

# ORCA – a Benchmark for Data Web Crawlers

Michael Röder  
DICE Group

Department of Computer Science  
Paderborn University, Germany;  
Institute for Applied Informatics  
Leipzig, Germany

Geraldo de Souza Jr.  
and Denis Kuchelev

and Abdelmoneim Amer Desouki  
DICE Group  
Department of Computer Science  
Paderborn University, Germany

Axel-Cyrille Ngonga Ngomo  
DICE Group

Department of Computer Science  
Paderborn University, Germany;  
Institute for Applied Informatics  
Leipzig, Germany

**Abstract**—The number of RDF knowledge graphs available on the Web grows constantly. Gathering these graphs at large scale for downstream applications hence requires the use of crawlers. Although Data Web crawlers exist, and general Web crawlers could be adapted to focus on the Data Web, there is currently no benchmark to fairly evaluate their performance. Our work closes this gap by presenting the ORCA benchmark. ORCA generates a synthetic Data Web, which is decoupled from the original Web and enables a fair and repeatable comparison of Data Web crawlers. Our evaluations show that ORCA can be used to reveal the different advantages and disadvantages of existing crawlers. The benchmark is open-source and available at <https://w3id.org/dice-research/orca>.

## I. INTRODUCTION

The number of RDF knowledge graphs (KGs) available on the Web has grown continuously over recent years.<sup>1</sup> These KGs are provided through means ranging from SPARQL endpoints over simple dump files to information embedded in HTML pages. Data Web crawlers are employed to discover and make use of such KGs [16]. With the growing amount of available KGs on the web, these crawlers become more important, e.g., in government-funded projects.<sup>2</sup> The efficiency and effectiveness of such crawlers are typically evaluated by crawling the Web for a set amount of time while measuring different performance indicators such as the number of requests performed by the crawler [12], [16]. While this kind of experiment can be performed for a crawler at a given point in time, the experiments are virtually impossible to repeat and thus, hard to compare with similar experiments. This is due to several factors, including primarily the fact that the Web is an ever-changing, evolving network of single, partly unreliable nodes. Another influence is the geographical location of the machine on which the crawler is executed. For example, geo-blocking can have an influence on the shape of the crawled network. Executing the same crawler on the same hardware might also lead to different evaluation results when various internet service providers offering different connections and bandwidths are used. In addition, the ground truth is not known in such experiments. Since the content of the complete Web is unknown, it is hard to measure the effectiveness of a crawler, i.e., its ability to retrieve relevant data.

<sup>1</sup>For an example, see <https://lod-cloud.net/>.

<sup>2</sup>For example, the projects OPAL (<http://projekt-opal.de/>) and LIMBO (<https://www.limbo-project.org/>) funded by the German government.

To overcome these limitations, we propose ORCA—a benchmark for Web Data Crawlers. The basic idea of ORCA is to alleviate the limitations of current benchmarking approaches by (1) generating a synthetic Data Web and (2) comparing the performance of crawlers within this controlled environment. The generation of the synthetic Web is based on statistics gathered from a sample of the real Data Web. The deterministic generation process implemented by our approach ensures that crawlers are benchmarked in a repeatable and comparable way.

This paper has four main contributions.

- 1) We provide an approach to generate a synthetic Data Web.
- 2) Based on this generator, we present our second contribution, ORCA, the first extensible FAIR benchmark for Data Web crawlers, which can measure the efficiency and effectiveness of crawlers in a comparable and repeatable way.
- 3) We present the first direct comparison of 2 Data Web crawlers in a repeatable setup.
- 4) We show how ORCA can be used to evaluate the politeness of a crawler, i.e., whether it abides by the Robots Exclusion Protocol [20].

The rest of the paper is organised as follows. Section II presents related work while Section III defines prerequisites. In Section IV, we describe the approach and its implementation. The experiments and their results are presented in Section V and discussed in Section VI. We conclude the paper with Section VII.

## II. RELATED WORK

We separate our overview of the related work into two parts. First, we present related publications regarding crawlers and their evaluations. Note that due to the limited space, we mainly focus on Data Web crawlers. Second, we present a brief overview of related work, with statistics regarding the Semantic Web.

*a) Crawlers and their Evaluation.*: The Mercator Web Crawler [12] is an example of a general Web crawler. The authors describe the major components of a scalable Web crawler and discuss design alternatives. The evaluation of the crawler comprises an 8-day run, which has been compared to similar runs of the Google and Internet Archive crawlers. As performance metrics, the number of HTTP requests performed

in a certain time period, and the download rate (in both documents per second and bytes per second) are used. [25] presents an evaluation framework for comparing topical crawlers, i.e., crawlers that are searching for web pages of a certain topic. It relies on a given topic hierarchy and the real web, which makes it susceptible for the aforementioned drawbacks. The BUBiNG crawler [6] has been evaluated relying on the real web as well as a simulation. This simulation was carried out by using a proxy that generated synthetic HTML pages. However, the authors do not give further details about the structure of the simulated web.

A crawler focusing on structured data is presented in [11]. It comprises a 5-step pipeline and converts structured data formats like XHTML or RSS into RDF. The evaluation is based on experiments in which the authors crawl 100k randomly selected URIs. To the best of our knowledge, the crawler is not available as open source project. In [13], [14], a distributed crawler is described, which is used to index resources for the Semantic Web Search Engine. In the evaluation, different configurations of the crawler are compared, based on the time the crawler needs to crawl a given amount of seed URIs. To the best of our knowledge, the crawler is not available as open-source project. In [4], the authors present the LOD Laundromat—an approach to download, parse, clean, analyse and republish RDF datasets. The tool relies on a given list of seed URLs and comes with a robust parsing algorithm for various RDF serialisations. In [9], the authors use the LOD Laundromat to provide a dump file comprising 650K datasets and more than 28 billion triples. Two open-source crawlers for the Data Web are LDSpider [16] and Squirrel [23]. Both are described in detail in Section V-A.

