



POLITECNICO
MILANO 1863

System and Methods for Big and Unstructured Data – Group 20

Alessio Buda
Leonardo Cesani
Fausto Lasca
Gabriele Munafò
Matteo Paraboschi



POLITECNICO
MILANO 1863

Project presentation & assumptions

Project presentation



Relational
tables



Graphs



Documents



Keys



Columns



Project presentation



Graph model



Documental model



Resilient Distributed
Dataset (RDD)

Project presentation



4.84 millions of publications



Different types of publication



Relevant information on each publication

Assumptions



Conference edition

vs

Journal edition

Organizations' affiliations

Unique attributes

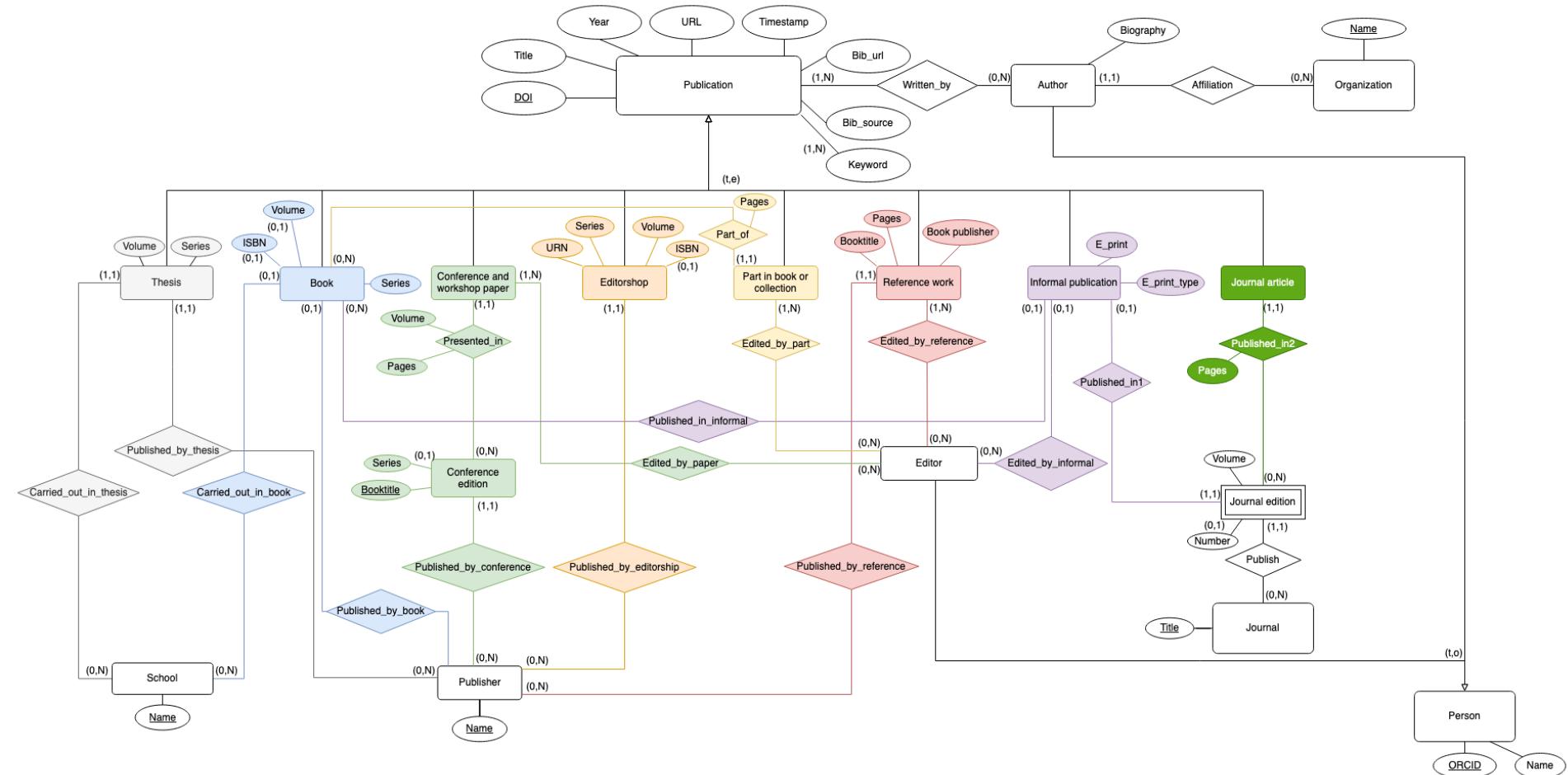


POLITECNICO
MILANO 1863

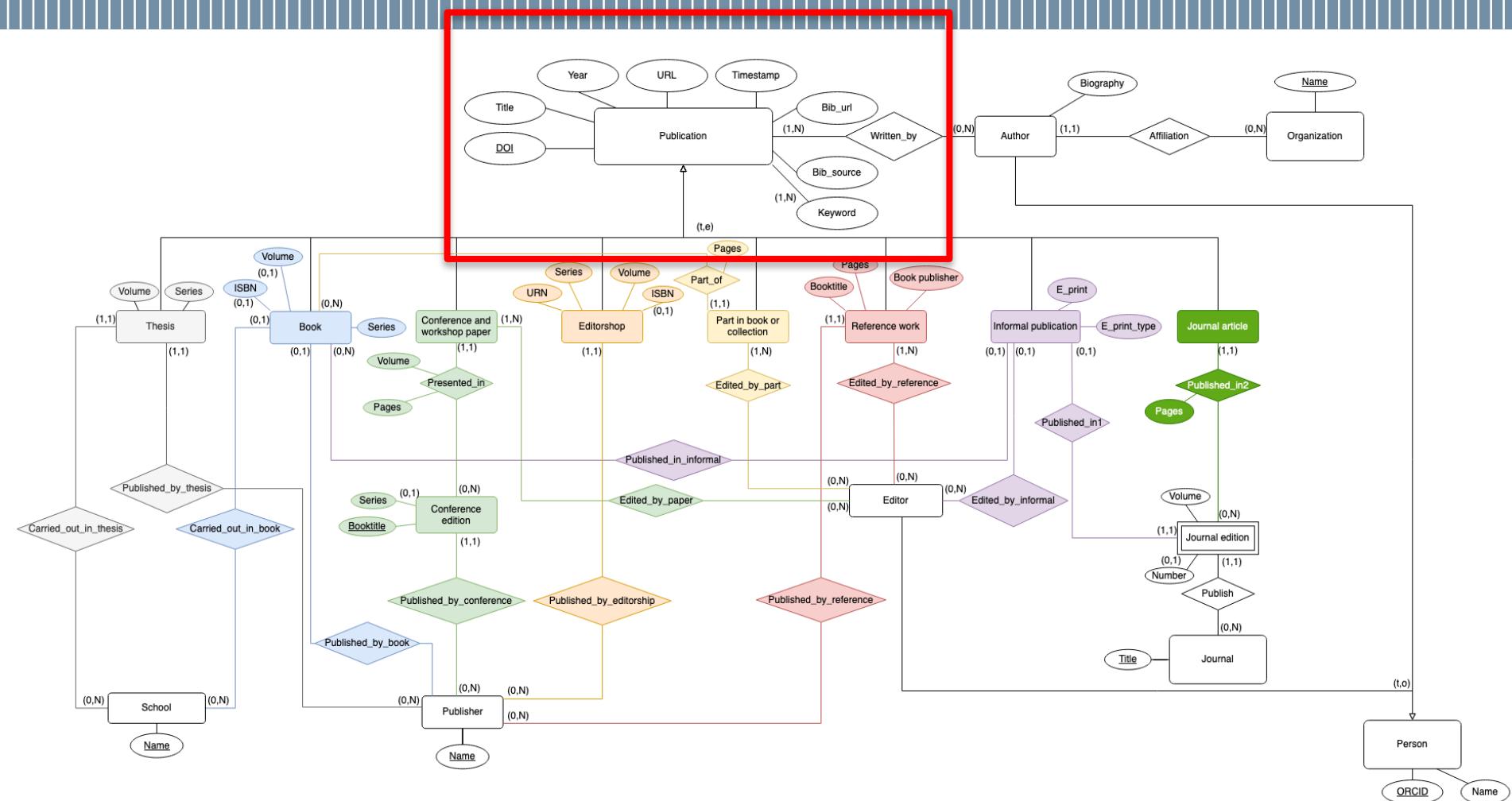
ER model

Conceptual data model

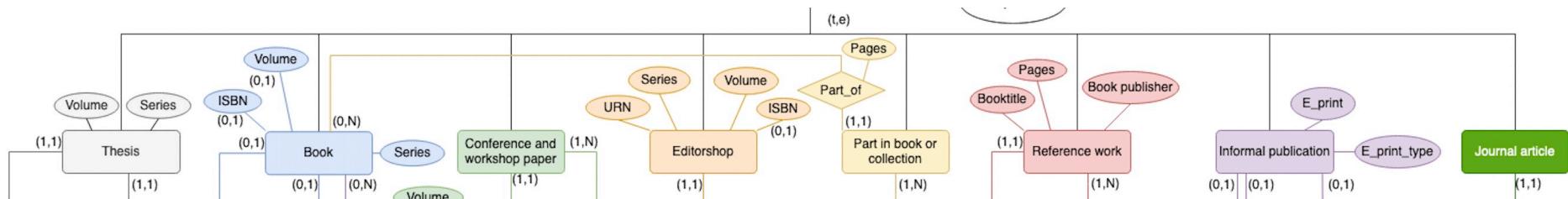
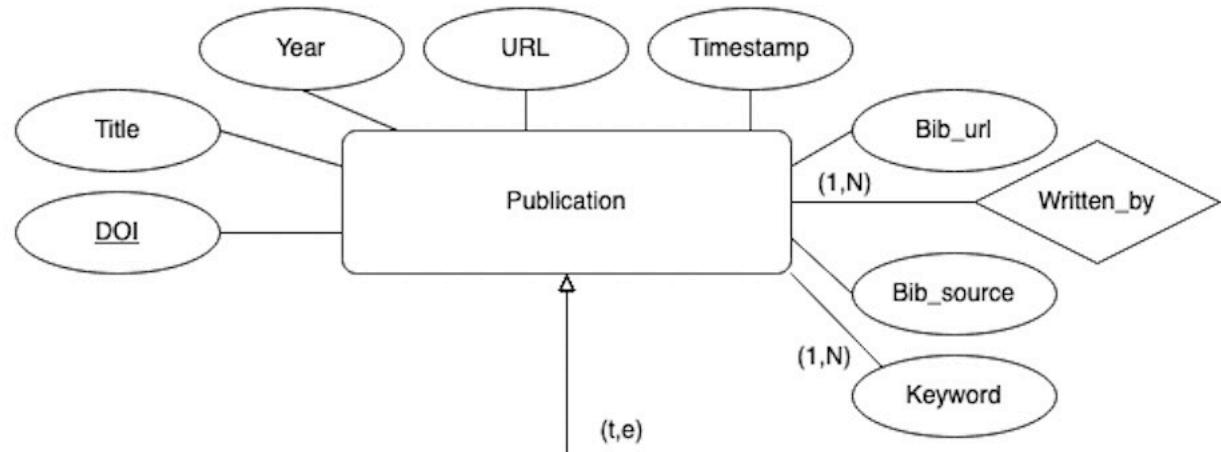
ER Model



