

Group A3D: SPARQL Queries and Analytics

Group members: Andrea Bruttomesso 2120933
Alessandro Corrà 2125034
Davide Seghetto 2122548
Andrea Stocco 2108885

Task

Provide at least 8 SPARQL queries over your RDF datasets. You may also perform advanced data analytics to uncover interesting insights from your datasets. Please submit a PDF that includes the SPARQL queries along with relevant plots or tables summarizing your analytics. For each query, provide a description that explains its purpose and overall objective.

1 Relationship between Nobel Prize winning ideas and published studies

```
1 PREFIX spif: <http://spinrdf.org/spif#>
2 PREFIX : <http://www.semanticweb.org/a3d/ontologies/2024/10/nobelOntology/>
3 PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
4
5 SELECT ?nobelTopic ?nobel (count(?paper) AS ?numPapers) WHERE {
6   {
7     SELECT ?paperTopic ?paper WHERE {
8       ?paper :hasAbstractTopics ?topics;
9       :hasYear ?year.
10      FILTER (?year = "2004"^^xsd:gYear)
11      ?paperTopic spif:split(?topics ",")
12    }
13  }
14  {
15    SELECT ?nobelTopic ?nobel WHERE {
16      ?nobel :hasMotivationTopics ?topics;
17      :hasYear ?year.
18      FILTER (?year = "2004"^^xsd:gYear)
19      ?nobelTopic spif:split(?topics ",")
20    }
21  }
22  FILTER (?nobelTopic = ?paperTopic)
23 }
24 GROUP BY ?nobelTopic ?nobel
25 ORDER BY DESC(?numPapers)
```

This query shows the topics present in both Nobel Prize motivations and paper abstracts. For a given year, it returns the number of paper in which a Nobel topic appears. This query can be used to find correlations between Nobel Prize topics and research papers.

Table 1 shows the output of this query on year 2004.

Table 1: Number of papers per Nobel topic in 2004

Nobel topic	Nobel Prize	Number of papers
protein	Chemistry 2004	28
development	Peace 2004	13
flow	Literature 2004	8
interaction	Physics 2004	4
discovery	Chemistry 2004	3
discovery	Physics 2004	3
degradation	Chemistry 2004	3
asymptotic	Physics 2004	2
forces	Economics 2004	1
cycles	Economics 2004	1
olfactory	Medicine 2004	1
organization	Medicine 2004	1

Considering the limited number of papers available, the topic “protein” appeared in 28 papers. The high number of papers mentioning this Nobel topic suggests that it was widely discussed or relevant in 2004.

We cannot conclude whether the research area of molecular biology was particularly active in that year, but in Section 2 we will further investigate this.

Unfortunately, this query is not always useful. In some cases, the main topics may include words like “method” and “analysis”, which are not informative enough to determine how extensively a specific topic was studied in a given year.

Due to the distribution of research papers in our dataset across different years, this query provides more meaningful results for years after 2000.

2 Most active research areas in a year

```

1 PREFIX : <http://www.semanticweb.org/a3d/ontologies/2024/10/nobel0ntology/>
2 PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
3 PREFIX skos: <http://www.w3.org/2004/02/skos/core#>
4
5 SELECT ?category (COUNT(?paper) as ?numPapers) WHERE {
6     ?paper :publishedIn ?venue;
7         :hasYear ?year.
8     ?venue :hasJournalCategory ?category.
9     ?category skos:broaderTransitive ?sub.
10    FILTER (?year = "2004"^^xsd:gYear)
11 }
12 GROUP BY ?category
13 ORDER BY DESC (?numPapers)

```

This query shows the number of papers published for each journal subcategory for a given year. It can be used to identify which research areas were particularly active in that year.

Table 2 continues the analysis started in the previous section.

Table 2: Number of papers for each journal subcategory in 2004

Journal Subcategory	Number of papers
Biochemistry Genetics Molecular Biology	418
Social Sciences	344
Decision Sciences	125
Arts Humanities	74
Business Management Accounting	68
Physics Astronomy	65
Neuroscience	55
Health Professions	34
Psychology	27
Earth Planetary Sciences	22
Economics Econometrics Finance	14
Materials Science	12
Environmental Science	12
Agricultural Biological Sciences	10
Energy	2
Pharmacology Toxicology Pharmaceuticals	2

Considering the limited number of papers in our dataset, molecular biology was the most active research area in 2004.

Building on the previous section, 28 out of 418 molecular biology papers focused “protein”.

This topic held central importance that year, which may explain why a Nobel Prize was awarded for it.