Apache Nutch [19] is an open-source Web crawler. However, the only available plugin for processing RDF stems from 2007, relies on an out-dated crawler version and was not working during our evaluation.<sup>3</sup>

*b) The Data Web.*: There are several publications analysing the Web of data that are relevant for our work, since we use their insights to generate a synthetic Data Web. The Linked Open Data (LOD) Cloud diagram project periodically generates diagrams representing the LOD Cloud and has grown from 12 datasets in 2007 to more than 1200 datasets in 2019.<sup>4</sup> These datasets are entered manually, require a minimum size and must be connected to at least one other dataset in the diagram. Other approaches for analysing the Data Web are based on the automatic gathering of datasets. LODStats [3], [7] collects statistical data about more than 9 000 RDF datasets gathered from a dataset catalogue.<sup>5</sup> In a similar way, [24] uses the LDSpider crawler [16] to crawl datasets in the Web. In [15], the authors gather and analyse 3.985 million open RDF documents from 778 different domains regarding their conformity to Linked Data best practices. The authors of [21] compare different methods to identify SPARQL endpoints on

the Web and suggest that most SPARQL endpoints can be found using dataset catalogues. [24] confirms this finding by pointing out that only 14.69% of the crawled datasets provide VoID metadata. In [5], the authors analyse the adoption of the different technologies, i.e., RDFa [1], Microdata [18] and Microformats and crawl 3 billion HTML pages to this end. None of these works target a benchmark for Data Web crawlers. We address this research gap with the work.

### III. PRELIMINARIES

*a) Data Web Crawler.*: Throughout the rest of the paper, we model a crawler as a program that is able to (1) download Web resources, (2) extract information from these resources and (3) identify the addresses of other Web resources within the extracted information. It will use these (potentially previously unknown) addresses to start with step 1 again in an autonomous way. A Data Web crawler is a crawler that extracts RDF triples from Web resources. Note that this definition excludes programs like the LOD Laundromat [4], which download and parse a given list of Web resources without performing the third step.

*b) Crawlable Graph.*: Let  $\mathbb{G} = (\mathbb{V}, \mathbb{E})$  be a directed graph.  $\mathbb{V} = \{v_1, v_2, \dots\}$  is the set of nodes of the graph.  $\mathbb{E} \subseteq \mathbb{V}^2$  is the set of edges of  $\mathbb{G}$ . Given an edge  $e = (u, v)$ , we call  $u$  the source node and  $v$  the target node of  $e$ . Let  $S \subseteq \mathbb{V}$  be a set of seed nodes. We call a graph  $\mathbb{G}$  *crawlable w.r.t.  $S$*  iff all nodes  $v \in \mathbb{V}$  can be reached from nodes  $u \in S$  in a finite number of steps by traversing the edges of the graph following their direction. A special case of crawlable graphs are graphs that are crawlable w.r.t. a singleton  $S = \{v_\epsilon\}$ . We call  $v_\epsilon$  an *entrance node* of such graphs.

*c) Data Web Analysis.*: The Data Web comprises servers of varying complexity. The types of nodes in this portion of the Web include simple file servers offering their data as dump files, Web servers able to dereference single RDF URIs or to serve HTML web pages with embedded structured data, and SPARQL endpoints that are able to handle complex queries. Hence, we define the different types of nodes in the synthetic Data Web that is to be generated and used for the benchmark: (1) **Dump file node.** This node comprises an HTTP server offering RDF data as a single dump file. In its current implementation, ORCA randomly chooses one of the following RDF serialisations: RDF/XML, Notation 3, N-Triples, or Turtle. Additionally, the file might be compressed with one of three available compression algorithms—ZIP, Gzip or bzip2.<sup>6</sup> (2) **Dereferencing node.** This node comprises an HTTP server and answers requests to single RDF resources by sending all triples of its RDF graph that have the requested resource as subject. The server offers all serialisations supported by Apache Jena.<sup>7</sup> When a request is received, the serialisation is chosen based on the HTTP `Accept` header sent by the

<sup>3</sup><https://issues.apache.org/jira/browse/NUTCH-460>

<sup>4</sup><https://lod-cloud.net/>

<sup>5</sup>The dataset catalogue is <http://thedatahub.org>.

<sup>6</sup>Details regarding the compressions can be found at <https://pkware.cachefly.net/Webdocs/APPNOTE/APPNOTE-6.3.5.TXT>, <https://www.gnu.org/software/gzip/> and <http://sourceware.org/bzip2/>, respectively.

<sup>7</sup><https://jena.apache.org/>

crawler. The complete list of serialisations supported by ORCA can be seen in Table II. (3) **SPARQL endpoint**. This node offers an API, which can be used to query the RDF data using SPARQL via HTTP.<sup>8</sup>(4) **RDFa**. This node offers HTML web pages via HTTP. The web pages contain structured data as RDFa. In its current version, the RDFa node relies on RDFa 1.0 and RDFa 1.1 test cases for HTML and XHTML of an existing RDFa test suite.<sup>9</sup> (5) **CKAN**. CKAN is a dataset catalogue containing meta data about datasets.<sup>10</sup> It offers human-readable HTML pages and an API that can be used to query the catalogue content.

d) *Robots Exclusion Protocol*.: The Robots Exclusion Protocol allows the definition of rules for bots like crawlers [20]. The draft of the standard defines two rules—allow and disallow. They allow or disallow access to a certain path on a domain, respectively. The rules are defined in a `robots.txt` file, which is typically hosted directly under the domain in which the rules have been defined. Although additional rules are not covered by the standard, the standard allows the addition of lines. Some domain owners and crawlers make use of a `Crawl-delay` instruction to define how much delay a crawler should have between its requests to this single Web server.<sup>11</sup>

#### IV. APPROACH

The main idea behind ORCA is to ensure the comparable evaluation of crawlers by creating a local, synthetic Data Web. The benchmarked crawler is initialised with a set of seed nodes of this synthetic cloud and asked to crawl the complete cloud. Since the cloud is generated, the benchmark knows exactly which triples are expected to be crawled and can measure the completeness of the crawl and the speed of the crawler. Since the cloud generation is deterministic, a previously used cloud can be recreated for benchmarking another crawler, ensuring that evaluation results are comparable if the experiments are executed on the same hardware.

