



**POLITECNICO**  
MILANO 1863

SCUOLA DI INGEGNERIA INDUSTRIALE  
E DELL'INFORMAZIONE

# Systems and Methods for Big and Unstructured Data Project

Author(s): **Simona Malegori**

**Nicole Perrotta**

**Michele Simeone**

**Alberto Pirillo**

**Simone Tognocchi**

Group Number: **38**

Academic Year: 2022-2023



# Contents

<b>Contents</b>	<b>i</b>
<b>1 First Delivery</b>	<b>1</b>
1.1 Problem description . . . . .	1
1.2 Hypothesis . . . . .	1
1.3 Data . . . . .	2
1.3.1 ER Diagram . . . . .	2
1.3.2 Data set description . . . . .	4
1.3.3 Data Pre-Processing . . . . .	5
1.4 Neo4j . . . . .	6
1.4.1 Data Upload . . . . .	6
1.4.2 Graph Diagram . . . . .	9
1.4.3 Queries . . . . .	12
1.4.4 Creation/Update Commands . . . . .	22
<b>2 Second Delivery</b>	<b>25</b>
2.1 Introduction . . . . .	25
2.2 Data . . . . .	25
2.2.1 Data Pre-Processing . . . . .	28
2.2.2 Data Completion . . . . .	28
2.3 MongoDB . . . . .	30
2.3.1 Data Upload . . . . .	30
2.3.2 Document Example . . . . .	30
2.3.3 Join Operation . . . . .	34
2.3.4 Queries . . . . .	34
2.3.5 Creation/Update Commands . . . . .	45
<b>3 References</b>	<b>55</b>



# 1 | First Delivery

## 1.1. Problem description

This project aims at building an Information System that manages a data set containing different type of scientific articles that can be used for: clustering with network and side information, studying influence in the citation network, finding the most influential papers and topics, modeling analysis.

The project is divided in the following steps.

At first it was made an ER Diagram that generalizes all the information gathered from different already existing data sets, then the most complete data set was chosen.

Afterwards the data set was pre-processed, transforming it from a JSON to a CSV format and then it was reduced in size.

After that it was uploaded on Neo4j and all the nodes, the relationships and the properties were edited to build the Graph Diagram.

At the end of the project 10 queries and 6 creation/update commands were initiated with a different level of complexity that was checked within the performance time.

## 1.2. Hypothesis

In order to model the database we made some assumptions:

- authors can have zero, one or more papers associated, assuming that authors are inserted in the database before their paper/s is/are;
- authors can have zero, one or more affiliated organizations, assuming that there can be authors that didn't provide their organization;
- papers have at least one author;
- papers have at least one field of study;
- not all papers have keywords associated, assuming that they may not have been

provided;

- papers can reference and be referenced by other papers;
- each paper has a venue, that is the place where it has been published/presented;
- a venue can host more than one paper;
- venues are of 4 types: Journal, Conference, Book and Patent;
- the volume  $n$  of a paper is the  $n$ -th published collection;
- the issue  $m$  of a paper is the  $m$ -th part of the volume in which it is published.

## 1.3. Data

### 1.3.1. ER Diagram

The ER Diagram of the chosen model is characterized by the following entities with the respective attributes:

- **Paper** is a scientific article that is associated with the following attributes:  
 $id$ ,  $title$ ,  $date$  that corresponds to the publication date,  $doi$  that is the Digital Object Identifier,  $volume$  that corresponds to the n-th published collection,  $issue$  that corresponds to the m-th part of the volume,  $language$ ,  $issn$  that is an identification code associated with the title of the publication,  $isbn$  that is a code that identifies printed or digital papers and it is used as inventory-tracking device,  $n\_citation$  that is the number of citations,  $page\_start$  and  $page\_end$  that are the starting and the ending point of the collection from which the paper was extracted,  $pdf\_url$  that is the source from which to recover the paper,  $abstract$  that is the summary of the paper,  $publisher$  and  $external\_url$  that corresponds to the sitography of the paper;
- **Author** with the attributes:  $id$ ,  $name$ ,  $surname$ ,  $email$ ,  $orcid$  that is a unique and persistent identification number and  $organization$ ;
- **Keyword** that represents a word that allows to define immediately the topic within the paper, its attributes are:  $id$  and  $name$ ;
- **FoS** that corresponds to the field of studies with the attributes:  $id$ ,  $name$ ,  $w$  that is the weight of the fields of study;
- **Venue** that is the collection from which the paper was extracted, with the attributes:  $id$  and  $name$ . The venue of the paper can be of different types, in fact there is a

total and exclusive generalization of the entity **Venue** that can be a:

- **Journal** that has also the attribute *addressee* that is the type of audience of the paper;
- **Book** that has also the attributes *category* and *edition*;
- **Conference** that has also the attributes *type* that can be physical and online, and *location* that has cardinality one only if the type is physical and it represents the place in which the conference takes place;
- **Patent** that has also the attributes *type* and *expiration*.

In the ER Diagram there are the following relationships with the respective cardinalities:

- **Writing** between the entities **Paper** and **Author**. The relationship means that a Paper can be written by at least 1 to a maximum of N authors and that an Author can write from 0 to N papers.
- **Containing** between the entities **Paper** and **Keyword**. The relationship means that a Paper can contain from 0 to N keywords and that a Keyword is contained into at least 1 to a maximum of N papers.
- **Dealing** between the entities **Paper** and **FoS**. The relationship means that a Paper can deal with at least 1 to a maximum of N field of studies and that a FoS can be dealt from 0 to N papers.
- **In** between the entities **Paper** and **Venue**. The relationship means that a Paper is extracted from exactly 1 venue and that a venue can be the collection in which at least one paper is contained.
- **Referencing** that is a relationship on the same entity of the **Paper**, in fact it contains the roles *referencer* and *referenced*. The relationship means that a Paper can or not reference other Papers and can be referenced or not from other Papers.

After all, in the ER Diagram there is an external constraint on the attributes of the entity **Conference**. In fact, the attribute *location* must have a cardinality of (1,1) if the type is *physical*.

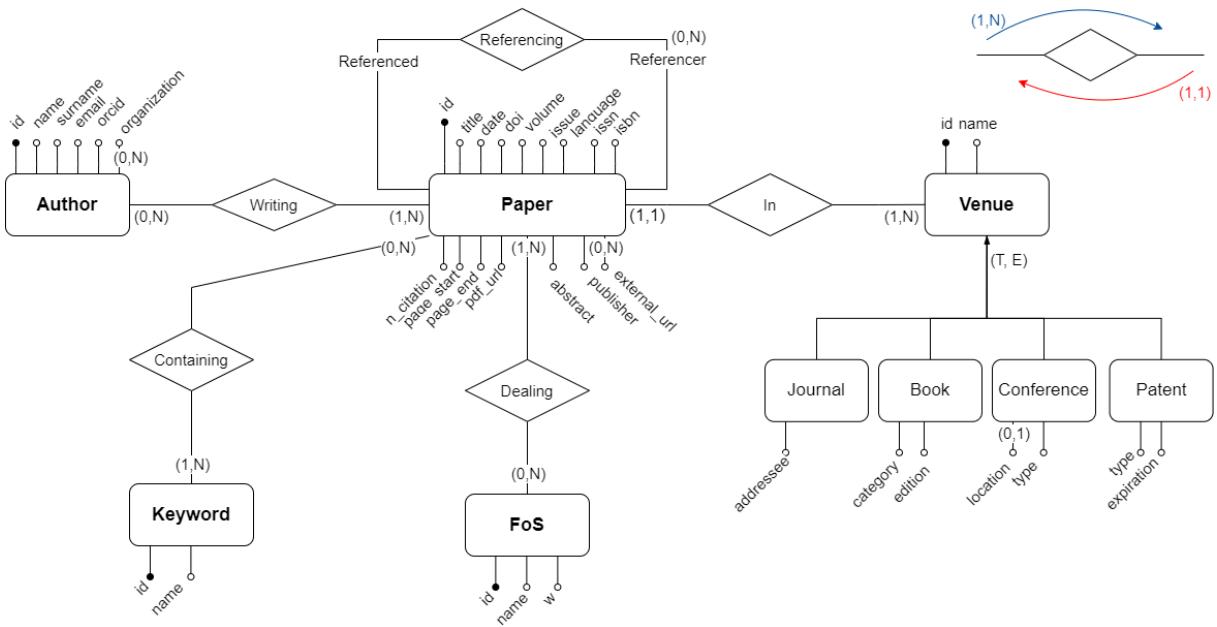


Figure 1.1: ER Diagram

### 1.3.2. Data set description

The data set that was chosen for the goal of the project contains 8335 papers, 730 venues, 27576 authors and 25005 field of studies. It is a reduced version of the more complex model on which the ER Diagram is based. In fact, it doesn't contain the entities Keyword and the sub-entities Journal, Book, Conference and Patent that, instead, are transformed into an attribute of the entity Paper that is called doc\_type. All the respective attributes of these entities are eliminated. Moreover, the attribute *date* of the Paper, that is generalized in the ER Diagram, becomes the year. Furthermore, in the reduced version of the data set, the relationship Referencing becomes a list of string that are the references of the Paper. Lastly, the chosen data set contains just the more significant attributes. Summarizing, the chosen data set is composed by the following entities and attributes:

- **Paper** with the attributes: id (integer type), title (string type), year (integer type), n\_citation (integer type), page\_start (integer type), page\_end (integer type), doc\_type (string type), publisher (string type), volume (integer type), issue (integer type), doi (string type), references (list of string type), abstract (string type);
- **Author** with the attributes: paper\_id (integer type), author\_name (string type), author\_id (integer type), author\_org (string type);
- **FoS** with the attributes: paper\_id (integer type), fos\_name (string type), fos\_weight (float type);

- **Venue** with the attributes: paper\_id (integer type), venue\_name (string type), venue\_id (integer type).

### 1.3.3. Data Pre-Processing

The original data set can be downloaded [here](#). Given that such data set was unnecessarily large for our purpose, we decided to reduce its size. We could have just cut it at a certain point, however we decided that it was better to work with consistent data, therefore we decided to carefully perform sub-sampling in an intelligent way. Such pre-processing was performed using multiple Python scripts and the Pandas library.

Here we provide a description of all the scripts used, together with the procedure to obtain the final data set starting from the initial one. All of these scripts can be found in the *scripts* folder of the *neo4j* section of the *GitHub repository* of the project.

The notebook **dataset\_exploration.ipynb** contains a short description of every operation for explanatory reasons. This notebook processes only a small chunk of the data set. The same operations are performed on the whole data set in the script **dataset\_preprocessing.py**.

Here is a summary of the operations performed:

- Removal of samples with Null and NaN values
- Removal of samples with an empty string in a field

The operations above are required to work with consistent data. Notice that we can afford to simply drop the samples that do not respect such conditions since we dispose of a very large data set.

The following operations are not required but were performed to reduce even further the size of the data set, with the objective of keeping only the "most important" samples.

We kept the samples:

- With a number of citations greater than a threshold
- With a reference count greater than a threshold

We converted the *indexed\_abstract* field from an inverted index to a string of text, to make it easier to query once inserted into the database. The field was also renamed to *abstract*.

We also processed the data set in order to remove some special characters not supported

by the import function of the database. An example of such characters is "\\". This operation is performed in the script **remove\_special\_characters.py**.

Given that many samples were removed from the data set, we also had to fix up the *references* field to keep only the valid ones. We consider a reference valid when it points to a sample of the data set that was kept. Otherwise, we say that such reference is invalid and we remove it from the data set. This operation is performed in the script **deprecated\_references.py**.

Lastly, when importing the data into the database, we realized that it was more practical to split the data set into multiple data sets to speed up the process and to produce cleaner code. This functionality was added in the **dataset\_preprocessing.py** script. We split the data set following the structure of ER diagram, ending up with one separate data set for each entity present in the original data set. Thus, we ended up with 4 data sets: **Paper**, **Author**, **Venue** and **Fos**.

To obtain the final data set starting from the downloaded one, run the scripts in this order:

1. **dataset\_preprocessing.py**
2. **deprecated\_references.py**
3. **remove\_special\_characters.py**

The input of the first script is the initial data set. Only the paper data set requires to be processed by the second script, then only the paper and the venue data sets require to be processed by the third script. At the end, you will obtain 4 data sets which are identical to the ones that we imported into the database.

## 1.4. Neo4j

### 1.4.1. Data Upload

To import the data into Neo4j, there is one last precaution needed: it is necessary to process the Paper data set to make the *references* field compatible with the import function of the database. To perform such action the script **preprocessing.py** can be used. The script is located in the *neo4j* folder of the *GitHub repository*. Now, the data sets can be correctly imported into the Neo4J database. In order to do that, you need to put the four data sets files in the program's import folder.

Thanks to the pre-processing that was performed, the four data sets are saved in the

simple CSV format and it is possible to use the LOAD CSV command to easily load all the entities and relationships.