3 papersPerTopic

```

1 PREFIX spif: <http://spinrdf.org/spif#>
2 PREFIX : <http://www.semanticweb.org/a3d/ontologies/2024/10/nobel0ntology/>
3
4 SELECT ?singleTopic (COUNT(?paper) AS ?numPapers) WHERE {
5     ?paper :hasAbstractTopics ?topics.
6     ?singleTopic spif:split(?topics ",")
7 }
8 GROUP BY ?singleTopic
9 ORDER BY desc(?numPapers)

```

4 sharedNobels

This query shows the number of Nobel Prizes shared by multiple laureates and the number of laureates sharing Nobel Prizes.

The query provides an interesting result: 242 out of 579 Nobel Prizes (41.8%) have been shared by multiple laureates, and 632 laureates have shared different Nobel Prizes. On average, a Nobel Prize is shared by more than two laureates (2.2 laureates per prize).

```

1 PREFIX : <http://www.semanticweb.org/a3d/ontologies/2024/10/nobel0ntology/>
2
3 SELECT ?share (COUNT(?nobel) AS ?numSharedNobels) (SUM(?share) AS
4     ?numLaureatesSharingNobels) WHERE {
5     ?nobel :hasPrizeShare ?share.
6 } GROUP BY ?share

```

For each value in the x axis, representing the number of people sharing a certain Nobel Prize, the following chart displays two bars. The first bar represents the total number of Nobel Prizes shared among the specified number of laureates. The second bar represents the total number of laureates who have shared these Nobel Prizes.

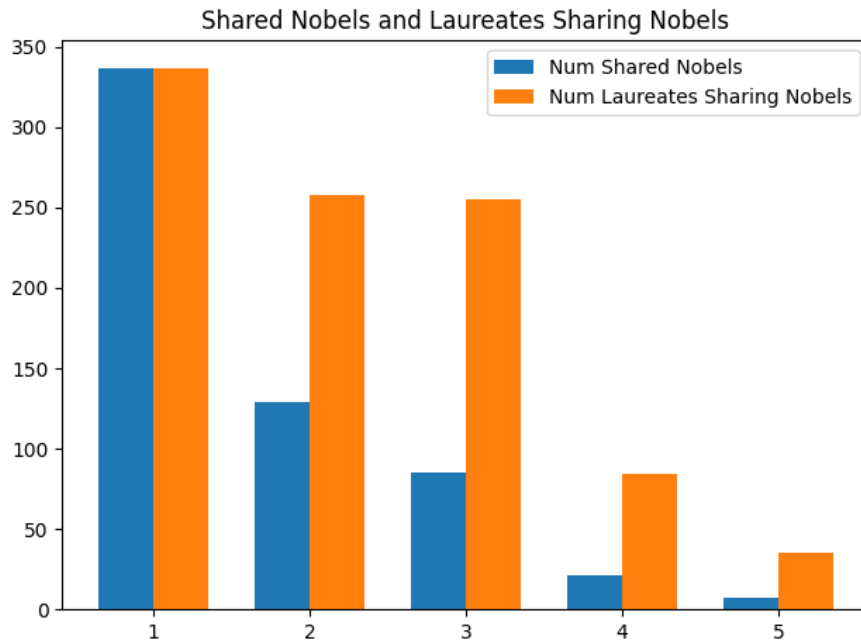


Figure 1: Graph showing the distribution of Nobel Prizes and laureates.

5 Collaborations among Nobel Laureates

```

1 PREFIX : <http://www.semanticweb.org/a3d/ontologies/2024/10/nobel0ntology/>
2 PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
3 PREFIX foaf: <http://xmlns.com/foaf/0.1/>
4
5 SELECT ?title (GROUP_CONCAT(?name; separator=", ") AS ?laureates) WHERE {
6     ?laureate rdf:type :Laureate .
7     ?paper rdf:type :Paper ;
8         :hasTitle ?title .
9     ?laureate :hasWritten ?paper .
10    ?laureate foaf:name ?name .
11 }
12 GROUP BY ?title
13 HAVING (COUNT(DISTINCT ?laureate) > 1)

```

The goal of this query was to explore whether Nobel laureates collaborate with each other by co-authoring scientific papers. As the table 3 shows, among approximately 53,000 papers and 904 laureates, only one paper was found to have been co-authored by multiple laureates.

