# Where is my URI?

Andre Valdestilhas[1]    Tommaso Soru[1]    Markus Nentwig[2]    Edgard
Marx[3]    Muhammad Saleem[1]    Axel-Cyrille Ngonga Ngomo[1,4]

[1]AKSW Group, University of Leipzig, Germany

[2]Database Group, University of Leipzig, Germany

[3]Leipzig University of Applied Sciences, Germany
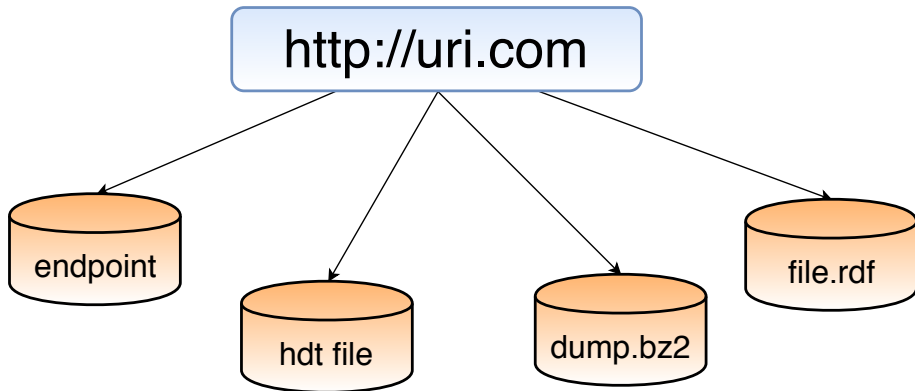
[4]Data Science Group, Paderborn University, Germany

May 14, 2018

# Outline

- Motivation
- Approach
- Use cases
- Experimental setup
- Statistics about Datasets
- Conclusions and Future work

**Where is my URI?**

http://uri.com

endpoint

hdt file

dump.bz2

file.rdf

# Motivation

## Data from the URI

- Problems to obtain data from a URI.

# Motivation

## Data from the URI

- Problems to obtain data from a URI.

## URI dataset

- Which dataset the URI was defined?

# Motivation

Example:



**Obvious** → http://dbpedia.org/resource/Leipzig → DBpedia

**Not so obvious** → http://citeseer.rkbexplorer.com/id/resource-CS116606 → km.aifb.kit.edu
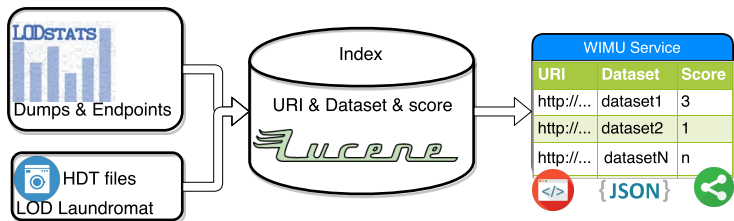
**No dataset** → http://quotationsbook.com/author/59 →

## Goal

Indexing URIs and their use in order to let Linked Data consumers find the respective RDF data source.

- Rank the datasets proportionally to the number of literals.



Keeping the provenance of the URI.

1. Retrieve list of datasets from sources.

2. Retrieve data from datasets.

```
SELECT ?s (count (?o) as ?c)
         WHERE {
?s [] ?o . FILTER(isliteral (?o))
         } GROUP BY ?s
```

3. Build three indexes from dump files, endpoints, and HDT files.

4. Merge the indexes into one.

5. Make the index available and browsable via a web application and an API service.

**LOD**STATS

Dumps & Endpoints

HDT files

LOD Laundromat

Index

URI & Dataset & score

*Lucene*

| WIMU Service | | |
|---|---|---|
| URI | Dataset | Score |
| http://... | dataset1 | 3 |
| http://... | dataset2 | 1 |
| http://... | datasetN | n |

</>  {JSON}

# Approach
The power of Literals

### The heuristic

- The raw data that can disambiguate a URI node in the most direct way.
- The semantic web architecture expects that datasets reusing a URI only refer to it without defining more literal values.

## Where is my URI?

| | | | | | |
|---|---|---|---|---|---|
| http://dbpedia.org/resource/Leipzig | | | 5 | | Find Datasets |

| Dataset | Literals | HDT | Original file | CBD |
|---|---|---|---|---|
| http://downloads.dbpedia.org/2016-10/core-i18n/en/infobox_properties_en.ttl.bz2 | 236 | - | | |
| http://gaia.infor.uva.es/hdt/dbpedia2015.hdt.gz | 165 | - | | |
| http://km.aifb.kit.edu/projects/btc-2009/btc-2009-chunk-061.gz | 142 | | | |
| http://data.dws.informatik.uni-mannheim.de/dbpedia/2014/en/infobox_properties_unredirected_en.nq.bz2 | 124 | | | |
| http://data.dws.informatik.uni-mannheim.de/dbpedia/2014/en/infobox_properties_en.nq.bz2 | 124 | | | |

```
curl "http://wimu.aksw.org/Find?top=5&uri=http://dbpedia.org/resource/Leipzig"

[{"dataset":"http://downloads.dbpedia.org/2016-10/core-i18n/en/infobox_properties_en.ttl.bz2","CountLiteral":"236"},
{"dataset":"http://gaia.infor.uva.es/hdt/dbpedia2015.hdt.gz","CountLiteral":"165"},
{"dataset":"http://download.lodlaundromat.org/a9dabf348fd6262edfbbcf7256b0f839?type=hdt","CountLiteral":"142"},
{"dataset":"http://download.lodlaundromat.org/dde1dcc095b38a1b65ebfbc7696d7998?type=hdt","CountLiteral":"124"}]
```