It should be noted that our approach relies on a local, synthetic Web comprising containerized servers. This differs to the proxy-based approach of [6] in which a single proxy implementation simulates the web. This allows us to use real-world implementations of servers (e.g., SPARQL endpoints) and eases future additions of further node types.

In the following, we describe the benchmark in detail. We begin by explaining the cloud generation in Section IV-A. An overview of the implementation and its details is given in Section IV-B.

##### A. Cloud Generation

Since the synthetically generated Data Web will be used to benchmark a Data Web crawler, we generate it as a crawlable

<sup>8</sup>In its current version, ORCA uses Virtuoso for this type of node.

<sup>9</sup><http://rdfa.info/test-suite/>

<sup>10</sup><https://ckan.org/>

<sup>11</sup>Examples are Bing (<https://blogs.bing.com/Webmaster/2012/05/03/to-crawl-or-not-to-crawl-that-is-bingbots-question/>) and Yandex (<https://yandex.com/support/Webmaster/controlling-robot/robots-txt.html>).

Table I  
CONNECTIVITY MATRIX  $\mathcal{C}$  USED FOR THE EXPERIMENTS.

from \ to	Deref.	Dump file	SPARQL	CKAN	RDFa
Deref.	1	1	1	1	1
Dump file	1	1	1	1	1
SPARQL	1	1	1	1	1
CKAN	0	1	1	1	1
RDFa	1	1	1	1	1

graph w.r.t. a set of seed nodes  $S$  as defined in Section III-0b. The generation of the synthetic Web can be separated into three steps—(1) Generating the single nodes of the cloud, (2) generating the node graph, i.e., the edges between the nodes, and (3) generating the RDF data contained in the single nodes.

a) *Node Generation*.: The set of nodes  $U$  is generated by virtue of types selected from the list of available types in Section III-0c. The number of nodes in the synthetic Web ( $\nu$ ) and the distribution of node types are user-defined parameters of the benchmark. The node generation process makes sure that at least one node is created for each type with an amount  $> 0$  in the configuration. Formally, let  $\Psi = \{\psi_1, \psi_2, \dots\}$  be the set of node types and  $\Psi_u \subseteq \Psi$  be the set of node types to be generated. To ensure that every type occurs at least once, the generation of the first  $|\Psi_u|$  nodes of  $U$  is deterministic and ensures every type in  $\Psi_u$  is generated. The remaining types are assigned using a seeded random model based on the user-defined distribution until  $|U| = \nu$ .

b) *Node Graph Generation*.: In the real-world Data Web, connections (i.e., edges) between instances of certain node types are unlikely. For example, an open data portal is very likely to point to dump files, SPARQL endpoints or even other open data portals. However, it is very unlikely that it points to a single RDF resource, i.e., to a server which dereferences the URI of the resource. To model this distribution, we introduce a connectivity matrix. Let  $\mathcal{C}$  be a  $|\Psi| \times |\Psi|$  matrix.  $c_{ij} = 1$  means that edges from nodes of type  $\psi_i$  to nodes of type  $\psi_j$  are allowed. Otherwise,  $c_{ij} = 0$ . An example of such a connectivity matrix is given in Table I and will be used throughout the paper. It allows all connections except the example mentioned above.

The algorithm that generates the node graph takes the matrix  $\mathcal{C}$ , the previously created list of typed nodes, and the user-configured average node degree as input. It starts with the first  $|\Psi_u|$  nodes of  $U$  and creates connections between them. For these initial nodes, all connections allowed in  $\mathcal{C}$  are created. This initial graph is extended step-wise by adding the other nodes from  $U$ . In each step, the next node from the list is added to the graph. The outgoing edges of the new node are added using a weighted sampling over the nodes that are permissible from the new node according to  $\mathcal{C}$ . Since the Web is known to be a scale-free network, the weights are the in-degrees of the nodes following the Barabási-Albert model for scale-free networks [2]. In the same way, a similar number of connections to the new node are generated.

After generating the node graph, a set of seed nodes  $S$  has

to be generated to make the graph crawlable as described in Section III-0b. This search is equivalent to the set cover problem [17]. Hence, searching for the smallest set of seed nodes would be NP-hard. Thus, we use a greedy solution which takes  $\mathbb{V}$  and  $\mathbb{E}$  of the generated node graph as input. It starts with defining all nodes as unmarked nodes. After this, the first unmarked node is added to  $S$ . A breadth-first search starts from this node and marks all reachable nodes. After that, the steps are repeated by choosing another unmarked node until all nodes are marked, i.e., all nodes are reachable from the nodes in  $S$ .

*c) RDF Data Generation.*: The benchmark can work with any RDF data generator. However, it comes with a simple generator that ensures the crawlability of the graph. The generator can be configured with three parameters: the average number of triples per graph ( $\tau$ ), the distribution of the sizes of the single graphs and the average degree of the RDF resources ( $d$ ) in the graph. In its current version, ORCA offers a simple approach that statically assigns the given average size to every RDF graph. However, this can be changed to use any other distribution.

