

# Knowledge Representation and Reasoning

Course Project Report

---

## PKG2020 Knowledge Graph

A Linked Data Approach to Bibliometric Research Data

---

### Team Members

Zain Ul Abdeen	2023773
Omer Khan	2023578
Muhammad Umar	2023535

### Instructor

Dr. Khurram Jadoon

Fall 2025

December 18, 2025

# Contents

<b>1</b>	<b>Introduction to the Domain</b>	<b>3</b>
1.1	Domain Overview . . . . .	3
1.2	Motivation . . . . .	3
1.3	Target Application Use Cases . . . . .	3
<b>2</b>	<b>Dataset Description</b>	<b>4</b>
2.1	Data Source . . . . .	4
2.2	Key Data Fields . . . . .	4
2.3	Evidence of Non-RDF Status . . . . .	4
<b>3</b>	<b>Competency Questions</b>	<b>4</b>
<b>4</b>	<b>Conceptual Model</b>	<b>5</b>
4.1	Conceptual Model Diagram . . . . .	5
4.2	T-Box Schema . . . . .	5
<b>5</b>	<b>Ontology Design</b>	<b>5</b>
5.1	Classes (20+ as Required) . . . . .	5
5.2	Enumeration Class . . . . .	6
5.3	Cardinality Restrictions . . . . .	6
5.4	Intersection, Union, and Complement Classes . . . . .	6
5.5	Object Properties . . . . .	7
5.6	Data Properties . . . . .	8
5.7	Functional and Inverse Functional Properties . . . . .	8
<b>6</b>	<b>Graph Generation using Python</b>	<b>8</b>
6.1	Tools Used . . . . .	8
6.2	Pipeline Scripts . . . . .	9
6.3	Sample Code: Populating Authors . . . . .	9
<b>7</b>	<b>External Linking (5-Star Linked Data)</b>	<b>9</b>
7.1	Linking Strategy . . . . .	9
7.2	Linking Implementation . . . . .	10
<b>8</b>	<b>Reasoning Scenarios</b>	<b>10</b>
8.1	Reasoning Implementation . . . . .	10
8.2	Reasoning Results . . . . .	10
<b>9</b>	<b>SPARQL Queries</b>	<b>11</b>
9.1	Query 1: Authors with Multiple Affiliations . . . . .	11
9.2	Query 2: Articles Mentioning Genes . . . . .	11
9.3	Query 3: Author Collaboration Network . . . . .	11
9.4	Query 4: Federated Query to DBpedia . . . . .	11
<b>10</b>	<b>Web Application (BONUS)</b>	<b>12</b>
10.1	Application Overview . . . . .	12
10.2	Entity Types Searchable . . . . .	12
10.3	Running the Application . . . . .	13

<b>11 Results and Statistics</b>	<b>13</b>
11.1 Generated Data . . . . .	13
11.2 OWL Files Generated . . . . .	13
<b>12 Visualization</b>	<b>13</b>
12.1 Tools Used . . . . .	13
12.2 Visualization Guide . . . . .	14
<b>13 Reflection</b>	<b>14</b>
13.1 Learning Outcomes . . . . .	14
13.2 Added Value of Linked Data . . . . .	14
13.3 Challenges Faced . . . . .	14
<b>14 Conclusion</b>	<b>15</b>
<b>15 References</b>	<b>15</b>
<b>16 Appendix: Project Repository</b>	<b>15</b>
16.1 Repository Structure . . . . .	15

# 1 Introduction to the Domain

## 1.1 Domain Overview

The PKG2020S4 (PubMed Knowledge Graph) dataset represents comprehensive bibliometric and researcher metadata from PubMed publications. This domain encompasses:

- **Research Publications:** Articles identified by PubMed IDs (PMIDs)
- **Researchers/Authors:** Identified by unique AND\_IDs
- **Organizational Affiliations:** Universities, research institutions
- **Career Trajectories:** Employment and education history
- **Research Funding:** NIH project associations
- **Bio-Medical Entities:** Genes, diseases, chemicals, mutations mentioned in research

