



POLITECNICO
MILANO 1863

SCUOLA DI INGEGNERIA INDUSTRIALE
E DELL'INFORMAZIONE

Systems and Methods for Big and Unstructured Data Project

Author(s): **Stefano Baroni**

Luca Bremilla

Alessia Menozzi

Caterina Giardi

Rocco Agnello

Group Number: **33**

Academic Year: 2022-2023

Contents

Contents	i
I Delivery 1	1
1 Introduction	3
2 Entities and Relations	5
2.1 Introduction	5
2.2 Description	6
3 Import	9
3.1 Choice of the dataset	9
3.2 Query import	11
3.3 Graph DB	13
4 Query	15
4.1 Creation/Update	15
4.1.1 Query 1	15
4.1.2 Query 2	15
4.1.3 Query 3	16
4.1.4 Query 4	17
4.1.5 Query 5	17
4.2 Information Retrieval	18
4.2.1 Query 1	18
4.2.2 Query 2	19
4.2.3 Query 3	20
4.2.4 Query 4	21
4.2.5 Query 5	21
4.2.6 Query 6	22
4.2.7 Query 7	23

4.2.8	Query 8	23
4.2.9	Query 9	24
4.2.10	Query 10	25
4.3	Performance	26
4.3.1	Creation/Update queries	26
4.3.2	Information Retrieval queries	27
II	Delivery 2	29
5	Introduction and description of the model	31
6	Import	35
6.1	Libraries used	35
6.2	Abstract and Authors	35
6.3	Metadata	36
6.4	Publication details	36
6.5	Sections	37
6.6	Bibliography	38
7	Queries in MongoDB	39
7.1	Data creation/update	39
7.1.1	Query 1	39
7.1.2	Query 2	41
7.1.3	Query 3	42
7.1.4	Query 4	44
7.1.5	Query 5	45
7.2	Queries	45
7.2.1	Filtering conditions	45
7.2.2	Filtering conditions and aggregation	47
7.2.3	Filtering conditions also on internal sub-documents	49
7.2.4	Filtering conditions through unwind	51
7.2.5	Filtering condition through refs	54
7.3	Performance of the queries	58
7.3.1	A closer look at performance with pipelines	59

Part I

Delivery 1

1 | Introduction

The aim of this project is to build a model for a database that contains bibliography concerning scientific papers.

First of all we have modeled the bibliographic database creating an E-R model to better model the data. Then we downloaded a json file containing all the information about the papers, authors, venue and the different fields of study.

We decided to decrease the dimension of the json file because Neo4J was becoming too slow and after the data import we worked on different queries to retrieve information and play with the nodes.

In this document we will further explain the procedures applied to import a given dataset on Neo4j, then we will discuss the structure of the dataset itself and finally we will test different types of queries and show the obtained results using screenshots.

2 | Entities and Relations

2.1. Introduction

We modeled the E-R model using erdplus.com. This model is based on the json file containing all the fields of the AMiner database given by the specifications of the project. It is modeled in order to be a general overview of what a bibliographic database might look like.

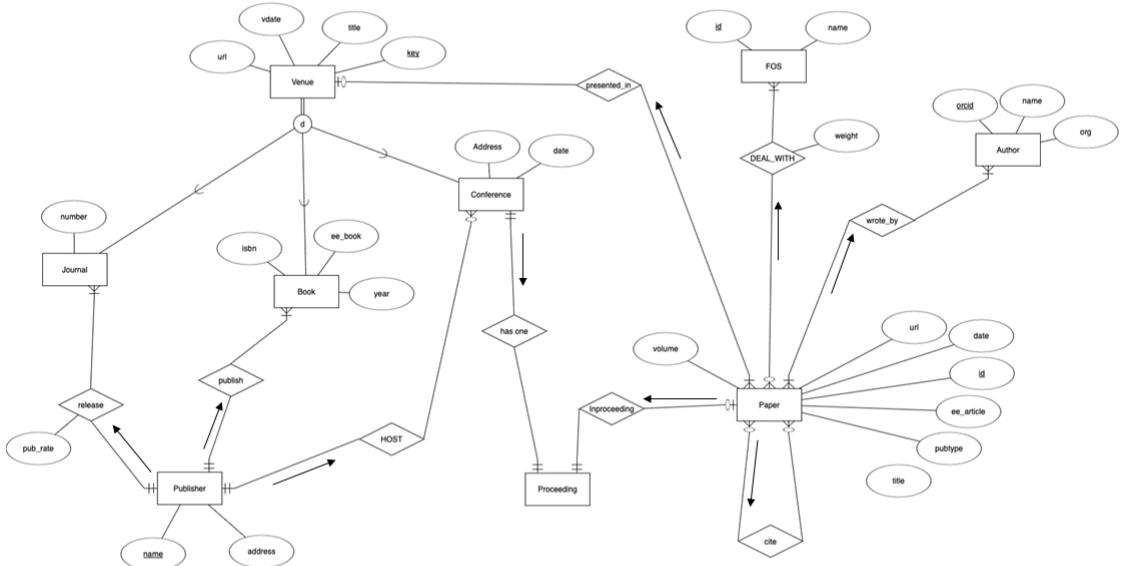


Figure 2.1: E-R Model

2.2. Description

We created nine entities:

- **Paper**, this entity is the core of our graph and that's why it contains the highest number of attributes: id (the primary key), title, date, URL, volume, pub_type (publication type f.i. informal, scientific, ...), and ee_article, that represents the position of the electronic edition. One paper can also contain citations of other papers and we represented it by inserting the *cite* relationship that connects two papers and has a 0-n cardinality for both.
- **Author**, characterized by attributes such as name, organization and orcid (primary key) that is an international unique identifier for scientific and academic authors. It is connected to the entity Paper through the relationship *wrote_by*: every author has written from 1 to n papers and every paper has from 1 to n authors.
- **FOS** (Field of Study) which is the topic of the paper, it has an id as primary key and a name. It's connected to it through the relationship *DEAL_WITH* that has cardinality 0-n on the paper's side and cardinality 1-n on the FOS's side. This relationship has the attribute *weight* that captures how much a topic is present inside the article.
- **Venue**, this is a supertype that includes 4 attributes: URL, vdate, title and the primary key. A Venue is connected to the Paper with the relationship *presented_in* that has cardinality 1-n on the Venue's side and 0-1 on the Paper's side. The attributes are inherited by the 3 disjoined subtypes.
 - **Journal**: since it can be compared to a periodic newspaper, we added the attribute *number*. Every journal must be connected to one Publisher using the relationship *release* that has the attribute *pub_rate* which describes the frequency of publication of every journal (weekly, daily, monthly...). On the other side, each publisher is connected to at least one journal.
 - **Book**: its attributes are the *year* in which the book was published, the *ISBN* code (an International Standard Book Number) and *ee_book*, which contains the position of the electronic edition. It's connected to the entity Publisher with the relationship *publish* and with the same cardinalities as the entity above.
 - **Conference**: its attributes are the date in which the conference took place and the address. This entity is connected to Publisher: a Publisher can host

from 0 to multiple Conference(s), and a Conference must be hosted by one Publisher.

- **Publisher**, it's the one that publishes books or journals, both containing one or more papers. Otherwise, it can host a conference, in which one or more papers are presented. Its attributes are name (primary key), and address.
- **Proceeding**, it represents a collection of papers presented in a conference. It's connected to the entity Conference with a one-to-one mandatory relation on both sides. It's also connected to the entity Paper through the relationship *inproceeding* which is optional: one proceeding can contain from 0 to multiple papers and one paper can be present in at most one conference. Note: This entity was inserted on the E-R diagram in order to make it the more general as possible but it doesn't appear on our specific dataset.

3 | Import

3.1. Choice of the dataset

We were given more choices about the dataset. So we have tried many of them and we then choose the best one.

First of all, we used the dataset by dblp. The file provided here is an xml. However, the bad thing about it is that a paper isn't linked to the papers cited, and we thought this is something not convenient for a project like this, because in bibliography the citations are one of the most important things.

For example, one paper in this file is:

```

1 <article mdate="2020-06-25" key="tr/meltdown/s18" publtype="informal">
2   <author>Paul Kocher</author>
3   <author>Daniel Genkin</author>
4   <author>Daniel Gruss</author>
5   <author>Werner Haas 0004</author>
6   <author>Mike Hamburg</author>
7   <author>Moritz Lipp</author>
8   <author>Stefan Mangard</author>
9   <author>Thomas Prescher 0002</author>
10  <author>Michael Schwarz 0001</author>
11  <author>Yuval Yarom</author>
12  <title>Spectre Attacks: Exploiting Speculative Execution.</title>
13  <journal>meltdownattack.com</journal>
14  <year>2018</year>
15  <ee type="oa">https://spectreattack.com/spectre.pdf</ee>
16 </article>
```

Moreover, as you can see, the authors are not identified by an id, so we would have to add the constraint that the key was their name and that it was not possible to have homonyms.

The other datasets tried are given by AMiner. They were all in json format. The citation data is extracted from DBLP (where we found the dataset above), ACM, MAG (Microsoft Academic Graph), and other sources. The first versions (till version 10) were too simple and spare. So, we opted for version 10. An example of a paper in that dataset is provided below. We now have references with other papers.

```

1 {
2     "abstract": "The purpose of this study is to develop a learning tool
3         for high school students studying the scientific aspects of
4         information and communication net- works. More specifically, we focus
5         on the basic principles of network proto- cols as the aim to develop
6         our learning tool. Our tool gives students hands-on experience to
7         help understand the basic principles of network protocols.",
8     "authors": ["Makoto Satoh", "Ryo Muramatsu", "Mizue Kayama", "
9         Kazunori Itoh", "Masami Hashimoto", "Makoto Otani", "Michio Shimizu",
10        "Masahiko Sugimoto"],
11     "n_citation": 0,
12     "references": ["51c7e02e-f5ed-431a-8cf5-f761f266d4be", "69b625b9-
13        ebc5-4b60-b385-8a07945f5de9"],
14     "title": "Preliminary Design of a Network Protocol Learning Tool
15         Based on the Comprehension of High School Students: Design by an
16         Empirical Study Using a Simple Mind Map",
17     "venue": "international conference on human-computer interaction",
18     "year": 2013,
19     "id": "00127ee2-cb05-48ce-bc49-9de556b93346"
20 }
```

Again, one thing not provided is the id of an author, but that's good enough.

However, to increase the number of nodes and edges, we should use version 12, where there are also “fos”, i.e. fields of study, publishers, id of authors and many other properties.