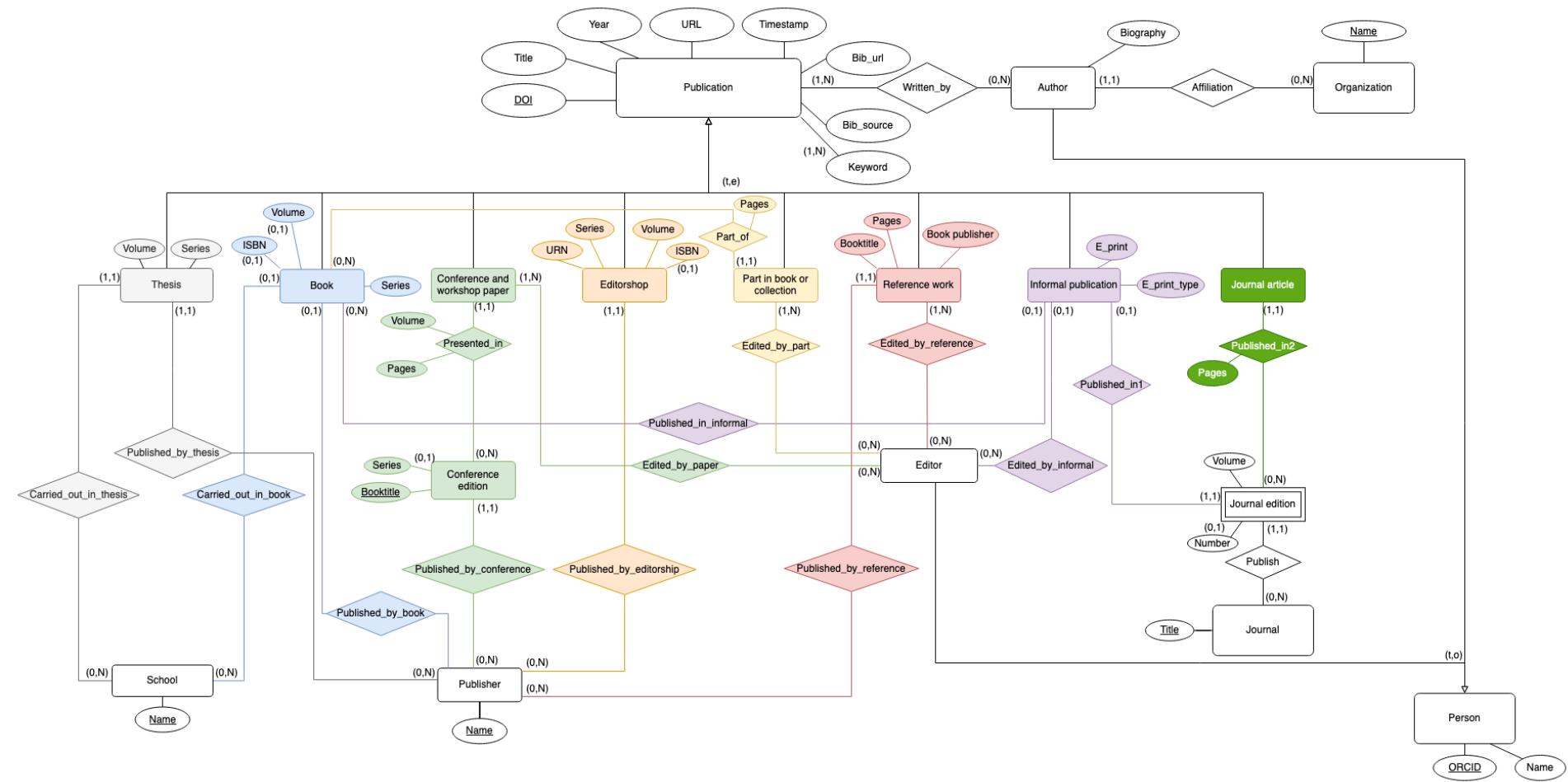
ER Model



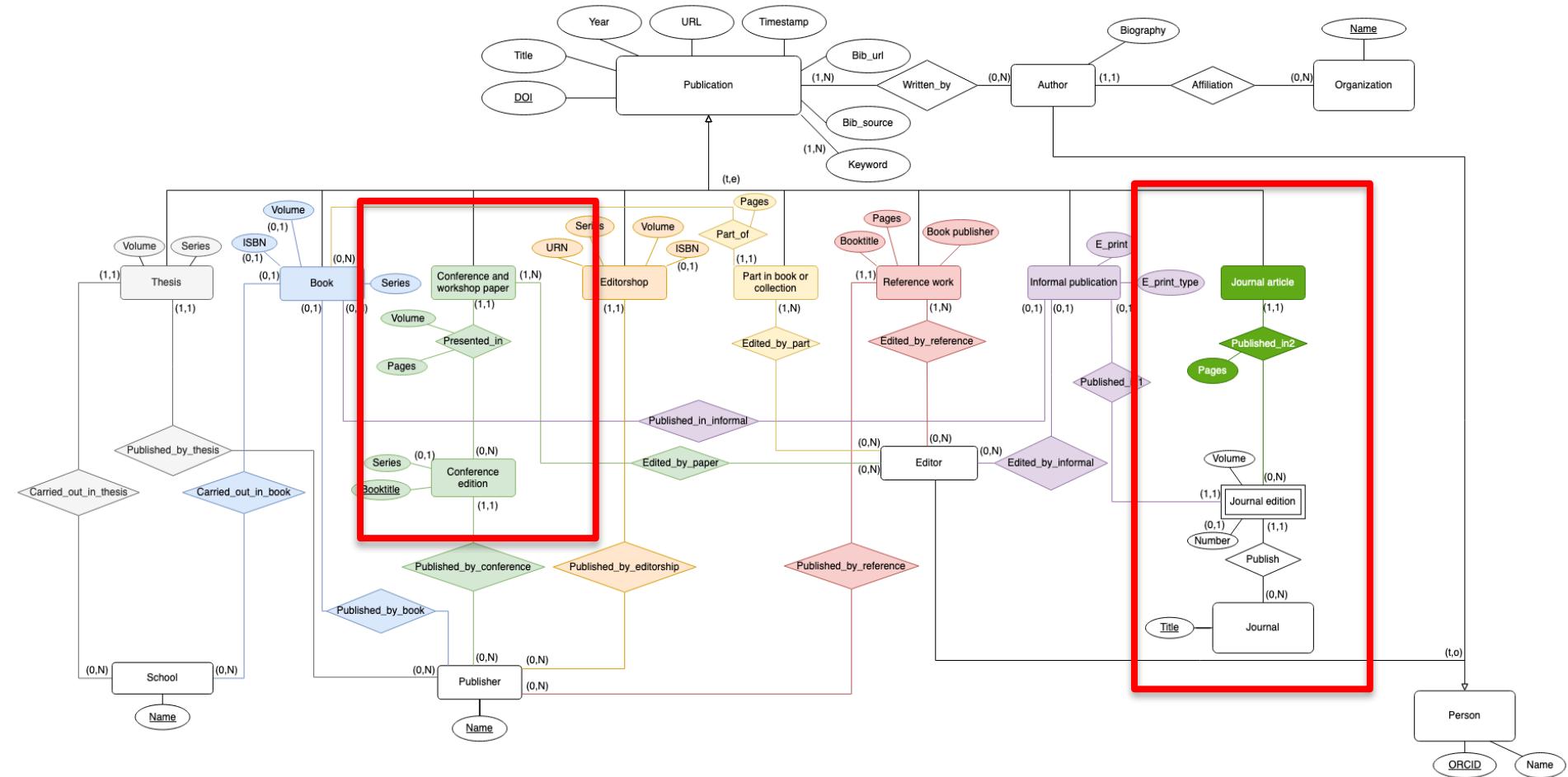
ER Model



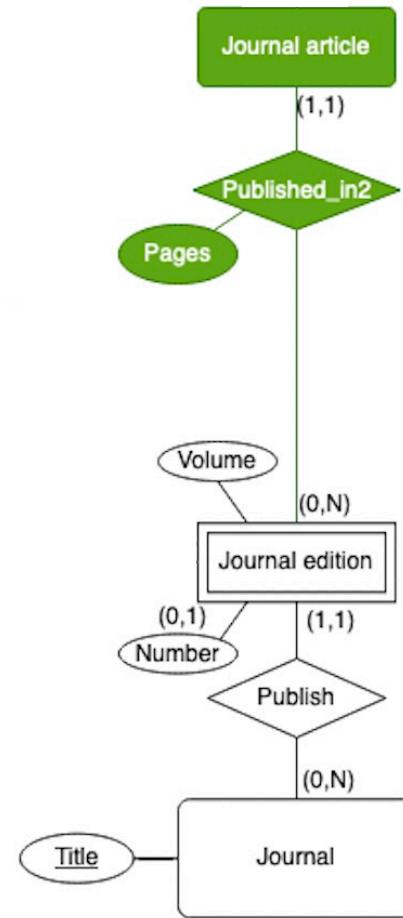
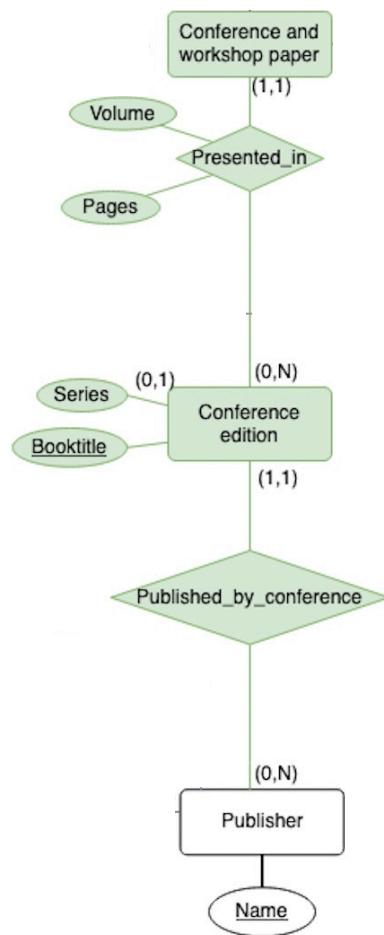
ER Model