Table 3: Paper co-authored by multiple Nobel Laureates

Title	Laureates
<i>Recursive Robust Estimation and Control Without Commitment</i>	Lars Peter Hansen, Thomas J. Sargent

This result suggests that collaboration between Nobel laureates is extremely rare. However, it is important to note that this outcome should not be taken as definitive, as the datasets used represent only a portion of all existing papers and laureates. Nevertheless, it provides an interesting insight into the rarity of such collaborations, offering a percentage-based perspective on how seldom laureates join forces to produce scientific work.

6 How fundings in R&D affect the possibility for a country to win a Nobel?

```

1 PREFIX : <http://www.semanticweb.org/a3d/ontologies/2024/10/nobel0ntology/>
2 PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
3
4 SELECT ?year ?topCountry (COUNT(DISTINCT ?laureate) AS ?numLaureates) (SUM(?
   fundingAmount) AS ?totalFunding) WHERE {
5   ?laureate rdf:type :Laureate ;
6             :hasWon ?nobelPrize ;
7             :bornIn ?city .
8   ?nobelPrize :hasYear ?year .
9   ?city :locatedIn ?topCountry .
10
11  OPTIONAL {
12    ?topCountry :hasFunded ?funding .
13    ?funding :hasYear ?year ;
14             :hasAmount ?fundingAmount .
15  }
16  { # Select country with most laureates
17    SELECT (?country AS ?topCountry) WHERE {
18      ?laureate rdf:type :Laureate ;
19                :bornIn ?city .
20      ?city :locatedIn ?country .
21    }
22    GROUP BY ?country
23    ORDER BY DESC(COUNT(DISTINCT ?laureate))
24    LIMIT 3
25  }
26 }
27 GROUP BY ?year ?topCountry
28 HAVING(SUM(?fundingAmount) > 0)
29 ORDER BY ?year ?topCountry

```

With this query, we identified the top three countries with the highest number of Nobel laureates born there, along with the annual amount of funding allocated to research and development (R&D) by these nations. To ensure data consistency, we focused exclusively on the years from 2000 to 2016.