3.2. Query import

We have used 3 different scripts to import the dataset. Obviously, we had to create the nodes (first script), then we had to delete some nodes that were "mistaken" (second and third scripts).

```

36    // Each paper has a list of "fos" (fields of study). We create each
37    author in that list and link it with the paper
38    FOREACH (f IN p.fos |
39      MERGE (fos:Fos {id: f.name})
40      MERGE (paper)-[:DEAL_WITH {weight:f.w}]->(fos)
41    )

```

Listing 3.1: Query used to import and create the database

Some papers don't have a publisher but that's provided anyways with empty string. So we need to delete it.

```

1 MATCH (p:Publisher)
2 WHERE p.id = ""
3 DETACH DELETE p

```

Listing 3.2: To delete the Publishers without name

In the same way, some paper don't have a venue but that's provided anyway with empty string. So we need to delete it.

```

1 MATCH (v:Venue)
2 WHERE v.raw=""
3 DETACH DELETE v

```

Listing 3.3: To delete the Venues without name

3.3. Graph DB

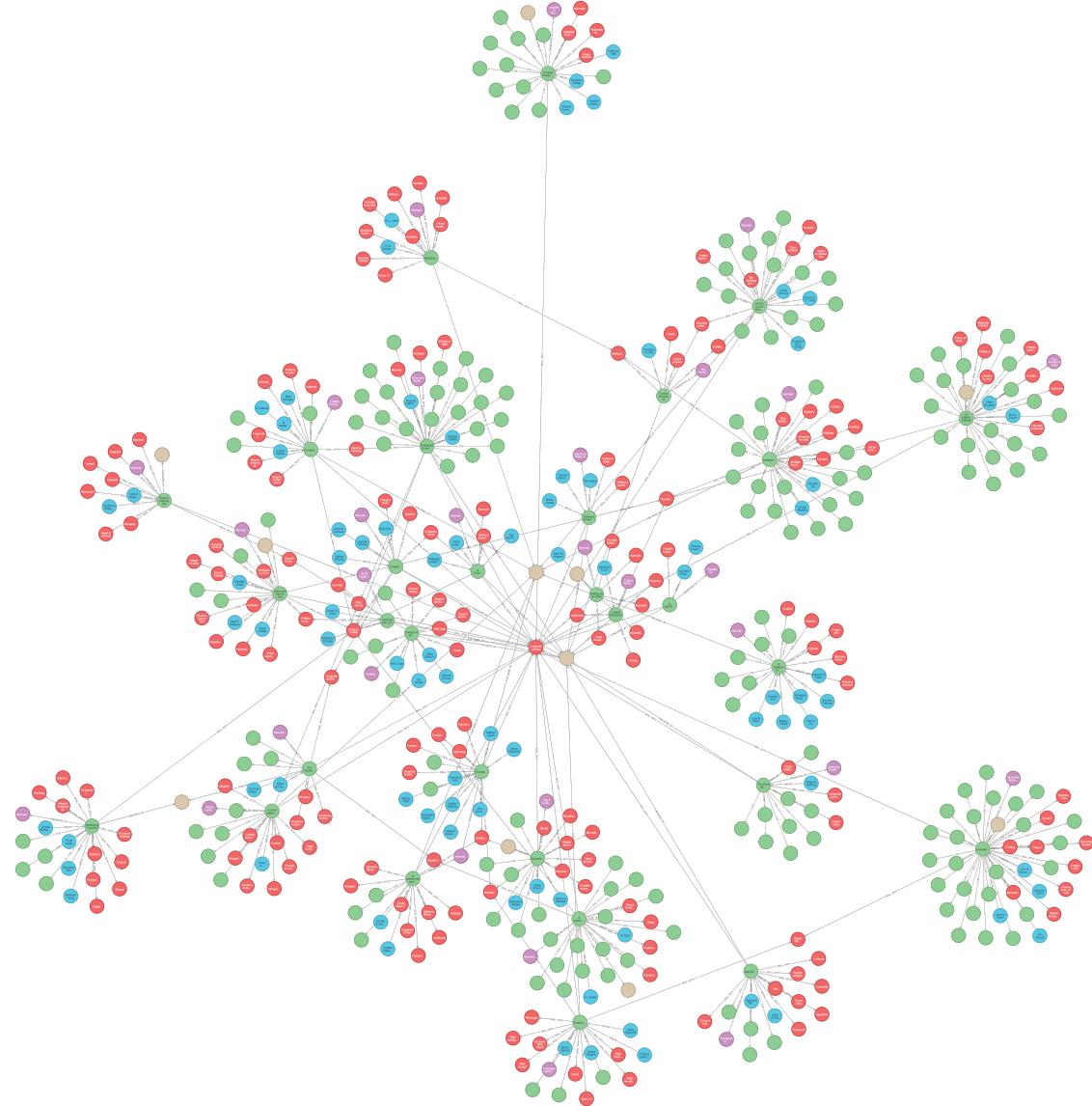


Figure 3.1: Graph Database

Description of the graph:

ENTITIES:

- **Venue**, whose attributes are id, raw(the name of the venue), numpages, type: J (for journal), C (for conference) and B (for Book)
- **Author**, whose attributes are id, name and org (optional)
- **Paper**, whose attributes are id, n_citations (number of citations included in that paper), title and year

- **FOS**, identified by an id which is their name
- **Publisher**, identified by an id which is their name

RELATIONS:

- CITE, links two papers that contains a citation on each other
- PUBLISHED_BY, links a paper to a publisher
- PRESENTED_IN, links a paper to a venue
- WROTE, links an author to a paper
- DEAL_WITH, links a paper to a field of study. This is the only relation that contains an attribute (weight). This attribute indicates how much that topic is present in the overall content of that specific paper.

4 | Query

4.1. Creation/Update

In the following queries we will create a new paper, associate it with an already existing author, and connect it with the fields of study that characterize the paper. These queries only represent an example of creation and update on a node, other queries regarding data creation have been used for the import.

4.1.1. Query 1

```
1      CREATE (p:Paper {title:"prova", id: 5000, year: 2022, n_citation: 0})
```

Listing 4.1: Creation of a new Paper

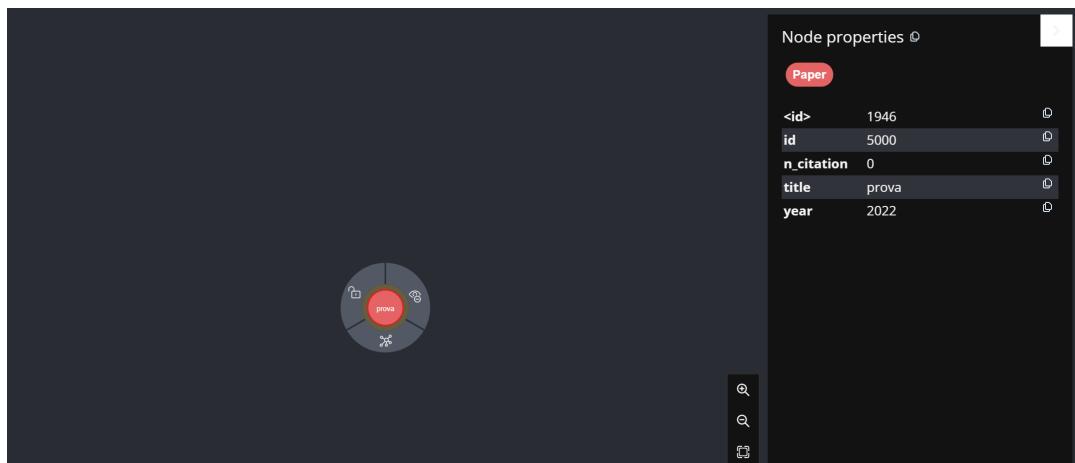


Figure 4.1: Result of new paper query

4.1.2. Query 2

```
1      MATCH (p:Paper), (fos:Fos)
2      WHERE p.title = 'prova' AND (fos.id = 'Mind map' OR fos.id = "Machine learning")
```

```
3     CREATE (p)-[r:DEAL_WITH {weight: 0.5}]->(fos)
```

Listing 4.2: Creation of a relationship between the new paper and its FOS



Figure 4.2: Result of creation query 2

4.1.3. Query 3

```
1     MATCH (a:Author), (p:Paper)
2     WHERE p.title = 'prova' AND (a.name = 'Rafael Alvarez' OR a.name = 'Masahiko Sugimoto')
3     CREATE (a)-[r:WROTE]->(p)
```

Listing 4.3: Creation of a relationship between the new paper and its authors

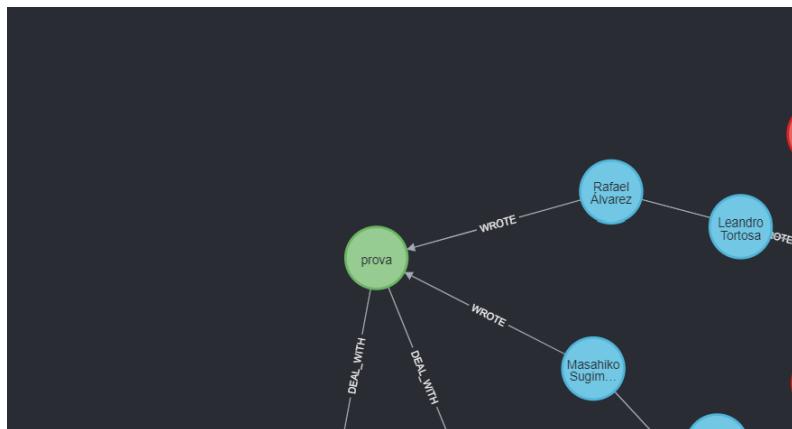


Figure 4.3: Result of creation query 3

4.1.4. Query 4

The first update query will perform the insertion of a new attribute `numpages` set to 20 in the case of the new paper we have just added to the database with the previous queries.

```

1   MATCH (p:Paper {title: 'prova'})
2   SET p.numpages = 20
3   RETURN p

```

Listing 4.4: Update Query 1

The screenshot shows two side-by-side views of the Neo4j Node Properties table for a node labeled "Paper".

Left View (Before Update):

<id>	1946	edit
id	5000	edit
n_citation	0	edit
title	prova	edit
year	2022	edit

Right View (After Update):

<id>	140	edit
id	5000	edit
n_citation	0	edit
numpage	20	edit
s		edit
title	prova	edit
year	2022	edit

Figure 4.4: Result of the update query, before and after

4.1.5. Query 5

If this update query finds a match based on the title of the paper, then the attribute `year` is incremented by 1, otherwise a new paper with that title and the year set on 2020 is created.

```

1   MERGE (p:Paper{title:"Fur Visualisation for Computer Game
Engines and Real-Time Rendering"})
2   ON CREATE SET p.year="2020"

```

```

3      ON MATCH SET p.year= p.year +1
4      return p

```

Listing 4.5: Update Query 2

Graph	"p"
Table	{"year":2016,"id":8763,"title":"Fur Visualisation for Computer Game Engines and Real-Time Rendering","n_citation":0}

Text	"p.year"
Code	2015

Figure 4.5: In the second row the non modified Paper, in the first one the modified one

4.2. Information Retrieval

In this section we will show queries that retrieve different kind of data. We tried to involve different nodes, using constraints and limitations.

Note that for some queries we have created datasets using other data different from those used in the graph reported in chapter 3.3. All queries would work if the database were created using all the data in the original file.

