



POLITECNICO
MILANO 1863

SCUOLA DI INGEGNERIA INDUSTRIALE
E DELL'INFORMAZIONE

Systems and Methods for Big and Unstructured Data Project

Author(s): **Balestrieri Niccolò - 10936955**

Bertogalli Andrea - 10702303

Tombini Nicolò - 10912627

Academic Year: 2023-2024

Contents

Contents	i
1 Introduction	1
2 Data wrangling / Data generation	3
3 Dataset	5
3.1 Neo4j	6
3.2 ElasticSearch	7
4 Queries	9
4.1 Neo4j	9
4.2 ElasticSearch	19
5 Extra	41
5.0.1 DISCLAIMER	41
5.1 WebApp structure	41
5.2 WebApp screens	42
5.2.1 Neo4j screens	42
5.2.2 ElasticSearch screens	45
List of Figures	47

1 | Introduction

In this report we will present a comprehensive analysis of the ArXiv dataset, a collaboratively funded, community-supported resource founded by Paul Ginsparg in 1991 and currently maintained and operated by Cornell University. Our project delves into the intricate web of academic knowledge encapsulated within ArXiv, utilizing advanced technologies to extract meaningful insights.

Our group strategically leveraged the power of two cutting-edge technologies to unravel the complexities inherent in the ArXiv dataset. The first technology used is Neo4j, a graph database management system. This choice was driven by the dataset's inherent graph-like structure, where nodes represent authors and articles, and relationships encapsulate the act of writing or not papers. Neo4j in handling complex relationships and its native support for graph structures made it the ideal choice for modeling and querying the interconnected web of academic contributions within the ArXiv repository.

Complementing our graph-oriented approach, we embraced Elasticsearch as our second technological choice. Elasticsearch, renowned for its robust full-text search capabilities, are well integrated with our project objectives. This technology addresses the need for powerful and efficient textual searches on the vast array of academic papers within the ArXiv dataset. Its flexibility and scalability make it a valuable asset in uncovering relevant information through textual queries, enhancing the overall usability and accessibility of the dataset.

In the following sections of this report, we will delve into our web application built on top of these two technologies, whose purpose is to show 20 queries (10 for each technology) in a user-friendly way.

In this chapter additional useful information are reported.

2 | Data wrangling / Data generation

The chosen dataset, in JSON format, presents several issues to address. Firstly, for computational and resource reasons, it was decided to limit the dataset to 30,000 randomly selected data points from the entire dataset.

Subsequently, cleaning was performed on the authors' names and dates. Moreover, upon examining a random abstract text, it was observed that the text is available in Latex format. Therefore, both the title and the abstract text were transformed to UTF-8 format using the `pylatexenc` library.

Before:

This result confirms that dominant contributions to the electric and magnetic polarizabilities may be represented in terms of two-photon couplings to the σ -meson having the predicted mass $m_\sigma=666$ MeV and two-photon width $\Gamma_{\gamma\gamma}=2.6$ keV.

After:

This result confirms that dominant contributions to the electric and magnetic polarizabilities may be represented in terms of two-photon couplings to the σ -meson having the predicted mass $m_\sigma = 666$ MeV and two-photon width $\Gamma_{\gamma\gamma} = 2.6$ keV.

Later on, we created the clean and well-formatted dataset in .csv format, which was then imported and utilized on Neo4j and ElasticSearch.

3 | Dataset

In general the dataset is a file JSON following this structure:

```

1 {
2   "id": "0704.0001",
3   "submitter": "Pavel Nadolsky",
4   "authors": "E. L. Berger, P. M. Nadolsky, C.-P. Yuan",
5   "title": "Paper title",
6   "comments": "37 pages, 15 figures; published version",
7   "journal-ref": "Phys.Rev.D76:013009,2007",
8   "doi": "10.1103/PhysRevD.76.013009",
9   "report-no": "ANL-HEP-PR-07-12",
10  "categories": "hep-ph",
11  "license": "",
12  "abstract": "Paper abstract",
13  "versions": [
14    0: {
15      "version": "v1"
16      "created": "Mon, 2 Apr 2007 19:18:42 GMT"
17    }
18  ],
19  "update_date": "2008-11-26"
20 }
```

After cleaning and parsing dataset from JSON into CSV format, we chose to use the same dataset for both technologies with slight modifications in terms of attributes. The details of these modifications will be explained in the following sections.

3.1. Neo4j

Regarding Neo4j, the dataset used considers authors and articles.

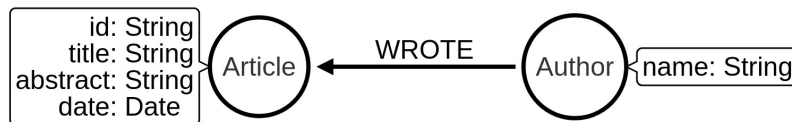


Figure 3.1: Neo4j DB structure.

As seen in the picture above, we have 2 types of nodes: **Article** and **Author**.

For **Article**, we have:

- **id**: article ID, String
- **title**: article title, String
- **abstract**: summary of the article content, String
- **date**: article publication date, Date

For **Author**, we have:

- **name**: author's name, String

Clearly, the relation **WROTE** is present if the author wrote the relative article.

3.2. Elasticsearch

Instead, for what concerns Elasticsearch, we started again from CSV dataset but we had to create a manual mapping, and for this purpose it was necessary to recover a sort of JSON file.

```
1 mapping = {
2     "settings": {
3         "number_of_shards": 3,
4         "number_of_replicas": 3
5     },
6     "mappings": {
7         "properties": {
8             "submitter": {"type": "keyword"},
9             "authors": {"type": "text"},
10            "title": {"type": "text"},
11            "abstract": {"type": "text"},
12            "comments": {"type": "text"},
13            "ref": {"type": "keyword"},
14            "categories": {"type": "keyword"},
15        }
16    }
17 }
```

Here we can spot these attributes:

- **submitter**: the author who submitted the journal, keyword
- **authors**: list of paper authors, text
- **title**: article title, text
- **abstract**: summary of the article content, text
- **comments**: comments related to an article, text
- **ref**: journal acronym, keyword
- **categories**: kind of categories related to an article, keyword

4 | Queries

4.1. Neo4j

Query 1. *Retrieves the articles from 2007 which have been written by at least 25 authors.*

```

1 MATCH (ar:Article)<-[:WROTE]-(au:Author)
2 WITH count(au) as number_of_authors, ar
3 WHERE date(ar.date).year = 2007 AND number_of_authors >= 25
4 MATCH p=(ar)<-[:WROTE]-(au:Author)
5 RETURN p;
```