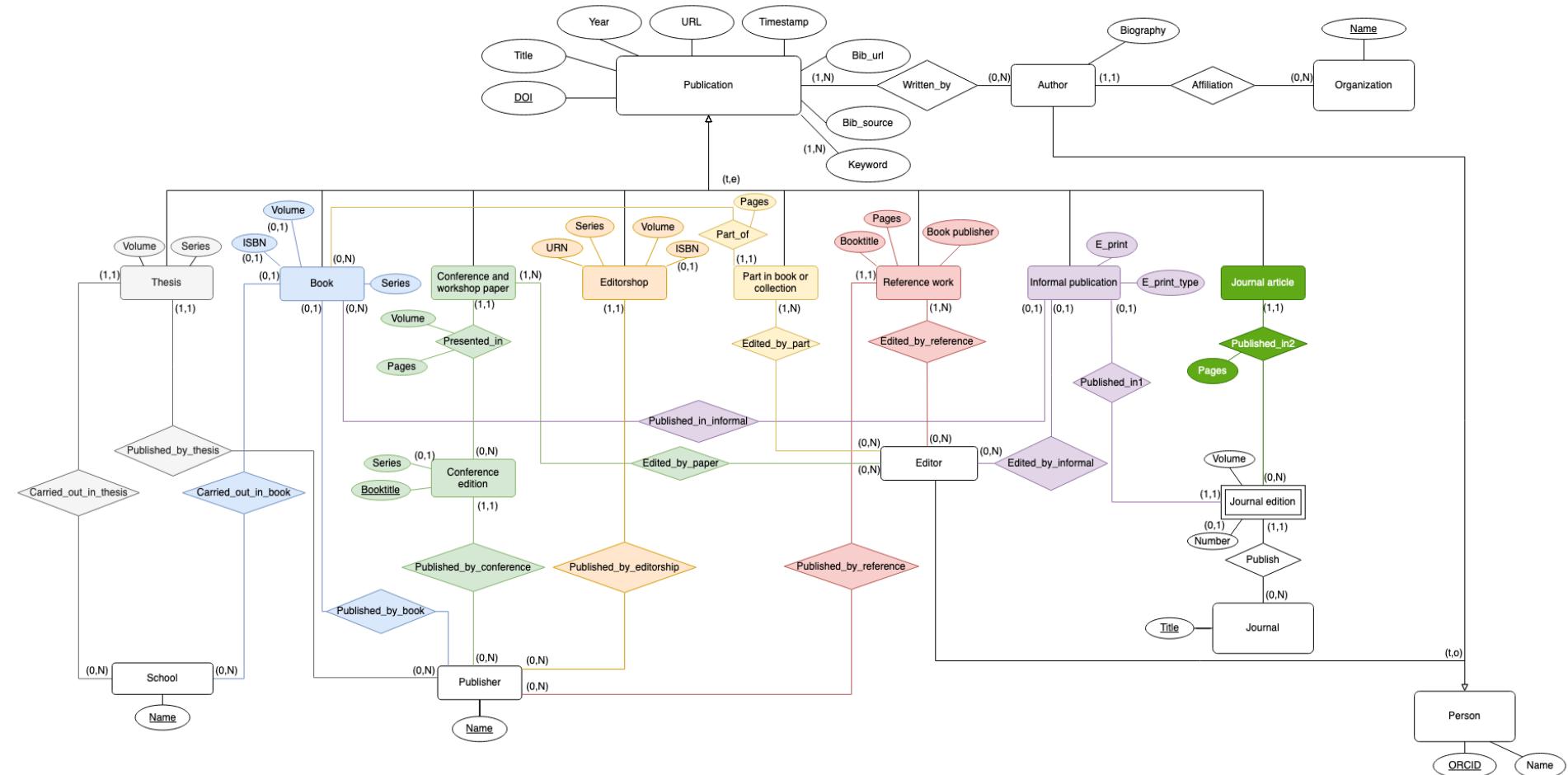
ER Model



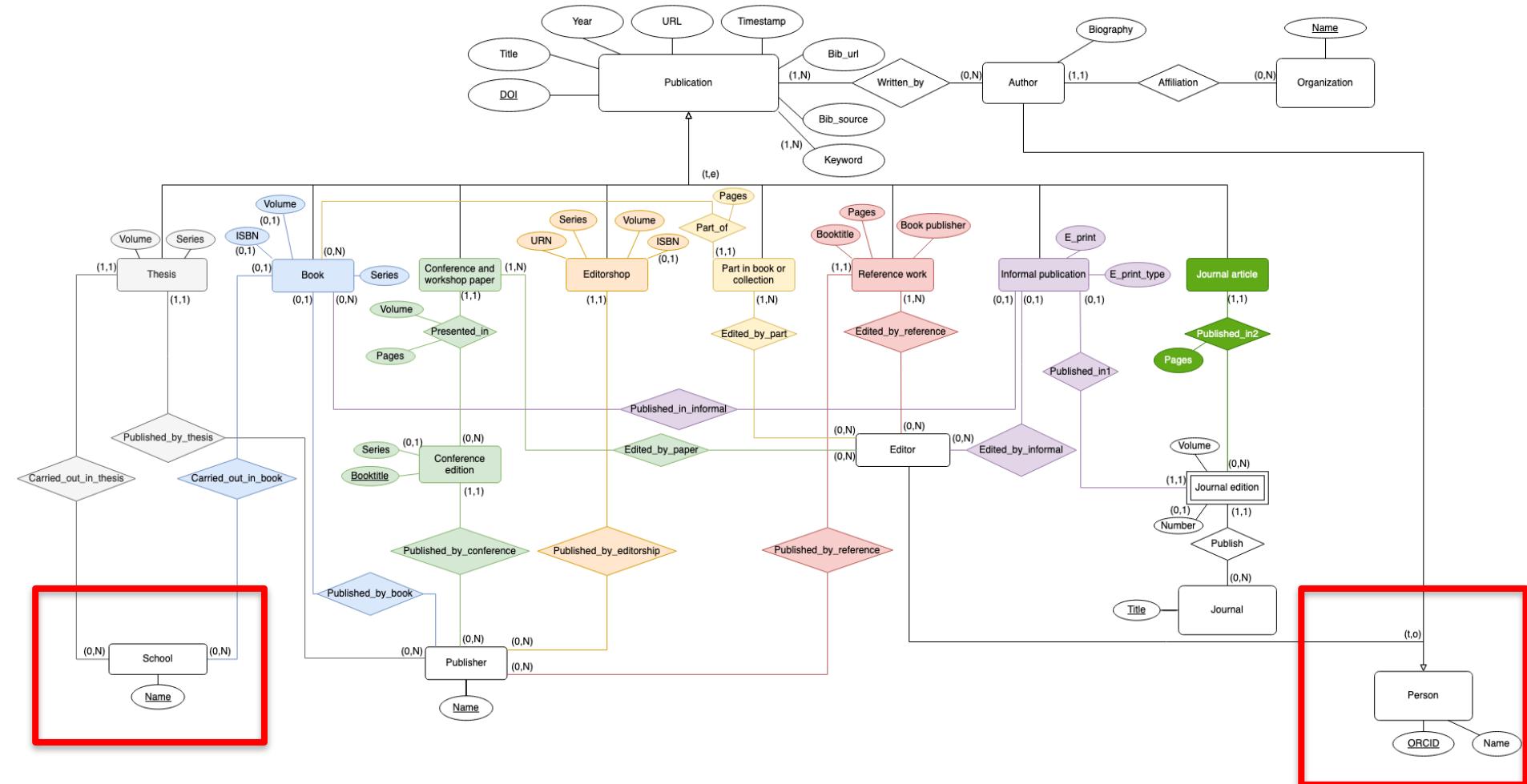
ER Model



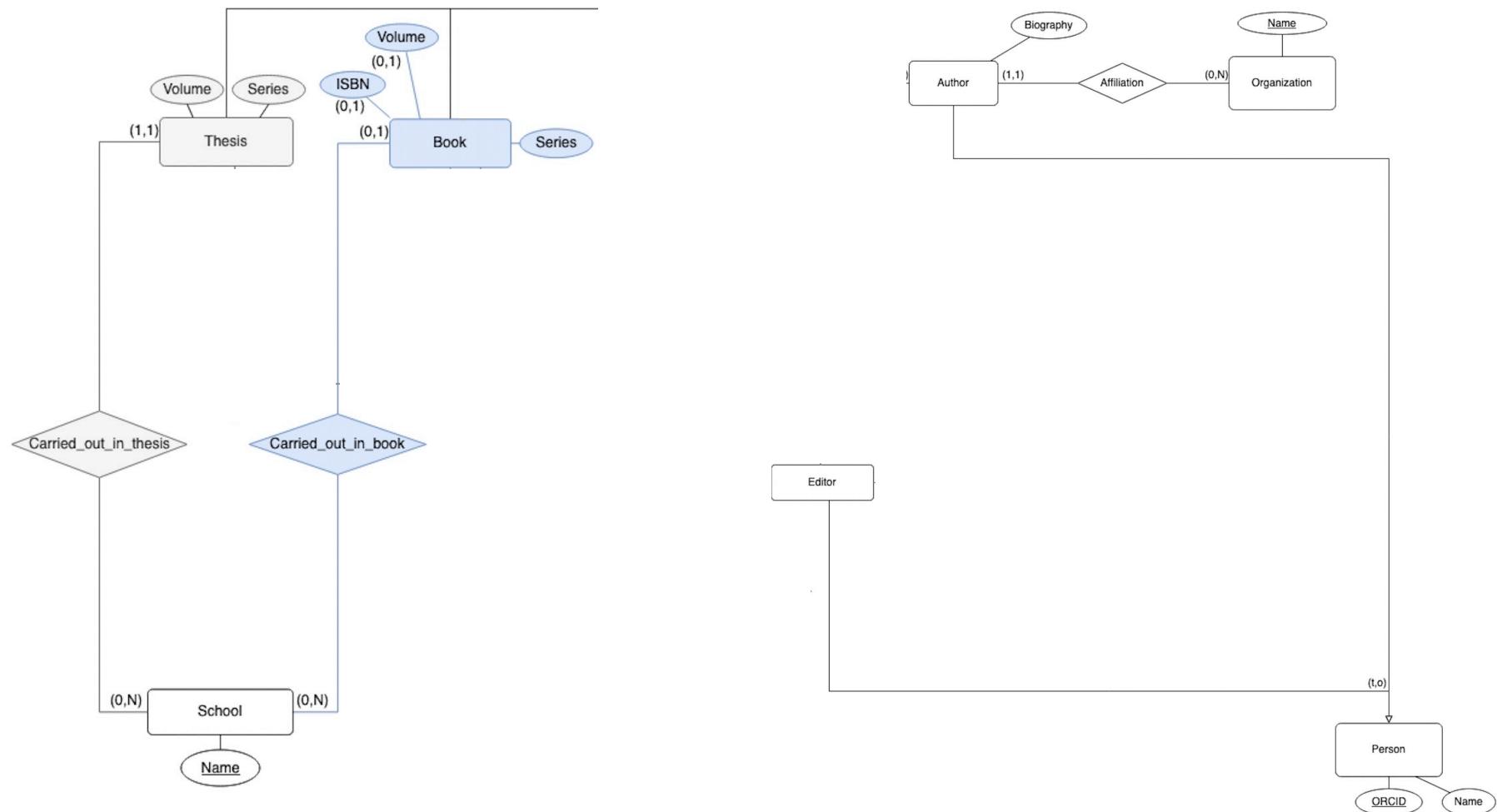
ER Model



ER Model



ER Model

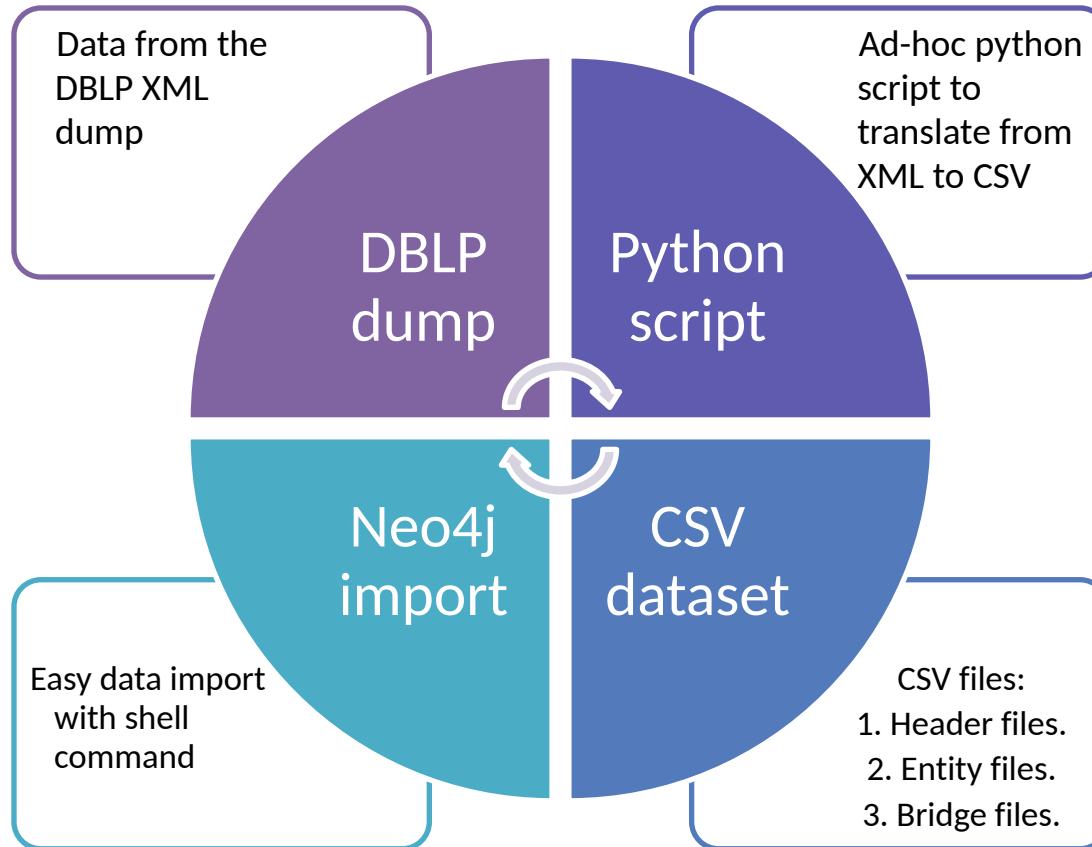




POLITECNICO
MILANO 1863

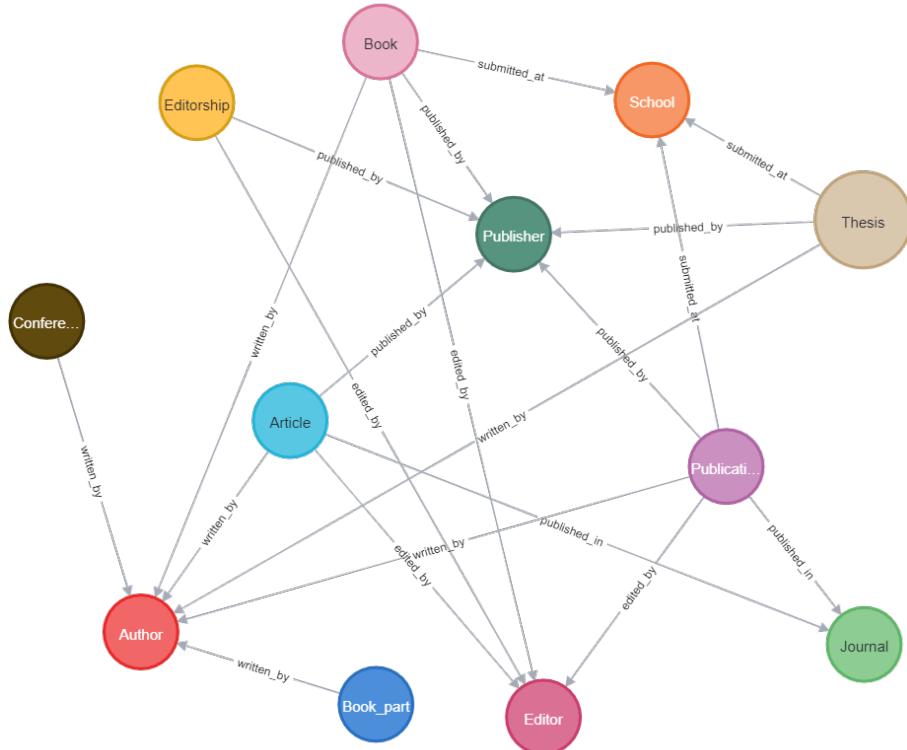
Neo4j

Neo4J : dataset generation and import

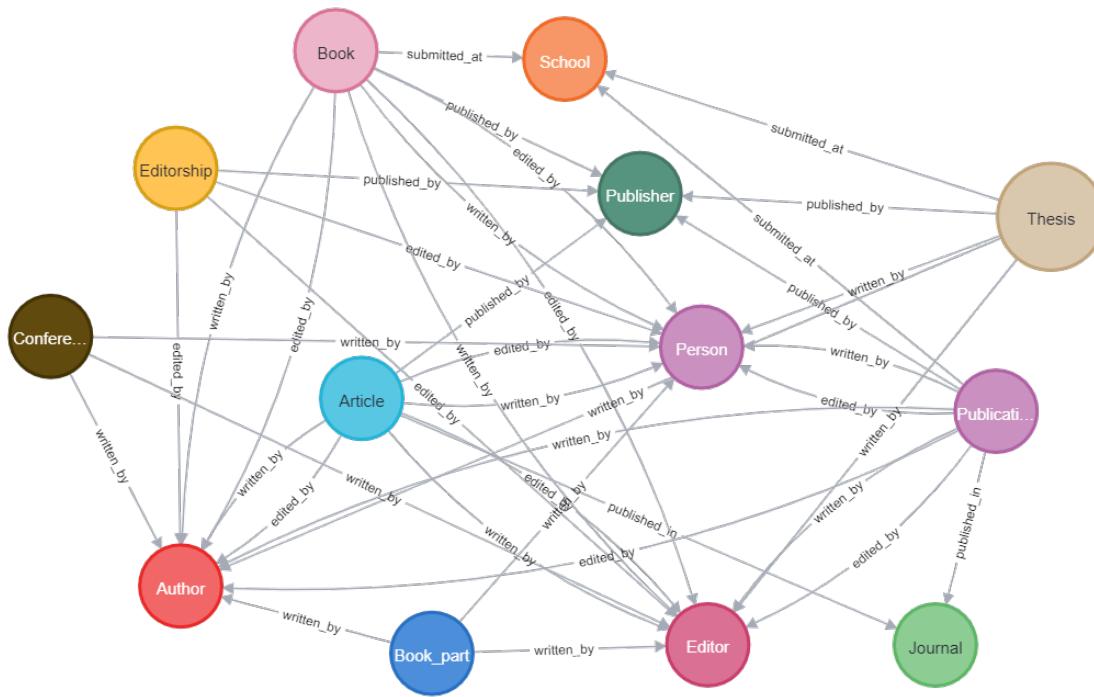


Neo4j: data structure

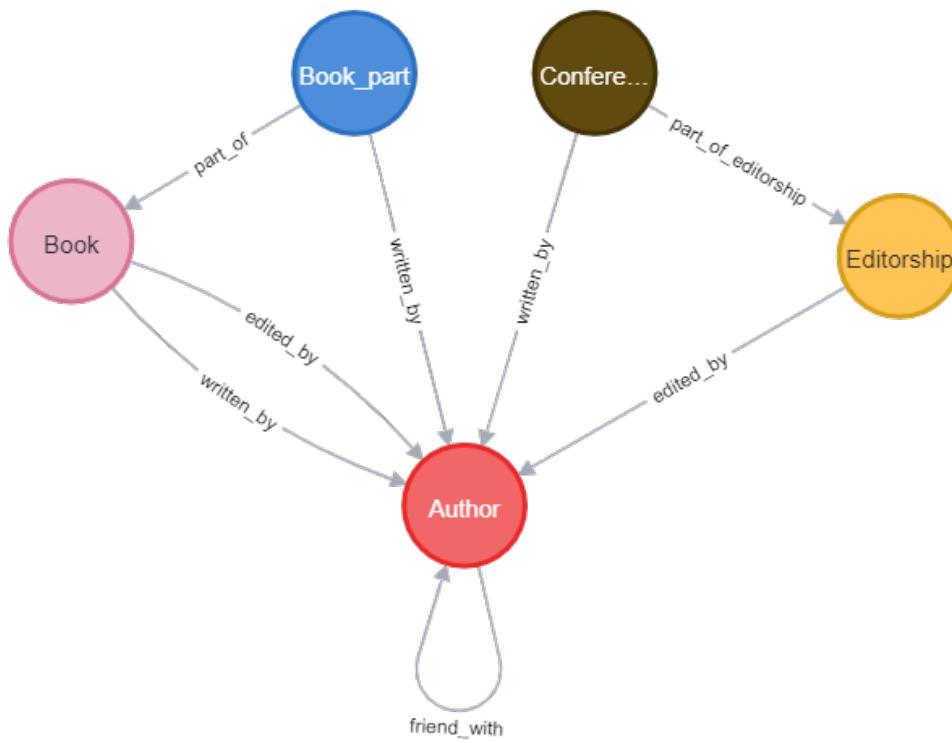
Neo4j: data structure



Neo4j: data structure

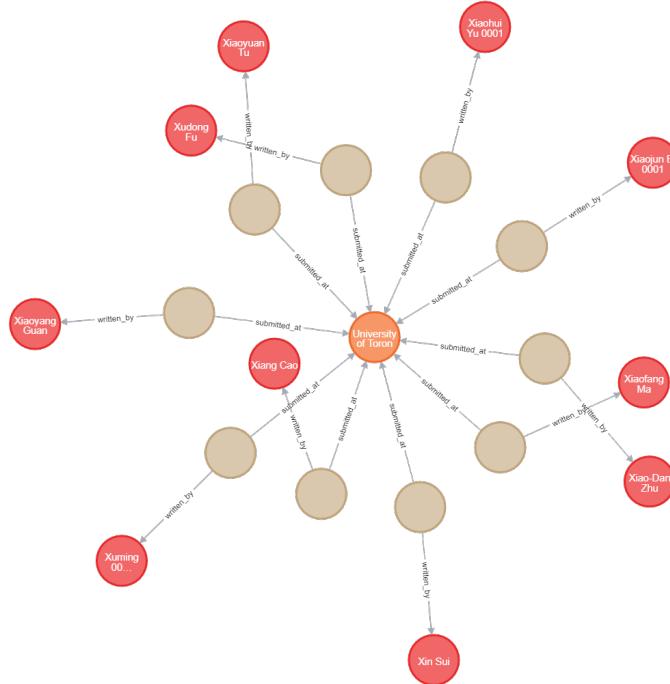


Neo4j: data structure



Neo4j: query

```
1 match (per:Person)←[]-(pub:Publication)-[s:submitted_at]→(school:School)
