# Knowledge Graph LAB Assignment Implement SDM UPC

Linhan Wang [*], Min Zhang[*]

{linhan.wang,min.zhang}@estudiantat.upc.edu

In this project we will be using GraphDB to create and manage knowledge graphs, and query and process them. We will use innovative methods to solve these problems. The problem list is divided into four parts: TBOX deninition, ABOX deninition, Create the final ontology, Querying the ontology.

## 1 B.1 Ontology creation

We used the same datasets(Scorpus) in the first lab and built our TBOX using grafo. And the visual representation of the TBOX can be seen in Figure 1, which was drawn by grafo.[1]

In the graph we defined superclasses and subclasses, which can demonstrate the ontology clearly. And in order to define the hierarchy, we made an assumption that reviewers are authors, so they can share some common properties with authors. So we defined superclasses as follows: Person is the superclass of Author, while Author is the superclass of Reviewer. Proceeding is the superclass of Proceeding_conference and Proceeding_workshop.

In this TBOX definition, there are three kinds of owl: owl:Class, owl:DatatypeProperty, owl:ObjectProperty. For each type, details are shown in Table 1.

## 2 B.2 ABOX deninition(programmatically)

We built the ABOX using Python's library RDFLib, and the ABOX was build in 2 steps: (1) Converting each row in csv file into RDF triples(Subject, predicate and object), (2) Serialize and output the ABOX into file.

The triples for each property and class are shown in Table 2.

When each csv is transformed into RDF triples, ABOX is constructed. We demonstrated this process by the example of Authors data conversion.

| | |
|---|---|
| owl:ObjectProperty | cites |
| | collection_of_conference |
| | collection_of_workshop |
| | has_keyword |
| | paper_belong_to_journal |
| | published_in_proceeding |
| | reviewed_by |
| | written_by |
| owl:DatatypeProperty | author_id |
| | author_name |
| | doi |
| | journal_name |
| | keyword |
| | title |
| | ...... |
| owl:Class | Author |
| | Conference |
| | Journal |
| | Keywords |
| | Paper |
| | ...... |