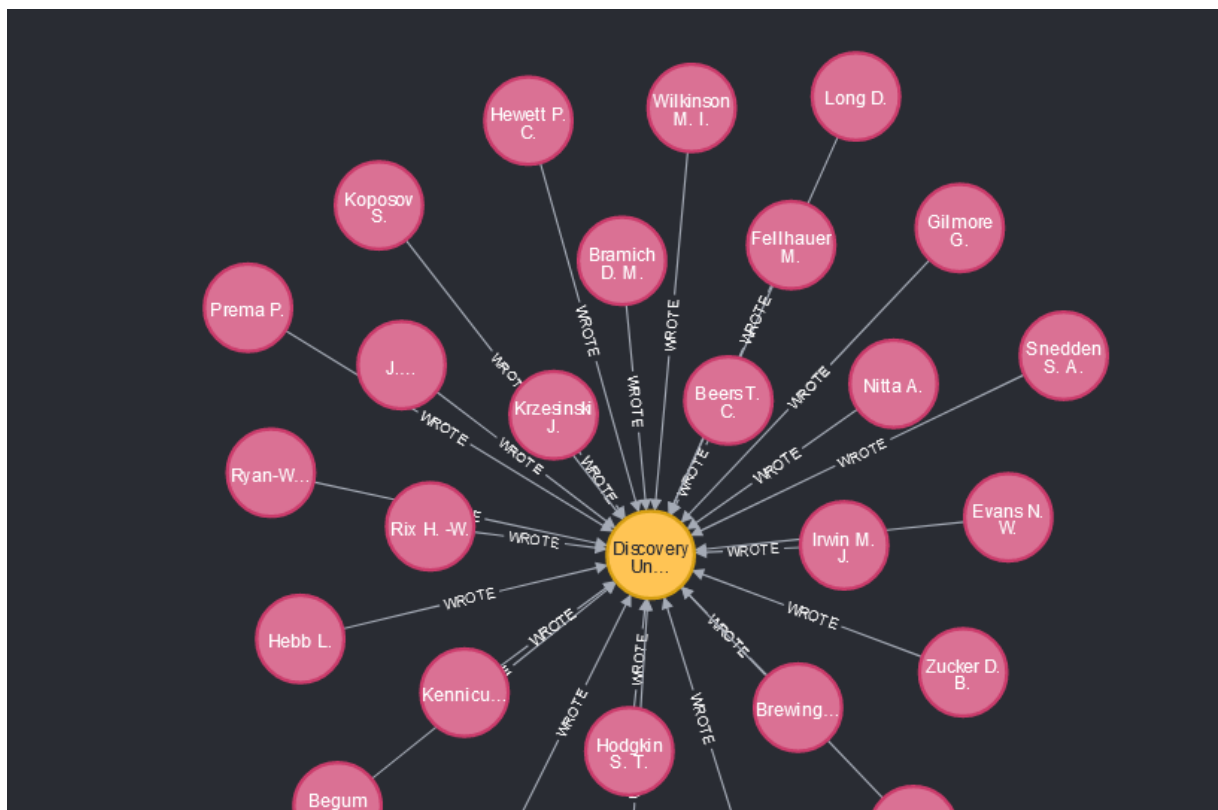


Figure 4.1: Partial outcome query 1

Query 2. *Retrieves the article with the highest number of authors.*

```

1 MATCH (ar:Article)<-[:WROTE]-(au:Author)
2 WITH ar, count(au) AS number_of_authors
3 WITH ar, number_of_authors
4 ORDER BY number_of_authors DESC
5 LIMIT 1
6 MATCH p=(ar)<-[:WROTE]-(Author)
7 RETURN p;
```

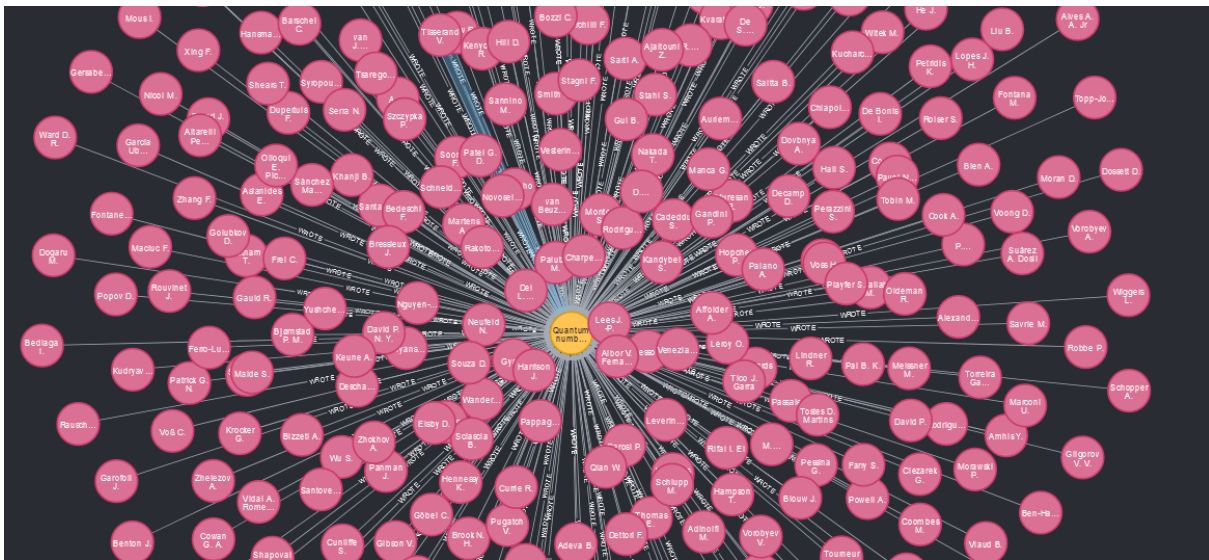


Figure 4.2: Partial outcome query 2

Query 3. *Returns all the papers that have Matteucci Matteo among the authors.*

```
1 MATCH (ar:Article)<-[:WROTE]-(au:Author)
2 WHERE au.name CONTAINS "Matteucci Matteo"
3 MATCH p=(ar)<-[:WROTE]-(Author)
4 RETURN p;
```

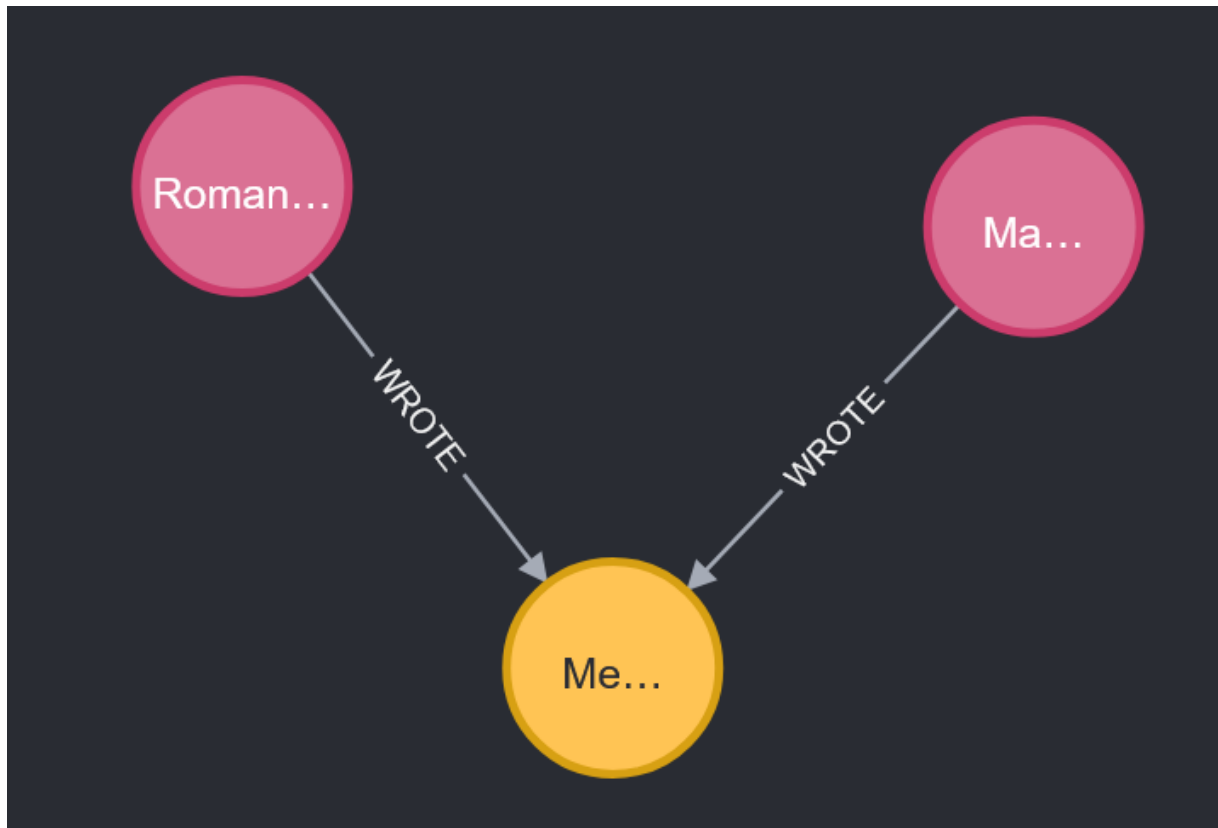


Figure 4.3: Query 3

Query 4. *Returns all the collaborators of Yann LeCun and their articles*

```
1 MATCH (ar:Article)<-[:WROTE]-(y:Author {name: "LeCun Yann"}),  
2 (au:Author)-[:WROTE]->(ar)  
3 MATCH p=(:Article)<-[:WROTE]-(au)  
4 RETURN p;
```