2 where per.author starts with 'X'
3 with school, count(s) as submissions, collect(pub) as publications, collect(per) as people
4 return school, people, publications, submissions order by submissions desc limit 1
```





POLITECNICO
MILANO 1863

MongoDB

MongoDB: Dataset generation

MongoDB: Dataset generation

We started by
defining the structure
of the document

Data structure

MongoDB: Data structure

```
▼ 0:
  _id: "63776ddd297f4328218ebc59"
  title: "On the Penrose process for rotating black holes"
  ▼ authors:
    ▼ 0:
      name: "Salinas Bauer"
      affiliation: "Geekol"
      email: "salinasbauer@geekol.com"
      ▶ bio:
        ▶ 1: {...}
        ▶ 2: {...}
        ▶ 3: {...}
        ▶ 4: {...}
        ▶ 5: {...}
    ▼ keywords:
      0: "Information Technology/Informatics"
      1: "Research Networking/GIRD"
      2: "Process automation"
    ▼ journal:
      name: "Optim. Lett."
      volume: 11
      number: 58
      date: "1987-05-\t\t\t\t\t23"
      pages: "197 - 388"
      ▶ abstract:
        " Penrose described a pr... of 1-1/sqrt{2} =29%.\\n"
    ▼ sections:
      ▶ 0: {...}
      ▶ 1: {...}
      ▶ 2:
        ▶ title: "Chain enumeration of \$k\$-s of classical\n types"
        ▼ paragraphs:
          ▶ 0: "Went here for brunch tod... or beer in the future."
        ▼ figure:
          [...]