Table 1: TBOX Definition

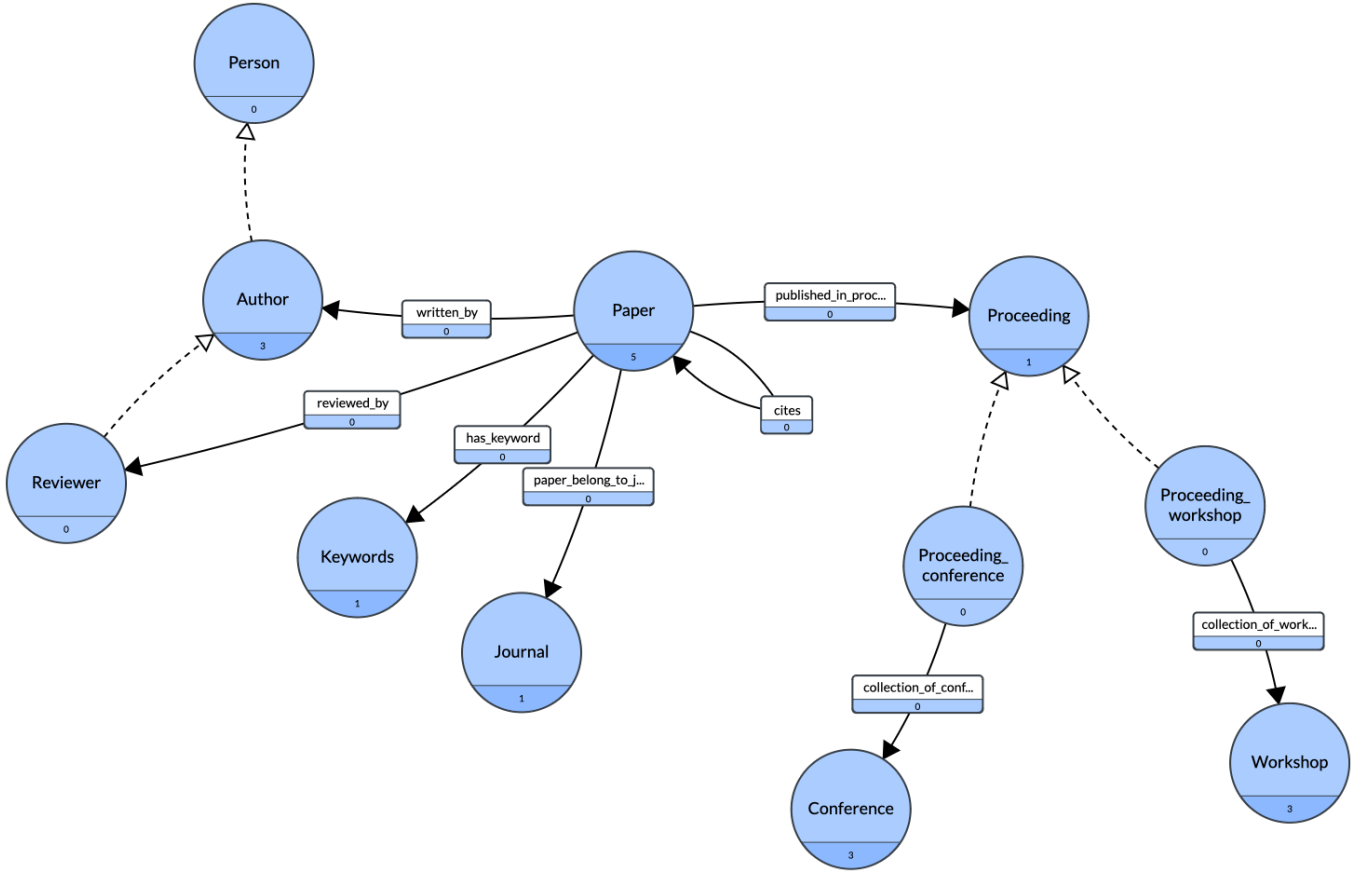[*]Universitat Politècnica de Catalunya (UPC), Erasmus Mundus Big Data Management and Analytics (BDMA)

Figure 1: Comprehensive view of the Knowledge Graphs.

| semantic | subject | predicate | object |
|---|---|---|---|
| cites | Paper | cites | Paper |
| has_keyword | Paper | has_keyword | Keyword |
| reviewed_by | Paper | reviewed_by | Reviewer |
| written_by | Paper | written_by | Author |
| ...... | ...... | ...... | ...... |
| author_id | Author | author_id_type | author_id |
| doi | Paper | doi_type | doi |
| keyword | Keyword | keyword_type | keyword |
| title | Paper | title_type | title |
| ...... | ...... | ...... | ...... |
| Author | Author | RDF.type | author_type |
| Journal | Journal | RDF.type | journal_type |
| Paper | Paper | RDF.type | paper_type |
| ...... | ...... | ...... | ...... |

Table 2: RDF triples

## 2.1 CSV to RDF Conversion Code for Authors

We used rdflib[2] to implement this conversion. The following Python code snippet reads data from a CSV file (`authors.csv`) and converts it into RDF triples, which are then added to an RDF graph. Each author in the CSV file is represented as a resource with properties such as type, name, ID, and affiliation.

```
1  # authors.csv
2  csv_data = []
3  with open("data/authors.csv", newline='') as csvfile
       :
4      reader = csv.DictReader(csvfile)
5      for row in reader:
6          csv_data.append(row)
7
8  for row in csv_data:
9      # author_type
10     author_uri = URIRef(EX[row["author_id"]])
11     author_type = URIRef("http://www.gra.fo/schema/
       untitled-ekg#Author")
12     abox_graph.add((author_uri, RDF.type,
       author_type))
13
14     # author_name
```