4.2.1. Query 1

This query returns the Venues of type C (Conference) and the Papers presented there. There is a constraint on the Papers: they must have been produced by at least 4 authors.

```

1      MATCH (v:Venue)-[:PRESENTED_IN]-(p:Paper)-[:WRITED_BY]-(a:Author)
2      WHERE v.type = "C"
3      WITH v.raw AS vtitle, p.title AS ptitle, collect(a.name) AS authors

```

```

4    WHERE size(authors) > 3
5    RETURN vtitle, ptitle, authors

```

Listing 4.6: Query 1

"vtitle"	"ptitle"	"authors"
"International Conference on Human-Computer Interaction"	"Preliminary Design of a Network Protocol Learning Tool Based on the Comprehension of High School Students: Design by an Empirical Study Using a Simple Mind Map"	["Kazunori Itoh", "Ryo Muramatsu", "Michio Shimizu", "Masahiko Sugimoto", "Makoto Otani", "Mizue Kayama", "Makoto Sato", "Masami Hashimoto"]
"Pattern Recognition and Machine Intelligence"	"Comparison of GARCH, Neural Network and Support Vector Machine in Financial Time Series Prediction"	["M. Nasser", "Altaf Hossain", "M. Mufakhkharul Islam", "Faisal Zaman"]
"International Conference on Computer Graphics Theory and Applications"	"COMPARING GNG3D AND QUADRIC ERROR METRICS METHODS TO SIMPLIFY 3D MESHES"	["José-Francisco Vicent", "Leandro Tortosa", "Antonio Zamora", "Rafael Alvarez"]
"International Conference on Mining Intelligence and Knowledge Exploration"	"Identifying Psychological Theme Words from Emotion Annotated Interviews"	["Ankita Brahmachari", "Priya Singh", "Dipankar Das", "Avadhesh Garg"]
"Conference of the International Speech Communication Association"	"Design of an audio-visual speech corpus for the czech audio-visual speech synthesis."	["Petr Cisar", "Zdenek Krnoul", "Milos Zelezny", "Jan Novák"]
"International Conference on Critical Infrastructure Protection"	"A Platform for Disaster Response Planning with Interdependency Simulation Functionality"	["José R. Martí", "Pranab Kini", "Abdullah Alsubaie", "Alberto Tofani", "Simone Palmieri", "Ting Fu Lin", "Antonio Di Pietro"]
"Language Resources and Evaluation"	"CleanEval: a Competition for Cleaning Web Pages"	["Marco Baroni", "Adam Kilgarriff", "Francis Chantree", "Serge Sharoff"]
"Operating Systems Design and Implementation"	"Leveraging legacy code to deploy desktop applications on the web"	["Jeremy Elson", "Jacob R. Lorch", "John R. Douceur", "John Howell"]

Figure 4.6: Result of query 1

4.2.2. Query 2

For every paper, the query returns the list of all FOS having a weight of at least 0.5 and having a number of citations that is less than 5.

```

1    MATCH(p:Paper)-[d:DEAL_WITH]->(f:Fos)
2    WHERE d.weight >= 0.5 and p.n_citation < 5
3    WITH p.title as title, collect(f.id) as topics, p.n_citation as citation
4    RETURN title, topics, citation

```

Listing 4.7: Query 2

"title"	"topics"	"citation"
"Preliminary Design of a Network Protocol Learning Tool Based on the Comprehension of High School Students: Design by an Empirical Study Using a Simple Mind Map"	["Mind map", "Communications protocol"]	1
"Further Results on Independence in Direct-Product Graphs."	["Direct product"]	1
"COMPARING GNG3D AND QUADRIC ERROR METRICS METHODS TO SIMPLIFY 3D MESHES"	["Quadric", "Polygon mesh"]	0
"Vectorial fast correlation attacks."	["Correlation attack"]	2
"Improved Secret Image Sharing Method By Encoding Shared Values With Authentication Bits"	["Peak signal-to-noise ratio", "Shamir's Secret Sharing", "Verifiable secret sharing", "Secret sharing", "Homomorphic secret sharing"]	0
"Formal agent-oriented ubiquitous computing: a computational intelligence support for information and services integration"	["Computational intelligence", "Context-aware pervasive systems", "Ubiquitous commerce", "Context awareness", "Utility computing", "Services computing", "Ubiquitous computing"]	0
"Fur Visualisation for Computer Game Engines and Real-Time Rendering"	["Visualization", "Rendering (computer graphics)", "Real-time rendering"]	0
"Identifying Psychological Theme Words from Emotion Annotated Interviews"	["Keyword spotting", "Sadness"]	0
"Multisymplectic Spectral Methods for the Gross-Pitaevskii Equation"	["Gross-Pitaevskii equation", "Spectral method", "Schrödinger equation", "Wave equation", "Nonlinear Schrödinger equation"]	2
"Speech training systems using lateral shapes of vocal tract and F1-F2 diagram for hearing-impaired children"	["Formant", "Vocal tract", "Vowel"]	1
"Software Evolution through Transformations."	["Software evolution"]	2

Figure 4.7: Result of query 2

4.2.3. Query 3

This query returns for every paper the list of papers that appear between his citations, but these papers must have at least a FOS in common.

For this particular query we created a graph using the data of a dataset in order to create the relation *CITE* (as represented below).

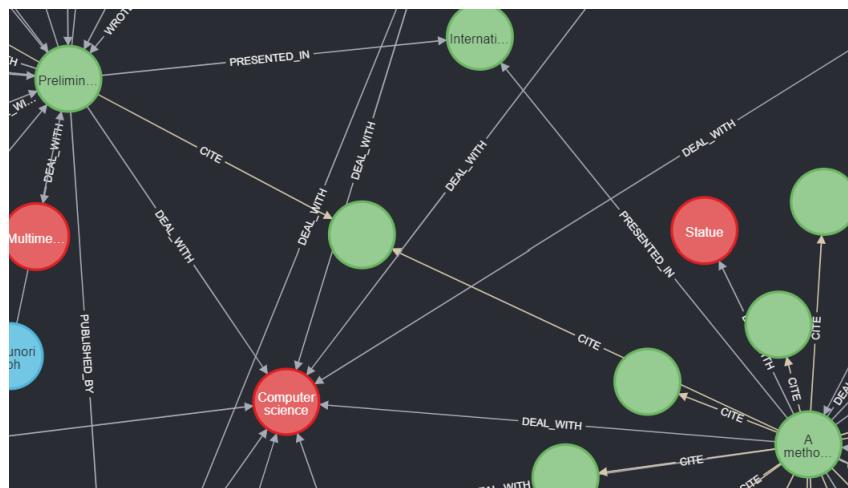


Figure 4.8: Example of the relation CITE

```

1 MATCH (f1:Fos) <- [:DEAL_WITH] -> (p1:Paper) -[:CITE] -> () <- [:CITE] -> (p2:Paper)
2 WITH p1.title AS original, collect(p2.title) as cited, f1.id AS name1, f2.id AS name2
  
```

```

3      WHERE name1 = name2
4      RETURN original, cited, name1

```

Listing 4.8: Query 3

"original"	"cited"	"name1"
"A methodology for the physically accurate visualisation of roman polychrome statuary"	["Preliminary Design of a Network Protocol Learning Tool Based on the Comprehension of High School Students: Design by an Empirical Study Using a Simple Mind Map"]	"Computer science"
"Preliminary Design of a Network Protocol Learning Tool Based on the Comprehension of High School Students: Design by an Empirical Study Using a Simple Mind Map"	["A methodology for the physically accurate visualisation of roman polychrome statuary"]	"Computer science"

Figure 4.9: Result of query 3

4.2.4. Query 4

The query 4 returns the papers written before the year 2000 and containing at least 6 citations. These papers must have been presented in a Venue where at least other 3 papers have been presented. In order to execute this query we created a graph "AD HOC" always based on the same dataset we imported.

```

1      MATCH (p1:Paper)-[:PRESENTED_IN]->(v:Venue)<-[r:PRESENTED_IN]-(:
2          Paper)-[:DEAL_WITH]->(fos:Fos)
3          WHERE p1.year<2000 and p1.n_citation>5 and fos.id="Machine learning"
4          WITH p1.title as paper, v.raw as raw, count(r) as
5          other_pap_presented
6          WHERE other_pap_presented>2
7          RETURN paper, raw, other_pap_presented

```

Listing 4.9: Query 4

"paper"	"raw"	"other_pap_presented"
"Cooperative learning over composite search spaces experiences with a multi-agent design system"	"National Conference on Artificial Intelligence"	3
"Search lessons learned from crossword puzzles"	"National Conference on Artificial Intelligence"	3

Figure 4.10: Result of query 4

4.2.5. Query 5

The query returns the list of authors that wrote a paper before the 2000 which was presented in a conference (for which we return the attribute *raw* which is the name of the

conference itself). This conference must have presented at least 5 papers (and returns the number of the corresponding papers, shown in a decreasing order).

```

1   MATCH (a:Author) -[:WRITED] ->(p1:Paper) -[r1:PRESENTED_IN] ->(v:Venue)
2     <- [r2:PRESENTED_IN] -(p2:Paper)
3     WHERE v.type = "C" and p1.year < 2000
4     WITH collect(distinct a.name) as authors, v.raw as raw, count(
5       distinct r2) AS num_papers
6     WHERE num_papers > 5
7     RETURN raw, num_papers, authors
8     ORDER BY num_papers DESC

```

Listing 4.10: Query 5

"raw"	"num_papers"	"authors"
"Conference of the International Speech Communication Association"	27	["Julian H. Page", "Mark Tatham", "Andy Tams", "Julia Hirschberg", "Richard Sproat", "David Yarowsky", "Jose Maria Elvira", "Javier Caminero", "Juan Carlos Torrecilla", "Mats Ljungqvist", "Anders Lindström", "Kjetil Gustafson", "Junichi Yamazaki", "Zhiyong Peng", "Makoto J. Hirayama", "Taro Sugahara", "Lisa-Jane Brown", "Sandra P. Whiteside", "John Locke", "Peter Jones", "Satoshi Nakamura", "Kiyohiro Shikano", "Shigeki Sagayama", "Hiroaki Hattori", "Klaus Beulen", "Hermann Ney", "Lutz Wellinger", "Annemie Vorstermans", "Nick Cremelie", "Jean-Pierre Martens", "Yoshinori Sagisaka", "Nobuyoshi Kaiki", "Katsuhiko Mimura", "Owen P. Kennedy", "Douglas J. Nelson", "Li-Qun Xu", "Tie-Cheng Yu", "G. D. Tattersall", "Laure Charonnat", "Michel Glutton", "Joel Crestel", "Herve Chuberre", "Arun C. Surendran", "Chin-Hui Lee", "Shuji Doshita", "Tatsuya Kawahara", "Toru Ogawa", "Shigeyoshi Kitazawa", "Cheol-Jun Hwang", "Hyun-Yeoil Chung", "Shi-wook Lee"]
"Parallel and Distributed Processing Techniques and Applications"	8	["I. M. Ikram", "Paul S. Wang", "Hong H. Ong"]
"International Conference on Human-Computer Interaction"	6	["Ann Doubleday", "Alistair G. Sutcliffe", "Ian Bennett", "Michele Ryan", "Anthony Savidis", "Pier Luigi Emilian", "Constantine Stephanidis"]
"International Joint Conference on Artificial Intelligence"	6	["K. Sparck Jones", "B. K. Bog", "Minoru Shigenaga", "Yoshihiro Sekiguchi", "Stephan Busemann", "P. Bonzon"]

