



**Building a Linked Data Knowledge Graph  
for the Scholarly Publishing Domain**

Markus Kaindl

Senior Manager Semantic Data &  
SN SciGraph Product Owner

LDBC Technical User Community (TUC)  
Meeting; Munich, September 1<sup>st</sup> 2017

**SPRINGER NATURE**

# Agenda

## Intro

- Springer Nature SciGraph
- Linked Open Data Publishing

## Status

- SN SciGraph Hack Day
- Analytics Dashboards

## Data

- Roadmap EOY and beyond

**Intro:**

- **Springer Nature SciGraph**
- **Linked Open Data Publishing**

**#1**



**Intro:**

- **Springer Nature SciGraph**
- **Linked Open Data Publishing**

**#1.1**



A world-leading  
research, educational  
and professional  
publisher

Formed in **May 2015** through the **merger** of Nature Publishing Group, Palgrave Macmillan, Macmillan Education and Springer Science+Business Media

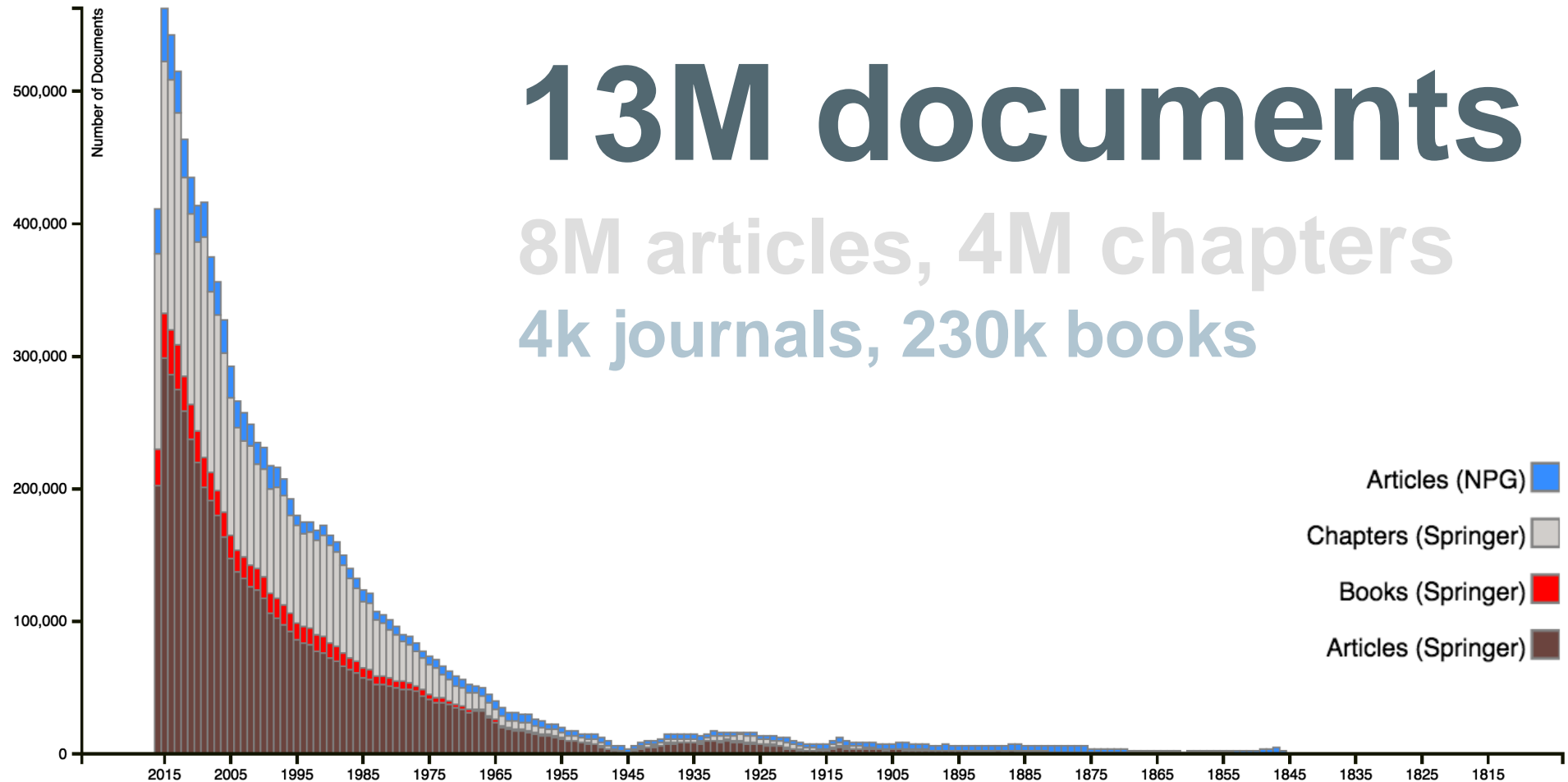
# [Pre-Merger] Springer Science + Business Media brands