2

```python
15    author_name = URIRef(urllib.parse.quote(row["
      author_name"]))
16    author_name_type = URIRef("http://www.gra.fo/
      schema/untitled-ekg#author_name")
17    abox_graph.add((author_uri, author_name_type,
      author_name))
18
19    # author_id
20    author_id = URIRef(row["author_id"])
21    author_id_type = URIRef("http://www.gra.fo/
      schema/untitled-ekg#author_id")
22    abox_graph.add((author_uri, author_id_type,
      author_id))
23
24    # author_affiliation
25    author_aff = URIRef(urllib.parse.quote(row["
      author_affiliation"]))
26    author_aff_type = URIRef("http://www.gra.fo/
      schema/untitled-ekg#author_affiliation")
27    abox_graph.add((author_uri, author_aff_type,
      author_aff))
```

Listing 1: CSV to RDF Conversion Code for Authors

The code performs the following steps:

1. **Reading the CSV File**: Opens the `authors.csv` file and reads its contents using a dictionary reader.

2. **Processing Each Row**: For each row, it creates RDF triples for the author's type, name, ID, and affiliation.

3. **Defining the Author Resource**: Converts the author ID into a URI and adds a triple to specify that it is of type `Author`.

4. **Adding the Author Name**: Encodes the author name as a URI and links it to the author resource.

5. **Adding the Author ID**: Directly links the author ID to the author resource.

6. **Adding the Author Affiliation**: Encodes the author affiliation as a URI and links it to the author resource.

7. **Adding RDF Triples to the Graph**: Add each triple to the RDF graph.

## 2.2 ABOX Step 2 - Output ABOX

After converting all the csv files into rdf, using serialize function to output ABOX in turtle format.

## 3 B.3 Create the final ontology

After generating the TBOX and ABOX, we created the GraphDB repository and imported both files into it. When importing .owl

| Label | Count |
|---|---|
| Number of Classes | 20 |
| Number of Main Properties | 23 |
| Number of Instances of Main Classes | 2719 |
| Number of Triples | 31011 |

Table 3: Summary of Knowledge Graph Statistics

| Class | Instances | Triples |
|---|---|---|
| Paper | 483 | 19896 |
| Author | 1942 | 10819 |
| Conference | 178 | 178 |
| Journal | 116 | 118 |

Table 4: Statistics for individual main Classes

file, TBOX is imported; when importing .ttl file, ABOX is imported.

We configured the ruleset of our repository to RDFS (Optimized), as shown in Figure 2. According to GraphDB, this ruleset supports the standard model-theoretic semantics of RDFS.



Figure 2: Repository-Setting

From Table3 and Table4, we can see the statistics information after of GraphDB.

From Figure 3 and Figure 4, we can see the overview of repository and graph in GraphDB. And from Figure 5, we can see the hierarchy of the GraphDB.

## 4 B.4 Querying

In this part, we listed all the 6 queries, as well as the screenshot of query results. Moreover, we also downloaded all the query results as CSV files in the query_results folder attached to our submission.
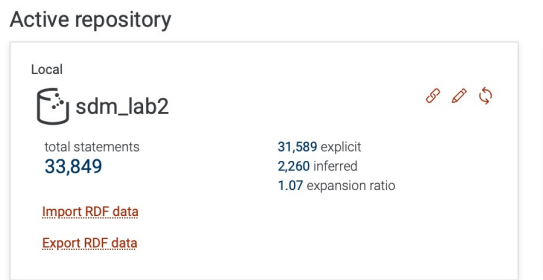
Figure 3: Repository Overview