Figure 4.11: Result of query 5

4.2.6. Query 6

The query returns the name of the FOS that has more papers published on the journals and the corresponding number of paper that deal with that field.

```

1   MATCH (fos:Fos) <- [r:DEAL_WITH] - (p:Paper) -[:PRESENTED_IN] -> (:Venue {
2     type: "J"})
3     WITH fos, count(r) as num_papers
4     ORDER BY num_papers desc
5     LIMIT 1
6     RETURN fos.id, num_papers

```

Listing 4.11: Query 6

"fos.id"	"num_papers"
"Computer science"	6

Figure 4.12: Result of query 6

4.2.7. Query 7

Since the dataset has several different "Springer" publishers from different cities, this query aims to catalog all of them, counting the number of papers they published. Moreover, it retrieves a list of the authors that wrote those articles, created using the function `collect()`. They are finally shown in decreasing order

```

1   MATCH (pub:Publisher)-[r1:PUBLISHED_BY]-(p:Paper)-[:WROTE]-(a:
2     Author)
3   WITH pub, count(distinct r1) as num_papers, collect(distinct a.name)
4     as authors
5   WHERE substring(pub.id, 0, 8)="Springer"
6   RETURN pub.id, num_papers, authors
7   ORDER BY num_papers desc

```

Listing 4.12: Query 7

"pub.id"	"num_papers"	"authors"
"Springer, Berlin, Heidelberg"	8	[{"Fencol C. C. Yung", "Prudence W. H. Wong", "Mihai Burcea", "M. Mufakhkh", "arul Islam", "Faisal Zaman", "Altaf Hossain", "M. Nasser", "Pranab Kini", "José R. Martí", "Alberto Tofani", "Simone Palmieri", "Antonio Di Pietro", "Abdullah Alsubaie", "Ting Fu Lin", "Yutaka Nakamura", "Hiroshi Furukawa", "C. M. Schober", "A. L. Islas", "Jill-Jenn Vie", "Michel Abdalla", "Phan Cong Vinh", "Mizue Kayama", "Makoto Satoh", "Makoto Otani", "Michio Shimizu", "Masami Hashimoto", "Ryo Muramatsu", "Masahiko Sugimoto", "Kazunori Itoh"}]
"Springer, Cham"	2	[{"Adam Jurczyński", "Dominik Szajerman", "Avdhesh Garg", "Priya Singh", "Ankita Brahmachari", "Dipankar Das"]]
"Springer, London"	1	[{"Güzin Ulutas", "Mustafa Ulutas", "Vasif V. Nabihev"}]

Figure 4.13: Result of query 7

4.2.8. Query 8

This query returns the authors that have written at least two different papers published by two different publishers. The authors, limited to 2, are then returned by order of the

total citations of all the papers in which they appears as authors. The result given is performed on an "AD HOC" dataset.

```

1   MATCH (p1:Publisher)-[:PUBLISHED_BY]-(:Paper)-[:WROTE]-(a:Author)
2     -[:WROTE]->(:Paper)-[:PUBLISHED_BY]->(p2:Publisher)
3   WHERE p1<>p2
4   MATCH (a)-[:WROTE]->(p:Paper)
5   WITH a, SUM(p.n_citation) as numCitations
6   ORDER BY numCitations DESC
7   LIMIT 2
8   RETURN (a)
```

Listing 4.13: Query 8

```

| "a"
|
| {"org": "Anyang University", "name": "Mohammad Hasmat Ullah", "id": 2203502152}
|
| {"name": "Hafner Cd", "id": 2983371265}
```

Figure 4.14: Result of query 8

4.2.9. Query 9

This query returns all the papers written by at least two authors (and also return them). One of these authors must have written at least one paper concerning the FOS "Artificial Intelligence" with weight ≥ 0.5 , the other author must have written one paper about "Management" with weight ≥ 0.5 .

```

1   MATCH (p1:Paper)-[d1:DEAL_WITH]->(f1:Fos),
2     (p2:Paper)-[d2:DEAL_WITH]->(f2:Fos),
3     (a1:Author)-[:WROTE]->(p:Paper)-[:WROTE]-(a2:Author)
4   WHERE
5     (a1)-[:WROTE]->(p1) AND
6     (a2)-[:WROTE]->(p2) AND
7     f1.id = "Artificial intelligence" AND f2.id = "Management" AND
8     d1.weight >= 0.5 AND d2.weight >= 0.5
9   RETURN (a1)-[]->(p)-[]-(a2)
```

Listing 4.14: Query 9 less performing version

The following version is more efficient. That's because the above version builds a Cartesian product between disconnected patterns.

```

1      MATCH
2      (:Fos {id: "Management"})-<-[d2:DEAL_WITH]-(p2:Paper)-<[:WROTE]-(a1:
3          Author)-[:WROTE]->(p:Paper)-<[:WROTE]->(a2:Author)-[:WROTE]->(p1:Paper)
4          )-[d1:DEAL_WITH]->(:Fos {id: "Artificial intelligence"})
5          WHERE d1.weight >=0.5 AND d2.weight >=0.5
6          RETURN (a1)-[]->(p)->[]-(a2)

```

Listing 4.15: Query 9 more performing version

```

["(a1)-[]->(p)->[]-(a2)"]

[[{"org": "Sapienza", "University of Rome", "name": "Benjamin Habegger", "id": 1998301270}, {}, {"year": 2006, "id": 12412048, "title": "Biosequence Use Cases in MoBIOS SQL.", "n_citation": 0}, {"year": 2006, "id": 12412048, "title": "Biosequence Use Cases in MoBIOS SQL.", "n_citation": 0}, {}, {"name": "Irene Pollach", "id": 2675219732}]

[[{"org": "Sapienza", "University of Rome", "name": "Benjamin Habegger", "id": 1998301270}, {}, {"year": 2002, "id": 12464753, "title": "The online university: the students perspective.", "n_citation": 3}, {"year": 2002, "id": 12464753, "title": "The online university: the students perspective.", "n_citation": 3}, {}, {"name": "Irene Pollach", "id": 2675219732}]]
```

Figure 4.15: Result of query 9

4.2.10. Query 10

This query finds the shortest Path from an author whose first name is "Michel" to the topic "Artificial intelligence" with the constraint that the author can't have written any Paper concerning Machine Learning. Since in our database the number of nodes is limited I haven't put any upper bound, but in other cases it would have been necessary in order not to result in very long execution times.

```

1      MATCH (FirstAuthor:Author),
2          (Topic:Fos{id:"Artificial intelligence"}) ,
3          p=shortestPath((FirstAuthor)-[*]-(Topic))
4          WHERE substring(FirstAuthor.name, 0, 6)="Michel" and not (
5              FirstAuthor)-[:WROTE]-()-[:DEAL_WITH] -(:Fos{id:"Machine learning"})
6          RETURN p,FirstAuthor.name

```

Listing 4.16: Query 10

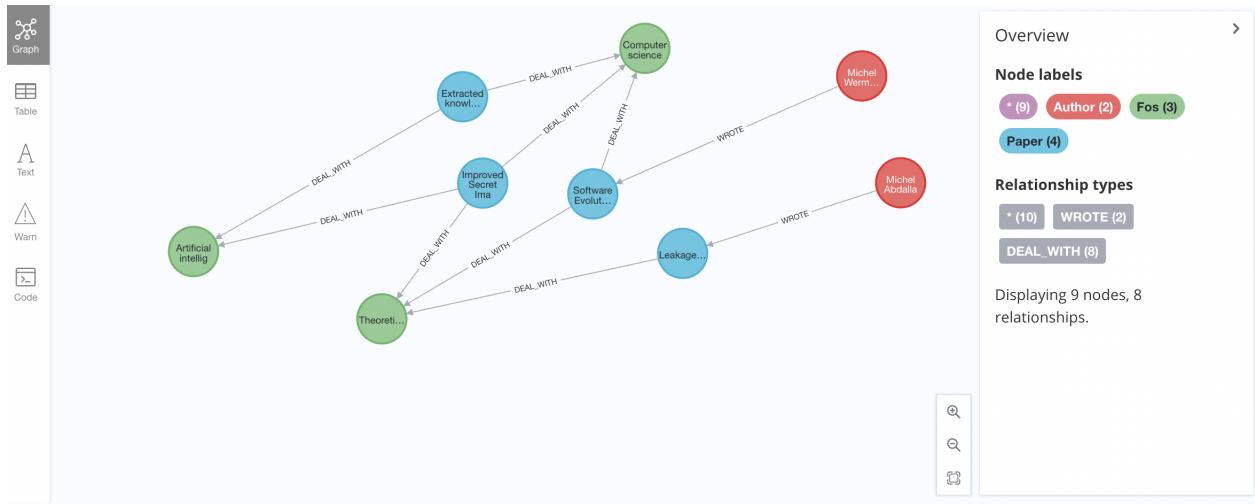


Figure 4.16: Result of query 10

4.3. Performance

Here we can see the performance of the queries described in a particular execution on a local database of 7926 nodes.

4.3.1. Creation/Update queries

Query	Time (ms)
Query 1	7
Query 2	48
Query 3	17
Query 4	13
Query 5	16

Table 4.1: Table of execution times of data creation/upload

4.3.2. Information Retrieval queries

Query	Time (ms)
Query 1	57
Query 2	63
Query 3	126
Query 4	32
Query 5	33
Query 6	21
Query 7	36
Query 8	21
Query 9 version 1	31
Query 9 version 2	25
Query 10	88

Table 4.2: Table of execution times of information retrieval queries

Part II

Delivery 2

5 | Introduction and description of the model

```

_id: 5411
> authors: Array
  title: "COMPARING GNG3D AND QUADRIC ERROR METRICS METHODS TO SIMPLIFY 3D MESHE..."
  n_citation: 4
  doc_type: "Conference"
  publisher: "PoliPrint, Milano"
  doi: ""
> venue: Object
  abstract: "Abstract of COMPARING GNG3D AND QUADRIC ERROR METRICS METHODS TO SIMPL..."
> metadata: Array
> references: Array
> bibliography: Array
> sections: Array

```

Figure 5.1: Example of a paper seen on MongoDB

Each paper contains information regarding:

- **Title:** contains the title of the paper;
- **Author:** the array *authors* contains all the paper authors and their details such as
 - *name*, which contains the name and surname of the author;
 - *email*, which contains the author's email;
 - *affiliation*, which contains the University affiliated to the author;
 - *bio*, which should contain a brief description of what the author does for a living, as well as their qualifications and education details (in our json file we use the package lorem to create random texts);

```

✓ authors: Array
  ✓ 0: Object
    name: "Altaf Hossain"
    id: 2300589394
    email: "Altaf.Hossain@mail.com"
    bio: "Bio of Altaf Hossain: ullam ea quo suscipit ex officiis expedita conse..."
    affiliation: "Department of Statistics , Rajshahi University , Rajshahi , Bangladesh"
  ✓ 1: Object
    name: "Faisal Zaman"
    id: 2308774408
    email: "Faisal.Zaman@mail.com"
    bio: "Bio of Faisal Zaman: labore ullam illum voluptatibus amet repellendus ..."
    affiliation: "Department of System Design and Informatics, Kyushu Institute of Techn..."
  > 2: Object
  > 3: Object