1. Clear the Database:

```
1 MATCH (x) DETACH DELETE x;
```

2. Load the Papers:

```
1 LOAD CSV WITH HEADERS FROM "file:///paper_dataset2.csv" AS csvLine
2 CREATE (p:Paper {id: toInteger(csvLine.id), title: csvLine.title,
    year: toInteger(csvLine.year), doi: csvLine.doi, volume: csvLine
    .volume, issue: csvLine.issue, abstract:csvLine.abstract,
    n_citation:toInteger(csvLine.n_citation), page_start:toInteger(
    csvLine.page_start), page_end:toInteger(csvLine.page_end),
    publisher:csvLine.publisher, doc_type:csvLine.doc_type})
```

3. Load the Authors:

```
1 LOAD CSV WITH HEADERS FROM "file:///author_dataset.csv" AS csvLine
2 CREATE (a:Author {id: toInteger(csvLine.author_id), name: csvLine.
    author_name, organization: csvLine.author_org})
```

4. Load the Fos:

```
1 LOAD CSV WITH HEADERS FROM "file:///fos_dataset_0.csv" AS csvLine
2 CREATE (f:Fos {weight: toFloat(csvLine.fos_weight), name: csvLine.
    fos_name})
```

5. Load the Venue:

```
1 LOAD CSV WITH HEADERS FROM "file:///venue_dataset.csv" AS csvLine
2 MERGE (v:Venue {id: toInteger(csvLine.venue_id), name: csvLine.
    venue_name})
```

6. Create the Writing relationship between papers and authors:

```
1 LOAD CSV WITH HEADERS FROM "file:///author_dataset.csv" AS csvLine
2 MATCH (p:Paper),(a:Author)
3 WHERE (toInteger(csvLine.paper_id)=p.id) and (toInteger(csvLine.
    author_id) = a.id)
4 CREATE (p)-[w:writing]->(a)
```

7. Create the Referring relationship between papers and papers:

```
1 LOAD CSV WITH HEADERS FROM "file:///paper_dataset2.csv" AS csvLine
2 UNWIND split(csvLine.references, ':') as ref
3 MATCH (p1:Paper),(p2:Paper)
4 WHERE (toInteger(csvLine.id)=p1.id) and (toInteger(ref)=p2.id)
```

```
5 CREATE (p1)-[r:referencing]->(p2)
```

8. Create the `in` relationship between papers and venues:

```
1 LOAD CSV WITH HEADERS FROM "file:///venue_dataset.csv" AS csvLine
2 MATCH (p:Paper),(v:Venue)
3 WHERE (toInteger(csvLine.paper_id)=p.id) AND (toInteger(csvLine.
    venue_id) = v.id)
4 CREATE (p)-[i:in_]->(v)
```

9. Create the `dealing` relationship between papers and FoS:

```
1 LOAD CSV WITH HEADERS FROM "file:///fos_dataset_0.csv" AS csvLine
2 MATCH (p:Paper),(f:Fos)
3 WHERE (toInteger(csvLine.paper_id)=p.id) AND (csvLine.fos_name = f.
    name) AND (toFloat(csvLine.fos_weight) = f.weight)
4 CREATE (p)-[d:dealing]->(f)
```

### 1.4.2. Graph Diagram

Entities:

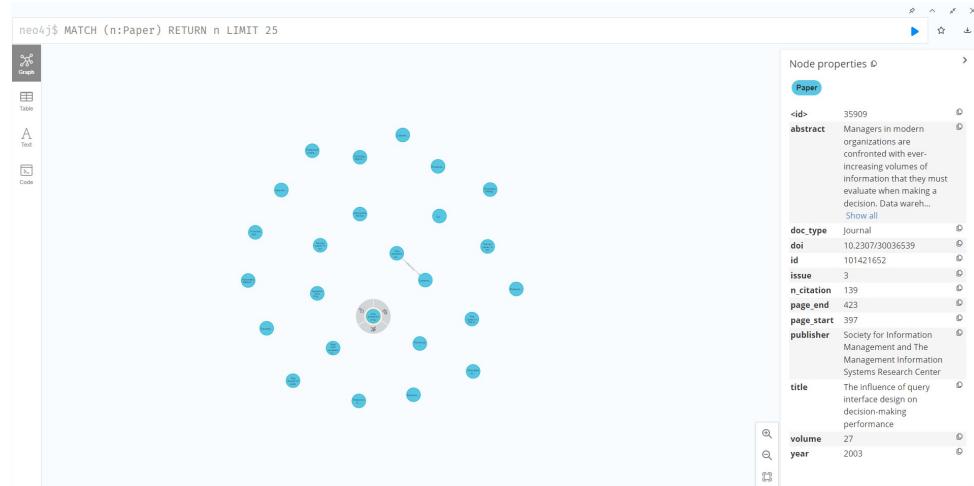


Figure 1.2: Paper

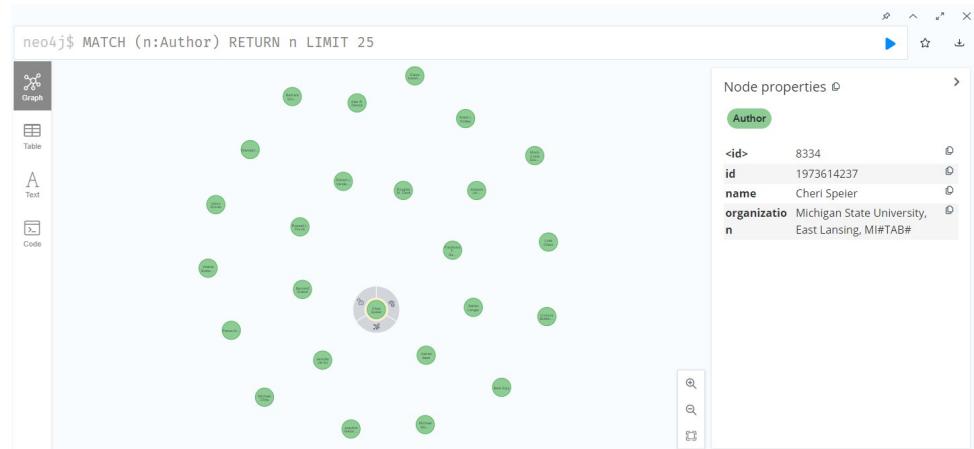


Figure 1.3: Author



Figure 1.4: FoS

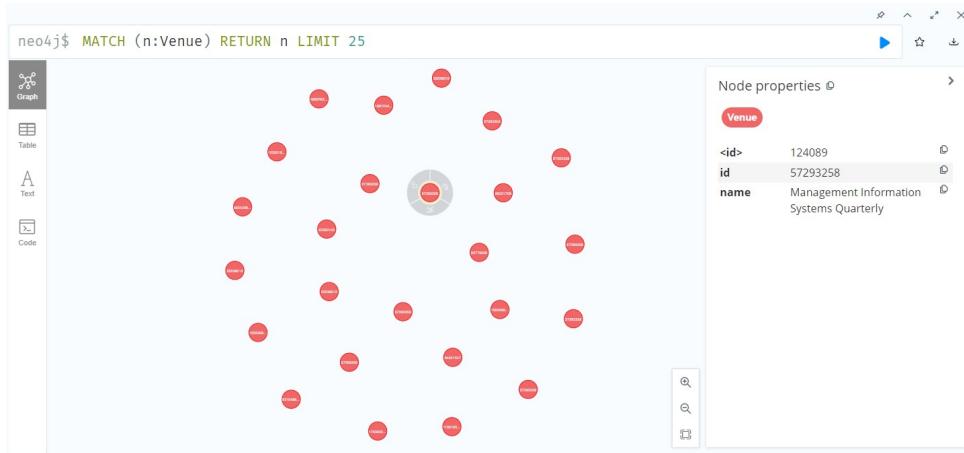


Figure 1.5: Venue

Relationships:

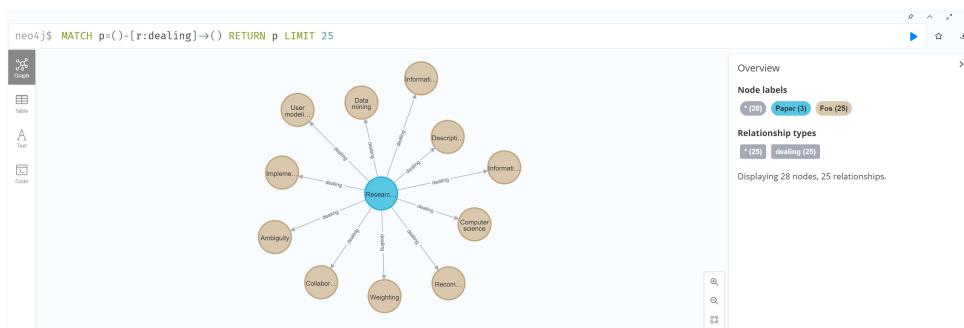


Figure 1.6: Paper-Dealing-&gt;FoS

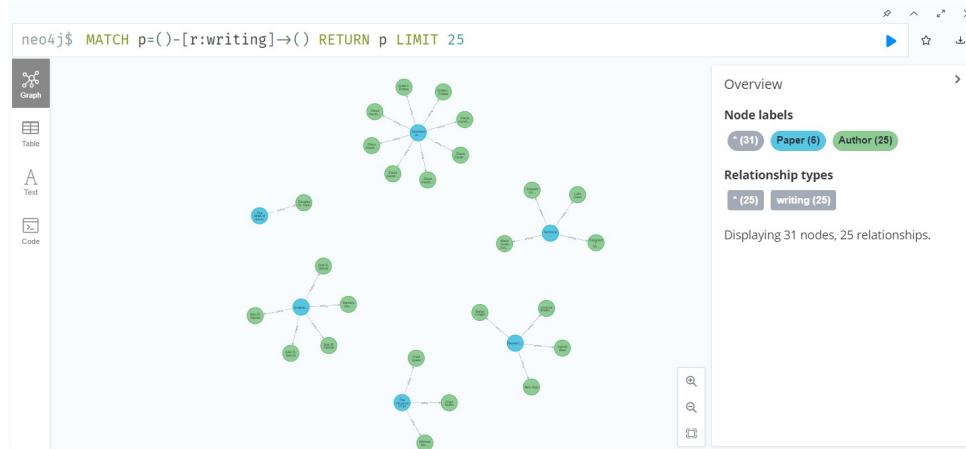


Figure 1.7: Paper-Writing-&gt;Author

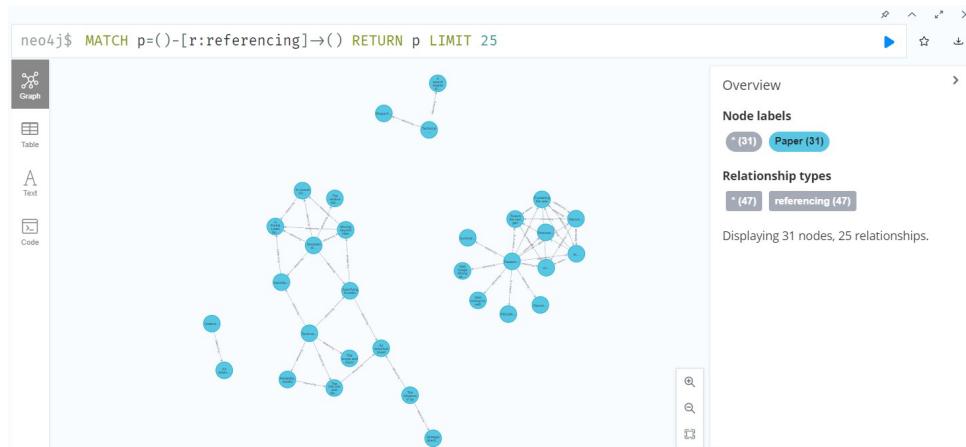


Figure 1.8: Paper-Referencing-&gt;Paper

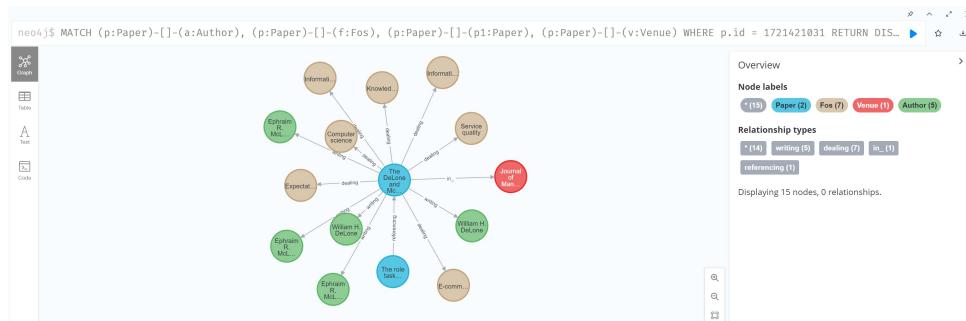


Figure 1.9: All the relationships

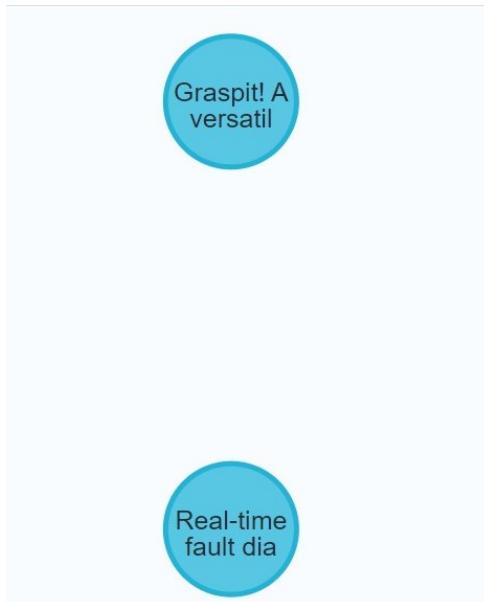
### 1.4.3. Queries

- Find papers of a determined venue, written after a certain date (Execution time 454ms):

```

1 MATCH (p: Paper)-[i:in_]->(v: Venue), (p: Paper)-[d:dealing]->(f:
  Fos)
2 WHERE (p.year > 2000) AND (v.name="IEEE Robotics & Automation
  Magazine")
3 RETURN p

```

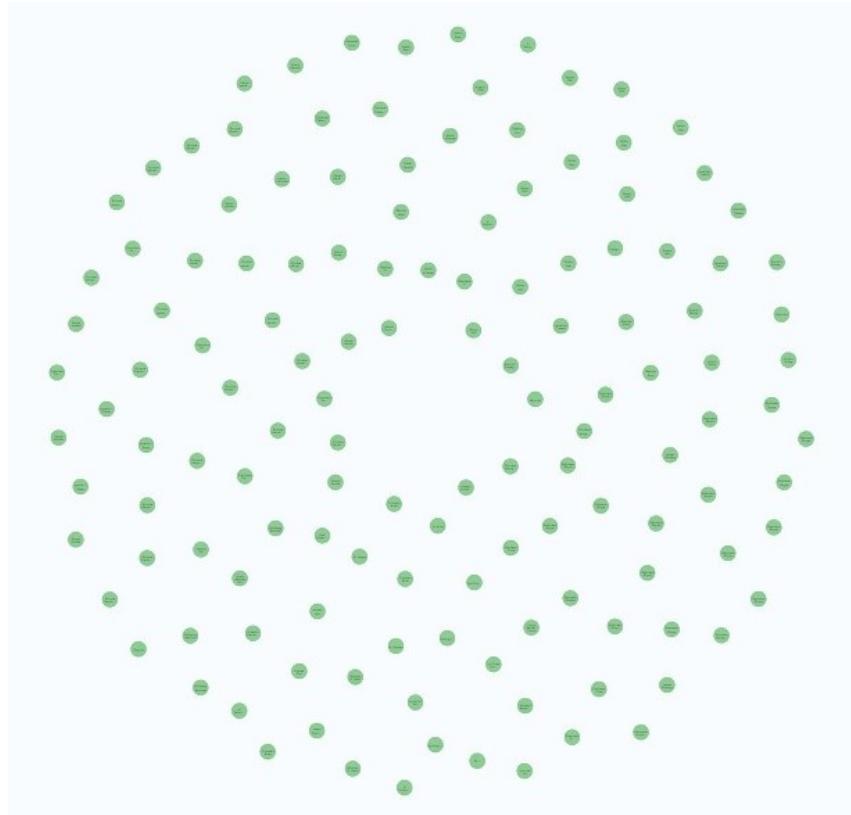


- Find papers in a set of venues, with a determined type (Execution time 2ms):