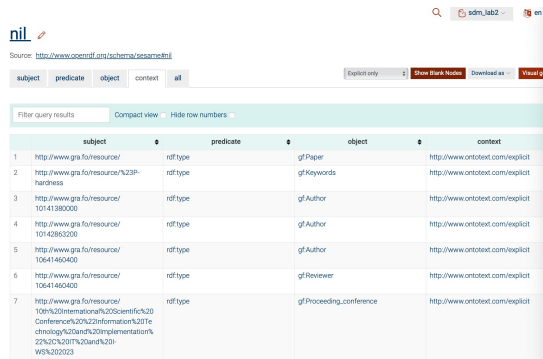

Figure 6: Class relationship
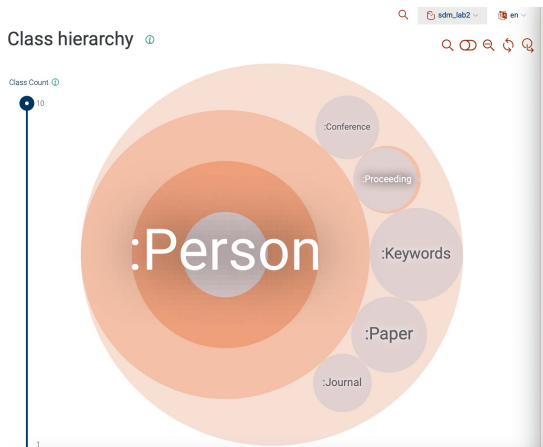


Figure 4: Graph Overview

```
8 }
```

Listing 2: SPARQL Query to Select Distinct Authors and Their Names



Figure 5: Class hierarchy



Figure 7: Result of Finding all Authors

## 4.1   1. Find all Authors.

The following Cypher query finds the top 3 most cited papers of each conference:

```
1 PREFIX gf: <http://www.gra.fo/schema/untitled-ekg#>
2 PREFIX ex: <http://www.gra.fo/resource/>
3
4 SELECT DISTINCT ?author ?authorName
5 WHERE {
6   ?author a gf:Author .
7   ?author gf:author_name ?authorName .
```

## 4.2   2. Find all properties whose domain is Author.

```
1 PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
2 PREFIX gf: <http://www.gra.fo/schema/untitled-ekg#>
3
4 SELECT DISTINCT ?property
5 WHERE {
6   ?property rdfs:domain gf:Author .
7 }
```

Listing 3: SPARQL Query to Select Distinct Properties of Authors

## 4.3   3. Find all properties whose domain is either Conference or Journal.

4

Figure 8: Result of Finding all properties whose domain is Author

```
1  PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
2  PREFIX gf: <http://www.gra.fo/schema/untitled-ekg#>
3
4  SELECT DISTINCT ?property
5  WHERE {
6    {
7      ?property rdfs:domain gf:Conference .
8    }
9    UNION
10   {
11     ?property rdfs:domain gf:Journal .
12   }
13 }
```



Figure 9: Result of Finding all properties whose domain is either Conference or Journal.

## 4.4  4. Find all the papers written by a given author - Amarilli 20A that where published in database conferences.

For this query, our query design is as follows: Firstly, we identify all papers in the database field by searching for all papers from conferences with the keyword "database". Secondly, regarding the specified author condition, we identify authors

from the previous step who have published more than two papers, and then sort these authors by the number of papers published, identifying the author with the highest number of publications. The third step is to find all papers published by the specified author (i.e., the author with the most publications in the database field) at database conferences.

```
1  PREFIX schema: <http://www.gra.fo/schema/untitled-
       ekg#>
2  PREFIX gf: <http://www.gra.fo/resource/>
3
4  SELECT ?paper ?paper_title
5  WHERE {
6    ?conference rdf:type schema:Conference .
7    FILTER regex(str(?conference), "Database", "i")
8    ?proceeding schema:collection_of_conference ?
       conference ;
9             rdf:type schema:Proceeding_conference
       .
10   ?paper schema:published_in_proceeding ?proceeding
       ;
11        rdf:type schema:Paper ;
12        schema:written_by ?author ;
13        schema:title ?paper_title .
14   ?author schema:author_name <https://gra.fo/
       Amarilli%20A.> .
15 }
```



Figure 10: Result of Finding all the papers written by author Amarilli 20A. that where published in database conferences.