```

Figure 5.2: Example document - author array

- **Number of citations**: contains the number of papers that are citated in the actual paper (this value can also be checked by counting the number of elements in the array *References*);
- **Doctype**: it explains the kind of publication in which the paper first appears;
- **Publisher**: it indicates the name of the publisher of the paper or the publisher of the journal of the paper;
- **Digital Object Identifier (doi)**: it's a digital identifier;
- **Abstract**: it should contain the abstract of the paper, but we set this field using the pattern: "Abstract of" + title of the paper + random words using *lorem.words()*;
- **Venue**: this object describes where the paper was presented and can be either a Conference, a Repository or a Journal (*doc_type*);

```

✓ venue: Object
  raw: "Pattern Recognition and Machine Intelligence"
  id: 1136274694
  type: "C"
  volume: 3
  number: 17
  date: "2009-03-16"
  pages: 4

```

Figure 5.3: Example document - venue object

- **Metadata**: it contains an array with all the keywords (i.e. the topics discussed) and the corresponding weight (the importance they have in the meaning of the overall document);

- **References:** it's an array that contains all the Ids of the paper citated in the current paper;
- **Bibliography:** it's an array used to complete the *references* array by giving a title to the referenced Id;

```

< references: Array
  0: 511601
  1: 672965
  2: 681007
  3: 512839
  4: 671169
> venue: Object
  abstract: "Abstract of Comparison of GARCH, Neural Network and Support Vector Mac..."
< metadata: Array
  < 0: Object
    keyword: "Autoregressive-moving-average model"
    weight: 0.49811
  < 1: Object
    keyword: "Computer science"
    weight: 0.41114
  > 2: Object
  > 3: Object
  > 4: Object
  > 5: Object
  > 6: Object
  > 7: Object
  > 8: Object
< bibliography: Array
  0: "An OSF/1 UNIX for Massively Parallel Multicomputers. Id: 511601"
  1: "Using Argumentation to Overcome Hypertext's HCI Failings. Id: 672965"
  2: "Improving User Stereotypes through Machine Learning Techniques Id: 681..."
  3: "Use of spatial adaptation for image rendering based on an extension of..."
  4: "Virtual coiling of intracranial aneurysms based on dynamic path planni..."
```

Figure 5.4: Example document - references, bibliography and metadata arrays

- **Sections** is a multiple array; every element is an object that contains: the field *title*, an array *textbyparagraph* containing the paragraphs of the section (before beginning with the subsections) again created with `lorem.words()`, another array *subsections* that has a title and an array of paragraphs, and the array *figures* that contains all the urls and captions of the figures of that section.

```
✓ sections: Array
  ✓ 0: Object
    title: "Title of section 1"
    ✓ textbyparagraph: Array
      ✓ 0: Object
        paragraph: "modi rem eos quis voluptates obcaecati saepe necessitatibus ipsum dolo..."
  ✓ figures: Array
    ✓ 0: Object
      url: "https://image.com/?paperid=1688&section=1&image=1"
      caption: "Caption of image 1.0"
  ✓ subsections: Array
    ✓ 0: Object
      title: "Title of subsection 1"
      ✓ textbyparagraph: Array
        ✓ 0: Object
          paragraph: "error asperiores exercitationem repudiandae corporis minus tenetur ut ..."
        ✓ 1: Object
          paragraph: "dolorem sapiente possimus quo ullam eveniet minima recusandae inventor..."
      > 1: Object
      > 2: Object
      > 3: Object
    > 1: Object
    > 2: Object
    > 3: Object
```

Figure 5.5: Example document - sections array

6 | Import

For the import phase, we used the same dataset of the first delivery, but this time we had to add some new fields to the structure of a paper. In fact, they were not present in the json file already used. Examples went from the simpler ones, like the absence of the *email* and *bio* for the authors, to the more difficult ones, like the generation from scratch of the *sections*, with their title, text by paragraph, subsections and list of figures.

To cope with these problems we wrote a Python script which reads the file of the original dataset, modifies it and deploys another json file with all the changes required, both in adding new fields and in deleting not interesting/useful ones. In the following sections we will present the issues encountered, together with the solutions we came up with.

6.1. Libraries used

To work on the raw dataset, we imported the following libraries:

- **json** for the methods of read/write operations on a .json file;
- **lorem**¹ to fill with random strings the new fields (wg. the paragraphs of the sections);
- **random** for the functions of picking a random number over an interval of numbers;
- **date** to work with dates, in particular to assign a date to a paper (the year was already provided);
- **urllib** to generate a random URL for the images.

6.2. Abstract and Authors

For the abstract, we simply inserted a new field for each paper and put some words from the *lorem* library.

¹<https://pypi.org/project/python-lorem/>

```

1 data[i]['abstract'] = "Abstract of " + data[i]['title'] + ": " + lorem.
    words(20)

```

In particular, `data[i]` is the array containing the papers extracted from the raw dataset.

As for the authors, we did not have neither the mail nor the bio. For having a more consistent database, we reused the data we have on the author's name and we created a mail like `name.surname@mail.com` using the method `.split()`. Next, for the bio we adopted the same way as for the abstract for the paper. Finally, we assigned an *affiliation* to the authors which did not have it, taken randomly from a vector of affiliations we created (`vect_org[]` in the code).

```

1 for author in data[i]['authors']:
2     author['email'] = author['name'].split()[0] + '.' + \
3                     author['name'].split()[len(author['name'].split()) - 1] + '@mail.com'
4     author['bio'] = 'Bio of ' + author['name'] + ':' + lorem.words(20)
5     if 'org' in author.keys():
6         author['affiliation'] = author['org']
7     else:
8         author['affiliation'] = vect_org[random.randint(0, len(vect_org) - 1)]

```

6.3. Metadata

For the required field *metadata*, we used the *fos* that the dataset already provided, renaming the name of the *fos* in the *keyword* of the metadata.

```

1 data[i]['metadata'] = []
2 if 'fos' in data[i]:
3     for j in range(len(data[i]['fos'])):
4         data[i]['metadata'].append({'keyword': str(data[i]['fos'][j]['name']), 'weight': data[i]['fos'][j]['w']})

```

6.4. Publication details

For the publication details, we partially used the field *venue*, assuming that each venue produces a journal, with the name of its raw, in which these papers are collected. For the *date*, we only had the attribute *year* of the paper; anyway, for the same reason of giving a minimum of consistency to data, we extracted a random date in that year. Then, for the integer values of *number*, *volume* and *pages* we assigned a random one. Finally, we set a fixed publisher randomly to some papers that didn't have it.

```

1 data[i]['venue']['date'] = date.fromordinal(
2     random.randint(date(day=1, month=1, year=data[i]['year']).toordinal()
3         (), date(day=31, month=12, year=data[i]['year']).toordinal()))
4     ).isoformat()
5 data[i]['venue']['volume'] = random.randint(1, 3)
6 data[i]['venue']['number'] = random.randint(1, 20)
7 data[i]['venue']['pages'] = random.randint(1, 20)
8 if data[i]['publisher'] == "" and i%2 == 0:
9     data[i]['publisher'] = "PoliPrint, Milano"

```

6.5. Sections

As said before, we had nothing about sections in the dataset, so we inserted them after creating from the ground up. A section has been modeled as an array composed of:

- A title;
- An array with the texts of the paragraphs;
- An array with a list of the figures, which are objects consisting in:
 - URL of the image, which has as values passed the id of the paper, the section in which it is present and the number of the image;
 - Caption of the image;
- An array with the subsections, which are objects consisting in:
 - Title of the subsection;
 - Array with the text of the various paragraphs;

```

1 data[i]['sections'] = []
2 for j in range(random.randint(1, 5)):
3     data[i]['sections'].append({'title': "Title of section " + str(j+1),
4                                 'textbyparagraph': [],
5                                 'figures': []})
6     data[i]['sections'][j]['subsections'] = []
7     for p in range(random.randint(1, 5)):
8         data[i]['sections'][j]['subsections'].append({'title': "Title of
9             subsection " + str(p+1),
10            ,
11            'textbyparagraph': []})
12            for k in range(random.randint(1, 5)):
13                data[i]['sections'][j]['subsections'][p]['textbyparagraph'].append({'paragraph': lorem.words(80)})

```

```

12
13     for p in range(random.randint(0,4)):
14         data[i]['sections'][j]['textbyparagraph'].append({'paragraph':
15             lorem.words(80)})
16
17     for p in range(random.randint(0, 3)):
18         data[i]['sections'][j]['figures'].append({"url": 'https://image.
com/?' + urllib.parse.urlencode( {'paperid': data[i]['id'], 'section':
: j+1, 'image': p+1}),
"caption": 'Caption of image ' + str(j+1) + '.' + str(p+1)})
```

6.6. Bibliography

For the *bibliography*, we created a *references* array with some random ids of the papers present in the dataset (excluding its own). *vect_id* is the array containing all the ids of the dataset.

```

1 data[i]['references'] = []
2     for j in range(random.randint(1, 6)):
3         number = vect_id[random.randint(0, n_doc-1)]
4         while number in data[i]["references"]:
5             number = vect_id[random.randint(0, n_doc-1)]
6         data[i]["references"].append(number)
```

7 | Queries in MongoDB

7.1. Data creation/update

Here are the five data creation/update commands.

7.1.1. Query 1

Simple query for the creation of a new document. In this example only some fields are filled for simplicity.

```

1 db.papers.insertOne({
2   "_id" : 1,
3   "title" : "New Paper",
4   "n_citation" : 0,
5   "doc_type" : "Conference",
6   "authors": [
7     {
8       "name": "Caterina Giardi",
9       "id": 1,
10      "email": "caterina.giardi@mail.polimi.it",
11      "bio": "I'm a Master Degree student at Politecnico di Milano",
12      "affiliation": "Politecnico di Milano"
13    },
14    {
15      "name": "Stefano Baroni",
16      "id": 2,
17      "email": "stefano.baroni@mail.polimi.it",
18      "bio": "I'm a Master Degree student at Politecnico di Milano",
19      "affiliation": "Politecnico di Milano"
20    },
21    {
22      "name": "Alessia Menozzi",
23      "id": 3,
24      "email": "alessia.menozzi@mail.polimi.it",
25      "bio": "I'm a Master Degree student at Politecnico di Milano",
26      "affiliation": "Politecnico di Milano"
27    }
28  }
29}
```

```
27     {
28         "name": "Rocco Agnello",
29         "id": 4,
30         "email": "rocco.agnello@mail.polimi.it",
31         "bio": "I'm a Master Degree student at Politecnico di Milano",
32         "affiliation": "Politecnico di Milano"
33     },
34     {
35         "name": "Luca Brembilla",
36         "id": 5,
37         "email": "luca.brembilla@mail.polimi.it",
38         "bio": "I'm a Master Degree student at Politecnico di Milano",
39         "affiliation": "Politecnico di Milano"
40     ],
41     "venue" : {
42         "raw": "New Conference",
43         "id": 1,
44         "type": "C",
45         "volume": 4,
46         "number": 1,
47         "date": "2022-11-21",
48         "pages": 10
49     },
50     "abstract": "Abstract of the new paper",
51     "references": []
52 })
```

```

_id: 1
  title: "New Paper"
  n_citation: 0
  doc_type: "Conference"
  authors: Array
    > 0: Object
      name: "Caterina Giardi"
      id: 1
      email: "caterina.giardi@mail.polimi.it"
      bio: "I'm a Master Degree student at Politecnico di Milano"
      affiliation: "Politecnico di Milano"
    > 1: Object
    > 2: Object
    > 3: Object
    > 4: Object
  venue: Object
    raw: "New Conference"
    id: 1
    type: "C"
    volume: 4
    number: 1
    date: "2022-11-21"
    pages: 10
  abstract: "Abstract of the new paper"
  references: Array