Let  $\mathcal{G} = (R, P, L, T)$  be an RDF graph, where  $R = \{r_1, r_2, \dots\}$  is the set of URI resources of the graph,  $P = \{p_1, p_2, \dots\}$  is the set of properties,  $L$  is the set of external URI resources, i.e., resources belonging to a different RDF graph, with  $R \cap L = \emptyset$  and  $T = \{t_1, t_2, \dots\}$  being the set of triples of the RDF graph where each triple has the form  $t_j = \{(s_j, p_j, o_j) | s_j \in R, p_j \in P, o_j \in (R \cup L)\}$ .<sup>12</sup>  $T$  can be separated into two subsets  $T = T_i \cup T_o$ . The set of graph-internal triples  $T_i$  comprises triples with objects  $o_j \in R$ . In contrast, the set of outgoing triples  $T_o$  (a.k.a. link set) contains only triples with external resources as objects ( $o_j \in L$ ). Further, let  $d$  be the average node degree of the resources, i.e., the number of triples a resource is part of.

Like the node graph, each created RDF graph has to be crawlable w.r.t. a set of resources. For the RDF graphs, we implemented an algorithm based on the Barabási-Albert model for scale-free networks [2]. The algorithm guarantees that all resources within the generated RDF graph can be reached from the first resource it generates. As defined in Section III-0b, this resource can be used later on as entrance node by all other RDF graphs which must generate links to this graph.

Let  $\tau$  be the RDF graph size that has been determined based on the chosen parameters. Based on the previously created node graph, the number of outgoing edges  $\tau_o = |T_o|$  as well as their objects, i.e., the set of external URI resources  $L$ , are known. Algorithm 1 takes  $\tau_i = \tau - \tau_o$  together with the average degree  $d$  and a generated set of properties  $P$  as input to generate an initial version of graph  $\mathcal{G}$ . The loop (lines 4–13) adds new resources to the graph until the number of necessary triples has been reached. For each new resource  $r_n$ , a URI is generated (line 5) before it is connected to the existing resources of the graph. After that, the degree of the

---

### Algorithm 1: Initial RDF graph generation

---

**Input** :  $\tau_i, P, d$   
**Output**:  $\mathcal{G}$

```

1  $E, T_i \leftarrow \{\}$ 
2  $R \leftarrow \{r_1\}$ 
3  $d_r \leftarrow \text{drawDegree}(d)$ 
4 while  $|T_i| < \tau_i$  do
5    $r_n \leftarrow \text{genResource}(|R|)$ 
6   while  $(|T_i| < \tau_i) \&\& (\text{degree}(r_n) < d_r)$  do
7      $R_c \leftarrow \text{drawFromDegreeDist}(R, T_i)$ 
8     for  $r_c \in R_c$  do
9       if  $(\text{degree}(r_n) == 0) \parallel (\text{bernoulli}(\frac{0.5d_r-1}{d_r-1}))$  then
10         $T_i \leftarrow T_i \cup \{\text{genTriple}(r_c, \text{draw}(P), r_n)\}$ 
11       else
12         $T_i \leftarrow T_i \cup \{\text{genTriple}(r_n, \text{draw}(P), r_c)\}$ 
13    $R \leftarrow R \cup \{r_n\}$ 
14  $\mathcal{G} \leftarrow \{R, P, \emptyset, T_i\}$ 

```

---

new resource  $d_r$  is drawn from a uniform distribution in the range  $[1, 2d]$  (line 3). The  $d_r$  resources  $r_n$  will be connected to are chosen based on their degree, i.e., the higher the degree of a resource, the higher the probability that it will be chosen for a new connection. The result of this step is the set  $R_c$  with  $|R_c| = d_r$ . For each of these resources, a direction of the newly added triple is chosen. Since the graph needs to be crawlable, the algorithm chooses the first triple to be pointing to the newly resourced node. This ensures that all resources can be reached, starting from the first resource of the graph. For every other triple, the decision is based on a Bernoulli distribution with a probability of  $\frac{0.5d_r-1}{d_r-1}$  being a triple that has the new node as an object (line 9). Based on the chosen direction, the new triple is created with a property that is randomly drawn from the property set  $P$  (lines 10 and 12).

After the initial version of the RDF graph is generated, the outgoing edges of  $T_o$  are created. For each link to another dataset, a triple is generated by drawing a node from the graph as subject, drawing a property from  $P$  as predicate and the given external node as object. Both— $T_o$  and  $L$ —are added to  $\mathcal{G}$  to finish the RDF graph.

*d) URI Generation.*: Every resource of the generated RDF graphs needs to have a URI. To make sure that a crawler can use the URIs during the crawling process, the URIs of the resources are generated depending on the type of node hosting the RDF dataset. All URIs contain the host name. The dump file node URIs contain the file extension representing the format of the file (the RDF serialization and the compression). If a resource of the SPARQL node is used in another generated RDF graph (i.e., to create a link to the SPARQL node), the URL of the SPARQL API is used instead of a resource URI.

<sup>12</sup>Note that this simplified definition of an RDF graph omits the existence of literals and blank nodes.

In a similar way, the links to the CKAN nodes are created by pointing to the CKAN’s Web interface.

## B. Implementation

a) *Overview.*: ORCA is a benchmark built upon the HOBBIT benchmarking platform [22].<sup>13</sup> This FAIR benchmarking platform allows Big Linked Data systems to be benchmarked in a distributed environment. It relies on the Docker<sup>14</sup> container technology to encapsulate the single components of the benchmark and the benchmarked system.

We adapted the suggested design of a benchmark described in [22] to implement ORCA. The benchmark comprises a benchmark controller, data generators, an evaluation module, a triple store and several nodes that form the synthetic Data Web. The benchmark controller is the central control unit of the benchmark. It is created by the HOBBIT platform, receives the configuration defined by the user and manages the other containers that are part of the benchmark. Figure 1 gives an overview of the benchmark components, the data flow and the single steps of the workflow. The workflow itself can be separated into 4 phases—creation, generation, crawling and evaluation. When the benchmark is started, the benchmark controller creates the other containers of the benchmark.<sup>15</sup> During this creation phase, the benchmark controller chooses the types of nodes that will be part of the synthetic Data Web, based on the parameters configured by the user. The Docker images of the chosen node types are started together with an RDF data generator container for each node that will create the data for the node. Additionally, a node data generator, a triple store and the evaluation module are started. The node data generator will generate the node graph. The triple store serves as a sink for the benchmarked Linked Data crawler during the crawling phase while the evaluation module will evaluate the crawled data during the evaluation phase.