# [Pre-Merger] Macmillan Science & Education brands



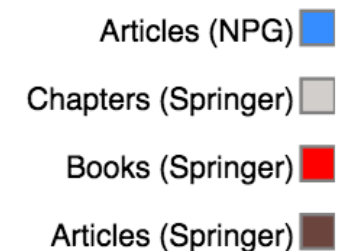
We publish a lot of science (since 1815)



# 13M documents

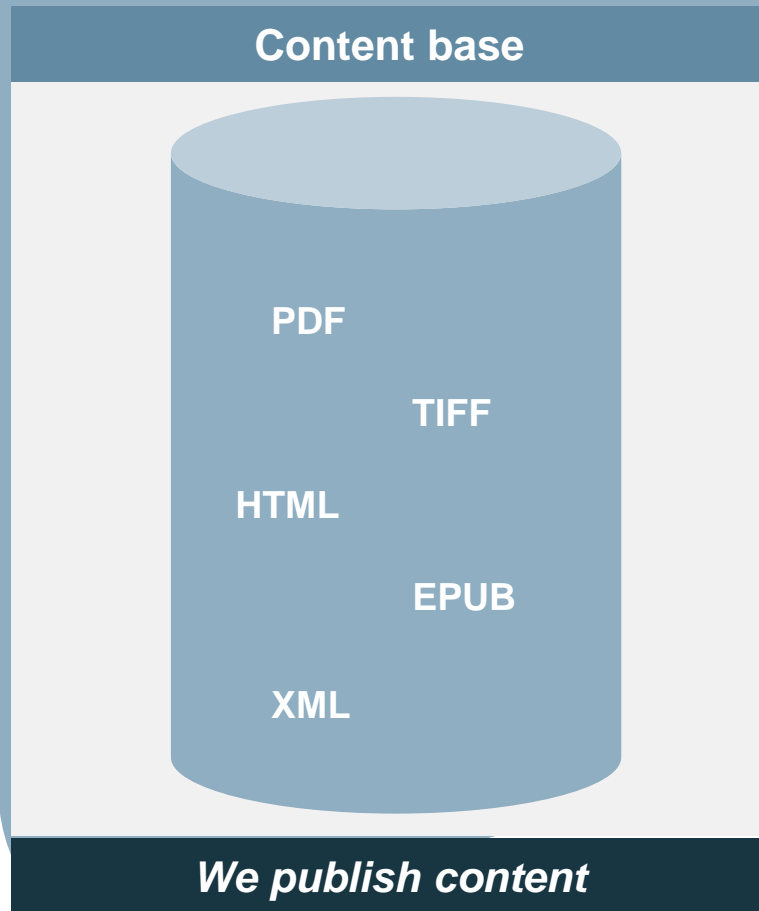
8M articles, 4M chapters

4k journals, 230k books



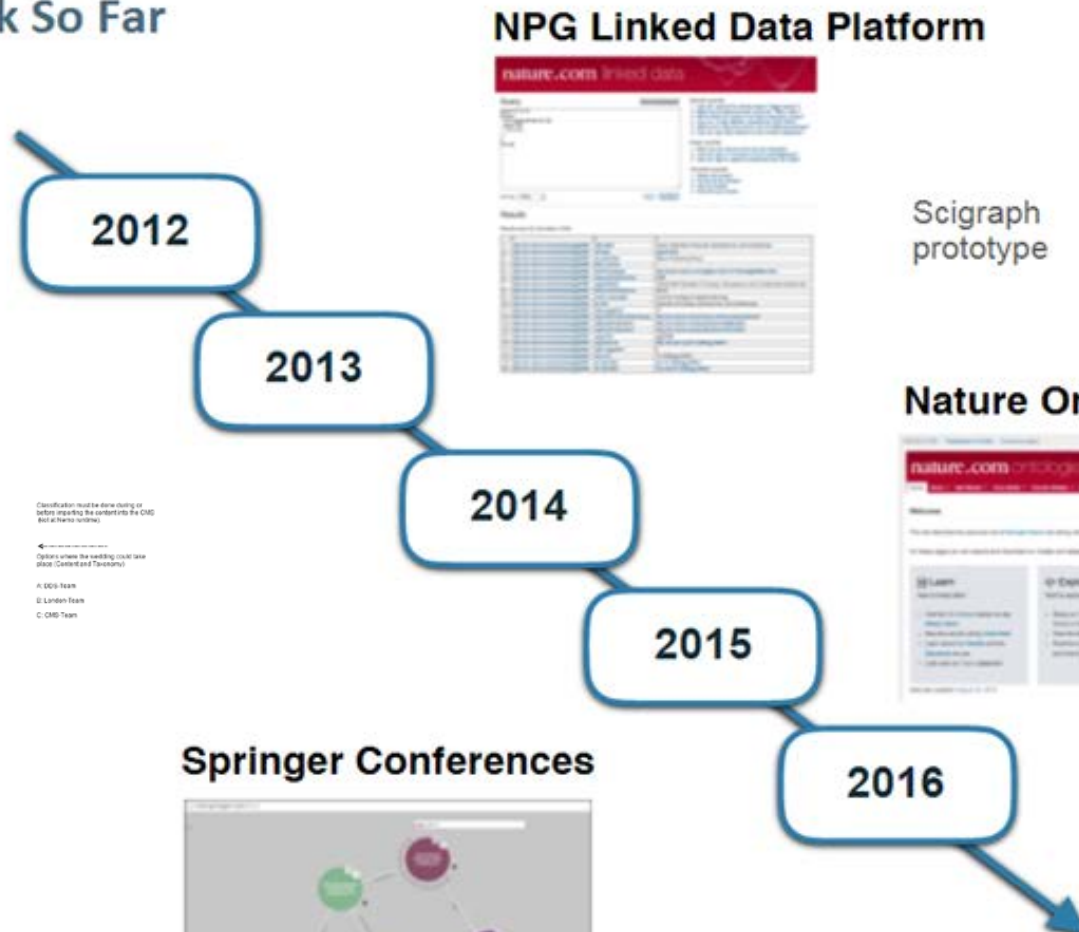


# From Content to Data



# Integration of Various Semantic Data Initiatives

## Our Work So Far



### NPG Linked Data Platform



CURI Semantic Annotation Project

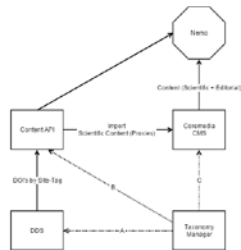
Scigraph prototype

Subject Pages

### Nature Ontologies Portal



### Linnaeus Project



Classification results are shown in a table or list view depending on the content type (e.g. article or reference).  
 Options shown in the sidebar could also be shown in the main content area.  
 A: EDS team  
 B: London team  
 C: CMS team