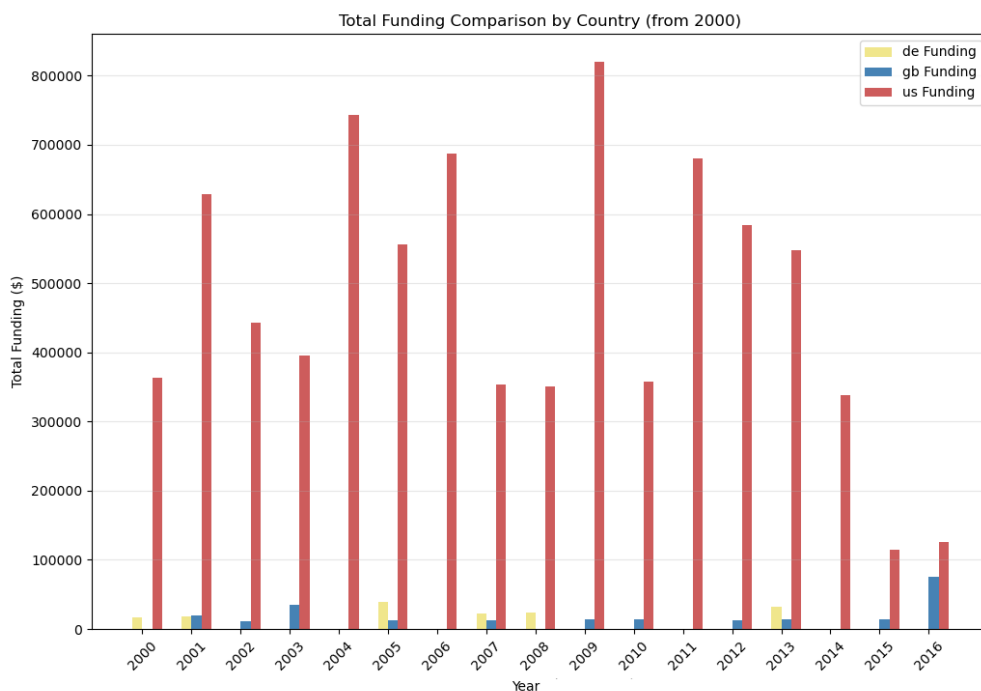


Figure 2: Funding Comparison by Country

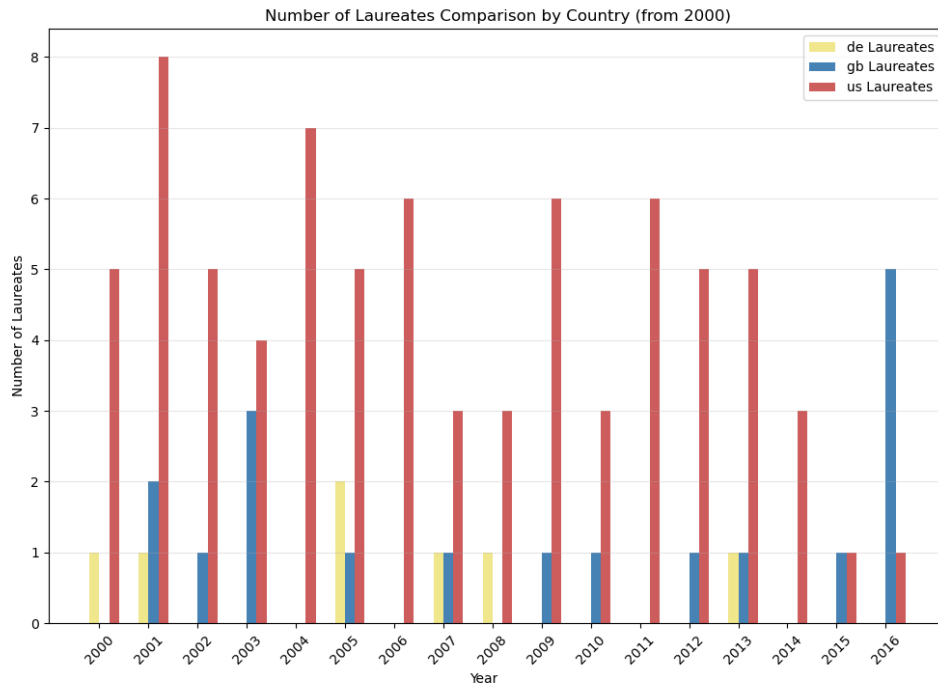


Figure 3: Laureates Comparison by Country

The graphs reveal a strong correlation between R&D funding and the number of Nobel laureates. In particular, the United States dominates both metrics, demonstrating how substantial investments in research directly contribute to significant achievements in this field, resulting in a higher number of laureates annually.

The situation in Great Britain, highlighted in the following plot, is particularly curious and further supports this observation:

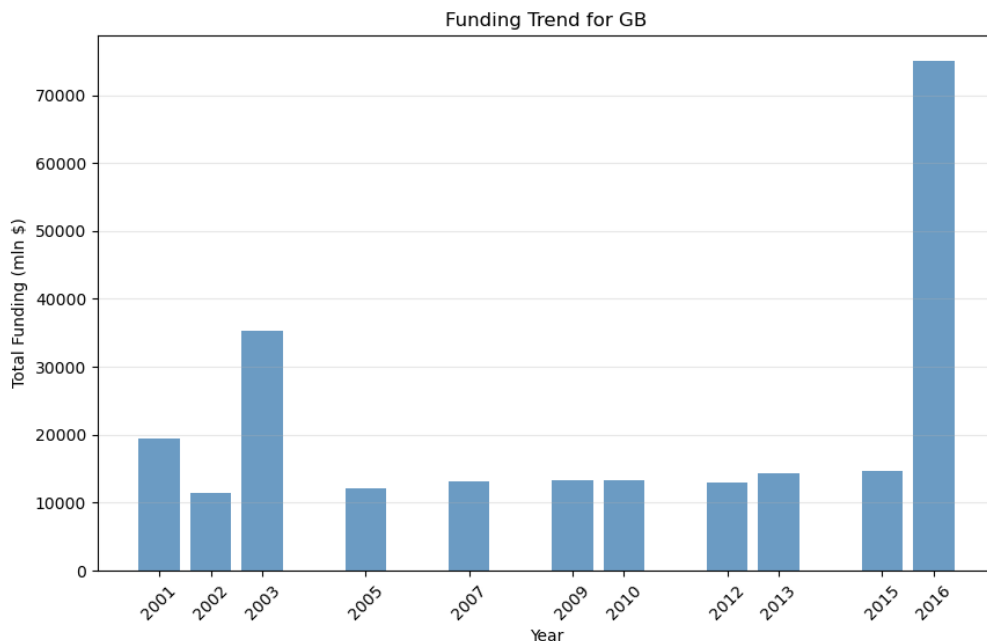


Figure 4: Great Britain R&D Funding Trend

It's clear that the trends in funding and the number of laureates mirror each other closely. From 2001 to 2003, we observe the same pattern in both metrics. Subsequently, a steady and low level of R&D funding still reflects the number of British Nobel laureates until 2016, when a sharp increase in Nobel prizes matches with a significant

rise in R&D investments.