        ▼ subsections:
          ▼ 0:
            title: "Green-Tao theorem in function fields"
            ▶ paragraphs: [...]
            ▶ figure:
              [...]
    ▼ bibliography:
      0: "63776dddcbf3d894a5f8aff3"
      1: "63777073d939a5250d1c80ef"
      2: "63776ddef403a87eea28accf"
```

MongoDB: Dataset generation

We started by
defining the structure
of the document

Data structure

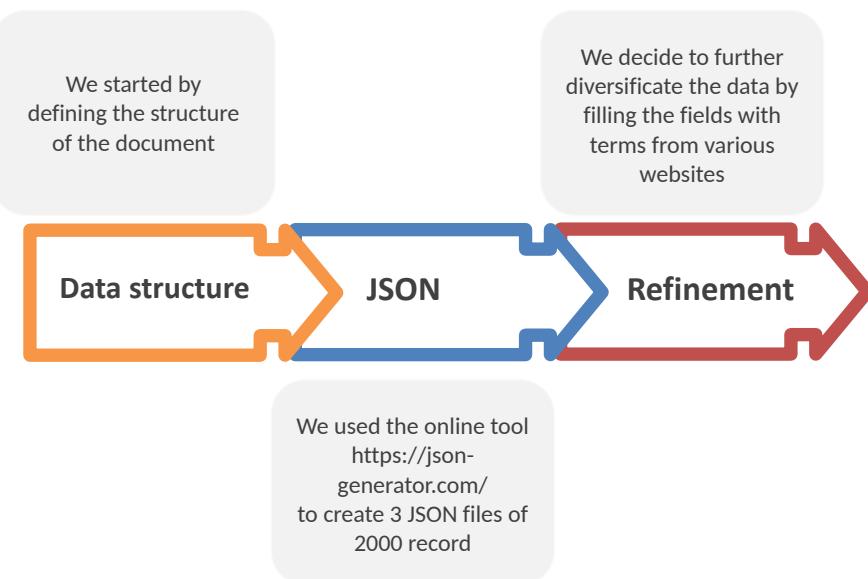
MongoDB: Dataset generation

We started by defining the structure of the document

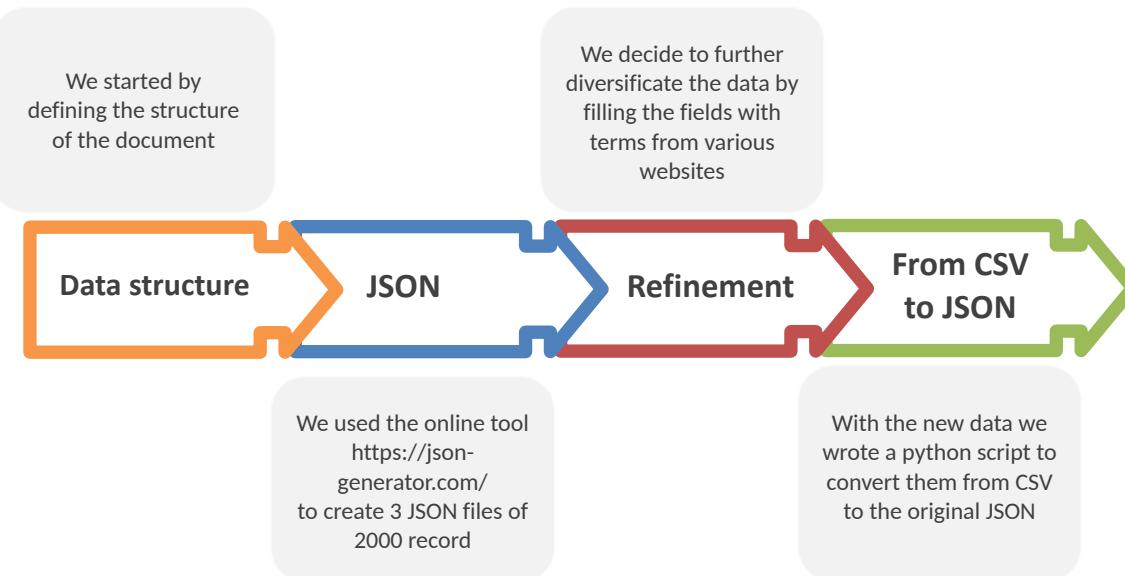


We used the online tool
<https://json-generator.com/>
to create 3 JSON files of 2000 record