2014

2015

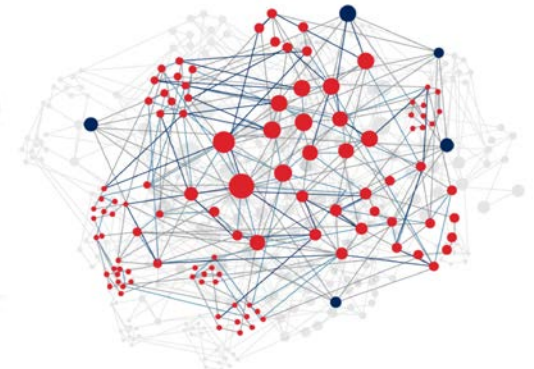
### Springer Conferences



Springer Protocols

2016

### SN SciGraph



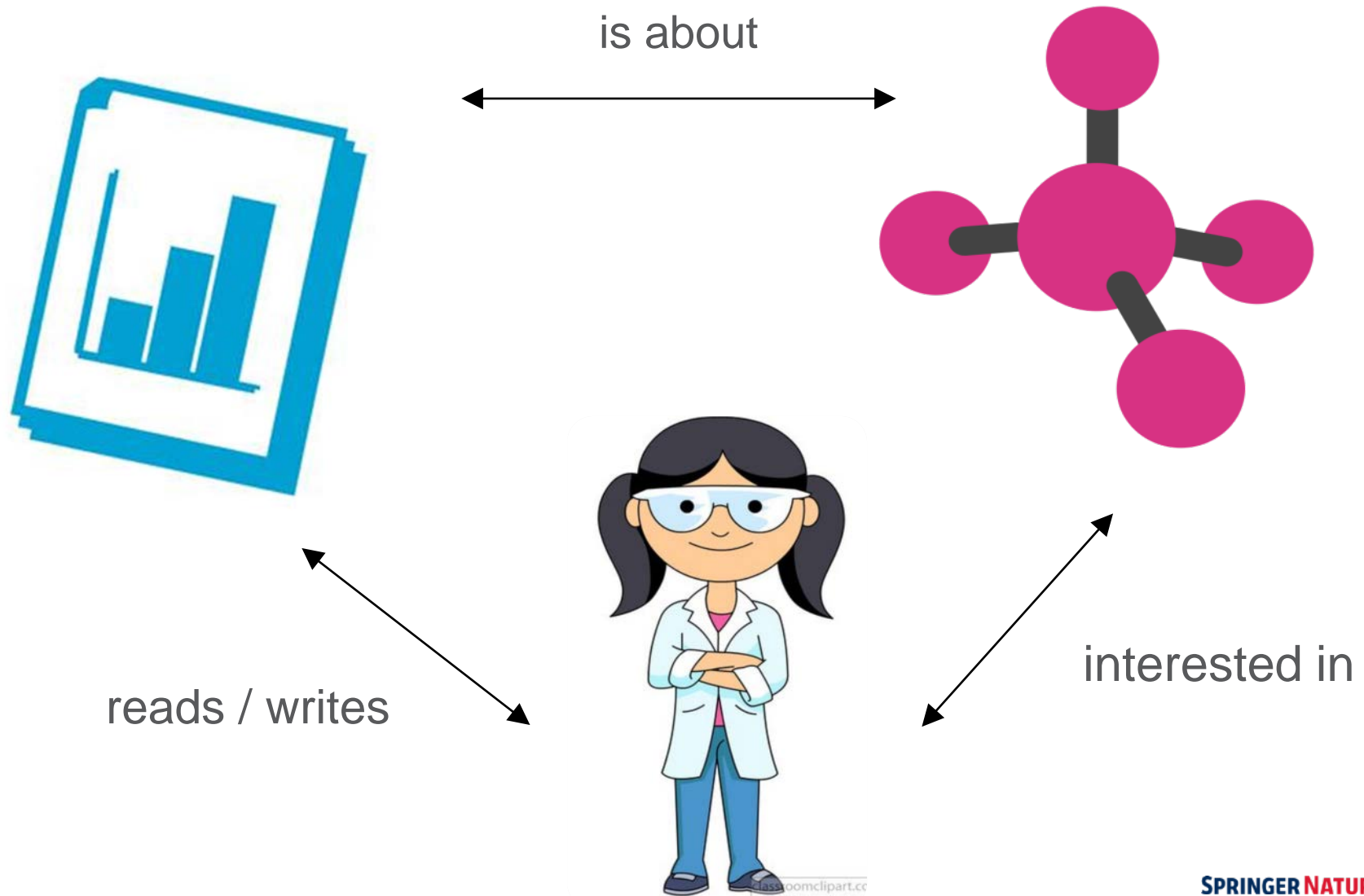
# Product Vision

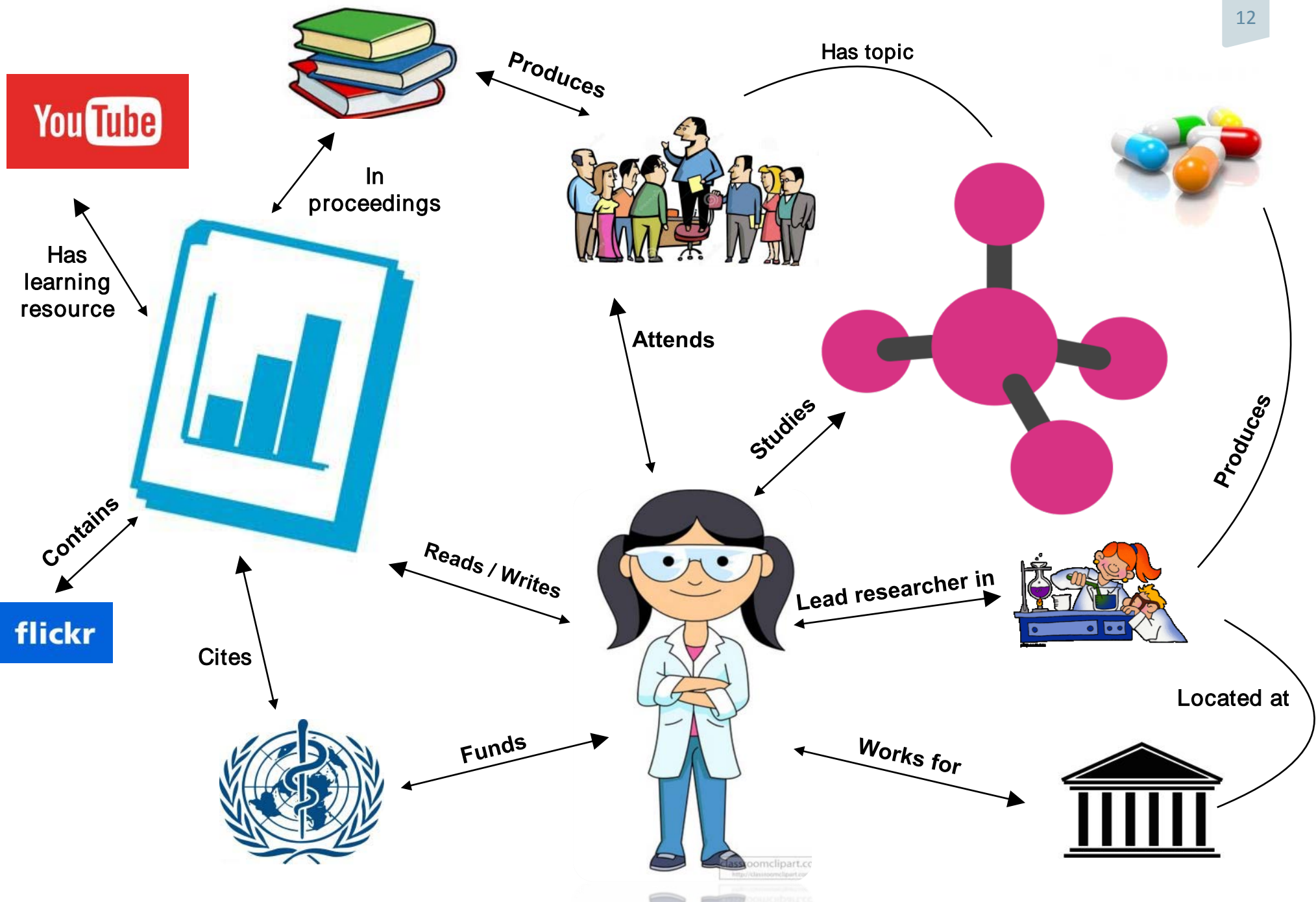
We create the largest state-of-the-art linked open data aggregation platform in the scholarly domain. In doing so, we increase content discoverability and provide data tools and services for researchers, authors, editors, librarians, data scientists, funders, conference organizers, and many others by adding value across all content types.