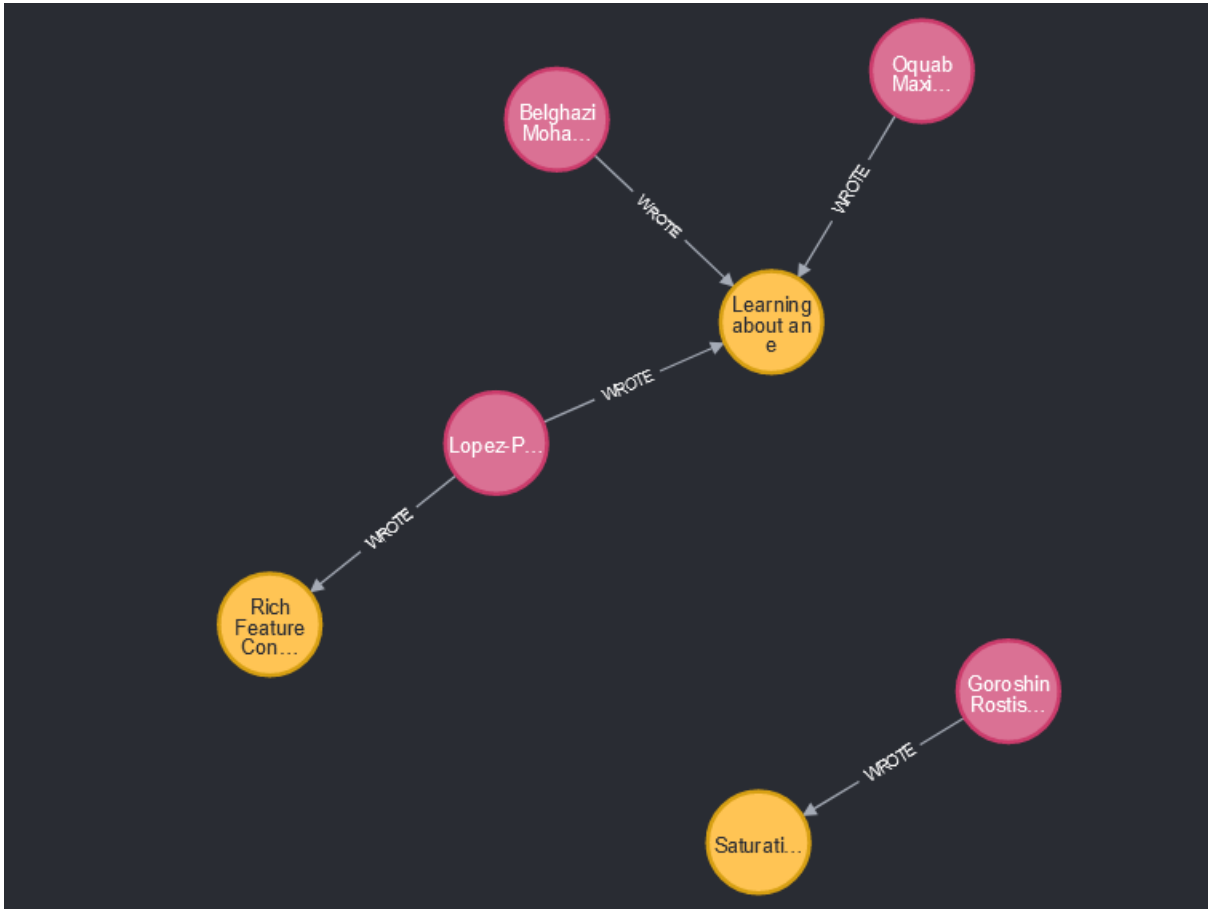


Figure 4.4: Query 4

Query 5. *Returns all the authors who wrote an article about Quantum physics in feb 2019.*

```
1 MATCH (a:Author)-[:WROTE]->(ar:Article)
2 WHERE date(ar.date).month = 2 AND date(ar.date).year = 2019
   AND ar.abstract CONTAINS 'Quantum'
3 MATCH p=(ar)<-[:WROTE]-(a)
4 RETURN p;
```

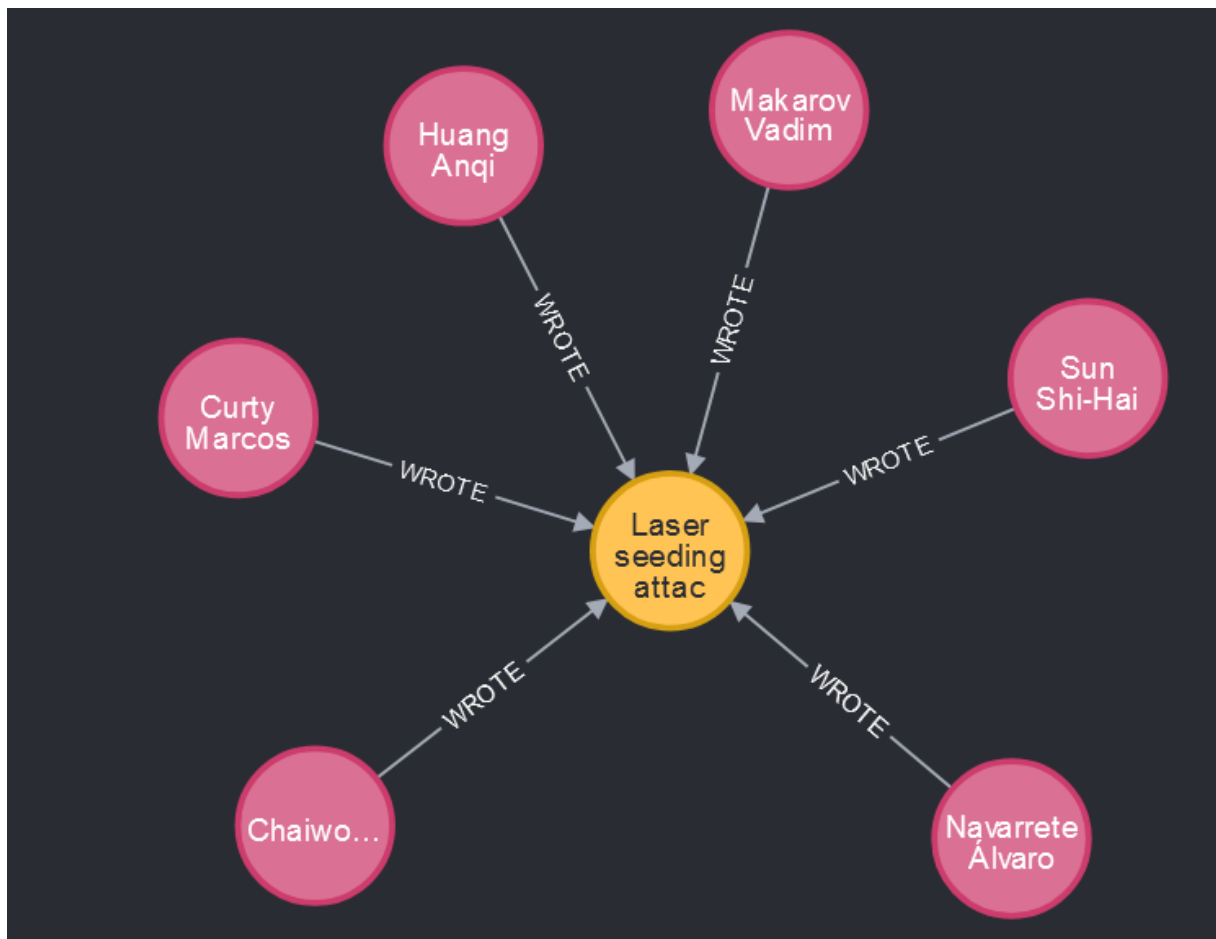


Figure 4.5: Query 5

Query 6. Gets all articles by authors who have written articles in at least three different years.

```

1 MATCH (ar:Article)<-[:WROTE]-(au:Author)-[:WROTE]->(other_ar:
   Article)
2 WITH au, count(DISTINCT date(other_ar.date).year) AS
   unique_years
3 WHERE unique_years >= 3
4 MATCH p=(ar)<-[:WROTE]-(au)
5 RETURN p
6 limit 250;