After the initial creation, the graph generation phase is started. This phase can be separated into two steps—initial generation and linking. During the first step, each RDF data generator creates an RDF graph for its Web node. In most cases, this is done using the algorithm described in Section IV-A. For data portal and RDFa nodes, the generation process differs. The portal nodes solely use the information to which other nodes they have to be linked to, i.e., each linked node is inserted as a dataset into the data portal node’s database. The RDFa node relies on an already existing test suite.<sup>16</sup> It generates an HTML page that refers to the single tests and to all other connected nodes using RDFa. The node data generator creates the node graph as described in Section IV-A. After this initial generation step, the node graph is sent to the benchmark controller and all RDF data generators (Step 1 in Figure 1). This provides the RDF data generators with the information to which other nodes their RDF graph

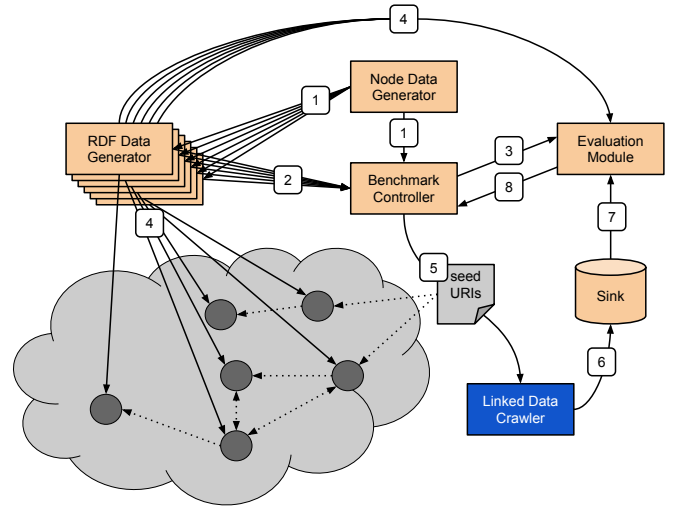


Figure 1. Overview of the Benchmark components and the flow of data. Orange: Benchmark components; Grey: Synthetic Data Web generated by the benchmark; Dark blue: The benchmarked crawler; Solid arrows: Flow of data; Dotted arrows: Links between RDF datasets; Numbers indicate the order of steps.

should be linked. Subsequently, the RDF data generators send their metadata to each other and the benchmark controller (Step 2). This provides the data generators with the necessary data to create links to the entrance nodes of other RDF datasets during the linking step. Additionally, the benchmark controller forwards the collected metadata to the evaluation module and the nodes in the cloud (Step 3).<sup>17</sup> At the end of the generation phase, the generated RDF graphs are forwarded to the single nodes and the evaluation module (Step 4). The generation phase ends as soon as all nodes have processed the received data.

After the generation phase is finished and the HOBBIT platform signals that the crawler has initialised itself, the benchmark controller submits the seed URIs to the crawler (Step 5). This starts the crawling process in which the crawler sends all downloaded RDF data to its sink (Step 6). When the crawler finishes its crawling—i.e., all given URIs and all URIs found in the crawled RDF data have been crawled—the crawler terminates and the crawling phase ends.

During the evaluation phase, the evaluation module measures the recall of the crawler by checking whether the RDF graphs generated by the data generators can be found in the sink (Step 7). The result of this evaluation is sent to the benchmark controller, which adds further data and results of the benchmarking process (Step 8). This can include data that has been gathered from the single nodes of the cloud, e.g., access times. After this, the final results are forwarded to the HOBBIT platform.

## V. EVALUATION

For evaluating Data Web crawlers, we use three different experiments. The first experiment uses all available node types

<sup>17</sup>The submission to the cloud nodes has been omitted in the figure.

<sup>13</sup><https://github.com/hobbit-project/platform>

<sup>14</sup><https://www.docker.com/>

<sup>15</sup>The benchmarked crawler is created by the HOBBIT platform as described in [22].

<sup>16</sup><http://rdfa.info/test-suite/>

to generate the synthetic Data Web and mainly focuses on the recall of the benchmarked crawlers. The second experiment uses a simpler Web to measure efficiency. The last experiment checks whether the crawlers abide by the Robots Exclusion Protocol.

For all experiments, the online instance of HOBBIT is used. It is deployed on a cluster with 3 servers that are solely used by the benchmark and 3 servers that are available for the crawler.<sup>18</sup>

#### A. Benchmarked crawlers

To the best of our knowledge, there are only two working open-source Data Web crawlers available—LDSpider and Squirrel. Other crawlers are either not available as open-source project or web crawlers without the ability to process RDF data. Table II shows the RDF serialisations supported by LDSpider and Squirrel in comparison to ORCA.

