahu, B., Dietze, S., Nunes, B. P., Casanova, M. A., Taibi, D., & Nejd, W. (2014). A Scalable Approach for Efficiently Generating Structured Dataset Topic Profiles. *Proceedings of the European Semantic Web Conference ESWC* (pp. 519–534). Springer. doi:doi:10.1007/978-3-319-07443-6\_35 doi:10.1007/978-3-319-07443-6\_35
- Gallego, M. A., Fernández, J. D., Martínez-Prieto, M. A., & de la Fuente, P. (2011). *An Empirical Study of Real-World SPARQL Queries*. USEWOD.
- Glaser, H., Millard, I., & Jaffri, A. (2008). RKBExplorer.com: A Knowledge Driven Infrastructure for Linked Data Providers. *Proceedings of the European Semantic Web Conference ESWC* (pp. 797–801). Springer. doi:10.1007/978-3-540-68234-9\_61
- Harris, S., Lamb, N., & Shadbolt, N. (2009). 4store: The Design and Implementation of a Clustered RDF Store. *Proceedings of the Scalable Semantic Web Systems Workshop (SWSS)*.
- Harris, S., Seaborne, A., & Prud'hommeaux, E. (2013, March). SPARQL 1.1 Query Language.
- Harth, A., Hose, K., Karnstedt, M., Polleres, A., Sattler, K.-U., & Umbrich, J. (2010). Data summaries for on-demand queries over linked data. *Proceedings of the International Conference on World Wide Web (WWW)* (pp. 411–420). doi:10.1145/1772690.1772733
- Harth, A., Umbrich, J., Hogan, A., & Decker, S. (2007). YARS2: A Federated Repository for Querying Graph Structured Data from the WHolst. *Proceedings of the International Semantic Web Conference (ISWC)* (pp. 211–224). Springer. doi:10.1007/978-3-540-76298-0\_16
- Hasnain, A., Kamdar, M. R., Hasapis, P., Zeginis, D., & Warren, C. N. Jr et al. (2014, October). Linked Biomedical Dataspace: Lessons Learned integrating Data for Drug Discovery. *Proceedings of the International Semantic Web Conference (In-Use Track)*.
- Hasnain, A., Mehmood, Q., Syeda, S. e Z., Saleem, M., Warren Jr, C., Zehra, D., ... Rebholz-Schuhmann, D. (2016). BioFed: Federated Query Processing over Life Sciences Linked Open Data. *Journal of Biomedical Semantics*.
- Heath, T., & Bizer, C. (2011). *Linked Data: Evolving the Web into a Global Data Space*. Morgan & Claypool.
- Holst, T., & Höfig, E. (2013). Investigating the Relevance of Linked Open Data Sets with SPARQL Queries. In *COMPSAC Workshops* (pp. 230–235).