```



Figure 4.6: Partial outcome query 6

Query 7. *Returns articles written by authors who have also published articles with 'Machine Learning' in the abstract.*

```

1 MATCH (ar:Article)<-[:WROTE]-(au:Author)-[:WROTE]->(ml:
   Article)
2 WHERE ml.abstract CONTAINS 'Machine Learning'
3 MATCH p=(ar)<-[:WROTE]-(au)
4 RETURN p;

```

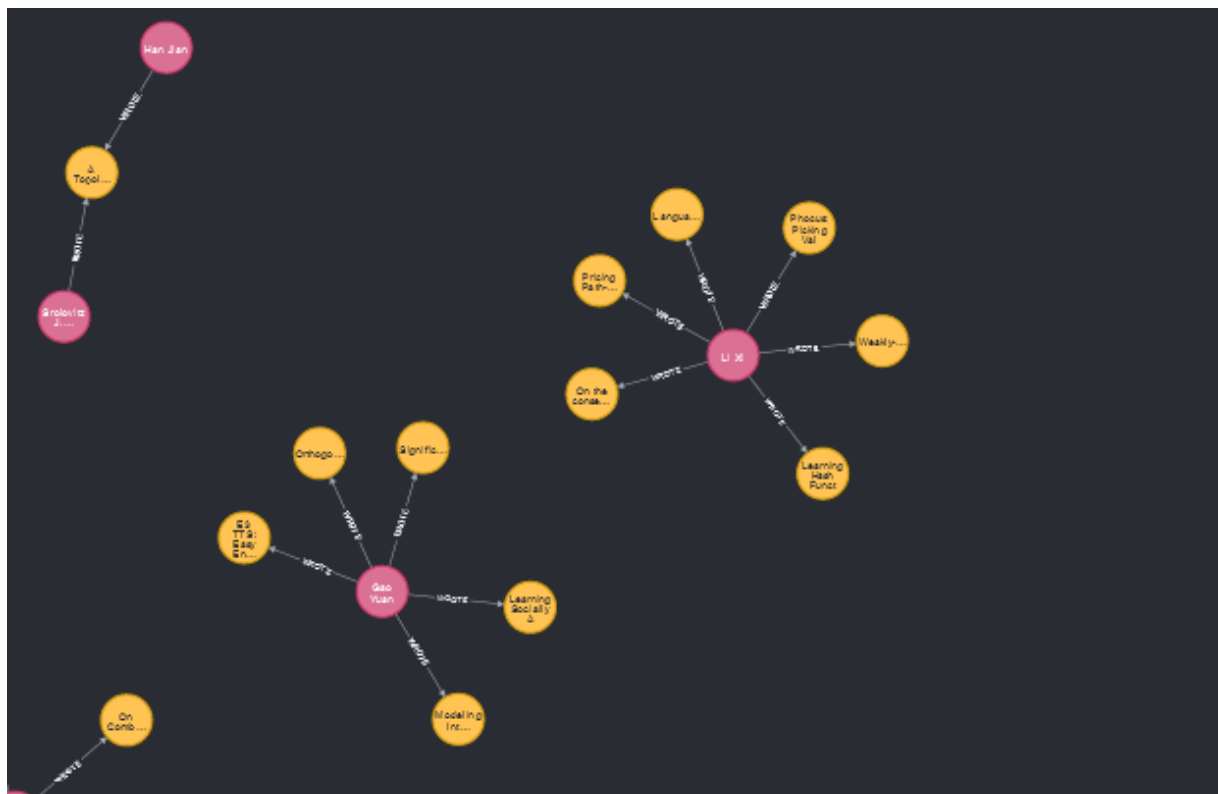


Figure 4.7: Partial outcome query 7

Query 8. *Retrieves articles with titles containing the word 'Graph' and lists their authors.*

```

1 MATCH (ar:Article)<-[:WROTE]-(au:Author)
2 WHERE ar.title CONTAINS 'Graph'
3 MATCH p=(ar)<-[:WROTE]-(au)
4 RETURN p
5 LIMIT 250;

```

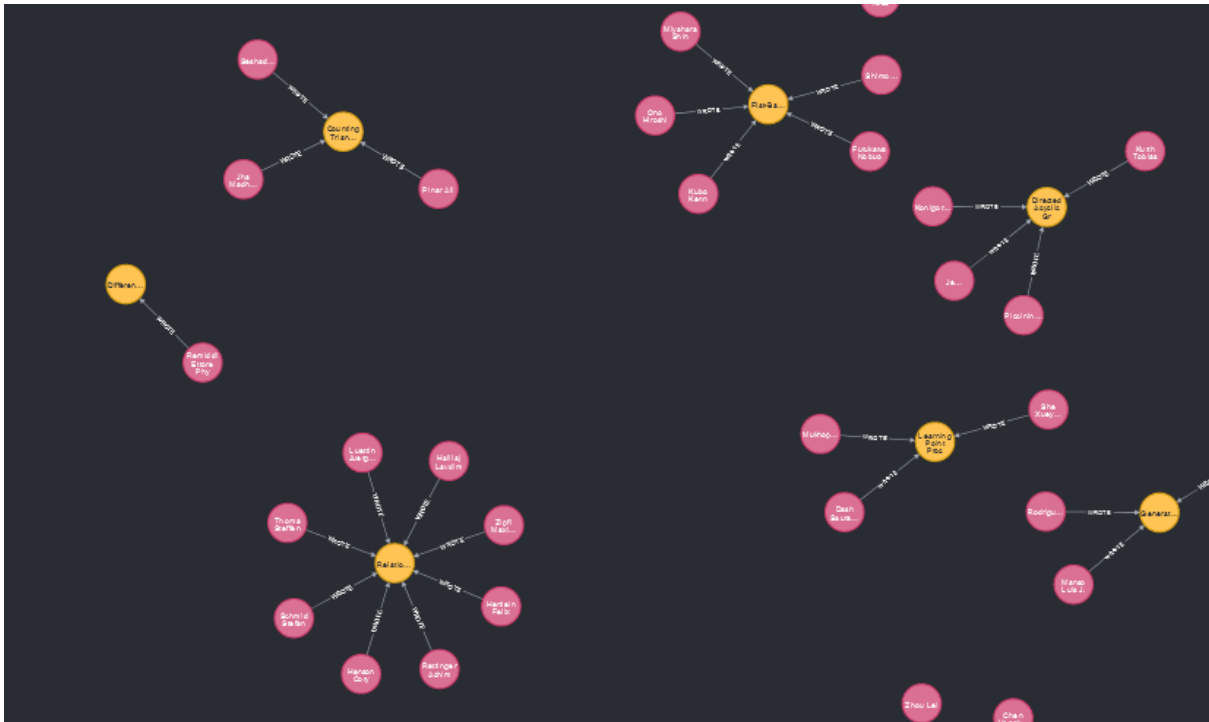


Figure 4.8: Partial outcome query 8

Query 9. *Retrivies articles written by authors who have not contributed to any articles in the year 2023 and returns the articles along with their author relationships.*

```

1 MATCH (ar:Article)<-[:WROTE]-(au:Author)
2 WITH ar, au
3 WHERE NOT (au)-[:WROTE]->(a:Article {date: '2023'})
4 MATCH p=(ar)<-[:WROTE]-(au:Author)
5 RETURN p
6 LIMIT 50;