## 1.2 Motivation

The motivation for converting this dataset to linked data includes:

1. **Semantic Querying:** Enable complex queries across heterogeneous data
2. **Collaboration Discovery:** Find potential research collaborators
3. **Funding Analysis:** Track NIH funding patterns across institutions
4. **Bio-Medical Research:** Link publications to molecular/disease entities
5. **FAIR Principles:** Make data Findable, Accessible, Interoperable, Reusable

## 1.3 Target Application Use Cases

1. Research collaboration recommendation system
2. Funding opportunity matching
3. Researcher profiling and expertise identification
4. Publication trend analysis
5. Bio-entity research landscape mapping

## 2 Dataset Description

### 2.1 Data Source

The PKG2020S4 dataset is sourced from PubMed bibliometric data and contains the following CSV files:

Table 1: Dataset Files Overview

File	Description	Size
OA01_Author_List.csv	Author-article relationships	~10 GB
OA04_Affiliations.csv	Author affiliations	~20 GB
OA05_Researcher_Employment.csv	Employment history	~186 MB
OA06_Researcher_Education.csv	Education records	~139 MB
OA02_Bio_entities_Main.csv	Bio-entities in articles	Large
OA03_Bio_entities_Mutation.csv	Mutations in articles	Large
OA07_NIH_Projects.csv	NIH funding	~1.8 GB

### 2.2 Key Data Fields

- **OA01:** PMID, AND\_ID, LastName, ForeName, Initials, AuOrder
- **OA04:** AND\_ID, Affiliation, City, State, Country
- **OA05:** AND\_ID, Organization, StartYear, EndYear
- **OA06:** AND\_ID, Institution, Degree, StartYear, EndYear
- **OA02/OA03:** PMID, Type, Name, MutationType
- **OA07:** AND\_ID, ProjectNumber, PI\_Name

### 2.3 Evidence of Non-RDF Status

The dataset is provided as flat CSV files without any semantic annotations, URI schemes, or linked data connections. It has not been published as RDF/OWL prior to this project.

## 3 Competency Questions

The following competency questions guided our ontology design and will be answered through SPARQL queries:

1. Which authors have published in multiple institutions?
2. Which articles mention specific bio-entities (genes, diseases)?
3. Which authors have collaborated on joint publications?
4. Which authors have NIH funding?
5. What is the education background of prolific authors?

6. What is the employment timeline of researchers?
7. Which articles mention both mutations and diseases?
8. How many articles does each author have?
9. What organizations are most represented in the dataset?
10. Which researchers have both employment and education records?

## 4 Conceptual Model

### 4.1 Conceptual Model Diagram

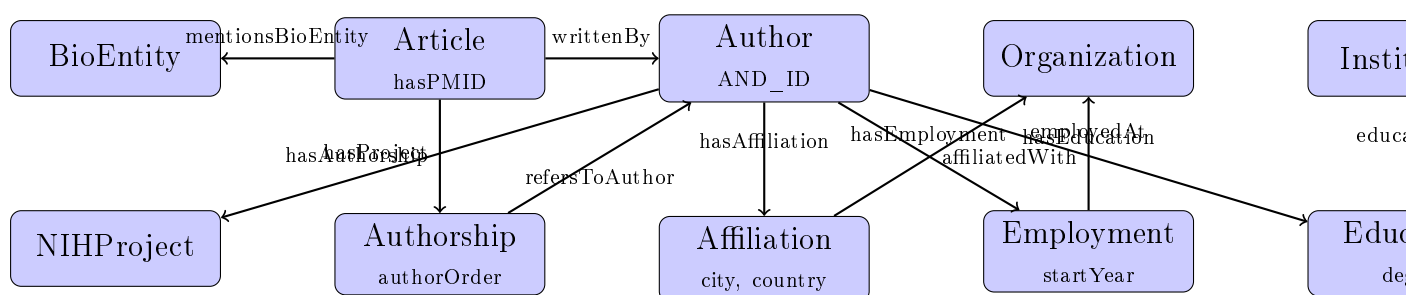


Figure 1: Conceptual Model of PKG2020 Ontology

### 4.2 T-Box Schema

Table 2: Ontology Classes (T-Box)