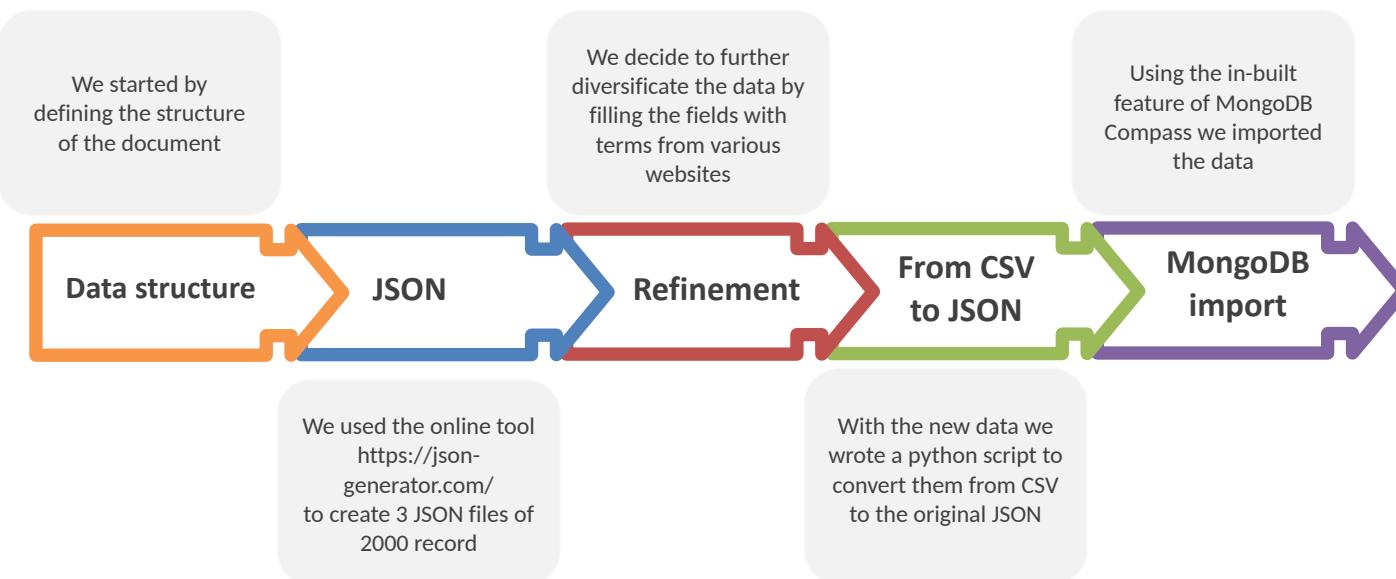
MongoDB: Dataset generation



MongoDB: Dataset generation



MongoDB: Dataset generation



MongoDB: Query

```
Atlas atlas-l7mawr-shard-0 [primary] journal_articles> db.articles.aggregate([
    {
        $match: {$or: [{"_id": ObjectId("63776ddde00236845a2fb904")}, {"_id": ObjectId("63776ddd2f7fbf02b6ecfdf3")}]}
    },
    {
        $lookup: {
            from: "articles",
            localField: "bibliography",
            foreignField: "_id",
            as: "bibliography"
        }
    },
    {
        $unwind: "$bibliography"
    },
    {
        $sort: {"bibliography.journal.date": -1}
    },
    {
        $project: {"title": 1, "bibliography.title":1, "bibliography.journal.date": 1}
    }
])|
```

```
< { _id: '63776ddd2f7fbf02b6ecfdf3',
  title: 'HD 77361: A new case of super Li-rich K giant with anomalous low 12C/13C\n  ratio',
  bibliography:
  [ '63777073bc1aeb06de1b7492',
    '63776dde07a64b8cb303d4e9',
    '63776dde6a9de277c1ce0ae1' ] }
{ _id: '63776ddde00236845a2fb904',
  title: '$F$-pure homomorphisms, strong $F$-regularity, and $F$-injectivity',
  bibliography:
  [ '63776dde84264e473ba7c413',
    '63776dde697c3f791c68642d',
    '63776dde6efbe362b4909aba',
    '637770ac8fdfb72787bb6363',
    '63776ddef5d187cd65b423b2' ] }
```

MongoDB: Query

MongoDB: Query

```
< { _id: '63776ddde00236845a2fb904',
    title: '$F$-pure homomorphisms, strong $F$-regularity, and $F$-injectivity',
    bibliography:
      { title: 'Computing the multifractal spectrum from time series: An algorithmic\n approach',
        journal: { date: 2020-04-11T22:00:00.000Z } }
    { _id: '63776ddde00236845a2fb904',
      title: '$F$-pure homomorphisms, strong $F$-regularity, and $F$-injectivity',
      bibliography:
        { title: 'Spin dimers under staggered and random field in\n Cu$_2$Fe$_2$O$_4$S$_{13}$',
          journal: { date: 2015-09-13T22:00:00.000Z } }
    { _id: '63776dd2f7fb02b6ecfdf3',
      title: 'HD 77361: A new case of super Li-rich K giant with anomalous low 12C/13C\n ratio',
      bibliography:
        { title: 'Phenomenology of ESR in heavy fermion systems: Fermi liquid and\n non-Fermi liquid regime',
          journal: { date: 2009-12-14T23:00:00.000Z } }
    { _id: '63776ddde00236845a2fb904',
      title: '$F$-pure homomorphisms, strong $F$-regularity, and $F$-injectivity',
      bibliography:
        { title: 'Quantum-optical state engineering up to the two-photon level',
          journal: { date: 2009-07-13T22:00:00.000Z } }
    { _id: '63776dd2f7fb02b6ecfdf3',
      title: 'HD 77361: A new case of super Li-rich K giant with anomalous low 12C/13C\n ratio',
      bibliography:
        { title: 'Bell\\'s Inequalities: Foundations and Quantum Communication',
          journal: { date: 2006-01-08T23:00:00.000Z } }
    { _id: '63776ddde00236845a2fb904',
      title: 'HD 77361: A new case of super Li-rich K giant with anomalous low 12C/13C\n ratio',
      bibliography:
        { title: 'Terrorism: Mechanisms of Radicalization Processes, Control of Contagion\n and Counter-Terrorist Measures',
          journal: { date: 2004-08-14T22:00:00.000Z } }
    { _id: '63776ddde00236845a2fb904',
      title: '$F$-pure homomorphisms, strong $F$-regularity, and $F$-injectivity',
      bibliography:
        { title: 'Incorporating characteristics of human creativity into an evolutionary\n art algorithm',
          journal: { date: 1990-10-30T23:00:00.000Z } }
    { _id: '63776ddde00236845a2fb904',
      title: '$F$-pure homomorphisms, strong $F$-regularity, and $F$-injectivity',
      bibliography:
        { title: 'On the magnetic equation of state in (2+1)-flavor QCD',
          journal: { date: 1973-03-24T23:00:00.000Z } } }
```