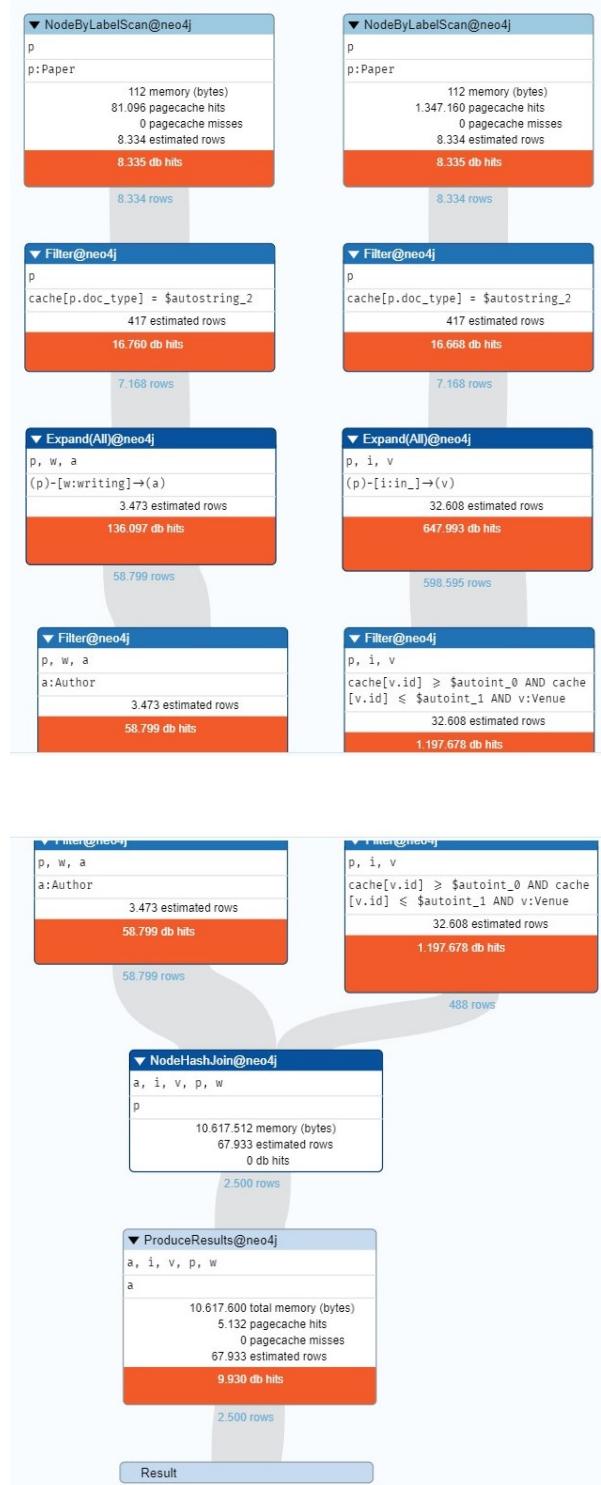
```

1 MATCH (p: Paper) -[w:writing]->(a:Author), (p: Paper) -[i: in_]->
  (v:Venue)
2 WHERE v.id >= 140000000 and v.id <= 140900000 AND p.doc_type='
  Journal'
3 RETURN a

```



Below we can see the output of the profile statement, that is used to track the query and the numbers of rows of each operation. It starts by scanning all the nodes with the papers, then it expands all the nodes with the 'writing' relationship on the left of the picture and the 'in' relationship on the right. Finally, it applies the filters on the venue id and on the paper's doc type and returns the results.



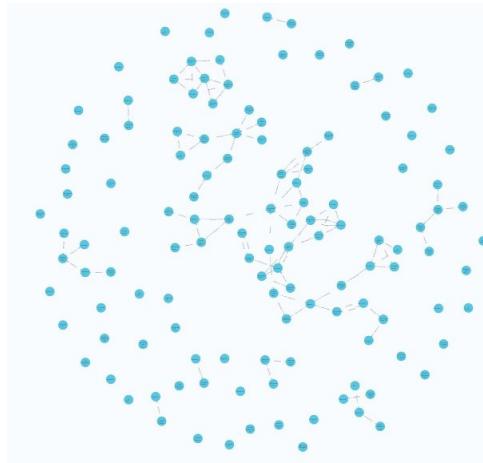
3. Find papers of a determined main argument, written after a certain date (Execution time 30ms):

```
1 MATCH (p: Paper)-[r:referencing]->(p2: Paper), (p2: Paper)-[d: dealing]->(f: Fos)
```

```

2 WHERE (f.weight >= 0.45) AND (p2.year >2010) AND (f.name="Artificial intelligence")
3 RETURN p

```



4. Count the authors that have written a famous paper of a determined argument (Execution time 38ms):

```

1 MATCH (p: Paper)-[w:writing]->(a:Author), (p: Paper)-[d: dealing]->(f: Fos)
2 WHERE f.name = "Machine learning" AND p.n_citation>5000
3 RETURN COUNT(a.id) AS num_aut

```

num_aut
1 2120

5. Count the papers divided per author written in a set of venue by a determined authors' organization (Execution time 103ms):

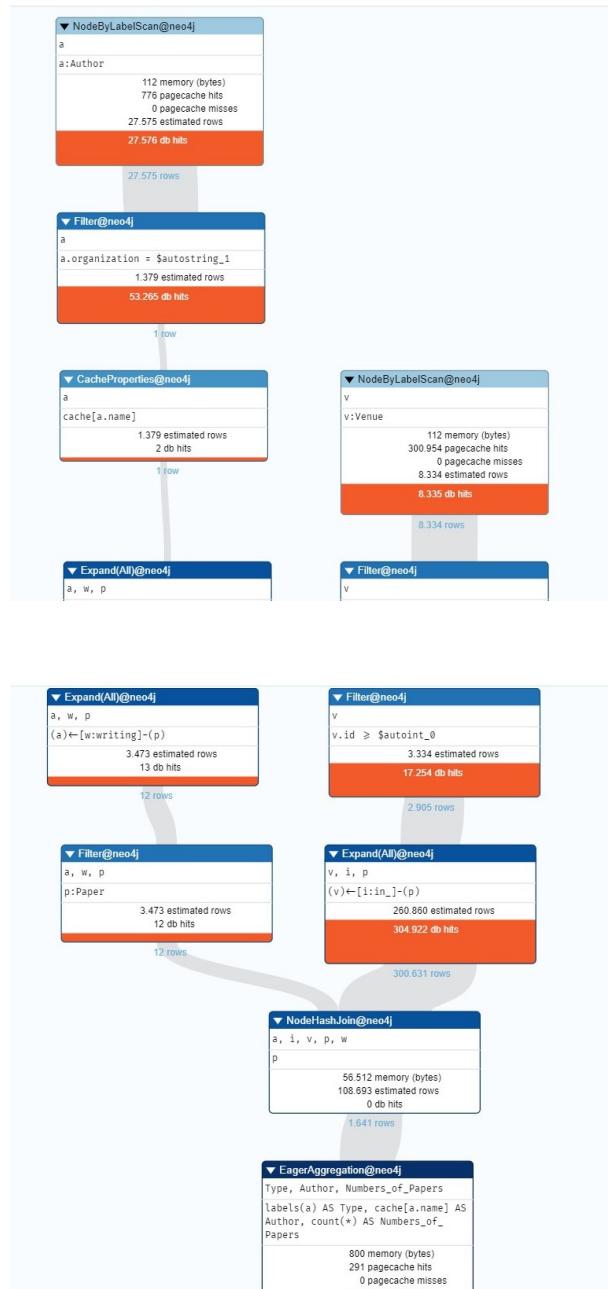
```

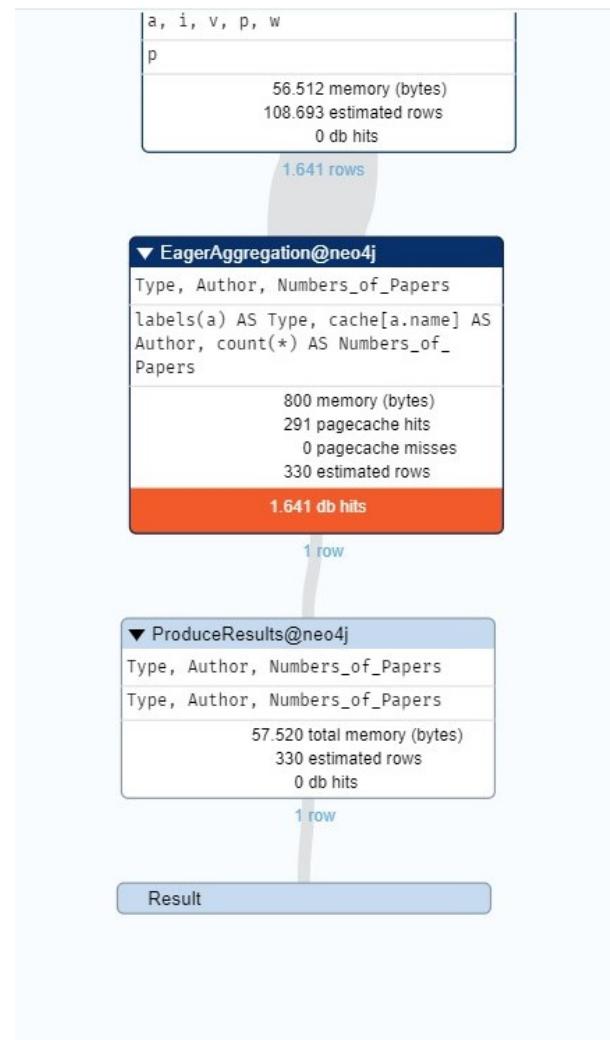
1 MATCH (p: Paper) -[w:writing]->(a:Author), (p: Paper) -[i: in_]->(v : Venue)
2 WHERE v.id >= 160000000 AND a.organization='Microsoft Research, One Microsoft Way, Redmond, WA 98052, USA. hshum@microsoft.com#TAB#
'
3 RETURN labels(a) AS Type, a.name AS Author, COUNT(*) AS Numbers_of_Papers

```

Type	Author	Numbers_of_Papers
["Author"]	"Heung-Yeung Shum"	1641

The profile statement scans all the authors, then it expands the 'writing' relationship and applies the organization's filter. At the same time it scans all the venues and filters their id. At the end the aggregation is applied and all the resulting values are returned.





6. Calculate the weight keywords' average of papers that have a determined publisher and a reference to a famous paper (Execution time 27ms):

```

1 MATCH (p: Paper) -[d: dealing]->(f: Fos), (p: Paper) -[r:
  referencing]->(p2: Paper)
2 WHERE p2.n_citation>2000 AND p.publisher='Society for Information
  Management and The Management Information Systems Research
  Center'
3 RETURN p.title, avg(f.weight)

```

p.title	avg(f.weight)
"Shackled to the status quo: the inhibiting effects of incumbent system habit, switching costs, and inertia on new system acceptance"	0.449569698283237
"Technostress: technological antecedents and implications"	0.44386502606134975
"Business intelligence in blogs: understanding consumer interactions and communities"	0.41906192304436024
"Web and wireless site usability: understanding differences and modeling use"	0.42455533500106846
"Competing perspectives on the link between strategic information technology alignment and organizational agility: insights from a mediation model"	0.4178663593780648
"Reliability, mindfulness, and information systems"	0.42862394017677813

7. Count the numbers of bilateral reference between two papers that have the same venue and a large number of citation (Execution time 119ms):

```

1 MATCH (p: Paper)-[r:referencing]->(p2: Paper), (p2: Paper)-[r2:
    referencing]->(p: Paper), (p: Paper)-[w:writing]->(a: Author),(p
    : Paper)-[i:in_]->(v: Venue), (p2: Paper)-[i2:in_]->(v2: Venue)
2 WHERE p2.n_citation>500 AND p.n_citation>500 AND v.name=v2.name AND
      a.organization="Royal Institute of Technology"
3 RETURN type(r) AS Relation, COUNT(*)/2 AS
      Num_of_bilateral_referencingsame_venue, a.organization AS
      organization
4 LIMIT 1

```

Relation	Num_of_bilateral_referencing_same_venue	organization
"referencing"	28900	"Royal Institute of Technology"

8. Count the total citation of a paper that have references to two papers that deal of different field of study (Execution time 1849ms):

```

1 MATCH (p: Paper)-[r: referencing]->(p2: Paper), (p: Paper)-[r2:
    referencing]->(p3: Paper), (p2: Paper)-[d: dealing]->(f: Fos), (
    p3: Paper)-[d2: dealing]->(f2: Fos)
2 WHERE p.page_end-p.page_start>10 AND p2 <> p3 AND f.name='
    Artificial Intelligence' AND f.weight>0.0 and f2.name = 'Machine
    learning' AND f2.weight>0.0
3 RETURN p.title AS Title, SUM(p.n_citation) AS Sum_n_citation
4 LIMIT 5

```

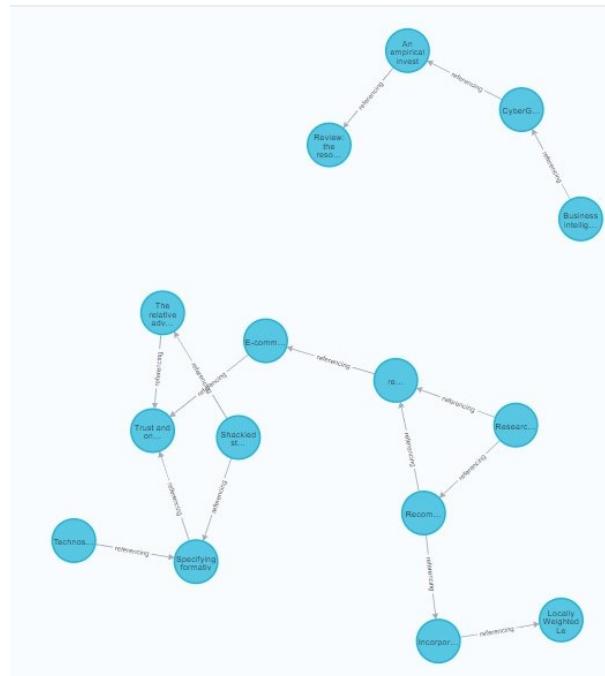
Title	Sum_n_citation
1 "The SUN Attribute Database: Beyond Categories for Deeper Scene Understanding"	890400
2 "Composition in Distributional Models of Semantics"	4053440
3 "Wikipedia-based semantic interpretation for natural language processing"	2170880
4 "Knowledge derived from wikipedia for computing semantic relatedness"	737760
5 "A survey of paraphrasing and textual entailment methods"	4324800