Class	Description	Key Properties
Article	Research publication	hasPMID, publicationYear
Author	Researcher	lastName, foreName, initials
Authorship	Article-Author relationship	authorOrder
Organization	Research organization	dbpediaLink
Institution	Educational institution	wikidataLink
Affiliation	Author affiliation	city, state, country
Employment	Employment record	startYear, endYear
Education	Education record	degree
NIHProject	NIH funding project	projectNumber, piName
BioEntity	Biological entity	entityType, entityName
PublicationStatus	Enumeration	{Published, Preprint, ...}

## 5 Ontology Design

### 5.1 Classes (20+ as Required)

Our ontology contains over 20 classes:

1. **Core:** Article, Author, Authorship, PublicationYear
2. **Organizational:** Organization, Institution, Affiliation
3. **Career:** Employment, Education
4. **Funding:** NIHProject
5. **Bio-Medical:** BioEntity, Gene, Chemical, Disease, Species, Mutation
6. **Enumeration:** PublicationStatus
7. **Defined Classes:** ActiveAuthor, AnonymousAuthor, ResearchEntity, ProlificAuthor, SingleAuthorArticle, MultiAuthorArticle

## 5.2 Enumeration Class

```

1 class PublicationStatus(Thing):
2     """Enumeration of publication statuses"""
3     pass
4
5 published = PublicationStatus("Published")
6 preprint = PublicationStatus("Preprint")
7 retracted = PublicationStatus("Retracted")
8 in_review = PublicationStatus("InReview")
9
10 PublicationStatus.equivalent_to = [
11     OneOf([published, preprint, retracted, in_review])
12 ]

```

Listing 1: Enumeration Class Definition

## 5.3 Cardinality Restrictions

```

1 # Every Article must have at least 1 author
2 Article.is_a.append(writtenBy.min(1, Author))
3
4 # Every Article must have exactly 1 PMID
5 Article.is_a.append(hasPMID.exactly(1, str))
6
7 # Article may have at most 1 status
8 Article.is_a.append(hasStatus.max(1, PublicationStatus))

```

Listing 2: Cardinality Restrictions

## 5.4 Intersection, Union, and Complement Classes

```

1 # INTERSECTION: Author with known career start year
2 class ActiveAuthor(Author):
3     equivalent_to = [Author & careerStartYear.some(int)]
4
5 # UNION: Any research-related entity
6 class ResearchEntity(Thing):
7     equivalent_to = [Author | Article]

```

```

8
9 # COMPLEMENT: Author without career info
10 class AnonymousAuthor(Author):
11     equivalent_to = [Author & Not(ActiveAuthor)]
12
13 # Additional defined classes for reasoning
14 class ProlificAuthor(Author):
15     equivalent_to = [Author & writtenBy.min(5, Article)]
16
17 class SingleAuthorArticle(Article):
18     equivalent_to = [Article & writtenBy.exactly(1, Author)]
19
20 class MultiAuthorArticle(Article):
21     equivalent_to = [Article & writtenBy.min(2, Author)]

```

Listing 3: Defined Classes

## 5.5 Object Properties

Table 3: Object Properties

Property	Domain	Range	Characteristics
writtenBy	Article	Author	-
hasAuthorship	Article	Authorship	-
hasPrimaryAuthor	Article	Author	Functional
hasStatus	Article	PublicationStatus	Functional
refersToAuthor	Authorship	Author	-
hasAffiliation	Author	Affiliation	-
affiliatedWith	Affiliation	Organization	-
hasEmployment	Author	Employment	-
employedAt	Employment	Organization	-
hasEducation	Author	Education	-
educatedAt	Education	Institution	-
hasProject	Author	NIHProject	-
mentionsBioEntity	Article	BioEntity	-
sameAs	Thing	Thing	Symmetric