```

Figure 7.1: Paper **after** the execution

7.1.2. Query 2

It is a simple update query in which the publisher of a specific paper, found by id, it's changed in "*PoliPrint, Milano*".

```

1 db.project.updateOne(
2   {"_id":8373},
3   {"$set": {"publisher":"PoliPrint, Milano"}}
4 )

```

```

_id: 8373
> authors: Array
  title: "Formal agent-oriented ubiquitous computing: a computational intelligen..."
  n_citation: 3
  doc_type: "Conference"
  publisher: "Springer, Berlin, Heidelberg"
  doi: "10.1007/978-3-642-28490-8_52"
> references: Array
> venue: Object
  abstract: "Abstract of Formal agent-oriented ubiquitous computing: a computationa..."
> metadata: Array
> bibliography: Array
> sections: Array

```

Figure 7.2: Paper **before** the changes

```

_id: 8373
> authors: Array
  title: "Formal agent-oriented ubiquitous computing: a computational intelligen..."
  n_citation: 3
  doc_type: "Conference"
  publisher: "PoliPrint, Milano" [highlighted]
  doi: "10.1007/978-3-642-28490-8_52"
> references: Array
> venue: Object
  abstract: "Abstract of Formal agent-oriented ubiquitous computing: a computationa..."
> metadata: Array
> bibliography: Array
> sections: Array

```

Figure 7.3: Paper **after** the changes

7.1.3. Query 3

In this update query we check for every keyword whose corresponding weight is equal to 0. If found the entry is removed from the metadata array.

```

1 db.papers.updateMany(
2   {
3     "metadata.weight": 0
4   },
5   {
6     $pull: {metadata: {weight:0}}
7 })

```

```

    ▼ metadata: Array
      ▼ 0: Object
        keyword: "Graph"
        weight: 0
      ▼ 1: Object
        keyword: "Discrete mathematics"
        weight: 0.45872
      ▼ 2: Object
        keyword: "Combinatorics"
        weight: 0.4515
      ▼ 3: Object
        keyword: "Direct product"
        weight: 0.59104
      ▼ 4: Object
        keyword: "Mathematics"
        weight: 0.42784
  
```

Figure 7.4: Before executing query 3

```

_id: 1388
> authors: Array
  title: "Further Results on Independence in Direct-Product Graphs."
  n_citation: 4
  doc_type: "Journal"
  publisher: "PoliPrint, Milano"
  doi: ""
> venue: Object
  abstract: "Abstract of Further Results on Independence in Direct-Product Graphs. ..."
  ▼ metadata: Array
    ▼ 0: Object
      keyword: "Discrete mathematics"
      weight: 0.45872
    ▼ 1: Object
      keyword: "Combinatorics"
      weight: 0.4515
    ▼ 2: Object
      keyword: "Direct product"
      weight: 0.59104
    ▼ 3: Object
      keyword: "Mathematics"
      weight: 0.42784
> references: Array
> bibliography: Array
> sections: Array
  
```

Figure 7.5: Result of Query 3

7.1.4. Query 4

The following update query removes the field doi when it is empty.

```

1 db.papers.updateMany({
2   "doi": ""
3 },
4 {
5   $unset: {doi: ""}
6 })

```

```

_id: 1388
> authors: Array
  title: "Further Results on Independence in Direct-Product Graphs."
  n_citation: 4
  doc_type: "Journal"
  publisher: "PoliPrint, Milano"
  doi: ""
> venue: Object
  abstract: "Abstract of Further Results on Independence in Direct-Product Graphs. : ..."
  > metadata: Array
    > 0: Object
      keyword: "Discrete mathematics"
      weight: 0.45872
    > 1: Object
      keyword: "Combinatorics"
      weight: 0.4515
    > 2: Object
      keyword: "Direct product"
      weight: 0.59104
    > 3: Object
      keyword: "Mathematics"
      weight: 0.42784
  > references: Array
  > bibliography: Array
  > sections: Array

```

Figure 7.6: Before executing update query 4

```

_id: 1388
> authors: Array
  title: "Further Results on Independence in Direct-Product Graphs."
  n_citation: 4
  doc_type: "Journal"
  publisher: "PoliPrint, Milano"
> venue: Object
  abstract: "Abstract of Further Results on Independence in Direct-Product Graphs. : ..."
  > metadata: Array
  > references: Array
  > bibliography: Array
  > sections: Array

```

Figure 7.7: Result of update query 4

7.1.5. Query 5

The last query is an example of delete. The query removes every paper that has a number of citation equals to 1.

```
1 db.papers.deleteMany({
2   "n_citation": 1
3 })
```

	_id Int32	authors Array	title String	n_citation Int32	doc_type String
1	1091	[] 8 elements	"Preliminary Design of a Ne...	6	"Conference"
2	1388	[] 1 elements	"Further Results on Independ...	4	"Journal"
3	1674	[] 2 elements	"A methodology for the physi...	1	"Conference"
4	1688	[] 4 elements	"Comparison of GARCH, Neural...	5	"Conference"
5	5411	[] 4 elements	"COMPARING GNG3D AND QUADRI...	4	"Conference"
6	5781	[] 2 elements	"Vectorial fast correlation ...	4	"Repository"
7	6522	[] 3 elements	"Improved Secret Image Sharin...	6	"Conference"
8	6762	[] 2 elements	"A Self-Stabilizing Algorith...	4	"Conference"
9	8373	[] 1 elements	"Formal agent-oriented ubiquit...	3	"Conference"
10	8763	[] 2 elements	"Fur Visualisation for Comp...	1	"Conference"
11	382732	[] 2 elements	"Bridging the Gap Between Sci...	4	"Conference"
12	385572	[] 3 elements	"A Generic Framework Based on...	3	"Conference"
13	385606	[] 3 elements	"Assessing the State of Too...	2	"Journal"

Figure 7.8: Before executing update query 5

	_id Int32	authors Array	title String	n_citation Int32	doc_type String
1	1091	[] 8 elements	"Preliminary Design of a Ne...	6	"Conference"
2	1388	[] 1 elements	"Further Results on Independ...	4	"Journal"
3	1688	[] 4 elements	"Comparison of GARCH, Neural...	5	"Conference"
4	5411	[] 4 elements	"COMPARING GNG3D AND QUADRI...	4	"Conference"
5	5781	[] 2 elements	"Vectorial fast correlation ...	4	"Repository"
6	6522	[] 3 elements	"Improved Secret Image Sharin...	6	"Conference"
7	6762	[] 2 elements	"A Self-Stabilizing Algorith...	4	"Conference"
8	8373	[] 1 elements	"Formal agent-oriented ubiquit...	3	"Conference"
9	382732	[] 2 elements	"Bridging the Gap Between Sci...	4	"Conference"
10	385572	[] 3 elements	"A Generic Framework Based on...	3	"Conference"
11	385606	[] 3 elements	"Assessing the State of Too...	2	"Journal"
12	385951	[] 3 elements	"A bimodal Korean address enc...	3	"Conference"
13	385984	[] 3 elements	"The Encourage Operators to ...	2	"Conference"

Figure 7.9: Result of update query 5

7.2. Queries

7.2.1. Filtering conditions

Query 1

This query matches two conditions that have to be both satisfied: having 564427 between the array of referenced IDs and having at least 5 sections. We show the title, the references array and the titles of the sections (to show that section 5 is present).

```

1 db.papers.find({
2   $and: [{references: { $in: [564427]}}, { "sections.4": {$exists: true
3 }}]
4 }, {
5   "title":1,
6   "references":1,
7   "sections.title":1
8 })

```

```

< { _id: 1091,
  title: 'Preliminary Design of a Network Protocol Learning Tool Based on the Comprehension of High School Students',
  references: [ 804213, 633195, 30346439, 608724, 621898, 564427 ],
  sections:
    [ { title: 'Title of section 1' },
      { title: 'Title of section 2' },
      { title: 'Title of section 3' },
      { title: 'Title of section 4' },
      { title: 'Title of section 5' } ] }
{ _id: 527403,
  title: 'From Sequential Processes to Grid Computation.',
  references: [ 659904, 703288, 792443, 473870, 378968, 564427 ],
  sections:
    [ { title: 'Title of section 1' },
      { title: 'Title of section 2' },
      { title: 'Title of section 3' },
      { title: 'Title of section 4' },
      { title: 'Title of section 5' } ] }
{ _id: 840705,
  title: 'The Role of Genetic Programming in Describing the Microscopic Structure of Hydrating Plaster',
  references: [ 460995, 510479, 418674, 499447, 796041, 564427 ],
  sections:
    [ { title: 'Title of section 1' },
      { title: 'Title of section 2' },
      { title: 'Title of section 3' },
      { title: 'Title of section 4' },
      { title: 'Title of section 5' } ] }
SMBUDProject> 
```

Figure 7.10: Result of query 1

Query 2

The following query filters the documents by checking if they have more than 18 pages and that the second section has at least 5 subsections. Once a document is found, the

query returns the title of the paper, the title of the fifth subsection and the number of pages.

```

1 db.papers.find({
2   "venue.pages": {$gt: 18},
3   "sections.1.subsections.4": {$exists: true}
4 },
5 {
6   "title": 1,
7   "sections[1].subsections[4].title": 1,
8   "venue.pages": 1
9 })

```

```

< { _id: 423593,
  title: 'Relative evaluation of informativeness in machine generated summaries.',
  venue: { pages: 19 } }
{ _id: 442490,
  title: 'Towards A Universal Tool For NLP Resource Acquisition.',
  venue: { pages: 19 } }
{ _id: 460995,
  title: 'Determinant of the Generalized Lucas RSFMLR Circulant Matrices in Communication',
  venue: { pages: 19 } }
{ _id: 628729,
  title: 'Evaluation of an alert system for selective dissemination of broadcast news',
  venue: { pages: 20 } }
{ _id: 505434,
  title: 'A knowledge-based approach to ontology learning and semantic annotation.',
  venue: { pages: 20 } }
{ _id: 687179,
  title: 'How clean is your software? The role of software validation in digital preservation research projects',
  venue: { pages: 19 } }
{ _id: 797742,
  title: 'On attributed relational graph matching using hopfield network',
  venue: { pages: 20 } }
{ _id: 804384,
  title: 'Some results on geometric independency trees.',
  venue: { pages: 20 } }
SMBUDProject>