9. Find the shortest path between two different paper that have a reference in common where the first is less famous than the second one (Execution time 61ms):

```

1 MATCH s = shortestPath(
2   (p: Paper)-[r:referencing*]->(p2: Paper)
3 )
4 WHERE p2.n_citation > 10*p.n_citation AND p <> p2
5 RETURN s
6 LIMIT 5

```

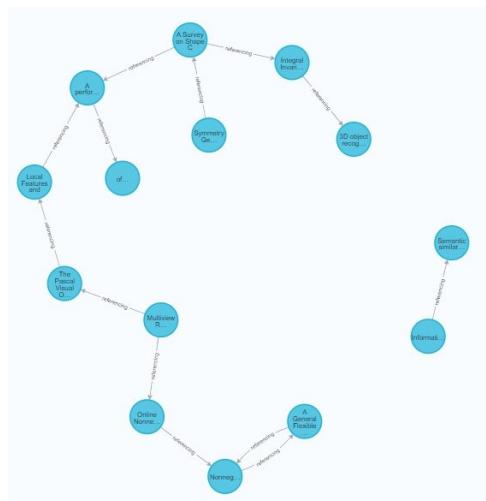


10. Find the shortest path between two different paper that have a reference in common different FoS (Execution time 42ms):

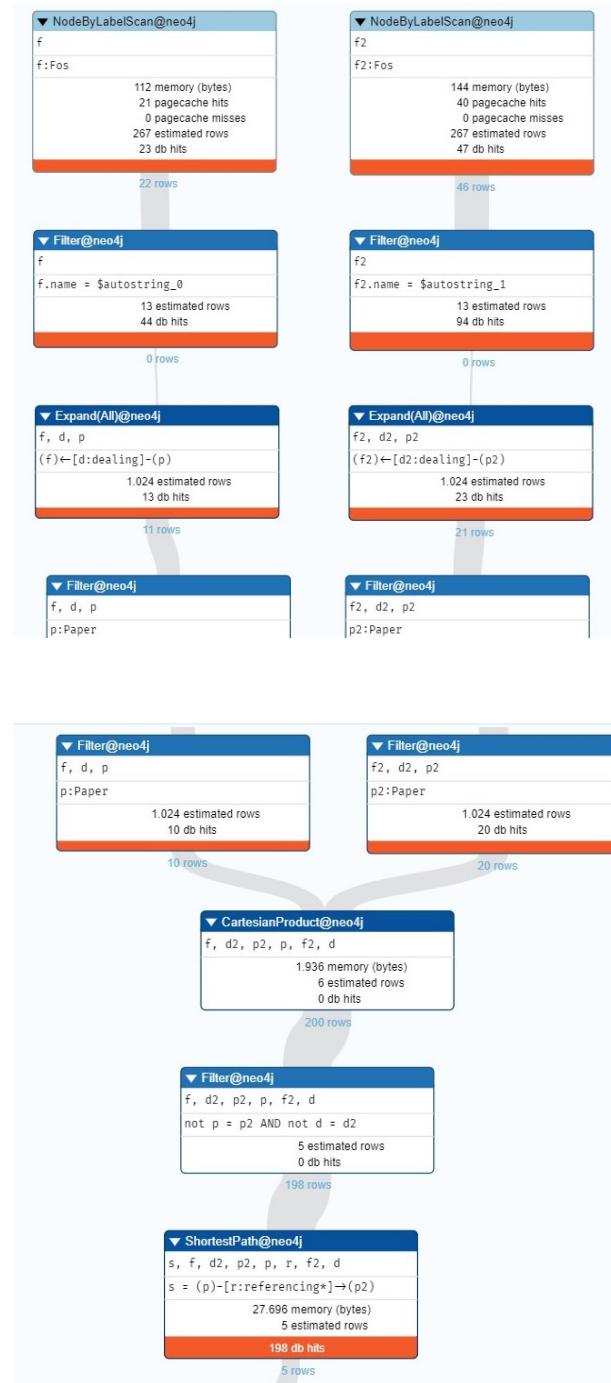
```

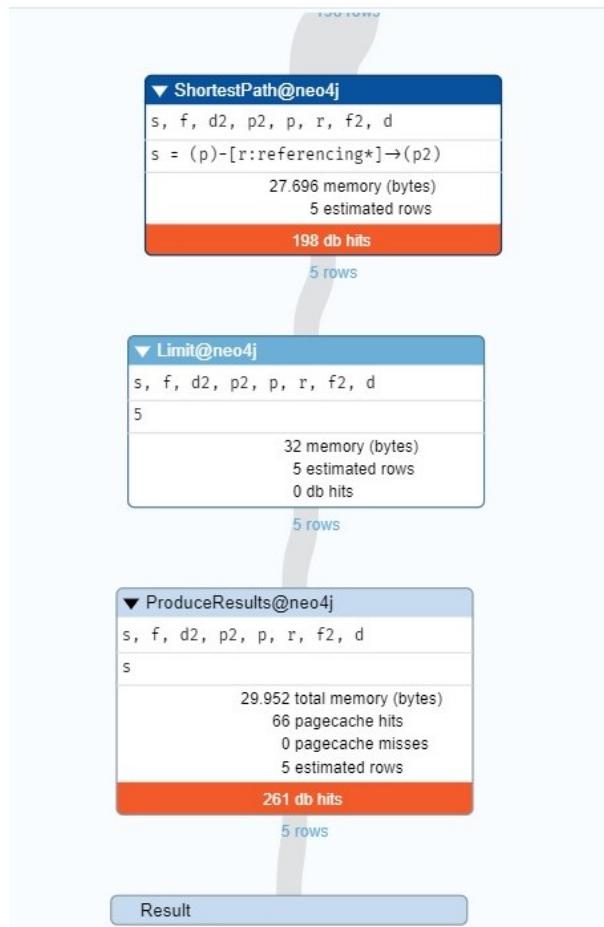
1 MATCH (p: Paper)-[d: dealing]->(f:Fos), (p2: Paper)-[d2:dealing]
      ->(f2:Fos), s = shortestPath( (p: Paper)-[r:referencing*]->(p2:
      Paper) )
2 WHERE f.name = "Artificial intelligence" AND f2.name = "Machine
      learning" AND p <> p2
3 RETURN s
4 LIMIT 5

```



In the following page we can see that the profile statement scans the FoS, applies the filters to them and expands the 'dealing' relationship. Then it checks that paper1 and paper2 are different and in the end it calls the function 'shortestPath' and returns the results.





#### 1.4.4. Creation/Update Commands

1. Insertion of a new author in the database:

```
1 CREATE (a:Author {id: 3000000000, name:"Mario Rossi", organization:
  "Politecnico di Milano"})
```

2. Insertion of a new field of study in the database:

```
1 CREATE (f:Fos {weight:0.57640203, name:"Ambient intelligence" })
```

3. Insertion of a new venue in the database:

```
1 CREATE (v:Venue { id:3000000001, name:"Journal of Cloud Computing"
  })
```

4. Insertion of a new paper in the database, paper's relationships with authors, venue, fields of study and other papers are created:

```
1 MATCH (a:Author) WHERE a.id=3000000000
2 MATCH (f:Fos) WHERE f.name="Cloud computing" AND f.weight=0.6410412
```

```
3 MATCH (v:Venue) WHERE v.id=3000000001
4 MATCH (ref:Paper) WHERE ref.doi="10.1109/JPROC.2015.2483592"
5 MATCH (cit:Paper) WHERE cit.doi="10.1007/s10845-008-0158-5"
6 CREATE (p:Paper { id: 3000000000, title: "Application of
    deterministic, stochastic, and hybrid methods for cloud provider
    selection", year: 2022, doi: "10.1186/s13677-021-00275-1",
    volume: "11" , issue: "1" , abstract: "Cloud Computing
    popularization inspired ... requests.", n_citation: 5,
    page_start: 5, page_end: 10, publisher: "SpringerOpen", doc_type
    : "Journal" })
7 CREATE (p)-[w:writing]->(a), (p)-[r:referencing]->(ref), (cit)-[r:
    referencing]->(p), (p)-[i:in_]->(v), (p)-[d:dealing]->(f)
```

##### 5. Update of the organization of an author:

```
1 MATCH (a:Author {name: 'Stefano Ceri'})
2 SET a.organization = "Politecnico di Milano"
```

##### 6. Update of the number of citations of a paper:

```
1 MATCH (p:Paper {title: "Markov localization for mobile robots in
    dynamic environments"})
2 SET p.n_citation = 728
```



# 2 | Second Delivery

## 2.1. Introduction

This project aims at building an Information System that manages a data set containing different type of scientific articles that can be used for: clustering with network and side information, studying influence in the citation network, finding the most influential papers and topics, modeling analysis.

For the same problem described in Chapter 1 (First Delivery), we modeled a different system following the subsequent steps.

At first the data set was pre-processed to add the fields that were missing in the original data set. Such fields were taken from different Kaggle data sets. Then it was reduced in size and after that it was uploaded on MongoDB.

At the end of the project 11 queries and 6 creation/update commands were initiated with a different level of complexity that was checked within the performance time.

## 2.2. Data

This is the structure of a paper inside the database, it contains different sub-documents: authors, fos (metadata), venue and sections.

```

1 {
2   "id": "...",
3   "title": "...",
4   "abstract": "...",
5   "authors": [
6     {
7       "name": "...",
8       "org": "...",
9       "email": "...",
10      "bio": "..."
11    }
12  ],
13 }
```

```

11 "fos": [
12   "name": "...",
13   "weight": ...
14 ],
15 "year": ...,
16 "page_start": ...,
17 "page_end": ...,
18 "doc_type": "...",
19 "publisher": "...",
20 "volume": ...,
21 "issue": ...,
22 "doi": "...",
23 "n_citation": ...,
24 "venue": {
25   "raw": "...",
26   "id": "...",
27 },
28 "references": [ ... ],
29 "sections": [
30   {
31     "id": ...,
32     "title": "...",
33     "text": "...",
34     "subsections": [
35       {
36         "id": ...,
37         "title": "...",
38         "text": ...
39       }
40     ],
41     "figures": [
42       {
43         "url": "...",
44         "caption": ...
45       }
46     ]
47   }
48 ]
49 }

```

This is the description of a paper:

**Paper** is a scientific article that is associated with the following fields:

- *id*: Int;
- *title*: String;
- *abstract*: that is the summary of the paper: String;

- *authors*:
  - *name*: String;
  - *organization*: String;
  - *email*: that is a personal contact of the author: String;
  - *bio*: a short introduction written by the author (it is added from other sources): String;
- *FOS (Field of Study)*:
  - *name*: String;
  - *w*: that is the weight of the fields of study: Float;
- *Publication Details*:
  - *year*: that corresponds to the publication year: Int;
  - *page\_start* and *page\_end*: that are the starting and the ending point of the collection from which the paper was extracted: Int;
  - *doc\_type*: that is the type of the paper, and can assume 3 values: *Journal, Conference, Patent*: String;
  - *publisher*: String;
  - *volume*: that corresponds to the n-th published collection: Int;
  - *issue*: that corresponds to the m-th part of the volume: Int;
  - *doi*: that is the Digital Object Identifier: String;
  - *n\_citation*: that is the number of citations: Int;
- *Venue*: that is the collection from which the paper was extracted, with the attributes:
  - *id*: String;
  - *raw*: that is the name of the collection: String;
- *References*:
  - *id*: String;

- *Sections*: every section has an id, a title and a content (textual), it could have some subsections too (with the same structure of a section). Furthermore every section has one or more images. Every image has an url and a caption:
  - *id*: Int;
  - *title*: String;
  - *text*: String;
  - *subsections*:
    - \* *id*: Int;
    - \* *title*: Int;
    - \* *text*: String;
  - *figures*:
    - \* *url*: String;
    - \* *caption*: String;

### 2.2.1. Data Pre-Processing

In this part of the project we used as a starting point the same data set used in the first delivery, which can be downloaded [here](#). However, the pre-processing this time was slightly different. Roughly the same operations of the first delivery were performed to clean and reduce the size of the data set. Those operations are performed with the script **dataset\_preprocessing.py**, which can be found in the *scripts* folder of the *mongodb* section of the repository. It is worth to notice that such script does not split the original data set into multiple ones like in the first delivery. Instead, a single data set is kept.

### 2.2.2. Data Completion

We decided to avoid working with synthetic data (i.e. generate random meaningless data) and instead we decided to add to our dataset the missing field by retrieving data from other data sets, trying to be as accurate as possible. The missing field w.r.t. the description in section 2.2 were:

- *authors.email*
- *authors.bio*

- *sections*

We generated the email of every author and we inserted their bio, that were taken from Kaggle's Goodread-Authors data set. This is performed with the script **update\_author.py**.

Furthermore we added the *Sections* part, using a Twitter dataset, that contains millions of tweets in different languages. For every paper we generated randomly a number of sections in the range from 1 to 3 and for each one from 0 to 2 subsections. For the title we used the text of one tweet and for the content the text of four tweets.

Then for the figures we used another data set (Train-GCC-training.tsv from the Google Conceptual Captions official website) in order to take the *image.url* and the *image.caption* fields. Every section contains a random number of figures between 1 and 3.

Those two operations above are executed in the notebook **sections\_preprocessing.py**, that can be found in the *mongodb* section of the GitHub repository.

The final data set consists of 8333 documents.

To obtain the final data set starting from the downloaded one, run the scripts in this order:

1. **dataset\_preprocessing.py**
2. **deprecated\_references.py**
3. **bio\_preprocessing.py**
4. **update\_author.py**
5. **sections\_preprocessing.py**

## 2.3. MongoDB

### 2.3.1. Data Upload

To import the data into a MongoDB collection, we rely on another Python script that exploits the PyMongo library, the official MongoDB driver for Python. We used the script `import_data.py`, that can be found in the `scripts` folder of the `mongodb` section of the GitHub repository. The script dumps the whole data set into a MongoDB collection, while preserving its complex and nested structure. It also performs some pre-processing on the `id` and on the `references` fields:

- The `id` field is renamed to `_id` in order to be used as an index inside of the database
- Both the `_id` field and the `references` field are converted to an ObjectId.