```

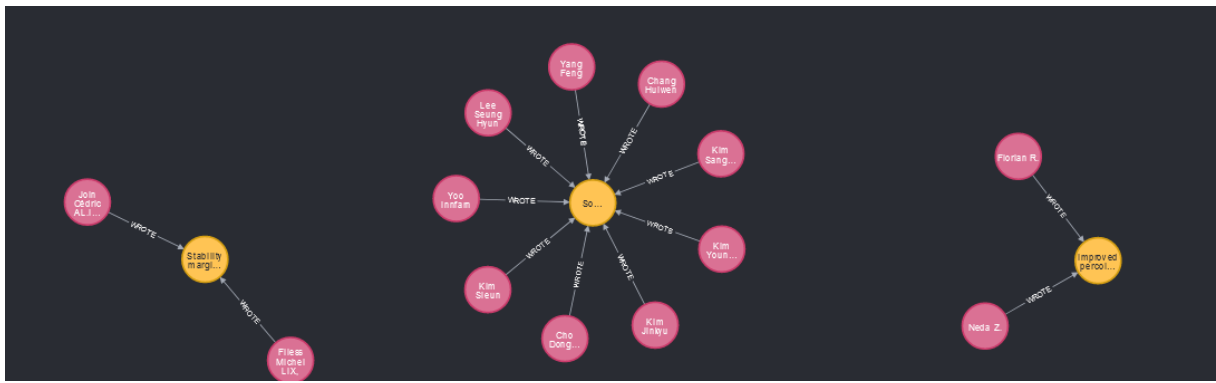


Figure 4.9: Query 9

Query 10. *Retrieves articles with the most recent publication date and their authors.*

```
1 MATCH (ar:Article)<-[:WROTE]-(au:Author)
2 WITH ar, date(ar.date) AS publication_date
3 ORDER BY publication_date DESC
4 LIMIT 1
5 MATCH p=(ar)<-[:WROTE]-(au)
6 RETURN p;
```

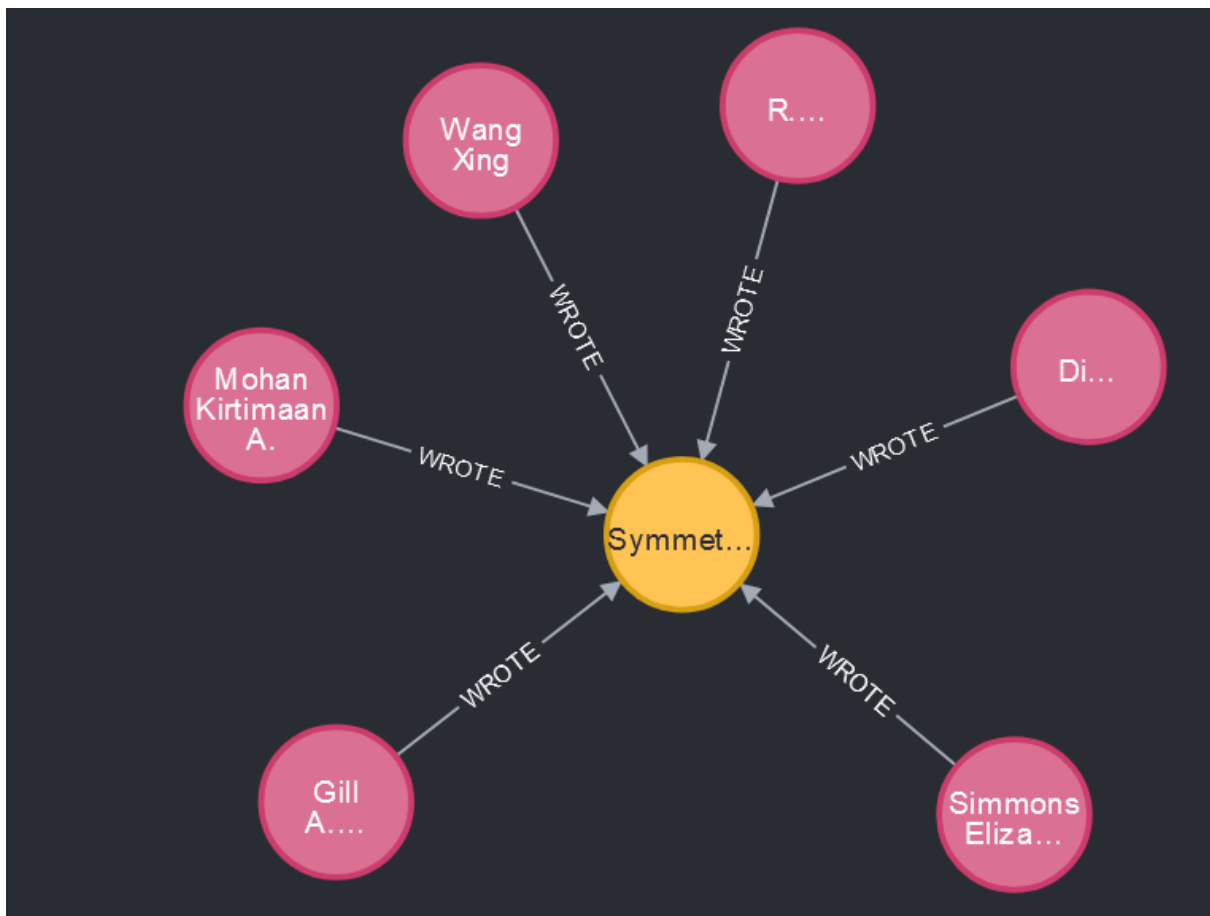


Figure 4.10: Query 10

4.2. Elasticsearch

Query 11. *Retrieves all the articles where the title contains the word "AI" but do not contains the word "ethic", prioritizing articles where the abstract contains the word "Neural".*

```

1 GET /arxiv/_search
2 {
3   "query": {
4     "bool": {
5       "must": [
6         { "match": { "title": "AI" } }
7       ],
8       "must_not": [
9         { "match": { "title": "Ethic" } }
10      ],
11      "should": [
12        { "match": { "abstract": "Neural" } }
13      ]
14    }
15  }
16 }

```

With this partial result:

```

1 "hits": [
2   {
3     "_index": "arxiv",
4     "_id": "13383",
5     "_score": 11.909199,
6     "_source": {
7       "submitter": "Maxime Gariel",
8       "authors": "Maxime Gariel, Brian Shimanuki, Rob
9         Timpe, Evan Wilson",
10      "title": "Framework for Certification of AI-Based
11        Systems",
12      "comments": "9 pages, 10 figures",
13      "ref": "",
14      "categories": "cs.LG cs.CV cs.SE",

```

```
13      "abstract": ""The current certification process
        for aerospace software is not adapted to "AI-
        based" algorithms such as deep neural networks.
        Unlike traditional aerospace software, the
        precise parameters optimized during neural
        network training are as important as (or more
        than) the code processing the network and they
        are not directly mathematically understandable.
        Despite their lack of explainability such
        algorithms are appealing because for some
        applications they can exhibit high performance
        unattainable with any traditional explicit line-
        by-line software methods. This paper proposes
        a framework and principles that could be used to
        establish certification methods for neural
        network models for which the current
        certification processes such as DO-178 cannot be
        applied. While it is not a magic recipe, it is
        a set of common sense steps that will allow the
        applicant and the regulator increase their
        confidence in the developed software, by
        demonstrating the capabilities to bring together
        , trace, and track the requirements, data,
        software, training process, and test results.""
14    }
15    },
16    ...
17  ]
```

Query 12. *Counts document by submitter.*

```
1 GET /arxiv/_search
2 {
3   {
4     "size": 0,
5     "aggs": {
6       "doc_by_submitter": {
7         "terms": {
8           "field": "submitter"
9         }
10      }
11    }
12  }
13 }
```

With this partial result:

```
1   "aggregations": {
2     "doc_by_submitter": {
3       "doc_count_error_upper_bound": 9,
4       "sum_other_doc_count": 29713,
5       "buckets": [
6         {
7           "key": "",
8           "doc_count": 186
9         },
10        {
11          "key": "EPTCS",
12          "doc_count": 46
13        },
14        ...
15      ]
16    }
17  }
```

Query 13. *Retrieves all the articles with 18 pages and 8 figures.*

```

1 GET /arxiv/_search
2 {
3   "query": {
4     "match": {
5       "comments": {
6         "query": "18 pages 8 Figures",
7         "operator": "and"
8       }
9     }
10  }
11 }
```

With this partial result:

```

1 "hits": [
2   {
3     "_index": "arxiv",
4     "_id": "1289",
5     "_score": 9.578911,
6     "_source": {
7       "submitter": "Evgeny Grines",
8       "authors": "Aleksei M. Arefev, Evgeny A. Grines,
9         Grigory V. Osipov",
10      "title": "Heteroclinic cycles and chaos in a system
11        of four identical phase oscillators with
12        global biharmonic coupling",
13      "comments": "18 pages, 8 figures",
14      "ref": "",
15      "categories": "nlin.CD math.DS",
16      "abstract": "abstract text"
17    }
18  },
19  {
20    "_index": "arxiv",
21    "_id": "6045",
22    "_score": 9.578911,
23    "_source": {
```

```
21         "submitter": "Fei Wu",
22         "authors": "Fei Wu, Chen Wu, Zhongzhou Ren",
23         "title": "Neutron stars including the effects of
                chaotic magnetic fields and the anomalous
                magnetic moments",
24         "comments": "18 pages, 8 figures",
25         "ref": "Chinese Physics C 41 (2017) 045102",
26         "categories": "astro-ph.HE nucl-th",
27         "abstract": "abstract text"
28     }
29 },
30 ...
31 ]
```