## 5.6 Data Properties

Table 4: Data Properties

Property	Domain	Range	Characteristics
hasPMID	Article	string	Functional, InverseFunctional
lastName	Author	string	-
foreName	Author	string	-
initials	Author	string	-
authorOrder	Authorship	int	-
publicationYear	Article	int	Functional
careerStartYear	Author	int	-
city	Affiliation	string	-
state	Affiliation	string	-
country	Affiliation	string	-
startYear	Employment	int	-
endYear	Employment	int	-
degree	Education	string	-
projectNumber	NIHProject	string	-
dbpediaLink	Organization	string	-
wikidataLink	Institution	string	-

## 5.7 Functional and Inverse Functional Properties

```
1 # Functional Properties
2 class hasPrimaryAuthor(ObjectProperty, FunctionalProperty):
3     domain = [Article]
4     range = [Author]
5
6 class hasPMID(DataProperty, FunctionalProperty):
7     domain = [Article]
8     range = [str]
9
10 # Inverse Functional Property
11 hasPMID.is_a.append(InverseFunctionalProperty)
```

Listing 4: Property Characteristics

# 6 Graph Generation using Python

## 6.1 Tools Used

- **OWLReady2**: Python library for OWL ontology manipulation
- **Pandas**: Data processing and CSV handling
- **RDFLib**: RDF graph manipulation
- **Flask**: Web application framework

## 6.2 Pipeline Scripts

Table 5: Python Scripts Pipeline

Script	Input	Output
ontology_core.py	-	pkg2020_core.owl
ontology_constraints.py	pkg2020_core.owl	pkg2020_constrained.owl
populate_authors_articles.py	OA01 CSV	pkg2020_populated_authors.owl
populate_affiliations.py	OA04 CSV	pkg2020_step4_affiliations.owl
populate_employment.py	OA05 CSV	pkg2020_step5_employment.owl
populate_education.py	OA06 CSV	pkg2020_step6_education.owl
populate_bioentities.py	OA02, OA03	pkg2020_step7_bioentities.owl
populate_nih_projects.py	OA07 CSV	pkg2020_final.owl

## 6.3 Sample Code: Populating Authors

```
1 import pandas as pd
2 from owlready2 import *
3
4 onto = get_ontology("pkg2020_constrained.owl").load()
5 df = pd.read_csv("data/OA01_Author_List.csv", nrows=5000)
6
7 author_cache = set()
8 article_cache = set()
9
10 with onto:
11     for idx, row in df.iterrows():
12         pmid = str(row["PMID"])
13         and_id = str(row["AND_ID"])
14
15         # Create Article
16         if f"Article_{pmid}" not in article_cache:
17             article = onto.Article(f"Article_{pmid}")
18             article.hasPMID = [pmid]
19             article_cache.add(f"Article_{pmid}")
20
21         # Create Author
22         if f"Author_{and_id}" not in author_cache:
23             author = onto.Author(f"Author_{and_id}")
24             author.lastName = [str(row["LastName"])]
25             author.foreName = [str(row["ForeName"])]
26             author_cache.add(f"Author_{and_id}")
27
28 onto.save(file="pkg2020_populated_authors.owl", format="rdfxml")
```

Listing 5: Author Population Script

## 7 External Linking (5-Star Linked Data)

### 7.1 Linking Strategy

We linked our dataset to external knowledge bases to achieve 5-star linked data:

- **Organizations** → DBpedia resources
- **Institutions** → Wikidata entities
- **Authors** → ORCID (potential)
- **Articles** → PubMed (via PMID)