Once the correct connection parameters and the correct file path have been specified, it is enough to run the script with Python to import the data. We will use `papers` as the name of the collection.

### 2.3.2. Document Example

Below we can see the general structure of a document, with the help of MongoDB Compass.

```
_id: ObjectId('101421652000000000000000')
title: "The influence of query interface design on decision-making performance"
> authors: Array
> venue: Object
year: 2003
n_citation: 139
page_start: 397
page_end: 423
doc_type: "Journal"
publisher: "Society for Information Management and The Management Information Syst..."
volume: 27
issue: 3
> flos: Array
doi: "10.2307/30036539"
> references: Array
abstract: "Managers in modern organizations are confronted with ever-increasing v..."
> sections: Array
```

Figure 2.1: General structure

We will now expand the structure of the fields with a complex type.

The *authors* field is an array of sub-documents in which every element represents an author that worked on the paper.

```

    ✓ authors: Array
      ✓ 0: Object
        name: "Luke Olsen"
        id: "2142664686"
        org: "Department of Computer Science, University of Calgary, Calgary, AB, Ca..."
        email: "lukeolsen@mit.edu"
        bio: "Lucie Dufresne was born in 1951 in Trois-Rivières between two rivers: ..."
      ✓ 1: Object
        name: "Faramarz F. Samavati"
        id: "1821145341"
        org: "Department of Computer Science, University of Calgary, Calgary, AB, Ca..."
        email: "faramarzf.samavati@hotmail.com"
        bio: "Ramón González Férriz es editor y periodista. Actualmente, es columnis..."
      ✓ 2: Object
        name: "Mario Costa Sousa"
        id: "2105740368"
        org: "Department of Computer Science, University of Calgary, Calgary, AB, Ca..."
        email: "mariocostasousa@mit.edu"
        bio: "Maja Ilisch, geboren 1975 in Dortmund, studierte Öffentliches Biblioth..."
      ✓ 3: Object
        name: "Joaquim A. Jorge"
        id: "2120678171"
        org: "Departamento de Engenharia Informática, Instituto Superior Técnico, Li..."
        email: "joaquima.jorge@gmail.com"
        bio: "Author of TINCTURE (<a target=_blank href="http://www.tincturestory..."
```

Figure 2.2: Structure of the *authors* field

The *venue* field is a sub-document with two internal fields.

```

    ✓ venue: Object
      raw: "Computers & Graphics"
      id: "94821547"
```

Figure 2.3: Structure of the *venue* field

The *fos* field is another array of sub-documents. Each element corresponds to one of the topics covered by the paper.

```

    ↘ fos: Array
      ↘ 0: Object
        name: "Computer vision"
        w: 0.427553326
      ↘ 1: Object
        name: "Artificial intelligence"
        w: 0
      ↘ 2: Object
        name: "Sketch recognition"
        w: 0.6052215
      ↘ 3: Object
        name: "Geometric modeling"
        w: 0.48624804600000004
      ↘ 4: Object
        name: "Human-computer interaction"
        w: 0.4427773
  
```

Figure 2.4: Structure of the *fos* field

The *references* field is implemented just as an array of ObjectIds.

```

    ↘ references: Array
      0: ObjectId('175816726800000000000000')
      1: ObjectId('191857022600000000000000')
      2: ObjectId('209868554100000000000000')
      3: ObjectId('211598118400000000000000')
      4: ObjectId('212745062100000000000000')
      5: ObjectId('213274786700000000000000')
  
```

Figure 2.5: Structure of the *references* field

Lastly, the *sections* field is an array of sub-documents. Each element of such array corresponds to one section of the paper. Every section can contain another array of sub-documents in the *subsections* field and another array of sub-documents in the *figures* field.

```

✓ sections: Array
  ✓ 0: Object
    id: 1
    title: "İlk önce kendi gücünün bilgisine sahip olmalısın; ikincisi, meydan oku..."
    text: "Soyumuz soylansın, Boyumuz Boylansın.. Al Yıldızlı bayrağımız KUDÜSTE,...""
  ✓ subsections: Array
    ✓ figures: Array
      ✓ 0: Object
        url: "http://lh6.ggpht.com/-IvRtNLNcG8o/TpFyrudaT6I/AAAAAAAAM6o/_11MuAAKalo/..."
        caption: "a very typical bus station    pop artist attends the 3rd annual at que..."
      ✓ 1: Object
        url: "http://lh6.ggpht.com/-IvRtNLNcG8o/TpFyrudaT6I/AAAAAAAAM6o/_11MuAAKalo/..."
        caption: "a very typical bus station    illustration of a map , its flag and a c..."
    ✓ 1: Object
      id: 2
      title: "Gurbete düşmüş bir insan, ne denli varlık içinde bir yaşam sürüyor olsam..."
      text: "Türk milleti vatanını koruyan ordumuzun yanındadır. #BarışPinarıViyadü...""
    ✓ subsections: Array
      ✓ 0: Object
        id: 1
        title: "Hep öлerek çоgaldık... Biz Oğuzun erleri. #DavanınGüçü"
        text: "Bizim milliyetçiliğimiz ayırcı değil birleştirici, çatışmacı değil ba..."
      ✓ 1: Object
        id: 2
        title: ""YPG silah ve malzemeleri bırakıp çekilsin, bu gece harekatı durduralım."
        text: "RT @Enesovvic: Ya siz kimi kimin toprağından kovuyorsunuz? Burası biz..."
    ✓ figures: Array
      ✓ 0: Object
        url: "http://lh6.ggpht.com/-IvRtNLNcG8o/TpFyrudaT6I/AAAAAAAAM6o/_11MuAAKalo/..."
        caption: "a very typical bus station    rock artist performs on stage at awards ..."
  
```

Figure 2.6: Structure of the *sections* field

### 2.3.3. Join Operation

To get a temporary collection in which every document also contains an array with inside all the other documents it references, we can use the `$lookup` operator.

```

1 db.papers.aggregate([
2   {
3     $lookup:
4       {
5         from: "papers",
6         localField: "references",
7         foreignField: "_id",
8         as: "refs"
9       }
10    }
11  ])

```

The referenced documents will be contained in the `refs` field.

### 2.3.4. Queries

In this section queries with a different level of complexity are presented, with a brief description and a figure that shows their results. Notice that, for some queries, that return several documents, the result is only partially shown.

1. Find papers of 2006 with issue equal to 3 and volume greater or equal to 5. Then show just the list of authors with their name and organization, and also the venue of the paper. Limit the result to 2:

```

1 db.papers.find(
2   {"year": 2006, "volume": {"$gte": 5}, "issue": 3},
3   {"authors.name":1, "authors.org":1, "venue":1}
4 ).limit(2)

```

(nReturned: 2, executionTimeMillis: 57, totalDocsExamined: 382)

```

< { _id: ObjectId("168383441500000000000000"),
  authors:
    [ { name: 'A.D. Murugan',
        org: 'Dept. of Electr. & Comput. Eng., Ohio State Univ., Columbus, OH, USA' },
      { name: 'H. El Gamal',
        org: 'Dept. of Electr. & Comput. Eng., Ohio State Univ., Columbus, OH, USA' },
      { name: 'Mohamed Oussama Damen', org: 'University of Waterloo' },
      { name: 'Giuseppe Caire', org: 'Institut Eurecom' } ],
  venue:
    { raw: 'IEEE Transactions on Information Theory',
      id: '4502562' } }
{ _id: ObjectId("196436990000000000000000"),
  authors:
    [ { name: 'Kenneth C. Barr',
        org: 'MIT Computer Science and Artificial Intelligence Laboratory, Cambridge, MA' },
      { name: 'Krsti Asanovic',
        org: 'MIT Computer Science and Artificial Intelligence Laboratory, Cambridge, MA' } ],
  venue: { raw: 'ACM Transactions on Computer Systems', id: '193109227' } }

```

2. Find 3 papers that were cited more than 500 times, published since 2015 in a Journal. Of them we retrieve only: title, publication year, number of citations, publisher, and doi:

```

1 db.papers.find(
2   {"n_citation": {"$gt": 500}, "year": {"$gte": 2015}, "doc_type":
3     → "Journal"}, 
4   {"title": 1, "year": 1, "n_citation": 1, "publisher": 1, "doi":
5     → 1}
6 ).limit(3)

```

(nReturned: 3, executionTimeMillis: 3, totalDocsExamined: 326)

```

< { _id: ObjectId("161299778400000000000000"),
  title: 'ORB-SLAM: A Versatile and Accurate Monocular SLAM System',
  year: 2015,
  n_citation: 619,
  publisher: 'IEEE',
  doi: '10.1109/TRO.2015.2463671' }
{ _id: ObjectId("188518597100000000000000"),
  title: 'Image Super-Resolution Using Deep Convolutional Networks',
  year: 2016,
  n_citation: 641,
  publisher: 'IEEE',
  doi: '10.1109/TPAMI.2015.2439281' }
{ _id: ObjectId("191065790500000000000000"),
  title: 'SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation',
  year: 2017,
  n_citation: 551,
  publisher: 'IEEE',
  doi: '10.17863/CAM.17966' }

```

3. Find 5 papers that contain the word "artificial" in their abstract.

To perform this query an index of type *text*, on the field *abstract* was previously created.

```

1 db.papers.find(
2   {"$text": {"$search": "artificial"}},
3   {"_id": 1, "title": 1, "abstract": 1}
4 ).limit(5)

```

(nReturned: 5, executionTimeMillis: 9, totalDocsExamined: 5)

```

< { _id: ObjectId("202911114700000000000000"),
  title: 'Logic and artificial intelligence',
  abstract: 'Nilsson, N.J., Logic and artificial intelligence, Artificial Intelligence 47 (1990) 31-56. The theoretical foundations of the logical approach to artificial intelligence',
{ _id: ObjectId("212351364800000000000000"),
  title: 'An artificial neural network (p,d,q) model for timeseries forecasting',
  abstract: 'Artificial neural networks (ANNs) are flexible computing frameworks and universal approximators that can be applied to a wide range of time series forecasting problems',
{ _id: ObjectId("202800282200000000000000"),
  title: 'An artificial bee colony algorithm for the capacitated vehicle routing problem',
  abstract: 'This paper introduces an artificial bee colony heuristic for solving the capacitated vehicle routing problem. The artificial bee colony heuristic is a swarm-based local search algorithm that uses the concept of the bee colony to find the optimal solution. The proposed algorithm is compared with other existing algorithms and it is shown that it is able to find better solutions in less time than the others',
{ _id: ObjectId("206693686000000000000000"),
  title: 'Applying artificial intelligence to virtual reality: Intelligent virtual environments',
  abstract: 'Research into virtual environments on the one hand and artificial intelligence and artificial life on the other has largely been carried out by two different groups of researchers. This paper attempts to bridge the gap between these two fields by presenting a survey of the applications of artificial intelligence in virtual reality. The survey covers a wide range of topics, including expert systems, hybrid intelligent systems, and nonlinear systems',
{ _id: ObjectId("212197026200000000000000"),
  title: 'A comparative survey of artificial intelligence applications in finance: artificial neural networks, expert system and hybrid intelligent systems',
  abstract: 'Nowadays, many current real financial applications have nonlinear and uncertain behaviors which change across the time. Therefore, the need to solve highly nonlinear problems has increased significantly. In this paper, we present a comparative survey of artificial intelligence applications in finance. We focus on three main areas: artificial neural networks, expert systems, and hybrid intelligent systems. We discuss the advantages and disadvantages of each approach and compare them based on various criteria. The results show that hybrid intelligent systems are more effective than individual approaches in solving complex financial problems'

```

4. Retrieve the sorted average number of pages, per doc\_type, of papers from a year after 2000, or with a number of citation greater or equal to 100:

```

1 db.papers.aggregate([
2   {"$match": {"$or": [{"n_citation": {"$gte": 100}}, {"year": {"$gte": 2000}}]}},
3   {
4     "$group": {
5       "_id": "$doc_type",
6       "averageNumberPages": { "$avg": {"$subtract": [
7         {"$page_end"}, {"$page_start"}]
8       }},
9       "$sort": {"averageNumberPages": -1}
10    }
11  ]

```

(nReturned: 3, executionTimeMillis: 195, totalDocsExamined: 8333)

```
< { _id: 'Patent', averageNumberPages: 23.428571428571427 }
  { _id: 'Journal', averageNumberPages: 23.040044649086088 }
  { _id: 'Conference', averageNumberPages: 16.00086281276963 }
```

5. Find the 7 papers that contain the most references, among the papers that cover a very relevant field of study (weight  $\geq 0.7$ ), and that were presented in a Conference:

```
1 db.papers.aggregate([
2   {"$match": { "$and": [{"fos.w": {"$gte": 0.7}}, {"doc_type": "Conference"}] }},
3   {
4     "$group": {
5       "_id": "$_id",
6       "n_refs": {"$push": {"$size": "$references"}},
7       "n_sections": {"$push": {"$size": "$sections"}}
8     },
9     {"$sort": {"n_refs": -1}},
10    {"$limit": 7}
11  ])
```

(nReturned: 7, executionTimeMillis: 14, totalDocsExamined: 8333)

```

< { _id: ObjectId("21318424030000000000000000"),
  n_refs: [ 16 ],
  n_sections: [ 2 ] }

{ _id: ObjectId("21350534600000000000000000"),
  n_refs: [ 11 ],
  n_sections: [ 1 ] }

{ _id: ObjectId("20175668840000000000000000"),
  n_refs: [ 10 ],
  n_sections: [ 3 ] }

{ _id: ObjectId("20448330800000000000000000"),
  n_refs: [ 9 ],
  n_sections: [ 3 ] }

{ _id: ObjectId("20248884880000000000000000"),
  n_refs: [ 8 ],
  n_sections: [ 2 ] }

{ _id: ObjectId("19894713220000000000000000"),
  n_refs: [ 8 ],
  n_sections: [ 2 ] }

{ _id: ObjectId("21562001890000000000000000"),
  n_refs: [ 7 ],
  n_sections: [ 3 ] }

```

6. Find the papers whose fos (field of studies) is "Computer engineering" or "Computer science", and that have at least one subsection:

```

1 db.papers.find(
2   {"$and": [ {"$or": [ {"fos.name": "Computer engineering"}, {
3     "fos.name": "Computer science"} ]}, {"sections.subsections": {
4       {"$exists": true}} } ],
5   {"title":1, "sections.title":1 }
6 )