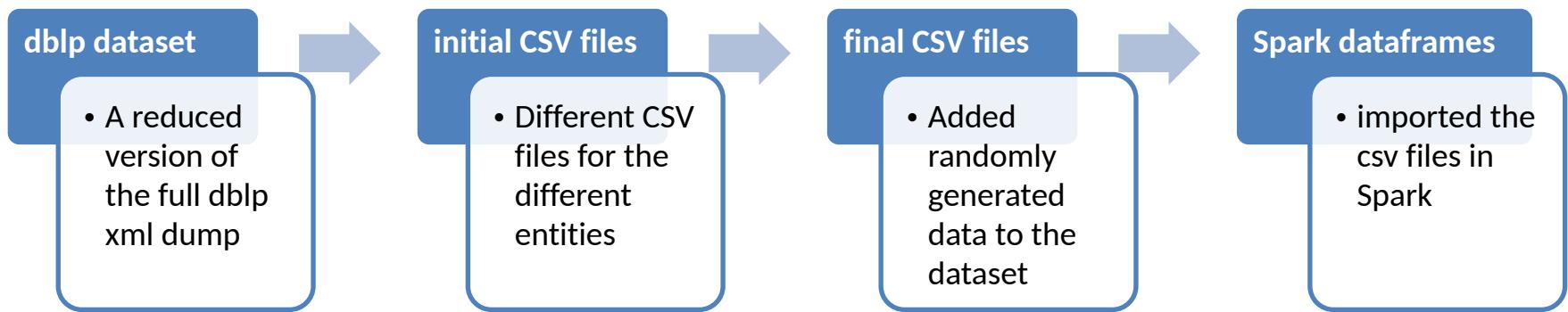


POLITECNICO
MILANO 1863

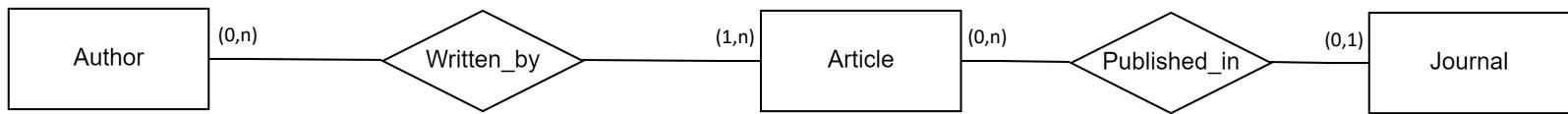
Apache Spark

Spark: dataset generation and import

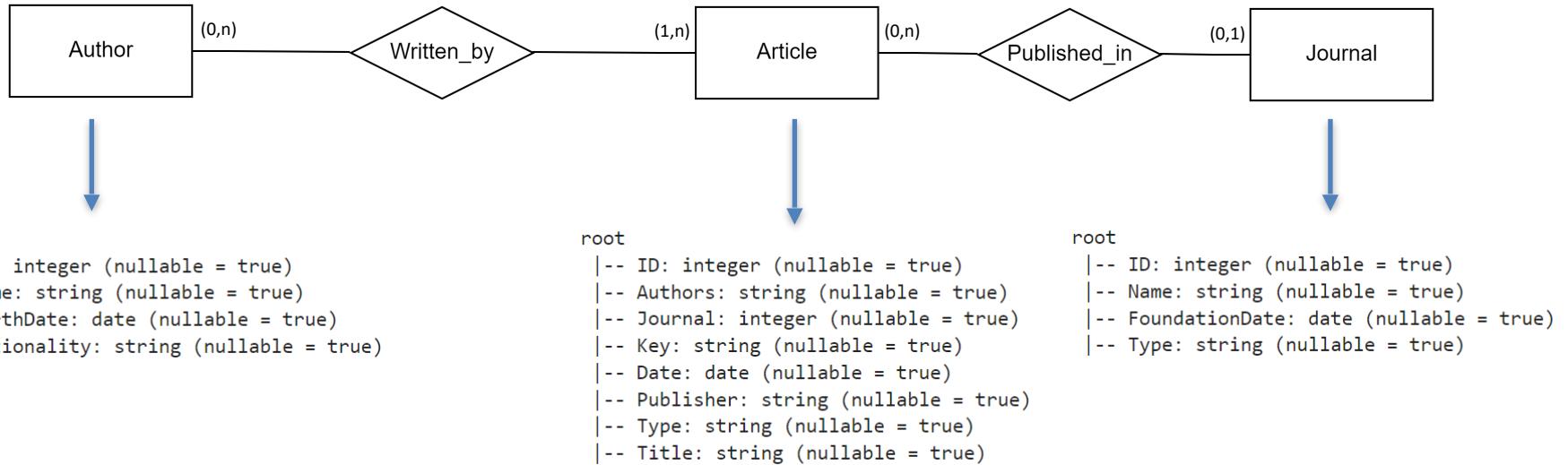
Similar process as for the Neo4j dataset



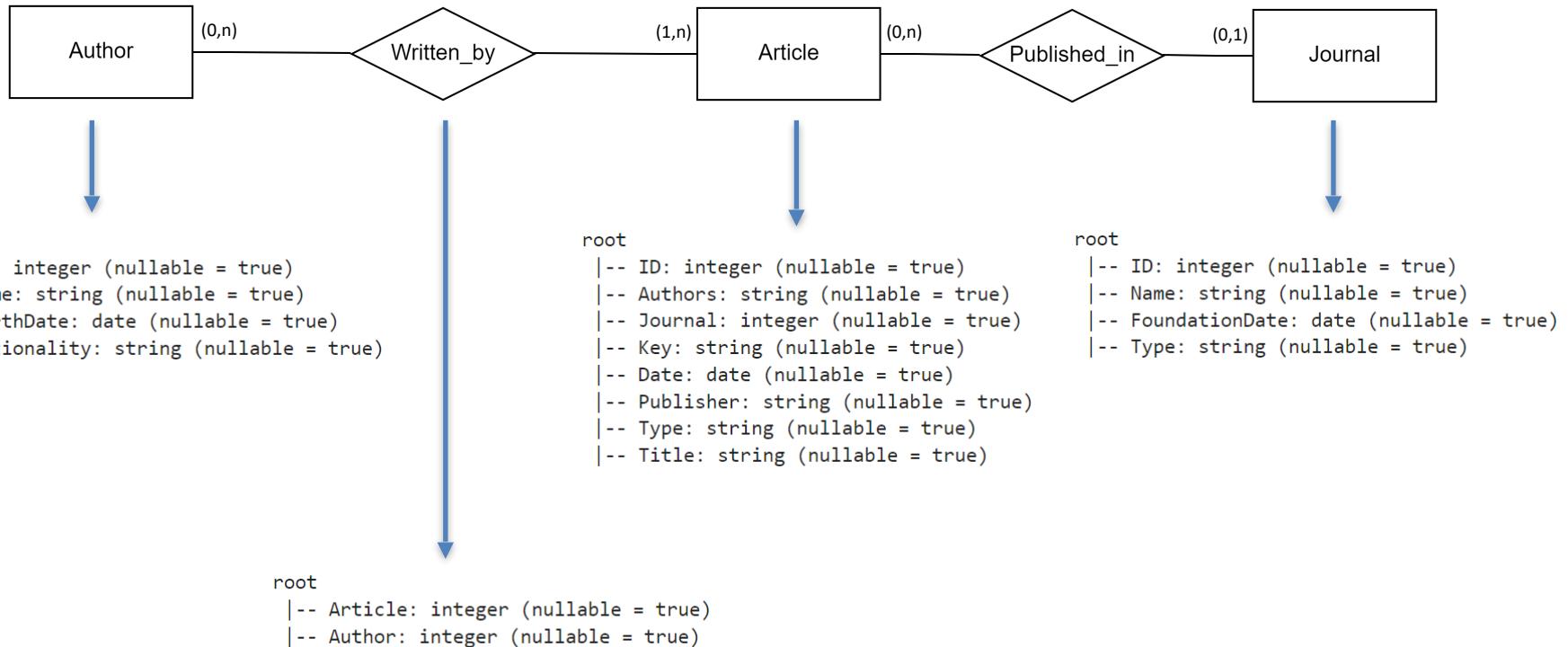
Spark: data structure



Spark: data structure



Spark: data structure



Spark: query

```
authors.filter(col("Name").like("L%")) \
    .join(written_by, written_by.Author == authors.ID, "inner") \
    .groupby(col("Name")).count() \
    .select(col("Name"), col("count").alias("#Articles")) \
    .sort(col("#Articles").desc()) \
    .show(10)
```

```
+-----+-----+
|      Name | #Articles |
+-----+-----+
| Lajos Hanzo |      124 |
|   Leo Storne |       41 |
| Lie-Liang Yang |      21 |
| Lothar Breuer |      17 |
| Lin Cai 0001 |      16 |
|     Lijun Ji |      16 |
|    Lian Zhao |      13 |
| Luca Benini |      13 |
|    Luxi Yang |      13 |
| Ludger Humbert |      12 |
+-----+-----+
only showing top 10 rows
```

The query returns the Authors with the name that start with an «L», ranked in descending order w.r.t. the number of articles they wrote