Query 14. *Counts all "hep-ph" articles for each submitter.*

```
1 GET /arxiv/_search
2 {
3   "size":0,
4   "query":{"
5     "match":{"
6       "categories":"hep-ph"
7     }
8   },
9   "aggs":{"
10     "hep_articles_per_submitter":{"
11       "terms":{"
12         "field":"submitter"
13       }
14     }
15   }
16 }
```

With this partial result:

```
1 "aggregations":{"
2   "hep_articles_per_submitter":{"
3     "doc_count_error_upper_bound":3,
4     "sum_other_doc_count":948,
5     "buckets":[
6       {
7         "key":"",
8         "doc_count":43
9       },
10      {
11        "key":"Adrian Signer",
12        "doc_count":4
13      },
14      ...
15    ]
16  }
17 }
```

Query 15. *Counts for each categories how many articles with the word "Quantum" in the title.*

```

1 GET /arxiv/_search
2 {
3   "size":0,
4   "query":{"
5     "match":{"
6       "title": "quantum"
7     }
8   },
9   "aggs":{"
10    "docs_per_category":{"
11      "terms":{"
12        "field": "categories"
13      }
14    }
15  }
16 }
```

With this partial result:

```

1 "aggregations":{"
2   "docs_per_category":{"
3     "doc_count_error_upper_bound":9,
4     "sum_other_doc_count":788,
5     "buckets":[
6       {
7         "key": "quant-ph",
8         "doc_count":388
9       },
10      {
11        "key": "cond-mat.mes-hall",
12        "doc_count":112
13      },
14      {
15        "key": "hep-th",
16        "doc_count":46
17      },
```

18
19
20
21

```
...  
]  
}  
}
```

Query 16. *Finds all articles with with "Neural Network" in the title, favouring articles containing "deep learning" in the abstract and favouring articles containing "cnn" in the abstract with double importance.*

```
1 GET /arxiv/_search
2 {
3   "query":{
4     "bool":{
5       "must":[
6         {
7           "match":{
8             "title":"neural network"
9           }
10        }
11      ],
12      "should":[
13        {
14          "match":{
15            "abstract":"deep learning"
16          }
17        },
18        {
19          "match":{
20            "abstract":{
21              "query":"CNN",
22              "boost":2.0
23            }
24          }
25        }
26      ]
27    }
28  }
29 }
```

With this partial result:

```
1 "hits":[
2   {
3     "_index":"arxiv",
```

```

4      "_id": "2833",
5      "_score": 31.632092,
6      "_source": {
7          "submitter": "Xingquan Liu",
8          "authors": "X. Zhang, Y. Huang, W. Lin, X. Liu,
9                      H. Zheng, R. Wada, A. Bonasera, Z.\n Chen
                      , L. Chen, J. Han, R. Han, M. Huang, Q. Hu,
                      Q. Leng, C. W. Ma, G. Qu, P.\n Ren, G.
                      Tian, Z. Xu, Z. Yang, and L. Zhang",
10         "title": "Determining impact parameters of
11                  heavy-ion collisions at low-intermediate
12                  incident energies using deep learning with
13                  convolutional neural network",
14         "comments": "21 pages, 16 figures, 1 table",
15         "ref": "",
16         "categories": "nucl-th",
17         "abstract": "abstract text."
18     }
19 },
20 {
21     "_index": "arxiv",
22     "_id": "15136",
23     "_score": 30.780548,
24     "_source": {
25         "submitter": "Tingle Li",
26         "authors": "Tingle Li, Qingjian Lin, Yuanyuan
27                     Bao, Ming Li",
28         "title": "Atss-Net: Target Speaker Separation
29                  via Attention-based Neural Network",
30         "comments": "Submitted to Interspeech 2020",
31         "ref": "",
32         "categories": "eess.AS cs.SD",
33         "abstract": "abstract text."
34     }
35 },
36 ...
37 ]

```

Query 17. *Counts articles about AI but not about ethic by category.*

```
1 GET /arxiv/_search
2 {
3   "size":0,
4   "query":{
5     "bool":{
6       "must":[
7         {
8           "match":{
9             "title":"AI"
10          }
11        }
12      ],
13      "must_not":[
14        {
15          "match":{
16            "title":"ethic"
17          }
18        },
19        {
20          "match":{
21            "abstract":"ethic"
22          }
23        }
24      ]
25    }
26  },
27  "aggs":{
28    "docs_per_category":{
29      "terms":{
30        "field":"categories"
31      }
32    }
33  }
34 }
```

With this partial result:

```
1 "aggregations":{
2   "docs_per_category":{
3     "doc_count_error_upper_bound":0,
4     "sum_other_doc_count":38,
5     "buckets":[
6       {
7         "key":"cs.AI",
8         "doc_count":6
9       },
10      {
11        "key":"cs.CL",
12        "doc_count":3
13      },
14      {
15        "key":"cs.CV",
16        "doc_count":3
17      },
18      ...
19    ]
20  }
21 }
```

Query 18. *Counts for each category how many articles have been submitted by each submitter.*

```

1 GET /arxiv/_search
2 {
3   "size":0,
4   "aggs":{
5     "by_categories":{
6       "terms":{
7         "field":"categories"
8       },
9       "aggs":{
10        "by_submitter":{
11          "terms":{
12            "field":"submitter"
13          }
14        }
15      }
16    }
17  }
18 }
```

With this partial result:

```

1 "aggregations":{
2   "by_categories":{
3     "doc_count_error_upper_bound":204,
4     "sum_other_doc_count":23623,
5     "buckets":[
6       {
7         "key":"astro-ph",
8         "doc_count":1064,
9         "by_submitter":{
10          "doc_count_error_upper_bound":3,
11          "sum_other_doc_count":1034,
12          "buckets":[
13            {
14              "key":"",
15              "doc_count":11
```



```
16         },
17         {
18             "key": "Charles Dermer",
19             "doc_count": 3
20         },
21         {
22             "key": "Alain Jorissen",
23             "doc_count": 2
24         },
25         ...
26     ]
27 }
28 }
29 ]
30 }
31 }
```

Query 19. *Counts for each submitter, for each categories the articles about AI favouring articles about transformers.*

```
1 GET /arxiv/_search
2 {
3   "size":0,
4   "query":{"
5     "bool":{"
6       "must":[
7         {
8           "match":{"
9             "title":"AI"
10          }
11        },
12        {
13          "match":{"
14            "abstract":"AI"
15          }
16        }
17      ],
18      "should":[
19        {
20          "match":{"
21            "abstract":"Transformers"
22          }
23        }
24      ]
25    }
26  },
27  "aggs":{"
28    "by_submitter":{"
29      "terms":{"
30        "field":"submitter"
31      },
32      "aggs":{"
33        "by_category":{"
34          "terms":{"
35            "field":"submitter"
```

```

36     }
37   }
38 }
39 }
40 }
41 }

```

With this partial result:

```

1  "aggregations":{
2    "by_submitter":{
3      "doc_count_error_upper_bound":0,
4      "sum_other_doc_count":43,
5      "buckets":[
6        "0":{
7          "key":"Abeba Birhane",
8          "doc_count":1,
9          "by_category":{
10             "doc_count_error_upper_bound":0,
11             "sum_other_doc_count":0,
12             "buckets":[
13               "0":{
14                 "key":"Abeba Birhane",
15                 "doc_count":1
16               }
17             ]
18           }
19         },
20         "1":{
21           "key":"Admela Jukan",
22           "doc_count":1,
23           "by_category":{
24             "doc_count_error_upper_bound":0,
25             "sum_other_doc_count":0,
26             "buckets":[
27               "0":{
28                 "key":"Admela Jukan",
29                 "doc_count":1