```

Figure 7.11: Result of query 2

7.2.2. Filtering conditions and aggregation

Query 1

Retrieve the maximum top 5 publishers per amount of citations of their most cited paper. The filter is applied on publisher, which has to exist and not to be equal to the empty string.

```

1 db.papers.aggregate([
2   {

```

```

3  /*We make sure that the publisher field is not empty and exists*/
4  "$match": {"publisher" : {"$exists" : true, "$ne" : ""}}
5 },
6 {
7 /*We group by the publisher and we record the biggest number of
8   citations in one of the papers the publisher appears in*/
9   "$group": {
10     "_id": "$publisher",
11     "maxNumOfCitations": { "$max": "$n_citation" }
12   }
13 },
14 { "$sort": { "maxNumOfCitations": -1 } },
15 /* Showing top 5*/
16 { "$limit": 5 }
17 ])

```

```

{ _id: 'Springer', maxNumOfCitations: 6 }
{ _id: 'Springer, Berlin, Heidelberg', maxNumOfCitations: 6 }
{ _id: 'American Medical Informatics Association',
  maxNumOfCitations: 6 }
{ _id: 'University of California, Los Angeles',
  maxNumOfCitations: 6 }
{ _id: 'Rinton Press, Incorporated', maxNumOfCitations: 6 }

```

Figure 7.12: Result of the query 1 with aggregation

Query 2

Retrieve the average authors for each publisher (the publisher must exist and not be the empty string) of a document of type *Conference*, then sort the results in descending order.

```

1 db.papers.aggregate([
2 /*Matching conditions/
3 { "$match":
4   { "$and":
5     [ { "doc_type": "Conference" }, { "publisher": { "$exists" :
6       true, "$ne": "" } } ] }
7 },
8 /*Group by publisher and we compute the average number of authors*/
9 { "$group":
10   { "_id": "$publisher",
11     "avg_authors": { "$avg": "$n_authors" }
12   }
13 },
14 { "$sort": { "avg_authors": -1 } }
15 ]

```

```

9           "avgAuthorPerPublisher": {"$avg": { $size: "$authors"
10      }}})
11  /*Descending order*/
12  { "$sort": { "avgAuthorPerPublisher": -1 }}])

```

```

_id: "International Society for Music Information Retrieval (ISMIR)"
avgAuthorPerPublisher: 10

_id: "Stud Health Technol Inform"
avgAuthorPerPublisher: 5.5

_id: "CEUR-WS.org"
avgAuthorPerPublisher: 5.333333333333333

_id: "Aka GmbH"
avgAuthorPerPublisher: 4

_id: "Springer, Vienna"
avgAuthorPerPublisher: 4

```

Figure 7.13: Top 5 results retrieved

7.2.3. Filtering conditions also on internal sub-documents

Query 1

Retrieve the papers written by "Pranay Chaudhuri" (at least, it may have other authors), they need to include the keyword "Computer science" and in the title they need to have the string "Algorithm". It returns only the title, the authors' name and the metadata of each document.

```

1 db.papers.find({
2 /*Filters the three conditions explained above */
3   "$and": [
4     {
5       "authors": {
6         "$elemMatch": {"name": "Pranay Chaudhuri"}
7       }
8     },
9     {
10       "metadata": {
11         "$elemMatch": {"keyword": "Computer science"}
12       }
13     },
14     {
15       "title": { $regex : "Algorithm" }
16     }
17   ]
18 },
19 /*Display the fields*/

```

```

20 {
21   "title": 1,
22   "authors.name": 1,
23   "metadata": 1
24 })

```

```

{ _id: 6762,
  authors: [ { name: 'Pranay Chaudhuri' }, { name: 'Hussein Thompson' } ],
  title: 'A Self-Stabilizing Algorithm for Finding the Cutting Center of a Tree.',
  metadata:
    [ { keyword: 'Computer science', weight: '0.44568' },
      { keyword: 'Parallel computing', weight: '0.46131' } ] }

```

Figure 7.14: A document retrieved by the query 1 with filters on sub-documents

Query 2

Retrieve the papers which have exactly 2 authors, the date of the venue in which the paper was presented is 2014 and in the metadata we have the keyword "Computer Science" with a weight greater or equal than 0.40. The query doesn't retrieve all the document but also some properties (i.e. title, authors, etc.).

```

1 db.papers.find({
2 /*Filters the three conditions explained above */
3   "$and": [
4     {
5       "authors": { "$size": 2 }
6     },
7     {
8       "metadata": {
9         "$elemMatch" : {
10           "keyword": { $regex: "Artificial intelligence" },
11           "weight": { "$gte": 0.40 }
12         }
13       }
14     },
15     {
16       "venue.date" : { "$regex" : "2014" }
17     }
18   ],
19   {
20   /*Display the fields*/
21     "title": 1,
22     "authors.name": 1,
23     "venue.raw": 1,

```

```

24     "venue.date": 1,
25     "n_citation": 1,
26     "metadata": 1
27 }
28 ).sort({"n_citations": -1}).limit(3)

```

```

{
  "_id: 1674,
  authors: [ { name: 'G. Beale' }, { name: 'G. Earl' } ],
  title: 'A methodology for the physically accurate visualisation of roman polychrome statuary',
  n_citation: 6,
  venue:
    { raw: 'International Conference on Virtual Reality',
      date: '2014-04-12' },
  metadata:
    [ { keyword: 'Statue', weight: 0.40216 },
      { keyword: 'Engineering drawing', weight: 0.43427 },
      { keyword: 'Virtual reconstruction', weight: 0 },
      { keyword: 'Computer science', weight: 0.42062 },
      { keyword: 'Visualization', weight: 0.4595 },
      { keyword: 'Polychrome', weight: 0.4474 },
      { keyword: 'Artificial intelligence', weight: 0.40496 } ] }
{ _id: 840705,
  authors:
    [ { name: 'Judith Ellen Devaney' },
      { name: 'John G. Hagedorn' } ],
  title: 'The Role of Genetic Programming in Describing the Microscopic Structure of Hydrating Plaster',
  n_citation: 3,
  venue:
    { raw: 'Genetic and Evolutionary Computation Conference',
      date: '2014-05-17' },
  metadata:

```

Figure 7.15: Result of query 2 with filters on sub-documents

7.2.4. Filtering conditions through unwind

Query 1

This query first selects only the papers whose type is *Journal* and then, through **\$unwind** it shows the top 10 keywords that are the most used (appear in the *metadata* array) in these kind of papers.

```

1 db.papers.aggregate([
2 {
3 /*We filter the doc_type*/
4   $match: { "doc_type": "Journal" }
5 },
6 {
7 /*We unwind metadata array*/

```

```

8     $unwind: { path: "$metadata" }
9 },
10 {
11 /*We count for each keyword the number of times it appears inside a
12   paper*/
13   $group: {
14     _id: "$metadata.keyword",
15     count: { $sum: 1 }
16   }
17 {
18 /*Descending order*/
19   $sort: { count: -1 }
20 },
21 {
22 /*Showing the top 10*/
23   $limit: 10
24 }
25 ])

```

```

{ _id: 'Computer science', count: 44 }
{ _id: 'Mathematics', count: 20 }
{ _id: 'Artificial intelligence', count: 12 }
{ _id: 'Discrete mathematics', count: 11 }
{ _id: 'Knowledge management', count: 10 }
{ _id: 'Mathematical optimization', count: 6 }
{ _id: 'Software engineering', count: 6 }
{ _id: 'Algorithm', count: 6 }
{ _id: 'Distributed computing', count: 6 }
{ _id: 'Combinatorics', count: 5 }

```

Figure 7.16: Result of filtering using unwind 1

Query 2

This query find the percentages of the presence inside the whole documents of the different types of papers (*Repository*, *Journal*, *Conference*) or when it's left blank.

```

1 db.papers.aggregate([
2 {
3 /*We group and count for each type the number of corresponding papers*/
4   $group: {

```

```
5      _id: '$doc_type',
6      counttypes: {
7          $sum: 1
8      }
9  },
10 }, {
11 /*We dinamically count the total number of papers and we push into an
   array the name of the type and the corresponding quantity*/
12 $group: {
13     _id: null,
14     total: {
15         $sum: '$counttypes'
16     },
17     Typequantity: {
18         $push: {
19             name: '$_id',
20             npapers: '$counttypes'
21         }
22     }
23 },
24 }, {
25 /*We unwind the array of types also adding an index field*/
26 $unwind: {
27     path: '$Typequantity',
28     includeArrayIndex: 'IndexType',
29     preserveNullAndEmptyArrays: false
30 }
31 }, {
32 /*We calculate the percentage and we project the result*/
33 $project: {
34     _id: '$Typequantity.name',
35     percentage: {
36         $multiply: [
37             {
38                 $divide: [
39                     '$Typequantity.npapers',
40                     '$total',
41                 ]
42             },
43             100
44         ]
45     }
46 }
47 }
```

```
48 ]
49 )
```



Figure 7.17: Result of unwind query 2

7.2.5. Filtering condition through refs

Query 1

It retrieves, for each paper which has at least three references, the percentage of them which were published by "*PoliPrint, Milano*". The results are shown in decreasing order with respect to the percentage.

It is a pipeline of several operators; the first **\$match** is useful to find all the papers with at least four references and it is done as first operation in order to have less document on which do the **\$lookup**. After the **\$lookup**, it finds the referenced papers with that publisher and soon after, through the first **\$project** it counts the total number of the references and the ones with the publisher required. Then, through another **\$project** it makes the division between them and multiplying it by 100 we obtain the percentages and finally we sort the results.

```

1 db.papers.aggregate([
2   /*Matches the papers with at least 3 references*/
3   {"$match": {"references.3": {"$exists": true}}}, 
4   /*Perform the lookup on the same collection papers */
5   {"$lookup": {
6     from: "papers",
7     localField: "references",
8     foreignField: "_id",
9     as: "refs"
10   }},
11  /*Matches the papers whose publisher is PoliPrint, Milano*/
12  {"$match": {"refs.publisher": "PoliPrint, Milano"}}, 
13  {"$project": {
14    /*Compute the total number of references(n_total_references) and the
      total references with the publisher matched before(n_refs_poliprint)
      */
15    "_id": 1,
16    "n_total_references": {"$size": "$refs"}, 
17    "n_refs_poliprint": {"$size": {"$filter": {"input": "$refs", "cond": {"$eq": ["$$this.publisher", "PoliPrint, Milano"]}}}}}, 
18  {"$project": {
19    /*Compute the percentages through a project*/
20    "_id": 1,
21    "percentage_of_poliprint_refs": {
22      "$trunc": [
23        {"$multiply": [
24          {"$divide": ["$n_refs_poliprint", "n_total_references"]}], 100}], 2]}}, 
23  /*Descending order*/
24  {"$sort": {"percentage_of_poliprint_refs": -1}}
25 ])

```

_id: 2532333
percentage_of_poliprint_refs: 83.33
_id: 1829660
percentage_of_poliprint_refs: 80
_id: 2394107
percentage_of_poliprint_refs: 80
_id: 2532194
percentage_of_poliprint_refs: 80
_id: 1987092
percentage_of_poliprint_refs: 80