```

(nReturned: 1789, executionTimeMillis: 31, totalDocsExamined: 8333)

```
< ( _id: ObjectId("101567523200000000000000"),
  title: 'Research-paper recommender systems: a literature survey',
  sections:
  [ { title: 'Kızı gekinlikle görünce ağlayan erkek net ılıktır kızlar sıktır edin bakmaz eve duygusal püşt.' },
    { title: 'Değişimi ve devrimi somuna kadar giticeğiz. Korku İmparatorluğu değil sevgiyeli egemen kılacağız. Kardeşçe beraber olacağız. #GameOfCells' },
    { title: 'Türkiye, Fırat Kalkanı Harekät'iyla birlikte gerçek anlamda bağımsızlığını kavuşturma sürecine girdi. O yuzden Fırat Kalkanı Harekätı, Türkiye'nin istiklal ve istikametini temsil ediyor.' },
  ],
  _id: ObjectId("101567523200000000000000"),
  title: 'The State of the Art in Text Filtering',
  sections: [ { title: 'Rus ve Kürtçe ortaklar karışımlı... Başsavcılık: Darbe girişimi engellendi - https://t.co/3cv3DAcrlp https://t.co/JUXLGThWjI' } ]
  ],
  _id: ObjectId("121952765400000000000000"),
  title: 'Business intelligence in blogs: understanding consumer interactions and communities',
  sections:
  [ { title: 'İlk önce kendin gücün bilgisine sahip olmalıdır; ikincisi, meydan okumaya cesaretin olmalı; ve üçüncü, yapacak inancı sahip olmalıdır. #DavanınGüçü' },
    { title: 'Gurbete düşgümüz bir insan, ne denli varlık içinde bir yaşam sürüyor olsa da doğup bulunduğu yeri arar. Denedim koynunda yattıkça benimsin ey güzeli toprak, nefer yapma' }
  ]
}
```

## 7. Find 10 papers with author affiliated with "IEEE", that have 3 sections and at least a figure's caption containing the word "bus":

```
1 db.papers.find(
2   { "$and": [ {"authors": { "$elemMatch": { "org": "IEEE" }}}, 
3     {"sections": { "$elemMatch": { "figures.caption": { "$regex": "/bus/}}}}, {"sections": { "$size": 3}}]}, 
4   {"title": 1, "year": 1, "n_citation": 1, "publisher": 1, "doi": 1, "sections": 1}
5 ).limit(10)
```

(nReturned: 10, executionTimeMillis: 21, totalDocsExamined: 7247)

```
< ( _id: ObjectId("200201661200000000000000"),
  title: 'Performance guarantees for Web server end-systems: a control-theoretical approach',
  year: 2002,
  n_citation: 539,
  publisher: 'IEEE Computer Society',
  doi: '10.1109/TI.980028',
  sections:
  [ { id: 1,
    title: 'RT @C66wCY4hr: #W68k12BmYjYSYj+7ye#LRSjnZE3LWNjw=: Eşgülörbornan Dış etmeden Yılmadan Dördüncüdürüm hak yolumuzda, Adım adım demokratik haklarımızın pevresinde yetişti. RT @C66wCY4hr: #W68k12BmYjYSYj+7ye#LRSjnZE3LWNjw=: Hayırlı olsun #WfHRCjE7c4rQExpxo2u00B6+OybzrSAFTCK0j0= Federasyonlarda yolunda bir adım daha ilerledik. #Fakultetlerde' },
    subsections:
    [ { id: 1,
      title: 'RT @izzet29723814: Hıighb Esnaf: Ticari itibarını, İşmini, Çocuğu gibi boyuttuğu firmasını, Çoklerinin yazılmasını istemey. Çeklerini ödeye.', 
      text: 'RT @Altinn_n: Demirtaşın çizgisini beğeneler bu teröristlere cesaret verdi. #Allahbelanızıversin #Hakkari pkk'nın yanında olanlar Allah RT @hulyayurt_:#' },
      figures:
      [ { url: 'http://lh6.ggpht.com/-IvRtNNGGBo/tPfryudax61/AAAAAAAAM6o/_1IMuAAKaiQ/IMG_3422.JPG?imgmax=800 https://lh6.googleusercontent.com/proxy/KSGMGN... Name: 5699, caption: 'a very typical bus station politician plays the piano at a charity concert. Name: 5699, type: object' } ],
      id: 2,
      title: 'RT @hulyayurt_: Geçenin özeti nedir biliyor musunuz? Kronik bir yalancının yüzüne yüzme milyonlarca insanın şahit olduğu "yalan söyleyorsun.' ,
      text: 'RT @Semihardic: AK Parti Sozcumuz Romercoşelik CHP, Kulliyeye giden CHP li İddiayıla ısrarla bu yalan siyasetini sürdürmeye devam etti. Net RT @hulyayurt_:_ Fransa' },
      subsections:
      [ { id: 1,
        title: 'RT @PUSAT01071: Vefatlarının yıl dönümünde Edebiyatımızın iki kıymetli değerli İsmi #AbdurrahimKarakoç ve #CahitZarifoğlu nu saygı ve Rahmetle', 
        text: 'RT @Semihardic: Parti Sozcumuz @semihardic Duyanın her tarafında bir takım Türkiye karşıtları kendi hükümetlerini işleyasette Türkiye'nin RT @TayfunMelekKar' },
        id: 3,
        title: 'RT @Nurrr1980: Sokakta top oynayan çocukların yerde 100 lira buluyor ve camiye bırakıyorlar siz de şubesinizin açık yanken eğilir RAMAZANı Bahtı.', 
        text: 'RT @Nurrr1980: RabbinizARLAN,fikrimiz zirkuların,kabimizin nuru Resulullah ,evelimiz ALLAH,rehberimiz Kelâullah,Kubbimiz hayırlara lisan RT @1968_simek: RT' },
        figures:
        [ { url: 'http://lh6.ggpht.com/-IvRtNNGGBo/tPfryudax61/AAAAAAAAM6o/_1IMuAAKaiQ/IMG_3422.JPG?imgmax=800 https://media.gettyimages.com/photos/roy-hodges... Name: 5700, caption: 'a very typical bus station football player , manager applauds his team du... Name: 5700, type: object' },
          { url: 'http://lh6.ggpht.com/-IvRtNNGGBo/tPfryudax61/AAAAAAAAM6o/_1IMuAAKaiQ/IMG_3422.JPG?imgmax=800 https://ak9.picdn.net/shutterstock/videos/6038... Name: 5701, caption: 'a very typical bus station footage of air bubbling up through water which... Name: 5701, type: object' },
          { url: 'http://lh6.ggpht.com/-IvRtNNGGBo/tPfryudax61/AAAAAAAAM6o/_1IMuAAKaiQ/IMG_3422.JPG?imgmax=800 https://i.pinimg.com/736x/b2/07/f7/b207f73b70d... Name: 5702, caption: 'a very typical bus station coaches name players following win over americ... Name: 5702, type: object' } ] ],
      id: 3,
      title: 'RT @OzlemYucel172: #Bismillah okunan #Ezanlar a handolsun uyanan gözlerle kalplerde şükürler olsun. Allâhim huzuruma 'geldim beni senden sevgi.', 
      text: 'RT @Semihardic: Vatan İhneeler le, Toplularla bir olup yurumek isteyenlerin değil.. Yan Koyumna bay Koymadan #şahit düşünlərinidir. RT @OzlemYucel172: Elif i' subsections:
      [ { id: 1,
        title: 'RT @hulyayurt_ : Medeniyetin başlığı () Fransa kadınlarla ayrı wagon verdi. Kadınlar için pembe otobüsler yaptığımızda tüm feministler, laikle.', 
        text: 'RT @Sozbirsamattir: 70 lerde olsak bir sur plak alırdım sana, 80 lerde açık hava sinemasına götürür, izledikten sonra muhalihci ısmarlardı RT @OzlemYucel172: subsections:
        [ { url: 'http://lh6.ggpht.com/-IvRtNNGGBo/tPfryudax61/AAAAAAAAM6o/_1IMuAAKaiQ/IMG_3422.JPG?imgmax=800 https://st2.depositphotos.com/5616164/8201/v/4... Name: 5703, caption: 'a very typical bus station gates of paradise on a white background. Name: 5703, type: object' } ] ] }
```

8. Retrieve for each venue the content of the papers, with more than 500 citations, and volume greater than 12, published in it.

This is done by first matching the conditions stated above, then unwinding on the sub-document "sections". In the end grouping by venue, and retrieving the sections, we obtain the following results (for each venue, an array with the contents collected from all the papers published in it):

```

1 db.papers.aggregate([
2   {"$match": {"$and": [{"n_citation": {"$gte": 500}}, {"volume": 
3     {"$gt": 12}}]}},
4   {"$unwind": {"path": "$sections"}},
5   {"$group": {"_id": "$venue", "content of papers": {"$push": 
6     "$sections"}}}
])

```

(nReturned: 923, executionTimeMillis: 126, totalDocsExamined: 8333)

```

{
  "_id": {
    "raw": "International Journal of Intelligent Systems",
    "id": "079500554"
  },
  "content of papers": [
    { "_id": 1,
      "title": "RT @1986susem: Bayatta en güzel şey,Kimine göre mutluluk,Kimine göre sevgidir,Kimine göre parıdır,Hastaya sorsan sağlık,Yalnızca sorsan yold...",
      "text": "RT @1986susem: 1/Bugün pazarRTesi,Zaferi kuvvetli olan kazanır,Malazgirt ve 30 Ağustos Zafer Bayramını Milletimizle kutuyoruz,Zafer inanancı RT #hanedan_ozmu",
      "subsections": [
        { "_id": 1,
          "title": "RT @YusufEkiyildirim: @HikmetYT1 @HikmetY04297948 @Hach11 @bestope @ERTErdogan @tosavumma @MnstafaSontop @SSadihBilic @sdhilic32 @mehmetcumcun @...",
          "text": "RT @GalatasaraySK: Aydınlanma, paçḍaplaǵma ve baǵımsızlıǵımızın simgesi Cumhuriyet'iminin ilanının 96. yıl dönümünde 29 Ekim Cumhuriyet Bay RT #seluk58",
          "figures": [
            { "url": "http://ih6.gppht.com/-IvRtNLNcG8o/TpFyrudaT6I/AAAAAAAAM6o/_1IMuAAKaiQ/IMG_3422.JPG?imgmax=800 https://ak9.picdn.net/shutterstock/videos/8879... Name: 23527, dtype: object" },
            { "url": "http://ih6.gppht.com/-IvRtNLNcG8o/TpFyrudaT6I/AAAAAAAAM6o/_1IMuAAKaiQ/IMG_3422.JPG?imgmax=800 https://media.gettyimages.com/photos/actor-jos... Name: 23528, dtype: object" }
          ]
        }
      ]
    }
  ],
  "content of papers": [
    { "_id": 1,
      "title": "Yeni tek başıma sinemaya gidiyorum sk",
      "text": "Sigara içen kız iticiliği der susarım.Dünyayı Big Mac yonetsinArkasını yamamaya çalışan knk sevgilisi iticiliği ???Yatak sal beni gün bitiyorrr",
      "subsections": [],
      "figures": [
        { "url": "http://ih6.gppht.com/-IvRtNLNcG8o/TpFyrudaT6I/AAAAAAAAM6o/_1IMuAAKaiQ/IMG_3422.JPG?imgmax=800 http://17.alamy.com/zooms/6f7d9ef8d6549e4ba8c... Name: 1607, dtype: object" }
      ]
    },
    { "_id": 2,
      "title": "'Son Dakika' Teror örgütü propagandası: şurpa, suyu ve suyuñun övme devleti alemin aşığılaşma suçlarından 27. dönem milletvekilleri Gülistan Kılıç Kocigitit, İ",
      "text": "Satıcı ol, alıcı ol, Kalıcı ol, bulucu ol, ama BÖLÜCÜ OLMA.. Davet et, bayret et, Af et, tövbe et, ama İHANET ETME.. #TümHainlerBirleşmiş https://t.co/Wkw",
      "subsections": [
        { "_id": 1,
          "title": "'Kırmızı listede aranan PKK'nın Sincar bölgesinde en üst düzey kadın yetkilisi Berat Arığın MIT'in operasyonu sonucunda yok edildi. https://t.co/w6FPO",
          "text": "Turkey, to collect securing the borders, fighting names regardless of ideology or against all terrorists. Our country is determined to clear Deash and",
          { "_id": 2,
            "title": "'QCokularının dağa kaçırılmasından HDP yi sorulmuş tutan Diyarbakır annelerinin, partinin İl binası önünde 3 Eylül de başlattığı oturma eylemi 72. günün",
            "text": "Twitter'dan skandal: '#TurkeyFightsISISandYPG' etiketine sansür Cumhurbaşkanı Recep Tayyip Erdoğan ile ABD Başkanı Donald Trump'ın görüşmesi öncesi bi",
            "figures": [
              { "url": "http://ih6.gppht.com/-IvRtNLNcG8o/TpFyrudaT6I/AAAAAAAAM6o/_1IMuAAKaiQ/IMG_3422.JPG?imgmax=800 http://i2.cdn.turner.com/money/dam/assets/1603... Name: 1608, dtype: object" },
              { "url": "http://ih6.gppht.com/-IvRtNLNcG8o/TpFyrudaT6I/AAAAAAAAM6o/_1IMuAAKaiQ/IMG_3422.JPG?imgmax=800 http://17.alamy.com/zooms/8d62bfa7ff584494bc79... Name: 1609, dtype: object" }
            ]
          }
        ]
      ]
    },
    { "_id": 1,
      "title": "Aylardır -- Öğretmenler mutsuz, --veliler umutsuz --Öğrenciler huzursuz. Benim çocukların #Dogakoleji nde okuması da, Ben de katılıyorum bu çiğdeja #Dogak",
      "text": "Çunku günden bizim Çunku Haber biziz Orneği olmayan bir hak arayışı bu #BasinYerDiyorBu məcədale Bu güzel aile Daha coook konusulur Çunku biz haklıyız ve",
      "subsections": [],
      "figures": [
        { "url": "http://ih6.gppht.com/-IvRtNLNcG8o/TpFyrudaT6I/AAAAAAAAM6o/_1IMuAAKaiQ/IMG_3422.JPG?imgmax=800 http://c8.alamy.com/comp/D41702/view-from-high... Name: 5078, dtype: object" },
        { "url": "http://ih6.gppht.com/-IvRtNLNcG8o/TpFyrudaT6I/AAAAAAAAM6o/_1IMuAAKaiQ/IMG_3422.JPG?imgmax=800 https://www.sanctuary-care.co.uk/sites/default... Name: 5079, dtype: object" }
      ]
    },
    { "_id": 2,
      "title": "RT @yusuf_oxzeli: Dokunmayın doğuya.. Elinizde çekin nefesinden #Kazdağın dokunuşuma https://t.co/FZDIO86Yi",
      "text": "RT @yusuf_oxzeli: İye başladığımız şartlarla, Sık ile yaptığıımız sözleşmeye göre emeklilik hakkını istiyorum #EmeklilikteYasaTakilanlar RT #aslandegirmen",
      "subsections": [],
      "figures": [
        { "url": "http://ih6.gppht.com/-IvRtNLNcG8o/TpFyrudaT6I/AAAAAAAAM6o/_1IMuAAKaiQ/IMG_3422.JPG?imgmax=800 http://www.secretsofparis.com/storage/newslett... Name: 5080, dtype: object" },
        { "url": "http://ih6.gppht.com/-IvRtNLNcG8o/TpFyrudaT6I/AAAAAAAAM6o/_1IMuAAKaiQ/IMG_3422.JPG?imgmax=800 http://www.brian-coffee-spot.com/wp-content/up... Name: 5081, dtype: object" }
      ]
    },
    { "_id": 3,
      "title": "RT @yusuf_oxzeli: Dokunmayın doğuya.. Elinizde çekin nefesinden #Kazdağın dokunuşuma https://t.co/FZDIO86Yi",
      "text": "RT @yusuf_oxzeli: İye başladığımız şartlarla, Sık ile yaptığıımız sözleşmeye göre emeklilik hakkını istiyorum #EmeklilikteYasaTakilanlar RT #aslandegirmen",
      "subsections": [],
      "figures": [
        { "url": "http://ih6.gppht.com/-IvRtNLNcG8o/TpFyrudaT6I/AAAAAAAAM6o/_1IMuAAKaiQ/IMG_3422.JPG?imgmax=800 http://www.secretsofparis.com/storage/newslett... Name: 5082, dtype: object" },
        { "url": "http://ih6.gppht.com/-IvRtNLNcG8o/TpFyrudaT6I/AAAAAAAAM6o/_1IMuAAKaiQ/IMG_3422.JPG?imgmax=800 http://www.brian-coffee-spot.com/wp-content/up... Name: 5083, dtype: object" }
      ]
    }
  ]
}