```

```

30         }
31     ]
32 }
33 },
34 "2":{
35     "key":"Aleksander Slominski",
36     "doc_count":1,
37     "by_category":{
38         "doc_count_error_upper_bound":0,
39         "sum_other_doc_count":0,
40         "buckets":[
41             "0":{
42                 "key":"Aleksander Slominski",
43                 "doc_count":1
44             }
45         ]
46     }
47 },
48 "3":{
49     "key":"Alexander Preuhs",
50     "doc_count":1,
51     "by_category":{
52         "doc_count_error_upper_bound":0,
53         "sum_other_doc_count":0,
54         "buckets":[
55             "0":{
56                 "key":"Alexander Preuhs",
57                 "doc_count":1
58             }
59         ]
60     }
61 },
62 "4":{
63     "key":"Ameet Deshpande",
64     "doc_count":1,
65     "by_category":{
66         "doc_count_error_upper_bound":0,

```

```
67         "sum_other_doc_count":0,  
68         "buckets":[  
69             "0":{  
70                 "key":"Ameet Deshpande",  
71                 "doc_count":1  
72             }  
73         ]  
74     },  
75     },  
76     ...  
77 ]  
78 }  
79 }
```

Query 20. *Retrieves Articles with Attention in Title, Excluding CNN in Abstract, and Emphasizing Bidirectional LSTM and Transformer.*

```
1 GET /arxiv/_search
2 {
3   "query":{
4     "bool":{
5       "must_not":[
6         {
7           "match":{
8             "abstract":"CNN"
9           }
10        }
11      ],
12      "must":[
13        {
14          "match":{
15            "title":"Attention"
16          }
17        }
18      ],
19      "should":[
20        {
21          "match":{
22            "abstract": {
23              "query": "Bidirectional LSTM",
24              "operator": "and",
25              "boost": 4.0
26            }
27          }
28        },
29        {
30          "match":{
31            "abstract":{
32              "query": "Transformer",
33              "boost": 2.0
34            }
35          }
```

```

36     }
37   ]
38 }
39 }
40 }

```

With this partial result:

```

1  "hits": [
2    {
3      "_index": "arxiv",
4      "_id": "7600",
5      "_score": 23.900904,
6      "_source": {
7        "submitter": "Jaeyoung Kim",
8        "authors": "Jaeyoung Kim, Mostafa El-Khamy,
9                  Jungwon Lee",
10       "title": "T-GSA: Transformer with Gaussian-
11               weighted self-attention for speech
12               enhancement",
13       "comments": "5 pages, Submitted to ICASSP 2020",
14       "ref": "",
15       "categories": "eess.AS cs.SD",
16       "abstract": "abstract text"
17     },
18     {
19       "_index": "arxiv",
20       "_id": "1253",
21       "_score": 21.699568,
22       "_source": {
23         "submitter": "Xiangyu Chen",
24         "authors": "Xiangyu Chen, Xintao Wang, Wenlong
25                   Zhang, Xiangtao Kong, Yu Qiao, \n Jiantao
26                   Zhou, and Chao Dong",
27         "title": "HAT: Hybrid Attention Transformer for
28                 Image Restoration",

```

```
24         "comments": "Extended version of HAT",
25         "ref": "",
26         "categories": "cs.CV",
27         "abstract": "abstract text"
28     }
29 },
30 ...
31 ]
```


5 | Extra

To enhance our project, we considered creating a web app that would simulate a user-friendly graphical interface for performing the desired queries.

5.0.1. DISCLAIMER

Due to temporal limit we had, we didn't manage error handling perfectly: so you have to write syntactically correct queries.

5.1. WebApp structure

As shown in the photo below, our application is built on a client-server architecture that leverages the HTTP protocol. The user initiates an HTTP request to the server (where a NodeJS instance is running), and the server responds with the correct HTML page along with the corresponding scripts and style sheets.

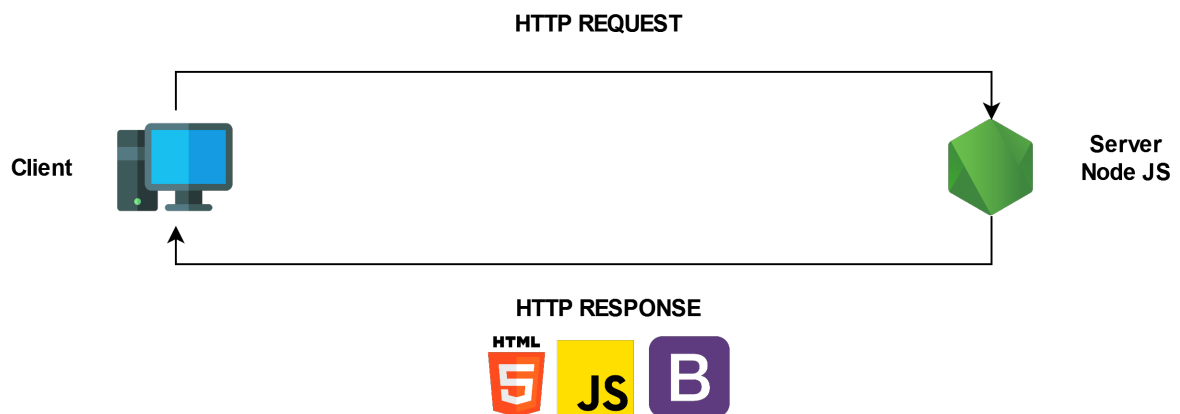


Figure 5.1: WebApp structure

Later, as we will see in the explanatory screenshots, the user selects the desired technology (Neo4j, ElasticSearch) and performs the query.

Subsequently, our server communicates with Neo4j or ElasticSearch, depending on the user's choice, to execute the requested query and return the result in the form of a JSON file, which will then be parsed on the client side.

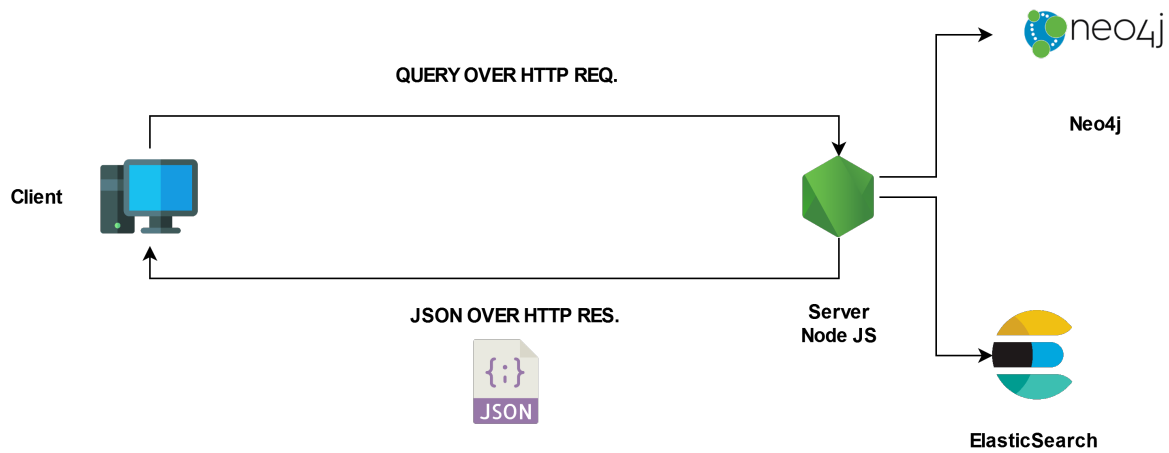


Figure 5.2: WebApp user interaction

5.2. WebApp screens

5.2.1. Neo4j screens

Initially, the website appears as follows: at the center, there is a canvas where the graph representing the result of the performed query will be drawn. In the toolbar, it's possible to select the desired technology, choose a query from the provided options, clear the input, and initiate the query. Finally, if the user wishes, they can write their own custom query in the lower bar.



Figure 5.3: Main screen for Neo4j

Clicking the yellow star reveals a screen with 10 queries that can be clicked to automatically insert them into the text field below.