## 4.5  5. Find the author that having the most citations in Computational conferences.

This query is designed to identify the author with the highest number of citations among those who have presented papers at conferences related to "Computational." The query first filters out conferences whose names contain "Computational," then

finds the proceedings of these conferences and the papers published in them. From these papers, it extracts the authors and their names, counts the citations for each author, sorts the results in descending order of citation count, and finally returns the author with the most citations.

```
1  PREFIX schema: <http://www.gra.fo/schema/untitled-
        ekg#>
2
3  SELECT ?author_name (COUNT(?citation) AS ?
        citation_count)
4  WHERE {
5    ?conference rdf:type schema:Conference .
6    FILTER regex(str(?conference), "Computational", "i
        ")
7    ?proceeding schema:collection_of_conference ?
        conference ;
8              rdf:type schema:Proceeding_conference .
9    ?paper schema:published_in_proceeding ?proceeding
        ;
10         rdf:type schema:Paper ;
11         schema:written_by ?author .
12   ?author schema:author_name ?author_name .
13   OPTIONAL {
14     ?paper schema:cites ?citation .
15   }
16 }
17 GROUP BY ?author_name
18 ORDER BY DESC(?citation_count)
19 LIMIT 1
```
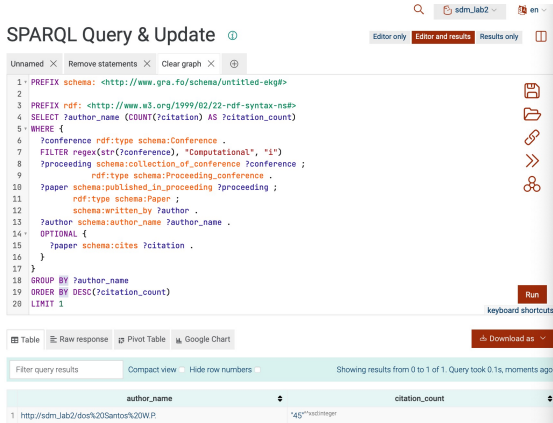


Figure 11: Result of Finding the author that having the most citations in Computational conferences.

## 4.6  6. Find the author that published most papers in journals.

This query is to identify the author who has published the most papers in journals. By matching journals and their related papers, the query retrieves each author's name and counts the number of papers they have published, ultimately returning the author with the highest paper count.

```
1  PREFIX schema: <http://www.gra.fo/schema/untitled-
        ekg#>
2  PREFIX gf: <http://www.gra.fo/resource/>
3
4  SELECT ?author_name (COUNT(?paper) AS ?paper_count)
5  WHERE {
6    ?journal rdf:type schema:Journal .
7    ?paper schema:paper_belong_to_journal ?journal ;
8          rdf:type schema:Paper ;
9          schema:written_by ?author .
10   ?author schema:author_name ?author_name .
11 }
12 GROUP BY ?author_name
13 ORDER BY DESC(?paper_count)
14 LIMIT 1
```
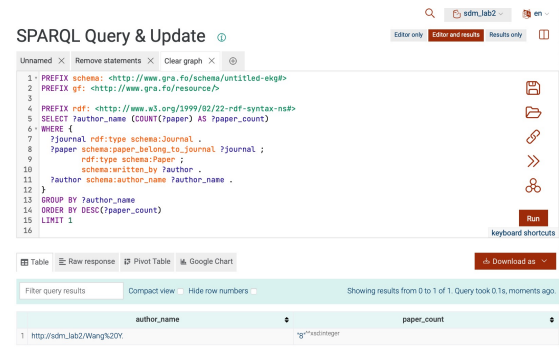


Figure 12: Result of Finding the author that published most papers in journals.

# References

[1] grafo. Design Knowledge Graphs, Simply, 2024.

[2] rdflib. rdflib 7.0.0 Documentation, 2024.