7 moreThanOneNobel

```
1 PREFIX spif: <http://spinrdf.org/spif#>
2 PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
3 PREFIX jur: <http://sweet.jpl.nasa.gov/2.3/humanJurisdiction.owl#>
4 PREFIX skos: <http://www.w3.org/2004/02/skos/core#>
5 PREFIX foaf: <http://xmlns.com/foaf/0.1/>
6 PREFIX : <http://www.semanticweb.org/a3d/ontologies/2024/10/nobelOntology/>
7 PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
8
9 SELECT ?laureate (COUNT(?nobel) AS ?numNobels) WHERE {
10     ?laureate :hasWon ?nobel.
11 }
12 GROUP BY ?laureate
13 HAVING (?numNobels > 1)
```

who won more than one nobel prize

8 papersPerVenue

This plot shows the number of papers published over the years by major venues (those with at least 800 papers published, according to our dataset).

In recent years, Bioinformatics could be considered one of the most influential venue due to its consistently higher number of papers published compared to others. IEEE venues, are the most prominent in the fields of information and tecnology.

For instance, on 2009, the research community focused more on the field of communications. That same year, the Physics Nobel Prize was awarded for "groundbreaking achievements concerning the transmission of light in fibers for optical communication".

```
1 PREFIX : <http://www.semanticweb.org/a3d/ontologies/2024/10/nobelOntology/>
2
3 SELECT ?venue ?year (COUNT(?paper) AS ?numPapers) WHERE {
4
5     # get the most important venues (the ones with at least 800 papers published)
6     {
7         SELECT ?venue (COUNT(?paper) AS ?totPapers) WHERE {
8             ?paper :publishedIn ?venue.
9         }
10        GROUP BY ?venue
11        HAVING (?totPapers > 800)
12        ORDER BY DESC (?totPapers)
13    }
14
15    # get the number of paper published in the most important venues for each year
16    ?paper :publishedIn ?venue;
17           :hasYear ?year.
18 }
19 GROUP BY ?venue ?year
20 ORDER BY ASC (?year)
```

9 papersPerCategory

The following query allows us to extract, for each year, the number of scientific articles published in each relevant category. The categories returned as results are the TopConcepts categories of our SKOS taxonomy, and they include in the count their various subcategories. For example, in the count of papers for the medicine category, articles belonging to subcategories like neuroscience are also included.

To obtain this data, the query uses two distinct subqueries. The first subquery extracts the number of articles published for each main category (TopConcept), while the second identifies the number of articles associated with the subcategories of each main category. The sum of the results of the two subqueries, aggregated by year and category, provides the total number of articles published for each category and for each year.

```

1 PREFIX skos: <http://www.w3.org/2004/02/skos/core#>
2 PREFIX : <http://www.semanticweb.org/a3d/ontologies/2024/10/nobel0ntology/>
3 # Extracts the number of papers we have for each category over the years --> the most
   studied research areas over the years
4 SELECT ?year ?category (SUM(?howmany) AS ?totalPapers) WHERE {
5     # Inner query to extract the number of papers published in journals that have at
       least one category that is a top concept of our skos scheme
6     {
7         SELECT ?year ?category (COUNT(DISTINCT ?paper) AS ?howmany) WHERE {
8             ?journal :hasJournalCategory ?category .
9             :journalCategoryScheme skos:hasTopConcept ?category .
10            ?paper :publishedIn ?journal ;
11                :hasYear ?year .
12        }
13        GROUP BY ?year ?category
14    }
15    UNION
16    # Inner query to extract the number of papers published in journals that have at
       least one category that is a subcategory of a top concept category
17    {
18        SELECT ?year ?category (COUNT(DISTINCT ?paper) AS ?howmany) WHERE {
19            ?journal :hasJournalCategory ?cat .
20            ?cat skos:broaderTransitive ?category .
21            ?paper :publishedIn ?journal ;
22                :hasYear ?year .
23        }
24        GROUP BY ?year ?category
25    }
26 }
27 GROUP BY ?year ?category
28 ORDER BY DESC (?totalPapers)

```


This approach offers a comprehensive view of the distribution of published articles over time, allowing us to identify which research areas, related to Nobel categories, have attracted the most attention from scholars over the years. Below is a plot representing the trend of the number of papers published over the years, divided by their respective categories. In the recent years, the most studied field is medicine which got a big leap around 2002, while in the nineties the most studied one was economics, which reached its peak in 2008, probability due to the economic crisis of that time.