```

9. Find the number of papers of the first 10 authors, affiliated to an University, that wrote most papers in Journals after year 2000.

This is done by first matching the conditions, and unwinding on the sub-documents "authors". Then by grouping by *authors.id*, the number of papers for each author can be computed. In the end the results are sorted, and the first 10 are shown.

```

1 db.papers.aggregate([
2   {"$match":{ "$and": [ {"doc_type": "Journal"}, {"year": {"$gt": 
→ 2000}}, {"authors.org": {"$regex": /University/}}]}},
3   {"$unwind": {"path": "$authors"}},
4   {"$group": {"_id": "$authors.id", "n_papers": {"$sum":1}}},
5   {"$sort": {"n_papers": -1}},
6   {"$limit": 10}
7 ])

```

(nReturned: 10, executionTimeMillis: 256, totalDocsExamined: 8333)

```

< { _id: '2141382980', n_papers: 26 }
  { _id: '2104129307', n_papers: 16 }
  { _id: '224175856', n_papers: 16 }
  { _id: '2149762431', n_papers: 13 }
  { _id: '2141728717', n_papers: 11 }
  { _id: '2231782831', n_papers: 10 }
  { _id: '737083156', n_papers: 10 }
  { _id: '2299437103', n_papers: 10 }
  { _id: '2469405535', n_papers: 10 }
  { _id: '2104966155', n_papers: 9 }

```

10. Find the referenced papers whose *fos* (field of studies) is "Data mining" and whose volume is 3, and that have at least one section:

```

1 db.papers.aggregate([
2   {"$match": { "$and": [ {"fos.name": "Data mining"}, {"volume": 
→ 3}, {"sections": {"$exists": true} } ] } },
3   {
4     "$lookup":
5       {
6         "from": "Project2",
7         "localField": "references",
8         "foreignField": "_id",

```

```

9      "as": "refs"
10     }
11   }
12 ])

```

(nReturned: 24, executionTimeMillis: 37, totalDocsExamined: 8333)

```

< { _id: ObjectId("101567523200000000000000"),
  title: 'Research-paper recommender systems: a literature survey',
  authors:
  [ { name: 'Joern Beel',
    id: '2032888927',
    org: 'Docer, Magdeburg, Germany#TAB#',
    email: 'joernbeel@yahoo.com',
    bio: 'For years, Terri Savelle Foy's life was average. She had no dreams to pursue. Each passing day was just a repeat of the day before. Finally, with a marriage in trou
    { name: 'Bela Gipp',
      id: '72611330',
      org: 'University of Konstanz, Konstanz, Germany#TAB#',
      email: 'belagipp@mit.edu',
      bio: 'Phil Nisman is the co-author of the #1 national best selling true crime book <a href="https://www.goodreads.com/book/show/28525592_Gitchie_Girl" title="Gitchie Gi
    { name: 'Stefan Langer',
      id: '2135709281',
      org: 'Otto-von-Guericke University, Magdeburg, Germany#TAB#',
      email: 'stefan.langer@mit.edu',
      bio: 'John "Red" Shea, 40, was a top lieutenant in the South Boston Irish mob run, led by James "Whitey" Bulger. An ice-cold enforcer with a red-hot temper, Shea was a le
    { name: 'Corinna Breitinger',
      id: '2063223331',
      org: 'Linnaeus University, Kalmar, Sweden#TAB#',
      email: 'corinna.breitinger@mit.edu',
      bio: 'Mary Kelly was an English crime writer best known for the Inspector Brett Nightingale series. Writing in the 1950s and 1960s, Kelly was celebrated for the sense of
  venue:
  { raw: 'International Journal on Digital Libraries',
    id: '1106155884' },
  year: 2016,
  n_citation: 106,
  page_start: 305,
  page_end: 338,
  doc_type: 'Journal',
  publisher: 'Springer Berlin Heidelberg',
  volume: 17,
  issue: 4,
  fosc:
  [ { name: 'Computer science', w: 0.405663 },
    { name: 'Information needs', w: 0.5010579360000001 },
    { name: 'Descriptive statistics', w: 0.4312128999999999 },
    { name: 'Implementation', w: 0.4495434 },
    { name: 'Information retrieval', w: 0.44728106300000003 } ],
  doi: '10.1007/s00799-015-0156-0',
  references:
  [ ObjectId('197104055000000000000000'),
    ObjectId('201245115200000000000000'),
    ObjectId('201944326400000000000000'),
    ObjectId('205011383800000000000000'),
    ObjectId('211285679700000000000000'),
    ObjectId('211518608700000000000000'),
    ObjectId('212451972500000000000000'),
    ObjectId('212533036900000000000000'),
    ObjectId('213933818200000000000000'),
    ObjectId('217196077000000000000000'),
    ObjectId('244249597300000000000000') ],
  abstract: 'In the last 16 years, more than 200 research articles were published about research-paper recommender systems. We reviewed these articles and present some descriptive sections:
  [ { id: 1,
    title: 'Kızı gelenlikle giরুনে আগুণ একে নেতৃত্বে কিছি সংক্ষিপ্ত বাকস এবং দৃঢ়ান্বল পৃষ্ঠা।',
    text: 'Dügünseme memoni aşıyorsun tak 100 k takipçin var. Keşke memom olsa.Emaneti çektiğim gorisini adalet duşunsun.İlk okul mezunu adam çaycılak yapıp 10 bin lira maaş al
  subsection:
  [ { id: 1,
    title: 'Sinema, duygular, duşler ve ıçgûdu dünyalarını anlatmak için en iyi araçtır. Luis Bunuel @mustafayalcin_ @talasbelediyesi #FestivalReyecaniTalasta',
    text: 'RT @yilmaazulu20: YARU Bİ ANGET YAPALIM DEDİK, DOSTA DÜŞMANA KARŞI.. KARŞIDAN TEK Bİ KİŞİ DÜŞTU.. ENGELİ YEDİ.. BİZİMKİLERE NE OLUYOR ALLAH'a RT @abdullahohuk11:
  figures:
  [ { url: 'http://lh6.ggpht.com/-IvRtNINcGBo/TpYryrdaT6I/AAAAAAAAM6o/_IMuAAKAkQ/IMG_3422.JPG?imgmax=800' https://thumbl.shutterstock.com/display_pic_wl... Name: 3, dtype: object' },
    caption: 'a very typical bus station - cybernetic scene isolated on white background - Name: 3, dtype: object' },

```

```

        {
            "url": "http://lh6.ggpht.com/-IvRtNLNgG8o/TpFyruTa6I/AAAAAAAAM6o/_1lMuAAKAlQ/IMG_3422.JPG?imgmax=800 https://media.gettyimages.com/photos/jayz-atte... Name: 4, dty: object' } ] },
            "_id": 2,
            "title": "Birliğimizdir ve devrinin sonuna kadar gideceğiz. Korku imparatorluğu değil sevgiyi egemen kılacajız. Kardeşce beraber olacağiz. #GameOfCels
            text: 'Birlikte yürüyecok daha çok yolumuz var. Aşkimiz, sevdamız, yârimiz, yaranimiz, vatanimiz var. #Erdoðanla2023Yolundaþimdi, Elini uzat, Başlasın en gacılı devir. Yet
            subsections: [],
            figures: [
                { "url": "http://lh6.ggpht.com/-IvRtNLNgG8o/TpFyruTa6I/AAAAAAAAM6o/_1lMuAAKAlQ/IMG_3422.JPG?imgmax=800 https://prismpub.com/wp-content/uploads/2016/1... Name: 5, dty: object' } ] },
            "_id": 3,
            "title": "Türkiye, Fırat Kalkanı Harekâti'yla birlikte gerekçanla bağımsızlığını kavuþma sürecine girdi. O yüzden Fırat Kalkanı Harekâti, Türkiye'nin istiklal ve istik
            text: 'Cumhurbaskaný Erdoðan, tüm Müslümanlar olarak, bir olup, bir duvarın tuþlaları gibi dayanışma içerisinde hareket edildiðinde, onumuzda hiçbir engelin dayanamayaca
            subsections: [
                { "_id": 1,
                    "title": "Hala su gereþi anlamak istemiyorlar. Türkiye'de siyaset sahnesinde 17 yıldır tek bir ADAM var O partiden gok ote lider farkı. Yani Recep Tayyip Erdoðan. I
                    text: 'Londra merkezi uluslararası haber ajansı Reuters da, Erdoðan in açıklamasını abonelelerine, Erdoðan, Türkiye nin Iran la petrol ve doğal gaz ticaretine devam
                figures: [
                    { "url": "http://lh6.ggpht.com/-IvRtNLNgG8o/TpFyruTa6I/AAAAAAAAM6o/_1lMuAAKAlQ/IMG_3422.JPG?imgmax=800 https://thumbl.shutterstock.com/display_pic_wi... Name: 6, dty: object' },
                    { "url": "http://lh6.ggpht.com/-IvRtNLNgG8o/TpFyruTa6I/AAAAAAAAM6o/_1lMuAAKAlQ/IMG_3422.JPG?imgmax=800 https://media.gettyimages.com/photos/bryan-mcc... Name: 7, dty: object' } ] },
                    "ref": [
                        { "_id": ObjectId("197104055000000000000000"),
                            "title": "Evaluating collaborative filtering recommender systems' },
                        { "_id": ObjectId("201245115200000000000000"),
                            "title": "Web mining for web personalization' },
                        { "_id": ObjectId("201944326400000000000000"),
                            "title": "PROGRESS IN DOCUMENTATION THE COMPLEXITIES OF CITATION PRACTICE: A REVIEW OF CITATION STUDIES' },
                        { "_id": ObjectId("205011383800000000000000"),
                            "title": "Evaluating recommender systems from the user's perspective: survey of the state of the art' },
                        { "_id": ObjectId("211285679700000000000000"),
                            "title": "Recommender systems: from algorithms to user experience' },
                        { "_id": ObjectId("211518608700000000000000"),
                            "title": "Recommender Systems Research: A Connection-Centric Survey' },
                        { "_id": ObjectId("212451972500000000000000"),
                            "title": "Personalisation and recommender systems in digital libraries' },
                        { "_id": ObjectId("212533036900000000000000"),
                            "title": "Explaining the user experience of recommender systems' },
                        { "_id": ObjectId("213933818200000000000000"),
                            "title": "Web Usage Mining as a Tool for Personalization: A Survey' },
                        { "_id": ObjectId("217196077000000000000000"),
                            "title": "Toward the next generation of recommender systems: a survey of the state-of-the-art and possible extensions' },
                        { "_id": ObjectId("244249597300000000000000") }
                    ]
                ]
            ]
        }
    ]
}

```

11. Find papers that are referenced by papers regarding machine learning, written by an IEEE author or published by "IEEE Computer Society" after 2010. Of the referenced papers retrieved with the join, show id, title and year of the ones published before 2005:

```

1 db.papers.aggregate([
2     { "$match": { "$and": [ {"year": { "$gte": 2010}}, {"$or": [
3         {"authors.org": "IEEE"}, {"publisher": "IEEE Computer
4         Society"}]}, {"fos.name": "Machine learning"}]} },
5     {
6         "$lookup": {
7             "from": "papersCollection",
8             "localField": "references",
9             "foreignField": "_id",
10            "pipeline": [ { "$match": { "$expr": { "$lte": ["$year",
11            2005]}}, {"$project": {"_id": 1, "title": 1, "year": 1}}}],
12            "as": "refs"
13        }
14    }
15 ])