**Service**: `https://wimu.aksw.org/Find`

| Parameter | Default | Description |
|-----------|---------|-------------|
| top | 0 | Top ocurrences of the datasets. |
| uri | - | URI expected to search. |
| link | - | URL from a linkset. |
| cbd | - | The URI that will origin the CBD |
| ds | - | URL to download the dataset. |

# Approach
Examples

## Single URI, top 5 datasets

```
https://wimu.aksw.org/Find?top=5&uri=http:
//dbpedia.org/resource/Leipzig
```

## Linkset

```
https://wimu.aksw.org/Find?link=http://www.linklion.org/
download/mapping/sws.geonames.org---purl.org.nt
```

## CBD

```
https://wimu.aksw.org/Find?cbd=http:
//dbpedia.org/resource/Leipzig&ds=http://download.
lodlaundromat.org/a9dabf348fd6262edfbbcf7256b0f839?type=hdt
```

# Approach

## Relevant points

- Rank the datasets from LODStats and LODLaundromat using a score function.
- is able to process linksets (more than one URI per request).



| Input |
|-------|
| URI |
| Linkset file |

| Output | | |
|--------|--------|------|
| **URI** | **Dataset** | **Rank** |
| http://... | dataset1 | 3 |
| http://... | dataset2 | 1 |

{JSON}

## Data quality in Link Repositories

- Regenerate mappings using the CBDs to reapply link discovery algorithms in order to validate the mappings.
- Part of LinkLion 2.0

### Federated Query Processing

- Query planning and Source (dataset) selection.
- WIMU will find relevant sources against the individual triple patterns of a given SPARQL query.
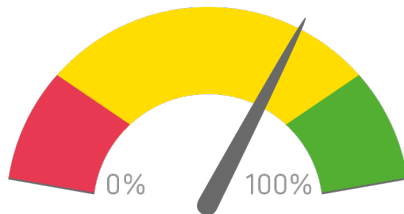


```
SELECT ?v1 ?v2 WHERE {
?uri <p1> ?v1. // Triple Pattern 1 (TP1)
?uri <p2> ?v2. // Triple Pattern 2 (TP2)
}
```

Source selection

```
WIMU source selection
TP1 = {S1,15}, {S2,9}, {S3,7}
TP2 = {S1,20}, {S2,11}, {S4,10}
Total selected sources = 6
```

S1   S2   S3   S4

RDF   RDF   RDF   RDF

# Experimental Setup

- Creation of the index: 3 days and 7 hours.
- Hardware: 64 CPU cores, 126 GB RAM, 2 TB hard disk.

## Statistics about Datasets
Not all datasets are ready to use

|  | LOD Laundromat | LODStats | Total |
|---|---|---|---|
| **URIs indexed** | 4,185,133,445 | 31,121,342 | 4,216,254,787 |
| **Datasets checked** | 658,206 | 9,960 | 668,166 |
| **Triples processed** | 19,891,702,202 | 38,606,408,854 | 58,498,111,056 |

### LODStats

- 60% offline, 14% empty.
- 8% triples with literals as objects are blank nodes.
- 35% online datasets present some error using Jena.
- 69.8% datasets with parser errors.
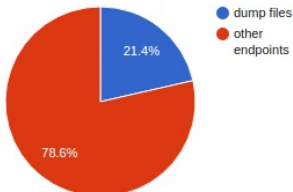
### LODLaundromat

- 2.3% parsing errors. 99% indexed by WIMU.
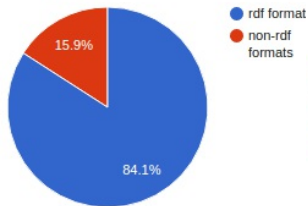
# Statistics about Datasets
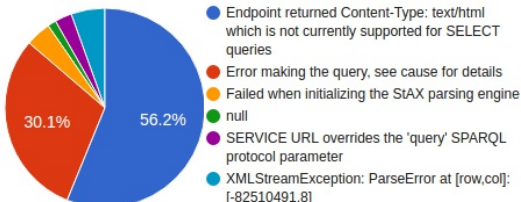


dumps by file extension

endpoints by types

dumps by format

errors JENA parsing

## Conclusions and Future work

- A regularly updated database index of more than 660K datasets from LODStats and LODLaundromat.
- An efficient service on the web that inform which dataset most likely defines a URI.
- Various statistics of datasets indexed from LODStats and LODLaundromat.
- **Future work**: Integrate the second version of LINKLION.
- http://www.linklion.org

# Thanks!
# Questions?

Contact e-mail: valdestilhas@informatik.uni-leipzig.de