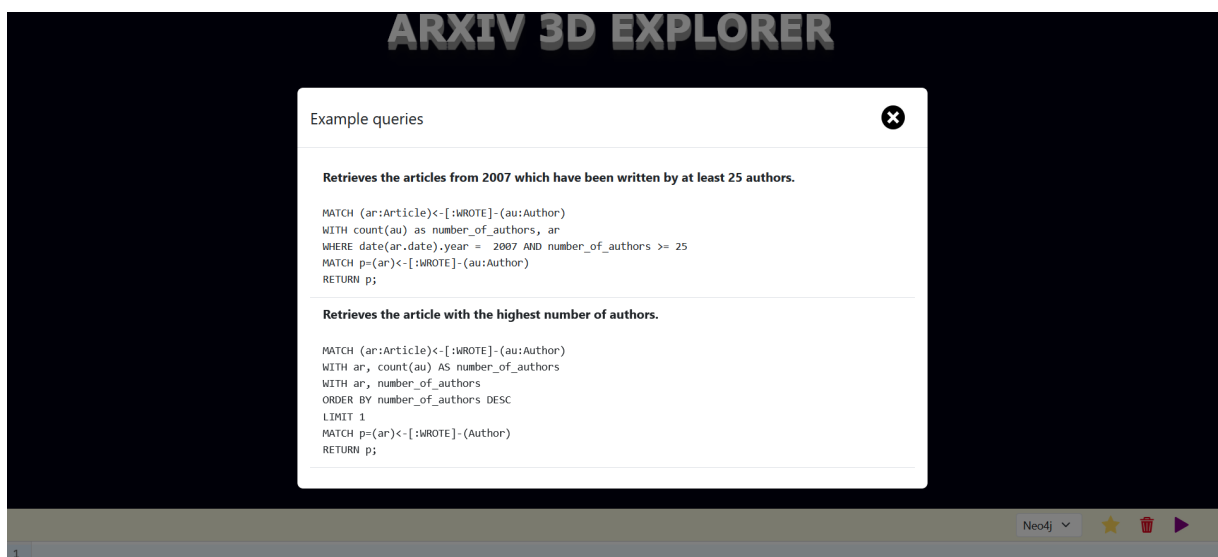


Figure 5.4: Queries screen for Neo4j

When the client receives the JSON file containing the query result, parsing is performed, and the corresponding graph is displayed. It's important to note that nodes represent authors and articles, while edges represent relationships between authors and articles.



Figure 5.5: Result screen for Neo4j

5.2.2. ElasticSearch screens

The main screen is almost identical to the one before, with the only difference being the presence of a container in the center where the query result will be displayed.



Figure 5.6: Main screen for ElasticSearch

Clicking the yellow star reveals a screen with 10 queries that can be clicked to automatically insert them into the text field below.

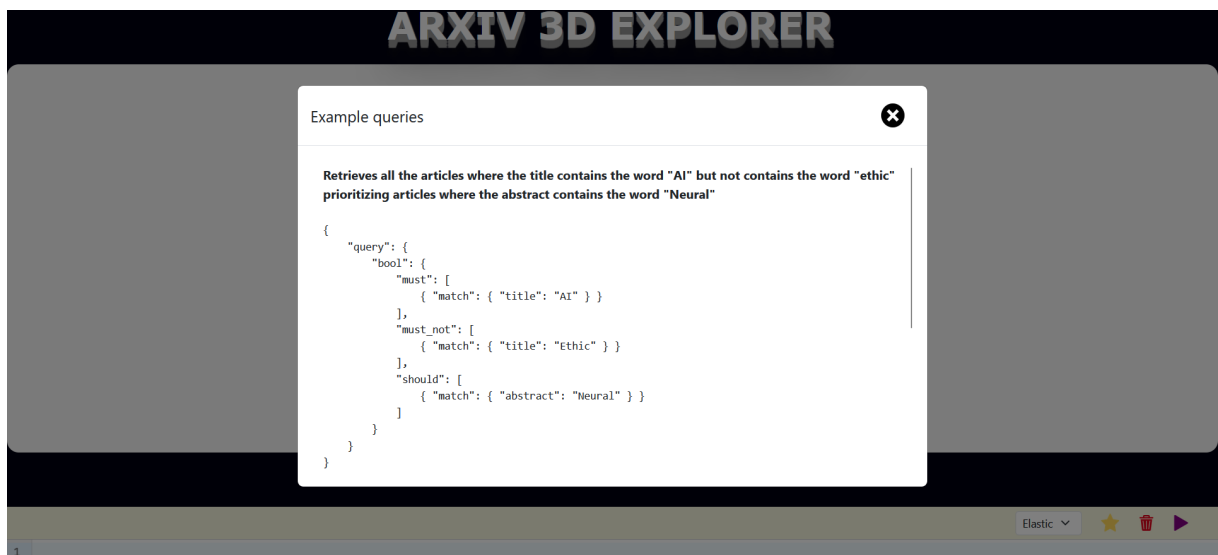


Figure 5.7: Queries screen for ElasticSearch

When the client receives the JSON file containing the query result, parsing is performed, and the corresponding parsed text is shown.

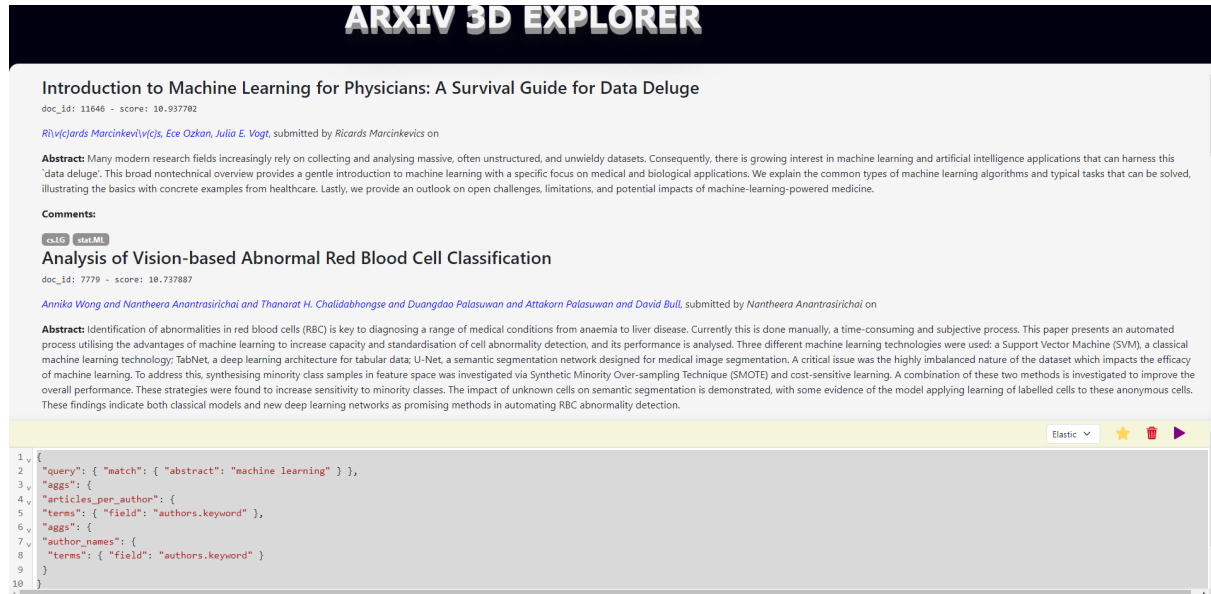


Figure 5.8: Result screen for ElasticSearch

List of Figures

3.1	Neo4j DB structure.	6
4.1	Partial outcome query 1	9
4.2	Partial outcome query 2	10
4.3	Query 3	11
4.4	Query 4	12
4.5	Query 5	13
4.6	Partial outcome query 6	14
4.7	Partial outcome query 7	15
4.8	Partial outcome query 8	16
4.9	Query 9	17
4.10	Query 10	18
5.1	WebApp structure	41
5.2	WebApp user interaction	42
5.3	Main screen for Neo4j	43
5.4	Queries screen for Neo4j	43
5.5	Result screen for Neo4j	44
5.6	Main screen for ElasticSearch	45
5.7	Queries screen for ElasticSearch	45
5.8	Result screen for ElasticSearch	46