LDSpider [16] is an open-source Linked Data crawler that has been used in several publications to crawl data from the Web [24].<sup>19</sup> Following the documentation, it is able to process triples serialised as RDF/XML, N3 and Turtle. Additionally, it supports Apache Any23, which can be used to parse RDFa, Microdata, Microformats and JSON-LD.<sup>20</sup> The crawler uses multiple threads to process several URIs in parallel. It offers different crawling strategies—a classic breadth-first strategy (BFS), and a load-balancing strategy (LBS). The latter tries to crawl a given number of URIs as fast as possible by sending parallel requests to different domains without overloading the server of a single domain. The crawler implements a static politeness strategy and offers the configuration of the delay that is inserted between two consecutive requests. The crawler can write its results either to files or send them to a triple store. For our experiments, we dockerized LDSpider and implemented a system adapter to make it compatible with the HOBBIT platform. We created several LDSpider instances with different configurations. LDSpider (T1), (T8), (T16) and (T32) use BFS and 1, 8, 16 or 32 threads, respectively. During our first experiments, we encountered issues with LDSpiders’ SPARQL client, which was not storing the crawled data in the provided triple store. To achieve a fair comparison of the crawlers, we extended our system adapter to implement our own SPARQL client, used LDSpider’s file sink to get the output of the crawling process, and sent file triples to the benchmark sink. These instances of LDSpider are marked with the addition “FS”. Additionally, we configured the LDSpider instance (T32,FS,LSB), which makes use of the load-balancing strategy to compare the two strategies offered by the crawler.

Squirrel [23] is an open-source, distributed Linked Data crawler, which uses Docker containers to distribute its components. It crawls resources in a similar way to the LSB strategy of LDSpider, by grouping URIs based on their domain

and assigning the created URI sets to its workers. Following the documentation, it supports all RDF serialisations implemented by Apache Jena. It uses Apache Any23 to parse Microdata, Microformats and the Semargl parser for RDFa.<sup>21</sup> Furthermore, it supports the crawling of HDT dump files [10], SPARQL endpoints and open data portals. The latter includes the crawling of CKAN portals and configuration of a scraper that can extract information from HTML pages following predefined rules. For compressed dump files, Squirrel implements a recursive decompression strategy for ZIP, Gzip, bzip2 and tar files. For our experiments, we implemented an adapter for the Squirrel crawler. Squirrel (W1), (W3), (W9) and (W18) are instances of the crawler using 1, 3, 9 or 18 worker instances, respectively.<sup>22</sup>

#### B. Data Web Crawling

The first experiment simulates a real-world Data Web and focuses on the effectiveness of the crawlers, i.e., the amount of correct triples they retrieve. To define the distribution of node types, we analyse the download URLs of the LODStats [3], [7] dump from 2016. Based on this analysis, the generated cloud comprises 100 nodes with 40% dump file nodes, 30% SPARQL nodes, 21% dereferencing nodes and 5% CKAN nodes. 4% are RDFa nodes to represent the 1 billion RDFa triples identified by [5] in comparison to the 28 billion RDF triples gathered from the semantic web by [9]. Following [15], the average degree of each node is set to 20. For each node, the RDF graph generation creates 1000 triples with an average degree of 9 triples per resource.<sup>23</sup> Based on our LODStats analysis, 30% of the dump file nodes use one of the available compression algorithms for the dump file. The usage of `robots.txt` files is disabled. Since LDSpider does not support the crawling of SPARQL endpoints, data catalogues like CKAN, or compressed dump files, we expect LDSpider to achieve a lower Recall than Squirrel. The results of the experiment are listed in Table III.<sup>24</sup>

#### C. Efficiency evaluation

The second experiment focuses on the efficiency of the crawler implementations. For this purpose, a synthetic Web comprising 200 dereferencing nodes is used since they offer to negotiate the RDF serialisation for transferring the data. This ensures that all crawlers can crawl the complete Web. The other parameters remain as before. For LDSpider, we use only the FS instances and configure its politeness to not insert any delay between two requests. We expect both crawlers to be able to crawl the complete cloud and that crawler instances

<sup>21</sup><https://github.com/semarglproject/semargl>

<sup>22</sup>Since the HOBBIT cluster assigns 3 servers for the benchmarked crawler, we use multiples of 3 for the number of workers.

<sup>23</sup>The dataset size is roughly the average document size from [15] (after excluding an outlier). The average degree is derived from the statistics of DBpedia, Freebase, OpenCyc, Wikidata and Yago from [8].

<sup>24</sup>The detailed results can be seen at <https://w3id.org/hobbit/experiments#1584544940477,1584544956511,1584544971478,15845403645660,1584545072279,1585230107697,1584962226404,1584962243223,1585574894994,1585574924888,1585532668155,1585574716469>.

<sup>18</sup>Each of the servers has 16 cores with Hyperthreading and 256 GB RAM. The details of the hardware setup that underlies the HOBBIT platform can be found at <https://hobbit-project.github.io/master#hardware-of-the-cluster>.

<sup>19</sup><https://github.com/ldspider/ldspider>

<sup>20</sup><https://any23.apache.org/>

Table II

THE RDF SERIALISATIONS AND COMPRESSIONS SUPPORTED BY ORCA AND THE TWO BENCHMARKED CRAWLERS. (✓) MARKS SERIALISATIONS IN ORCA THAT ARE NOT USED WITHIN HTML FILES OR FOR GENERATING DUMP NODES FOR THE SYNTHETIC LINKED DATA WEB. X MARKS SERIALISATIONS THAT ARE LISTED AS PROCESSIBLE BY A CRAWLER BUT WERE NOT WORKING DURING OUR EVALUATION.

	RDF Serialisations									Comp.			HTML			
	RDF/XML	RDF/JSON	Turtle	N-Triples	N-Quads	Notation 3	JSON-LD	TriG	TriX	HDT	ZIP	Gzip	bzip2	RDFa	Microdata	Microformat
ORCA	✓	(✓)	✓	✓	(✓)	✓	(✓)	(✓)	(✓)	-	✓	✓	✓	✓	-	-
LDSpider	✓	-	✓	X	✓	✓	✓	-	-	-	-	-	-	✓	✓	✓
Squirrel	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

Table III

RESULTS OF THE DATA WEB CRAWLING AND EFFICIENCY EXPERIMENTS.