- Isele, R., Umbrich, J., Bizer, C., & Harth, A. (2010). LDspider: An Open-source Crawling Framework for the Web of Linked Data. *Proceedings of the International Semantic Web Conference (ISWC) Posters & Demos*. CEUR.
- Jentzsch, A., Cyganiak, R., & Bizer, C. (2011, September). State of the LOD Cloud.
- Khatchadourian, S., & Consens, M. P. (2010). ExpLOD: Summary-Based Exploration of Interlinking and RDF Usage in the Linked Open Data Cloud. *Proceedings of the Extended Semantic Web Conference (ESWC)* (pp. 272–287). Springer. doi:10.1007/978-3-642-13489-0\_19
- Kinsella, S., Bojars, U., Harth, A., Breslin, J. G., & Decker, S. (2008). An Interactive Map of Semantic Web Ontology Usage. *Proceedings of the International Conference on Information Visualisation* (pp. 179–184). doi:10.1109/IV.2008.60
- Langegger, A., & Wöß, W. (2009). RDFStats – An Extensible RDF Statistics Generator and Library. In *DEXA Workshops* (pp. 79–83).
- Lebo, T., Sahoo, S., & McGuinness, D. (2013, April). PROV-O: The PROV Ontology.
- Lorey, J. (2014). Identifying and determining SPARQL endpoint characteristics. *International Journal on Semantic Web and Information Systems*, 10(3), 226–244. doi:10.1108/IJWIS-03-2014-0007
- Mäkelä, E. (2014). Aether – Generating and Viewing Extended VoID Statistical Descriptions of RDF Datasets. *Proceedings of the European Semantic Web Conference (ESWC)* (pp. 429–433). Springer.
- Mehdi, M., Iqbal, A., Hogan, A., Hasnain, A., Khan, Y., Decker, S., & Sahay, R. (2014). Discovering domain-specific public SPARQL endpoints: a life-sciences use-case. *Proceedings of the International Database Engineering & Applications Symposium (IDEAS)* (pp. 39–45).
- Mihindukulasooriya, N., Poveda-Villalón, M., García-Castro, R., & Gómez-Pérez, A. (2015). Loupe – An Online Tool for Inspecting Datasets in the Linked Data Cloud. *Proceedings of the International Semantic Web Conference (ISWC) Posters & Demos*. CEUR.
- Mountantonakis, M., Allocca, C., Fafalios, P., Minadakis, N., Marketakis, Y., Lantzaki, C., & Tzitzikas, Y. (2014). Extending VoID for Expressing Connectivity Metrics of a Semantic Warehouse. *Proceedings of the International Workshop on Dataset Profiling & federated Search for Linked Data (PROFILES)*.
- Omitola, T., Zuo, L., Gutteridge, C., Millard, I., Glaser, H., Gibbins, N., & Shadbolt, N. (2011). Tracing the provenance of Linked Data using void. *Proceedings of the International Conference on Web Intelligence, Mining and Semantics (WIMS)* (p. 17). doi:10.1145/1988688.1988709
- Paulheim, H., & Hertling, S. (2013). Discoverability of SPARQL Endpoints in Linked Open Data. *Proceedings of the International Semantic Web Conference (ISWC) Posters & Demos* (pp. 245–248). Springer.
- Prud'hommeaux, E., & Buil-Aranda, C. (2013, March). SPARQL 1.1 Federated Query.
- Qiu, D., & Srikant, R. (2004). Modeling and performance analysis of BitTorrent-like peer-to-peer networks. In *SIGCOMM* (pp. 367–378).
- Quilitz, B., & Leser, U. (2008). Querying Distributed RDF Data Sources with SPARQL. *Proceedings of the European Semantic Web Conference (ESWC)* (pp. 524–538). Springer. doi:10.1007/978-3-540-68234-9\_39
- Ratnasamy, S., Francis, P., Handley, M., Karp, R. M., & Shenker, S. (2001). *A scalable content-addressable network* (pp. 161–172). SIGCOMM;
- Ripeanu, M., Iamnitchi, A., & Foster, I. T. (2002). Mapping the Gnutella Network. *IEEE Internet Computing*, 6(1), 50–57. doi:10.1109/4236.978369
- Rowstron, A. I. T., & Druschel, P. (2001). Pastry: Scalable, Decentralized Object Location, and Routing for Large-Scale Peer-to-Peer Systems. *Proceedings of the IFIP/ACM International Conference on Distributed Systems Platforms (Middleware)* (pp. 329–350).
- Schwarte, A., Haase, P., Hose, K., Schenkel, R., & Schmidt, M. (2011). FedX: A Federation Layer for Distributed Query Processing on Linked Open Data. *Proceedings of the Extended Semantic Web Conference (ESWC)* (pp. 481–486). Springer. doi:10.1007/978-3-642-21064-8\_39

Stoica, I., Morris, R., Liben-Nowell, D., Karger, D. R., Kaashoek, M. F., Dabek, F., & Balakrishnan, H. (2003). Chord: A scalable peer-to-peer lookup protocol for internet applications. *IEEE/ACM Trans. Netw.*, *11*(1), 17–32. doi:10.1109/TNET.2002.808407

Umbrich, J., Hose, K., Karnstedt, M., Harth, A., & Polleres, A. (2011). Comparing data summaries for processing live queries over Linked Data. *World Wide Web Journal*, *14*(5-6), 495–544. doi:10.1007/s11280-010-0107-z

Verborgh, R., Hartig, O., De Meester, B., Haesendonck, G., De Vocht, L., & Vander Sande, M., ... de Walle, R. Van. (2014). Querying Datasets on the Web with High Availability. *Proceedings of the International Semantic Web Conference (ISWC)* (pp. 180–196). Springer. doi:doi:10.1007/978-3-319-11964-9\_12 doi:10.1007/978-3-319-11964-9\_12

Williams, G. T. (2013, March). SPARQL 1.1 Service Description.