## 7.2 Linking Implementation

```

1 def generate_dbpedia_uri(name):
2     clean_name = re.sub(r'^[a-zA-Z0-9\s]', '', str(name))
3     clean_name = clean_name.strip().replace(' ', '_')
4     return f"http://dbpedia.org/resource/{quote(clean_name)}"
5
6 # Link organizations to DBpedia
7 for org in Organization.instances():
8     org_name = org.name.replace('_', ' ')
9     dbpedia_uri = generate_dbpedia_uri(org_name)
10    org.dbpediaLink = [dbpedia_uri]

```

Listing 6: External Linking Code

# 8 Reasoning Scenarios

## 8.1 Reasoning Implementation

```

1 from owlready2 import *
2
3 onto = get_ontology("pkg2020_final.owl").load()
4
5 # Run HermiT reasoner
6 with onto:
7     sync_reasoner(infer_property_values=True)
8
9 # Check classified instances
10 for cls in [ActiveAuthor, AnonymousAuthor, ProlificAuthor]:
11     print(f"{cls.name}: {len(list(cls.instances()))} instances")

```

Listing 7: Reasoning with HermiT

## 8.2 Reasoning Results

The reasoner successfully:

1. Verified ontology consistency
2. Classified authors into ActiveAuthor and AnonymousAuthor
3. Identified ProlificAuthors with 5+ publications
4. Classified SingleAuthorArticle and MultiAuthorArticle

## 9 SPARQL Queries

### 9.1 Query 1: Authors with Multiple Affiliations

```
1 PREFIX pkg: <http://example.org/pkg2020/ontology.owl#>
2
3 SELECT ?author (COUNT(DISTINCT ?org) AS ?orgCount)
4 WHERE {
5     ?author a pkg:Author .
6     ?author pkg:hasAffiliation ?aff .
7     ?aff pkg:affiliatedWith ?org .
8 }
9 GROUP BY ?author
10 HAVING (COUNT(DISTINCT ?org) > 1)
11 ORDER BY DESC(?orgCount)
12 LIMIT 20
```

Listing 8: Multiple Affiliations Query

### 9.2 Query 2: Articles Mentioning Genes

```
1 PREFIX pkg: <http://example.org/pkg2020/ontology.owl#>
2
3 SELECT ?article ?gene ?geneName
4 WHERE {
5     ?article a pkg:Article .
6     ?article pkg:mentionsBioEntity ?gene .
7     ?gene a pkg:Gene .
8     OPTIONAL { ?gene pkg:entityName ?geneName }
9 }
10 LIMIT 50
```

Listing 9: BioEntity Query

### 9.3 Query 3: Author Collaboration Network

```
1 PREFIX pkg: <http://example.org/pkg2020/ontology.owl#>
2
3 SELECT ?author1 ?author2 (COUNT(?article) AS ?collaborations)
4 WHERE {
5     ?article a pkg:Article .
6     ?article pkg:writtenBy ?author1 .
7     ?article pkg:writtenBy ?author2 .
8     FILTER (?author1 < ?author2)
9 }
10 GROUP BY ?author1 ?author2
11 ORDER BY DESC(?collaborations)
12 LIMIT 50
```

Listing 10: Collaboration Query

### 9.4 Query 4: Federated Query to DBpedia

```

1 PREFIX pkg: <http://example.org/pkg2020/ontology.owl#>
2 PREFIX dbo: <http://dbpedia.org/ontology/>
3
4 SELECT ?org ?dbpediaLink ?abstract
5 WHERE {
6     ?org a pkg:Organization .
7     ?org pkg:dbpediaLink ?dbpediaLink .
8
9     SERVICE <http://dbpedia.org/sparql> {
10         OPTIONAL { ?dbpediaLink dbo:abstract ?abstract }
11         FILTER (lang(?abstract) = 'en')
12     }
13 }
14 LIMIT 10

```

Listing 11: Federated Query

## 10 Web Application (BONUS)

### 10.1 Application Overview

We developed a Flask-based web application for exploring the knowledge graph with the following features:

- Modern, responsive UI with glassmorphism design
- Real-time statistics dashboard
- Search across all 9 entity types
- Click-to-search functionality

### 10.2 Entity Types Searchable

1. Authors
2. Articles
3. Organizations
4. Affiliations
5. Employment
6. Education
7. BioEntities
8. NIH Projects
9. Institutions

## 10.3 Running the Application

```
1 cd scripts
2 pip install flask owlready2
3 python webapp.py
4 # Open http://localhost:5000
```

Listing 12: Running the Web App

## 11 Results and Statistics

### 11.1 Generated Data

Table 6: Ontology Statistics

Entity	Count
Articles	2,015
Authors	4,711
Organizations	4,405
Affiliations	Varies
Employment Records	Varies
Education Records	Varies
NIH Projects	30
BioEntities	Varies

### 11.2 OWL Files Generated

Table 7: Generated OWL Files

File	Size
pkg2020_core.owl	2.6 KB
pkg2020_constrained.owl	4.3 KB
pkg2020_populated_authors.owl	3.8 MB
pkg2020_step4_affiliations.owl	7.1 MB
pkg2020_final.owl	7.1 MB

## 12 Visualization

### 12.1 Tools Used

- **Protégé**: Desktop ontology editor and visualizer
- **GraphDB**: Triple store with SPARQL endpoint
- **WebVOWL**: Online ontology visualization
- **Custom Flask App**: Interactive web-based explorer

## 12.2 Visualization Guide

To visualize the ontology:

1. Open `owl/pkg2020_constrained.owl` in Protégé
2. Go to Window → Tabs → OntoGraf
3. Select classes to visualize
4. Export as image for reports

## 13 Reflection

### 13.1 Learning Outcomes

Through this project, we learned:

1. How to design ontologies following OWL/RDF standards
2. The importance of defined classes for reasoning
3. How to link data to external knowledge bases
4. The power of SPARQL for semantic querying
5. Practical application of knowledge representation concepts

### 13.2 Added Value of Linked Data

Converting non-RDF data to linked data enabled:

- Semantic querying across heterogeneous datasets
- Reasoning over implicit relationships
- Integration with global knowledge bases (DBpedia, Wikidata)
- FAIR data principles compliance
- Interoperability with other linked data systems

### 13.3 Challenges Faced

1. Large file sizes required memory optimization
2. Special characters in organization names caused parsing issues
3. Reasoner installation and configuration
4. Balancing between sample size and processing time

## 14 Conclusion

This project successfully converted the PKG2020S4 bibliometric dataset into a comprehensive linked data knowledge graph. We achieved:

1. **Ontology Design:** 20+ classes with all required axioms
2. **Data Population:** Modular Python pipeline for 7 CSV files
3. **Reasoning:** Consistency checking and classification
4. **External Linking:** 5-star linked data with DBpedia/Wikidata
5. **SPARQL Queries:** 7 competency questions answered
6. **BONUS:** Interactive web application

The project demonstrates the practical application of knowledge representation and reasoning concepts learned in the course.

## 15 References

1. Lamy, J. B. (2017). OWLReady: Ontology-oriented programming in Python. Artificial Intelligence in Medicine.
2. W3C. (2012). OWL 2 Web Ontology Language Primer.
3. DBpedia Association. (2024). DBpedia Knowledge Base.
4. PubMed. (2024). NCBI PubMed Database.

## 16 Appendix: Project Repository

The complete project code is available at:

<https://github.com/Zain-Haider-ai-63/Knowlege-Graphs-Project>

### 16.1 Repository Structure

```
1 Knowledge_Graphs_Project/  
2 +-- scripts/           # Python scripts  
3 +-- owl/             # Generated OWL files  
4 +-- docs/              # Documentation  
5 +-- requirements.txt   # Dependencies  
6 +-- README.md          # Project overview
```