```



```
{
  id: 0,
  title: 'RT @burus_huseyin: Örnek gösterdiğiniz ülkelerin yasaları bir kez olsun geriye iştekerlerké vatandasını mahdud etmiş mi? Neden hakkımız ol?',
  text: 'RT @Betzekeriya: ABD 2008 yılında 5510 sayılı yasa ile (28)e düşürecek ağılık sınırının altında emeklilik maasları alacak olan millet biziz.RT @burus_huseyin: @alc subsections:
  [ { id: 1,
    title: 'RT @aYzbFp3PymTnf2WJvtSM9Wt+b1JZmqMjej3ACRjXg=: Çalışma uzun soluklu olmasın. Çünkü #EYYinZamanıYok #RTErdoğan @tcbestepo @vedatbilgin',
    text: 'RT @burus_huseyin: Haklıyız, Varız.. Buradayız... Hakkımızı Alacağız! Sonuma, sonuç alana kadar #ETTiliyizVazgeçmeyorumRT @ArzuLast: Davama Sözüm Var Arka',
    id: 2,
    title: 'RT @One_minute: Millete dağıtılmış diye size emanet edilen A101 kartlarıyla alış veriş yapmak nasıl bir şeretsizliğin tezahürüdür Size d.',
    text: 'RT @MardinBuyuksehir: Mardin Büyükşehir Belediyesi, hava radar rölesi ve 1 Günde yapılan dekoratif ışıklandırma çalışmalarını tamamlandı..RT @malatyafilmfes figures:
  [ { url: 'http://lh6.ggpht.com/-IvRTNLNgG8o/TpFyrudaT6I/AAAAAAAAMGo/_11MuAAKAlq/IMG_3422.JPG?imgmax=800 https://afit4you.com/wp-content/uploads/2016/08... Name: 2827,
    caption: 'a very typical bus station a fit you - yoga for golfers Name: 2827, dtype: object' } ] },
  id: 3,
  title: 'RT @gibi_tokat: @HuzurIslandam_23 Kim ki eden adam be tivitle allahınızı sasırtıyo ekonomiyi olkeyi batırıyo şen kılıç çektiğ diyon allah ak..',
  text: 'RT @MardinBuyuksehir: Mardin Büyükşehir Belediyesi, Kızıltepe Yamanlar Mahallesi nde yol yapım çalışmalarına devam ediyor. https://t.co/OYZZ RT @Kusta_Luka: @cehher subsections:
  [ { id: 1,
    title: 'RT @biri_siz: Bir gün birisi için tweet bildirimini açarsam o gün burayı terk ederim...',
    text: 'RT @DeliKadirrr: GUNADIN! https://t.co/KPQbfBfworRT @Ryab6261098: bir dilek tutum ADI....SEN..... Huzurlu Matlu Akşamlar RT @Rya5626) figures:
  [ { url: 'http://lh6.ggpht.com/-IvRTNLNgG8o/TpFyrudaT6I/AAAAAAAAMGo/_11MuAAKAlq/IMG_3422.JPG?imgmax=800 https://media.musely.com/u/bf3ea812-a38c-4d8d-... Name: 2828,
    caption: 'a very typical bus station ... try adding a hint of lemon or lime to your... Name: 2828, dtype: object' },
    { url: 'http://lh6.ggpht.com/-IvRTNLNgG8o/TpFyrudaT6I/AAAAAAAAMGo/_11MuAAKAlq/IMG_3422.JPG?imgmax=800 https://chumbik.shutterstock.com/display_pic_wi... Name: 2829,
    caption: 'a very typical bus station fire in countryside or rural area that engulf... Name: 2829, dtype: object' },
    { url: 'http://lh6.ggpht.com/-IvRTNLNgG8o/TpFyrudaT6I/AAAAAAAAMGo/_11MuAAKAlq/IMG_3422.JPG?imgmax=800 https://chumbik.shutterstock.com/display_pic_wi... Name: 2830,
    caption: 'a very typical bus station vector illustration of a background for indepe... Name: 2830, dtype: object' } ] ],
  refs:
  [ { _id: ObjectId("20337730550000000000000000"),
    title: 'Automatic Facial Expression Analysis: A Survey',
    year: 2003 },
    { _id: ObjectId("215901723100000000000000000000"),
    title: 'Automatic analysis of facial expressions: the state of the art',
    year: 2000 },
    { _id: ObjectId("216335284800000000000000000000"),
    title: 'Multiresolution gray-scale and rotation invariant texture classification with local binary patterns',
    year: 2002 } ] ]
}
```

### 2.3.5. Creation/Update Commands

#### 1. Deletion of one paper from the database:

```
1 db.papers.deleteOne(
2   { "_id": ObjectId("1014216520000000000000000") }
3 )
```

#### 2. Deletion of some papers, based on some conditions:

```
1 db.papers.deleteMany(
2   {"$and" : [ {"n_citation" : {"$lt": 120} }, {"year": {"$lte": 2005} } ] }
3 )
```

#### 3. Insertion of a new paper:

```
1 db.papers.insertOne({
2   "title": "The influence of query interface design on
3     ↪ decision-making performance",
4   "authors": [
5     {
6       "name": "Cheri Speier",
7       "id": "1973614237",
8       "org": "Michigan State University, East Lansing, MI#TAB#",
9       "email": "cherispeier@polimi.it",
10      }
11  ]}
```

```
9      "bio": "Dr. Cheri Speier-Pero is the Ernst & Young Professor  
→ of Accounting and Information Systems, faculty director of the  
→ MS in Business Analytics program, and interim chairperson of the  
→ Department of Supply Chain Management in the Eli Broad College  
→ of Business. She joined the Broad College faculty in 1998, and  
→ has since served as associate dean for MBA/MS Programs, led  
→ executive development courses, and served on many different  
→ committees at the university, college- and department-levels."  
10     },  
11     {  
12         "name": "Greta L. Polites",  
13         "id": "2305721212",  
14         "org": "School of Management, Bucknell University,  
→ Lewisburg, PA#TAB#",  
15         "email": "gretal.polites@123mail.org",  
16         "bio": "Dr. Polites joined the Kent State faculty in Fall  
→ 2012. She earned her B.S., M.S., and M.B.A. degrees from the  
→ University of South Florida, and completed her Ph.D. in Business  
→ Administration at the University of Georgia.[. . .] She has  
→ published two papers in the field of invertebrate paleontology,  
→ and has two fossil mollusk species (Attiliosa gretae and Opalia  
→ politesae) named after her."  
17     }  
18 ],  
19 "venue": {  
20     "raw": "Management Information Systems Quarterly",  
21     "id": "57293258"  
22 },  
23     "year": 2003,  
24     "n_citation": 139,  
25     "page_start": 397,  
26     "page_end": 423,  
27     "doc_type": "Journal",  
28     "publisher": "Society for Information Management and The  
→ Management Information Systems Research Center",  
29     "volume": 27,  
30     "issue": 3,
```

```
31      "fos": [
32      {
33          "name": "Decision-making",
34          "w": 0.4957332
35      },
36      {
37          "name": "Knowledge management",
38          "w": 0.46034035100000004
39      },
40      {
41          "name": "Workload",
42          "w": 0.5390811560000001
43      },
44      {
45          "name": "Interface design",
46          "w": 0
47      },
48      {
49          "name": "Information system",
50          "w": 0.5624888
51      }
52  ],
53  "doi": "10.2307/30036539",
54  "references": [
55      ObjectId("151626165300000000000000"),
56      ObjectId("197873803500000000000000")
57  ],
```

```

58     "abstract": "Managers in modern organizations are confronted
→ with ever-increasing volumes of information that they must
→ evaluate when making a decision. Data warehousing and data
→ mining technologies have given managers a number of valuable
→ tools that can help them store, retrieve, and analyze
→ information contained in large databases; however, maximizing
→ user performance with these tools remains a challenge for
→ information systems professionals. [. . .] These results have
→ important implications for the design of managerial
→ decision-making systems, particularly in complex decision-making
→ environments.",
59     "sections": [
60         {
61             "id": 1,
62             "title": "Introduction,
63             "text": "Automated decision-making (ADM) is no longer
→ science fiction and systems are now making decisions that were
→ traditionally made by humans (Robert, Pierce, Marquis, Kim, &
→ Alahmad, 2020). Algorithms match passengers for a shared ride
→ and plan drivers' routes (Möhlmann & Zalmanson, 2017). In Hong
→ Kong, an algorithm organizes the underground's maintenance
→ schedule and assigns repair jobs to service technicians (Hodson,
→ 2014).",
64             "subsections": [],
65             "figures": [
66                 {
67                     "url": "http://lh6.ggpht.com/IMG_3422.JPG?imgmax=800",
68                     "caption": "Chat interaction from participants' point of
→ view"
69                 },
70                 {
71                     "url": "http://lh6.ggpht.com/IMG_3422.JPG?imgmax=800",
72                     "caption": "Interview with the system after round one"
73                 },
74                 {
75                     "url": "http://lh6.ggpht.com//IMG_3422.JPG?imgmax=800",

```

```

76             "caption": "Pictures of the anthropomorphism
77             ↳ manipulation"
78         }
79     ]
80   ]
81 })

```

#### 4. Insertion of more papers at once:

```

1 db.papers.insertMany(
2   [
3     {
4       "title": "Shackled to [...] acceptance",
5       "authors": [
6         {
7           "name": "Elena Karahanna",
8           "id": "270263451",
9           "org": "University of Georgia, Athens",
10          "email": "elenakarahanna@liberomail.com",
11          "bio": "Education: PhD, MIS, University of Minnesota,
12          ↳ 1993. [...]"
13        }
14      ],
15      "venue": {
16        "raw": "Management Information Systems Quarterly",
17        "id": "57293258"
18      },
19      "year": 2012,
20      "n_citation": 215,
21      "page_start": 21,
22      "page_end": 42,
23      "doc_type": "Journal",
24      "publisher": "Society for Information Management",
25      "volume": 36,
26      "issue": 1,
27      "fos": [

```

```
28         "w": 0.3945593
29     }
30 ],
31 "doi": "10.2307/41410404",
32 "references": [
33     ObjectId("175816726800000000000000"),
34     ObjectId("191857022600000000000000"),
35     ObjectId("209868554100000000000000")
36 ],
37 "abstract": "Given that adoption of a new system often implies
38 ← fully or partly replacing [...]",
39 "sections": [
40     {
41         "id": 1,
42         "title": "Introduction",
43         "text": "Warfarin, [...]",
44         "subsections": [
45             {
46                 "id": 1,
47                 "title": "Warfarin dose adjustment",
48                 "text": "Warfarin dose adjustment [...]"
49             },
50             {
51                 "id": 2,
52                 "title": "Time in therapeutic range (TTR)",
53                 "text": "Time in therapeutic range (TTR) has been widely
54 ← used. [...]"
55             }
56         ],
57         "figures": [
58             {
59                 "url": "http://lh6.ggpht.com/IMG_3422.JPG?imgmax=800",
60                 "caption": "PRISMA chart used for the selection of
61 ← articles"
62             },
63             ...
64         ]
65     }
```

```
62      },
63      {
64          "id": 2,
65          "title": "Method",
66          "subsections": [
67              {
68                  "id": 1,
69                  "title": "Search Strategy",
70                  "text": "We conducted an extensive literature search in
    ↳ a very systematic way [...]"
71              }
72          ],
73          "figures": [
74              {
75                  "url": "http://lh6.ggpht.com//IMG_3422.JPG?imgmax=800",
76                  "caption": "Poor hypertension control or age-related
    ↳ frailty"
77              }
78          ]
79      },
80      {
81          "id": 3,
82          "title": "Result",
83          "subsections": [
84              {
85                  "id": 1,
86                  "title": "Selection Criteria",
87                  "text": "1) Age: There is clear evidence on the benefits
    ↳ of warfarin [...]"
88              },
89              {
90                  "id": 2,
91                  "title": "Data Extraction",
92                  "text": "Two reviewers (TR and NKJ1) independently
    ↳ extracted data [...]"
93              }
94          . . .
```

```
95     ] ,  
96     "figures": [  
97         {  
98             "url": "http://lh6.ggpht.com/IMG_3422.JPG?imgmax=800",  
99             "caption": "Optimal anticoagulation"  
100        }  
101    ]  
102}  
103]  
104},  
105{  
106    "title": "The State of the Art in Text Filtering",  
107    "authors": [  
108        {  
109            "name": "Douglas W. Oard",  
110            "id": "7916806",  
111            "org": "University of Maryland, College Park, MD 20742,  
112            "U.S.A.",  
113            "email": "douglasw.oard@123mail.org",  
114            "bio": "Institute for Advanced Computer Studies.Information  
115            technology, user modeling, information retrieval. [...]"  
116        }  
117    ],  
118    "venue": {  
119        "raw": "User Modeling and User-adapted Interaction",  
120        "id": "160628929"  
121    },  
122    "year": 1997,  
123    "n_citation": 114,  
124    "page_start": 141,  
125    "page_end": 178,  
126    "doc_type": "Journal",  
127    "publisher": "Kluwer Academic Publishers",  
128    "volume": 7,  
129    "issue": 3,  
129    "fos": [  
130        {
```

```
130         "name": "Machine learning",
131         "w": 0.442438453
132     },
133     {
134         "name": "Computer science",
135         "w": 0.4239926
136     },
137 ],
138 "doi": "10.1023/A:1008287121180",
139 "references": [...],
140 "abstract": "This paper develops. [...] implications for
→ future research on text filtering.",
141 "sections": [
142     {
143         "id": 1,
144         "title": "Introduction",
145         "text": "Warfarin, a coumarin derivative oral anticoagulant
→ [...]",
146         "subsections": [
147             {
148                 "id": 1,
149                 "title": "Warfarin dose adjustment",
150                 "text": "Warfarin dose adjustment is based on regular
→ monitoring of international normalized ratio (INR). [...]"
151             },
152             {
153                 "id": 2,
154                 "title": "Time in therapeutic range (TTR)",
155                 "text": "Time in therapeutic range (TTR) has been widely
→ used to measure the quality of INR control [...]"
156             }
157         ],
158         "figures": [
159             {
160                 "url": "http://lh6.ggpht.com/IMG_3422.JPG?imgmax=800",
161                 "caption": "PRISMA chart used for the selection of
→ articles"
```

```

162         }
163     ]
164   }
165 ]
166 }]
167 )

```

5. Update of a paper, by modification of the title of a subsection:

```

1 db.papers.updateOne(
2   { "_id": ObjectId("1015675232000000000000000") , "authors.name":
3     "Joeran Beel" } ,
4   { "$set": { "sections.0.subsections.0.title": "Artificial
5     Intelligence: Machine Learning" } }
6

```

6. Update some papers, based on some conditions, modifying starting and ending pages:

```

1 db.papers.updateMany(
2   { "doc_type": "Conference" , "year": {"$gte": 2000} } ,
3   { "$set": { "page_start": 0 } },
4   { "$set": { "page_end": {"$subtract":
5     ["$page_end", "$page_start"] } } }
6

```

# 3 | References

- Aminer data set
- Bio data set
- Tweets data set
- Images data set
- Draw.io
- Neo4j
- MongoDB
- Overleaf
- Python