Zhao, B. Y., Huang, L., Stribling, J., Rhea, S. C., Joseph, A. D., & Kubiawicz, J. (2004). Tapestry: a resilient global-scale overlay for service deployment. *IEEE Journal on Selected Areas in Communications*, *22*(1), 41–53. doi:10.1109/JSAC.2003.818784

## ENDNOTES

- <sup>1</sup> Note that all prefixes used in this paper are listed in Table 15 of the Appendix.
- <sup>2</sup> Certain aspects of VoID may not be computable directly from a dataset, such as the author(s) of a dataset, how it is licensed, OpenSearch descriptions, etc. Likewise we do not include subjective criteria in the computable fragment – such as the categories of the dataset – even if candidates could be computed automatically (Fetahu et al., 2014).
- <sup>3</sup> <http://jena.apache.org/documentation/fuseki2/>; l.a. 2015/12/10
- <sup>4</sup> We say that an endpoint is *operational* if it can be accessed over HTTP through the SPARQL protocol and will return a valid non-empty response to the following query: `SELECT * WHERE ?s ?p ?o LIMIT 1`
- <sup>5</sup> <http://dbpedia.org/sparql>; l.a. 2015/12/10 (42 populations are returned at the time of writing).
- <sup>6</sup> The single server itself of course may be a distributed system, involving multiple replicated or clustered machines (Harris et al., 2009; Harth, Umbrich, Hogan, & Decker, 2007); however, this is all transparent from the perspective of the client, who sees one server
- <sup>7</sup> <http://DataHub.io/>; l.a. 2015/12/10.
- <sup>8</sup> <http://void.rkbexplorer.com/>; l.a. 2015/12/10
- <sup>9</sup> This does not imply that the endpoint is fully compliant with SPARQL 1.1; only that it supports a subset of features.
- <sup>10</sup> <http://ldf.fi/void-ext#>; denoted herein as e: .
- <sup>11</sup> We created a novel namespace (s:) available from <http://vocab.deri.ie/sad#>.
- <sup>12</sup> In fact, DBPEDIA contained no blank nodes, nor did any of the other datasets.
- <sup>13</sup> Many counts had the value of 1,996 or some other value close to a multiple of a thousand; these results were incorrect where some of the expected values were in the hundreds of thousands.
- <sup>14</sup> <http://DataHub.io>
- <sup>15</sup> These results correspond quite closely with those of Buil-Aranda et al. [10]. We believe that some Sesame endpoints may be within the Apache category since the default Sesame header is Apache-Coyote/1.1.
- <sup>16</sup> We do not explicitly type our entities with PROV-O classes simply to keep the data concise: memberships of the respective classes could be inferred from the domain/range of the PROV-O properties we use.
- <sup>17</sup> All such answers were generated from the SPORAL catalogue in March 2016.
- <sup>18</sup> The exact arguments used were `-s seeds.txt -n -o output.nq -b 0 -any23 -bl .xxx`, indicating to accept all formats supported by any23, to follow redirects but not links (i.e., download seeds), and to not blacklist any file extensions.
- <sup>19</sup> However, the cost of such an approach would be a prohibitively large number of requests if there are a large number of partitions.

## APPENDIX

### Prefixes

In *Table 15*, we list all if the prefixes used in the paper.

**Table 15.** IRI prefixes used in the paper

Prefix	IRI
bp:	<a href="http://www.biopax.org/release/biopax-level3.owl#">http://www.biopax.org/release/biopax-level3.owl#</a>
dcat:	<a href="http://purl.org/dc/terms/">http://purl.org/dc/terms/</a>
dbo:	<a href="http://dbpedia.org/ontology/">http://dbpedia.org/ontology/</a>
e:	<a href="http://ldf.fi/void-ext#">http://ldf.fi/void-ext#</a>
f:	<a href="http://xmlns.com/foaf/0.1/">http://xmlns.com/foaf/0.1/</a>
mo:	<a href="http://purl.org/ontology/mo/">http://purl.org/ontology/mo/</a>
p:	<a href="http://www.w3.org/ns/prov#">http://www.w3.org/ns/prov#</a>
s:	<a href="http://vocab.deri.ie/sad#">http://vocab.deri.ie/sad#</a>
sp:	<a href="http://spinrdf.org/spin#">http://spinrdf.org/spin#</a>
rs:	<a href="http://www.w3.org/2000/01/rdf-schema#">http://www.w3.org/2000/01/rdf-schema#</a>
v:	<a href="http://rdfs.org/ns/void#">http://rdfs.org/ns/void#</a>
x:	<a href="http://www.w3.org/2001/XMLSchema#">http://www.w3.org/2001/XMLSchema#</a>