Crawler	Data Web		Efficiency			
	Micro Recall	Runtime (in s)	Micro Recall	Runtime (in s)	CPU (in s)	RAM (in GB)
LDSpider (T8)	0.00	67	-	-	-	-
LDSpider (T16)	0.00	73	-	-	-	-
LDSpider (T32)	0.00	74	-	-	-	-
LDSpider (T1,FS)	0.31	1798	1.00	1847	627.2	1.8
LDSpider (T8,FS)	0.30	1792	1.00	1717	658.9	5.2
LDSpider (T16,FS)	0.31	1858	1.00	1753	1677.1	1.6
LDSpider (T32,FS)	0.31	1847	1.00	1754	1959.1	4.0
LDSpider (T32,FS,LBS)	0.03	66	0.01	56	-	-
Squirrel (W1)	0.98	6663	1.00	12051	1096.7	3.5
Squirrel (W3)	0.98	2686	1.00	4096	992.0	7.7
Squirrel (W9)	0.98	1412	1.00	1500	652.0	16.6
Squirrel (W18)	0.97	1551	1.00	893	424.0	24.0

with more threads or workers will crawl faster. The results of the experiment are listed in Table III.<sup>25</sup>

1500491 Total CPU usage (in ms) 652125.0

#### D. Robots Exclusion Protocol check

In the third experiment, we evaluate whether the crawlers follow the rules defined in a server’s `robots.txt` file. To this end, we configure ORCA to generate a smaller Web comprising 25 dereferencing nodes. Each of the nodes copies 10% of its RDF resources and marks the copies disallowed for crawling using the `disallow` instruction in its `robots.txt` file. The average number of requested disallowed resources (RDR) is used as metric. Additionally, we define a delay of 10 seconds between two consecutive requests using the `Crawl-delay` instruction in the same file. The average node degree of the nodes is configured as 5 while the average resource degree is set to 6. We calculate the crawl delay fulfilment (CDF) which we define as the requested delay

<sup>25</sup>Due to the lack of space we omit the standard deviations of the measures. While repeating the experiments, the measures turned out to be stable with standard deviations of  $\sim 2\%$  for the RAM,  $\sim 5\%$  for runtime and CPU time. The detailed results can be found at <https://w3id.org/hobbit/experiments#1608306725630,1608306734009,1609758770866,1604685169129,1608052657254,1609763993514,1609779299724,1609792299538,1608038447931>.

Table IV

RESULTS FOR A DATA WEB WITH `ROBOTS.TXT` FILES INCLUDING `DISALLOW` AND `CRAWL-DELAY` RULES. CDF = CRAWL DELAY FULFILMENT; RDR = REQUESTED DISALLOWED RESOURCES.

Crawler	CDF			RDR	Runtime (in s)
	Min	Max	Avg		
LDSpider (T32,FS,BFS)	0.052	0.122	0.089	0.0	224
LDSpider (T32,FS,LBS)	0.002	0.007	0.004	0.0	43
Squirrel (W18)	0.697	0.704	0.699	0.0	2384

divided by the average delay measured by the server. Table IV shows the results of this experiment.<sup>26</sup>

## VI. DISCUSSION

The experiment results give several insights. As expected, none of the instances of LDSpider were able to crawl the complete synthetic Linked Data Web during the first experiment. Apart from the expected reasons previously mentioned (i.e., the missing support for SPARQL, CKAN nodes and compressed dump files), we encountered two additional issues. First, as mentioned in Section V-A, the SPARQL client of LDSpider did not store all the crawled triples in the provided

<sup>26</sup>The detailed results can be seen at <https://w3id.org/hobbit/experiments#1575626666061,1575592492658,1575592510594>.

triple store. This leads to the different recall values of the LDSpider instances with and without the "FS" extension. Second, although we tried several content handler modules and configurations, LDSpider did not crawl dump files provided as N-Triples. In comparison, the Squirrel instances crawl the complete cloud, except for some triples of RDFa and CKAN nodes.

The second experiment reveals that overall, LDSpider is more resource-efficient than Squirrel. In nearly all cases, LDSpider crawls the Web faster and uses less resources than the Squirrel instances. Only with 9 or more workers Squirrel is able to crawl faster. For the size of the graph, the number of threads used by LDSpider do not seem to play a major role when employing the BFS strategy. It could be assumed that the synthetic Web, with 200 nodes, provides only rare situations in which several nodes are crawled by LDSpider in parallel. However, this assumption can be refuted since Squirrel achieves lower runtimes. Therefore, the load-balancing strategy of Squirrel seems to allow faster crawling of the Web than the BFS of LDSpider. However, the LDSpider (T32,FS,LBS) instance implementing a similar load-balancing strategy aborts the crawling process very early in all three experiments. Therefore, a clearer comparison of both strategies is not possible.

The third experiment shows that both crawlers follow the Robots Exclusion Protocol as both did not request disallowed resources. However, Squirrel seems to insert delays between its requests following the `Crawl-delay` instruction—although it reaches on average only 69.9% of the delay the server asked for—while LDSpider does not take this instruction into account and solely relies on its static politeness strategy with a configurable default delay of 500ms.

## VII. CONCLUSION

In this paper, we present ORCA—the first extensible FAIR benchmark for Data Web crawlers, which measures the efficiency and effectiveness of crawlers in a comparable and repeatable way. Using ORCA, we compared two Data Web crawlers in a repeatable setup. We showed that ORCA revealed strengths and limitations of both crawlers. Additionally, we showed that ORCA can be used to evaluate the politeness of a crawler, i.e., whether it abides by the Robots Exclusion Protocol. Our approach will be extended in various ways in future work. First, we will include HTML pages with Microdata, Microformat or JSON-LD into the benchmark. A similar extension will be the addition of further compression algorithms to the dump nodes (e.g., `tar`), as well as the HDT serialization [10]. The generation step will be further improved by adding literals and blank nodes to the generated RDF KGs and altering the dataset sizes. A simulation of network errors will round up the next version of the benchmark.