Figure 7.18: Top 5 results retrieved

Query 2

For this query we are using a lookup on two different collections. Since we decided at the beginning to create only one collection for the entire db, we manually created a json called *country* that contains the fields : *country*, *affiliation* and *continent*. We took some of the affiliations and added the country and the continent in order to perform aggregations. So this query performs the **\$lookup** between the collection paper and Country and returns, for each country listed in *country.json*, the list of the authors whose affiliation belongs to that country and the title of the paper they wrote.

```
Country.json prova.json +  
1 {"_id": 1, "country": "Italy", "affiliation": "Università della Calabria", "continent": "Europe"}  
2 {"_id": 2, "country": "Italy", "affiliation": "Università degli Studi di Modena e Reggio Emilia", "continent": "Europe"}  
3 {"_id": 3, "country": "Italy", "affiliation": "Polytechnic University of Turin#TAB#", "continent": "Europe"}  
4 {"_id": 4, "country": "Italy", "affiliation": "University of Milano-Bicocca", "continent": "Europe"}  
5  
6 {"_id": 5, "country": "Germany", "affiliation": "Cognitive Interaction Technology, Center of Excellence and Applied Informatics, Faculty of Technology, Bielefeld University, Bielefeld, Germany"}  
7 {"_id": 6, "country": "Germany", "affiliation": "Department of Information Systems and Business Administration, Johannes Gutenberg-Universität Mainz, Jakob-Welder-Weg 9, Mainz 55128, Germany"}  
8  
9 {"_id": 7, "country": "Bangladesh", "affiliation": "Department of Computer Science & Engineering, BUET, Dhaka, Bangladesh", "continent": "Asia"}  
10 {"_id": 8, "country": "Bangladesh", "affiliation": "Department of Statistics , Rajshahi University , Rajshahi , Bangladesh", "continent": "Asia"}  
11  
12 {"_id": 9, "country": "Mexico", "affiliation": "New Mexico State University,#TAB#", "continent": "America"}  
13 {"_id": 10, "country": "Mexico", "affiliation": "UNAM, Circuito Exterior s/n, 04510, Coyocan, Mexico D.F., Mexico#TAB#", "continent": "America"}  
14  
15 {"_id": 11, "country": "USA", "affiliation": "University of Southern California", "continent": "America"}  
16 {"_id": 12, "country": "USA", "affiliation": "University of Massachusetts,Boston", "continent": "America"}  
17 {"_id": 13, "country": "USA", "affiliation": "North Dakota State University,.", "continent": "America"}  
18 {"_id": 14, "country": "USA", "affiliation": "College of Computing, Georgia Institute of Technology Atlanta", "continent": "America"}  
19 {"_id": 15, "country": "USA", "affiliation": "Department of Radiology, The University of Chicago Chicago IL", "continent": "America"}  
20 {"_id": 16, "country": "USA", "affiliation": "Oakland University", "continent": "America"}  
21  
22 {"_id": 17, "country": "Spain", "affiliation": "Dept. of Computer Science, Uni. Politécnica de Catalunya, Spain#TAB#", "continent": "Europe"}  
23 {"_id": 18, "country": "Spain", "affiliation": "University of malaga", "continent": "Europe"}  
24  
25 {"_id": 19, "country": "Canada", "affiliation": "Univ. of Western Ontario#TAB#", "continent": "America"}  
26 {"_id": 20, "country": "Canada", "affiliation": "Electrical and Computer Engineering Department, University of Waterloo , Waterloo , Ontario , Canada", "continent": "America"}  
27 {"_id": 21, "country": "Canada", "affiliation": "School of Computer Science Carleton University Canada", "continent": "America"}  
28 {"_id": 22, "country": "Canada", "affiliation": "School of Electrical Eng. and Computer Science, University of Ottawa, Canada", "continent": "America"}  
29  
30 {"_id": 23, "country": "UK", "affiliation": "Computer Laboratory, University Of Cambridge", "continent": "Europe"}  
31 {"_id": 24, "country": "UK", "affiliation": "New England Aquarium", "continent": "Europe"}  
32 {"_id": 25, "country": "UK", "affiliation": "Unconventional Computing Group, University of the West of England, Bristol, UK BS16 10Y#TAB#", "continent": "Europe"}  
33  
34 {"_id": 26, "country": "Japan", "affiliation": "Corporate Research and Development Center, Toshiba Corporation, Kawasaki, Japan", "continent": "Asia"}  
35 {"_id": 27, "country": "Japan", "affiliation": "Institute of Industrial Science, University of Tokyo, Meguro-Ku, Japan", "continent": "Asia"}  
36 {"_id": 28, "country": "Japan", "affiliation": "Department of Human and Artificial Intelligence System, Graduate School of Engineering, University of Fukui, Japan#TAB#", "continent": "Asia"}  
37 {"_id": 29, "country": "Japan", "affiliation": "Department of Computer Science, Faculty of Engineering, Yamamanishi University, Kofu, Japan", "continent": "Asia"}  
38 {"_id": 30, "country": "Japan", "affiliation": "Nippon Hoso Kyokai", "continent": "Asia"}  
39  
40 {"_id": 31, "country": "France", "affiliation": "ESIEE-Paris", "continent": "Europe"}  
41 {"_id": 32, "country": "France", "affiliation": "UNIVERSITÉ LAVAL", "continent": "Europe"}  
42 {"_id": 33, "country": "Denmark", "affiliation": "University of Southern Denmark", "continent": "Europe"}  
43 {"_id": 34, "country": "Portugal", "affiliation": "University of Lisbon", "continent": "Europe"}  
44 {"_id": 35, "country": "UK", "affiliation": "##University of Glasgow##", "continent": "Europe"}  
45 {"_id": 36, "country": "Japan", "affiliation": "Shinshu University", "continent": "Asia"}  
Line: 25:42 JSON Tab Size: 4    
Line: 25:42 JSON Tab Size: 4  
```

Figure 7.19: Screenshot of the Country.json file created

```
1 db.papers.aggregate([{
2   /* needed for the lookup because authors is an array*/
3   $unwind: {
4     path: '$authors',
5   }
6 }, {
7   /*lookup between the two collections based on the field affiliation*/
8   $lookup: {
9     from: 'Country',
10    localField: 'authors.affiliation',
11    foreignField: 'affiliation',
12    as: 'affiliationDescription'
```

```
13  }
14 }, {
15 /*Since not all the affiliations are present in the json document, I
   only consider the papers whose authors found a match*/
16 $match: {
17   'affiliationDescription.0': {
18     $exists: true
19   }
20 }
21 }, {
22 /*Group by country and push into an array the informations we need to
   project*/
23 $group: {
24   _id: '$affiliationDescription.country',
25   authorsnames: {
26     $push: {
27       name: '$authors.name',
28       affiliation: '$authors.affiliation',
29       PaperTitle: '$title'
30     }
31   }
32 }
33 }, {
34   $project: {
35     _id: 1,
36     authorsnames: 1
37   }
38 }])
```

```

MONGOSH
{
  "_id": [ "Portugal" ],
  "authorsnames": [
    { "name": "João Ricardo Silva",
      "affiliation": "University of Lisbon",
      "PaperTitle": "Assigning Deep Lexical Types" },
    { "name": "António Branco",
      "affiliation": "University of Lisbon",
      "PaperTitle": "Assigning Deep Lexical Types" } ] }
{
  "_id": [ "Spain" ],
  "authorsnames": [
    { "name": "Gemma C. Garriga",
      "affiliation": "Dept. of Computer Science, Uni. Politècnica de Catalunya, Spain#TAB#",
      "PaperTitle": "On mining closed sets in multi-relational data" },
    { "name": "Amine Kerzazi",
      "affiliation": "university of malaga",
      "PaperTitle": "A Semantic Mediation Architecture for RDF Data Integration." },
    { "name": "Ismael Navas Delgado",
      "affiliation": "university of malaga",
      "PaperTitle": "A Semantic Mediation Architecture for RDF Data Integration." },
    { "name": "José Francisco Aldana Montes",
      "affiliation": "university of malaga",
      "PaperTitle": "A Semantic Mediation Architecture for RDF Data Integration." },
    { "name": "Othmane Chniber",
      "affiliation": "university of malaga",
      "PaperTitle": "A Semantic Mediation Architecture for RDF Data Integration." } ] }

```

Figure 7.20: Result of lookup query 2

7.3. Performance of the queries

Here's are the performance for a particular execution of the queries on a database hosted on the local machine (obviously times change if repeated, but we can see here how they differ in a simple way). The database contains 1000 documents.

Query	Time (ms)
Filtering conditions - query 1	3
Filtering conditions - query 2	8
Filtering conditions and aggregation - query 1	19
Filtering conditions and aggregation - query 2	10
Filtering conditions on sub-documents - query 1	3
Filtering conditions on sub-documents - query 2	10
Filtering conditions through unwind - query 1	9
Filtering conditions through unwind - query 2	11
Filtering condition through refs - query 1	174
Filtering condition through refs - query 2	605

Table 7.1: Table of execution times of the queries

7.3.1. A closer look at performance with pipelines

Let's take a closer look at performance for the query 7.2.2. We can write the same query in other two different way, using two stages for the filters, instead of using one with a logical and. So, the first way can be:

```

1 db.papers.aggregate([
2   { "$match": { "publisher": { "$exists": true, "$ne": "" } } },
3   { "$match": { "doc_type": "Conference" } },
4   { "$group": {
5     "_id": "$publisher",
6     "avgAuthorPerPublisher": {
7       "$avg": { $size: "$authors" } }
8   }},
9   { "$sort": { "avgAuthorPerPublisher": -1 } }
10 ])
11

```

and the second one:

```

1 db.papers.aggregate([
2   { "$match": { "doc_type": "Conference" } },
3   { "$match": { "publisher": { "$exists": true, "$ne": "" } } },
4   { "$group": {
5     "_id": "$publisher",
6     "avgAuthorPerPublisher": {
7       "$avg": { $size: "$authors" } }
8   }}
9

```

```

9     },
10    { "$sort": { "avgAuthorPerPublisher": -1 } }
11 ])

```

Reminding aggregations are pipelines, i.e. operations included in the query are performed sequentially, we expect different performance. For instance, for the second query, if we exclude nearly all the documents matching only the papers that are presented in a conference, we have then to apply only to these papers the filter on the publisher (a paper can be presented in a conference even without a specified publisher).

The execution times of the first query for three different executions are 8, 10 and 13. For the second one we get 11, 9 and 7. So the second one is more efficient ($31 > 27$).

To see the reason, we can see how many papers there are with the filters. For the Conference filter we get

```
1 db.papers.find({ "doc_type": "Conference" }).count()
```

equals to 613 out of 1000, while we get

```
1 db.papers.find({ "publisher" : { "$exists" : true, "$ne" : "" } }).count()
```

equals to 759, so the first filter is more restrictive. Thus, the query executing it firstly is more efficient (the second query seen above).

Another aspect that we can highlight is the time for the queries that performs the lookup operation. That operation is like a join for the relational databases, so it's obviously slower than the others. That's also the reason why in MongoDB are used sub-documents.