*Ali Hasnain is a PhD Candidate and Research Assistant at the Insight Centre for Data Analytics, National University of Ireland Galway. He received his Master from the Royal Institute of Technology, Stockholm Sweden in Engineering and Management of Information Systems. He received another Master Degree in Project Management and Operational Development from the same Institute. His research interests are mainly around the Semantic Web and Linked Data areas: working with SPARQL Query Federation, Cataloguing, Linking, Indexing and Searching, etc. He is on the Reviewing Committee board of several international conference including VOILA at ISWC 2016 and KESW 2016.*

*Qaiser Mehmood is a Research Assistant at Insight Centre for Data Analytics (NUIG). He completed a Master degree in "Computer Engineering" from Mid Sweden University, Sundsvall Sweden. His research interests are mainly around Dataset Profiling, Semantic Models, Provenance, Data Integration and Querying.*

*Syeda Sana E Zainab is a Master Student of College of Engineering and Informatics, National University of Ireland, Galway. She received her Bachelor's degree in Software Engineering and pursued her Master degree in the field of Semantic Web/Linked data. Her research interests are mainly in the Semantic Web area: working with SPARQL query federation and various Linked data visualisation approaches. She is now working as Research Assistant at INSIGHT Centre for Data Analytics, Galway and associated with Healthcare Life Sciences Linked Data research group.*

*Aidan Hogan is an Assistant Professor at the Department of Computer Science, Universidad de Chile and an Associate Researcher at the Center for Semantic Web Research. He received his PhD from the National University of Ireland, Galway in February 2011 while working with the DERI Galway research group. He continued as a PostDoc in Galway until moving to Chile in December 2013. His research interests centre around the Semantic Web area: working with lots of diverse structured Web data, be it indexing, querying, reasoning, searching, etc. He has won three best reviewer awards (SWJ, EKAW, ISWC), a best poster (ESWC) and a best evaluation paper (ISWC) and picked up a couple of other best paper nominations. He is on the editorial board of the Semantic Web Journal.*

# International Journal on Semantic Web and Information Systems

Volume 12 • Issue 3 • July-September 2016 • ISSN: 1552-6283 • eISSN: 1552-6291

**An official publication of the Information Resources Management Association**

## MISSION

The **International Journal on Semantic Web and Information Systems (IJSWIS)** is an archival journal that publishes high quality original manuscripts in all aspects of Semantic Web that are relevant to computer science and information systems communities. IJSWIS is an open forum aiming to cultivate the Semantic Web vision within the information systems research community. The main focus is on information systems discipline and working towards the delivery of the main implications that the Semantic Web brings to information systems and the information/knowledge society.

## SUBSCRIPTION INFORMATION

IJSWIS is published Quarterly: January-March; April-June; July-September; October-December by IGI Global. Full subscription information may be found at [www.igi-global.com/IJSWIS](http://www.igi-global.com/IJSWIS). The journal is available in print and electronic formats.

Institutions may also purchase a site license providing access to the full IGI Global journal collection featuring more than 100 topical journals in information/computer science and technology applied to business & public administration, engineering, education, medical & healthcare, and social science. For information visit [www.igi-global.com/isj](http://www.igi-global.com/isj) or contact IGI at [eresources@igi-global.com](mailto:eresources@igi-global.com).

## CORRESPONDENCE AND QUESTIONS

### EDITORIAL

Miltiadis D. Lytras, Editor-in-Chief • [IJSWIS@igi-global.com](mailto:IJSWIS@igi-global.com)

### SUBSCRIBER INFO

#### **IGI Global • Customer Service**

701 East Chocolate Avenue • Hershey PA 17033-1240, USA

**Telephone:** 717/533-8845 x100 • **E-Mail:** [cust@igi-global.com](mailto:cust@igi-global.com)

*The International Journal on Semantic Web and Information Systems* is indexed or listed in the following.

ACM Digital Library; Bacon's Media Directory; Burrelle's Media Directory; Cabell's Directories; Compendex (Elsevier Engineering Index); CSA Illumina; Current Contents®/Engineering, Computing, & Technology; DBLP; DEST Register of Refereed Journals; Gale Directory of Publications & Broadcast Media; GetCited; Google Scholar; INSPEC; Journal Citation Reports/Science Edition; JournalTOCs; Library & Information Science Abstracts (LISA); MediaFinder; Norwegian Social Science Data Services (NSD); Science Citation Index Expanded (SciSearch®); SCOPUS; The Index of Information Systems Journals; The Standard Periodical Directory; Thomson Reuters; Ulrich's Periodicals Directory; Web of Science