## ACKNOWLEDGEMENTS

This work has been supported by the German Federal Ministry of Transport and Digital Infrastructure (BMVI), for Economic Affairs and Energy (BMWi), and for Education and

Research (BMBF) within the projects OPAL (no. 19F2028A), LIMBO (no. 19F2029I), RAKI (no. 01MD19012D) and DAIKIRI (no. 01IS19085B).

## REFERENCES

- [1] Adida, B., Herman, I., Sporny, M., Birbeck, M.: RDFa 1.1 Primer – third edition. W3C Note, W3C (March 2015), <http://www.w3.org/TR/rdfa-primer/>
- [2] Albert, R., Barabási, A.L.: Statistical mechanics of complex networks. *Reviews of modern physics* 74(1), 47 (2002)
- [3] Auer, S., Demter, J., Martin, M., Lehmann, J.: Lodstats – an extensible framework for high-performance dataset analytics. In: *Knowledge Engineering and Knowledge Management*. pp. 353–362. Springer (2012)
- [4] Beek, W., Rietveld, L., Bazoobandi, H.R., Wielemaker, J., Schlobach, S.: Lod laundromat: A uniform way of publishing other people's dirty data. In: *The Semantic Web – ISWC 2014*. pp. 213–228. Springer (2014)
- [5] Bizer, C., Eckert, K., Meusel, R., Mühleisen, H., Schuhmacher, M., Völker, J.: Deployment of rdfa, microdata, and microformats on the web—a quantitative analysis. In: *International Semantic Web Conference*. pp. 17–32. Springer (2013)
- [6] Boldi, P., Marino, A., Santini, M., Vigna, S.: Bubing: Massive crawling for the masses. *ACM Trans. Web* 12(2) (Jun 2018)
- [7] Ermilov, I., Lehmann, J., Martin, M., Auer, S.: Lodstats: The data web census dataset. In: *The Semantic Web – ISWC 2016*. Springer (2016)
- [8] Färber, M., Bartscherer, F., Menne, C., Rettinger, A.: Linked data quality of DBpedia, Freebase, OpenCyc, Wikidata, and YAGO. *Semantic Web* 9(1), 77–129 (2018)
- [9] Fernández, J.D., Beek, W., Martínez-Prieto, M.A., Arias, M.: Lod-a-lot. In: *The Semantic Web – ISWC 2017*. pp. 75–83. Springer (2017)
- [10] Fernández, J.D., Martínez-Prieto, M.A., Gutiérrez, C., Polleres, A., Arias, M.: Binary rdf representation for publication and exchange (hdt). *Web Semantics: Science, Services and Agents on the World Wide Web* 19, 22–41 (2013)
- [11] Harth, A., Umbrich, J., Decker, S.: Multicrawler: A pipelined architecture for crawling and indexing semantic web data. In: *The Semantic Web - ISWC 2006*. pp. 258–271. Springer Berlin Heidelberg (2006)
- [12] Heydon, A., Najork, M.: Mercator: A scalable, extensible Web crawler. *World Wide Web* 2, 219–229 (1999)
- [13] Hogan, A.: Exploiting RDFS and OWL for Integrating Heterogeneous, Large-Scale, Linked Data Corpora. *Nat. University of Ireland, Galway* (2011), <http://aidanhogan.com/docs/thesis/>
- [14] Hogan, A., Harth, A., Umbrich, J., Kinsella, S., Polleres, A., Decker, S.: Searching and browsing linked data with SWSE: The semantic web search engine. *Web Semantics: Science, Services and Agents on the World Wide Web* 9(4), 365–401 (2011)
- [15] Hogan, A., Umbrich, J., Harth, A., Cyganiak, R., Polleres, A., Decker, S.: An empirical survey of linked data conformance. *Journal of Web Semantics* 14, 14–44 (2012)
- [16] Isele, R., Umbrich, J., Bizer, C., Harth, A.: LDspider: An open-source crawling framework for the Web of Linked Data. In: *Proceedings of the ISWC 2010 Posters & Demonstrations Track: Collected Abstracts*. vol. 658, pp. 29–32. CEUR-WS (2010)
- [17] Karp, R.M.: *Reducibility Among Combinatorial Problems* (1972)
- [18] Kellogg, G.: Microdata – second edition. W3C Note, W3C (December 2014), <https://www.w3.org/TR/microdata-rdf/>
- [19] Khare, R., Cutting, D., Sitaker, K., Rifkin, A.: Nutch: A flexible and scalable open-source web search engine. *Oregon State Uni.* 1 (2004)
- [20] Koster, M., Illyes, G., Zeller, H., Harvey, L.: Robots Exclusion Protocol. Internet-draft, Internet Engineering Task Force (IETF) (July 2019), <https://tools.ietf.org/html/draft-rep-wg-topic-00>
- [21] Paulheim, H., Hertling, S.: Discoverability of SPARQL Endpoints in Linked Open Data. In: *Proceedings of the ISWC 2013 Posters & Demonstrations Track*. vol. 1035. CEUR-WS.org (2013)
- [22] Röder, M., Kuchelev, D., Ngonga Ngomo, A.C.: HOBBIT: A platform for benchmarking Big Linked Data. *Data Science* (2019)
- [23] Röder, M., de Souza Jr, G., Ngonga Ngomo, A.C.: Squirrel – crawling rdf knowledge graphs on the web. In: *The Semantic Web – ISWC 2020*. Springer International Publishing (2020)
- [24] Schmachtenberg, M., Bizer, C., Paulheim, H.: Adoption of the linked data best practices in different topical domains. In: *The Semantic Web – ISWC 2014*. pp. 245–260. Springer International Publishing (2014)
- [25] Srinivasan, P., Menczer, F., Pant, G.: A general evaluation framework for topical crawlers. *Information Retrieval* 8(3), 417–447 (2005)