# SciGraph

# Three areas of knowledge we care about





## ETL Architecture: main features [in evolution]

### Tech stack

- > Airflow framework (Airbnb)\*
- > Amazon S3 to make backups
- > GraphDB triplestore (staging and presentation)
- > Elasticsearch and Kibana\*

### Components & Principles

- > Graph must be **'ephemeral'**
- > Data sources versioning algorithm
- > Identity Persistence service
- > Validation via SHACL (TopBraid API)\*

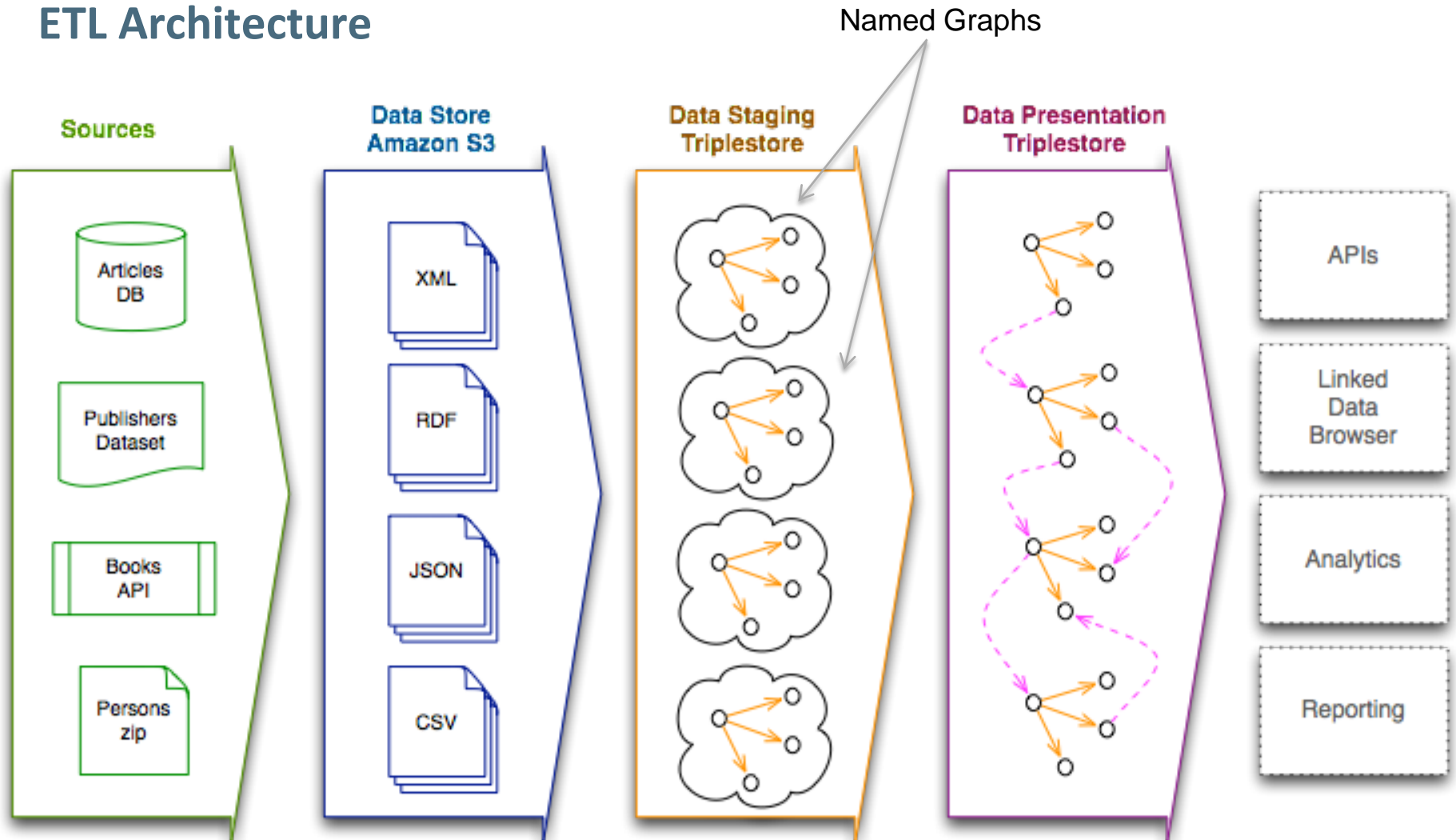


**Airflow**



\*Open Source

# ETL Architecture



- \* Versioning service
  - \* (md5 checksum, timestamps, origin version, etc...)

- \* Extraction
- \* Validation
- \* Identity Persistence
- \* Updating / Replacing named graphs

- \* Integration (union graph)
- \* Inference

**Intro:**

- Springer Nature SciGraph
- **Linked Open Data Publishing**

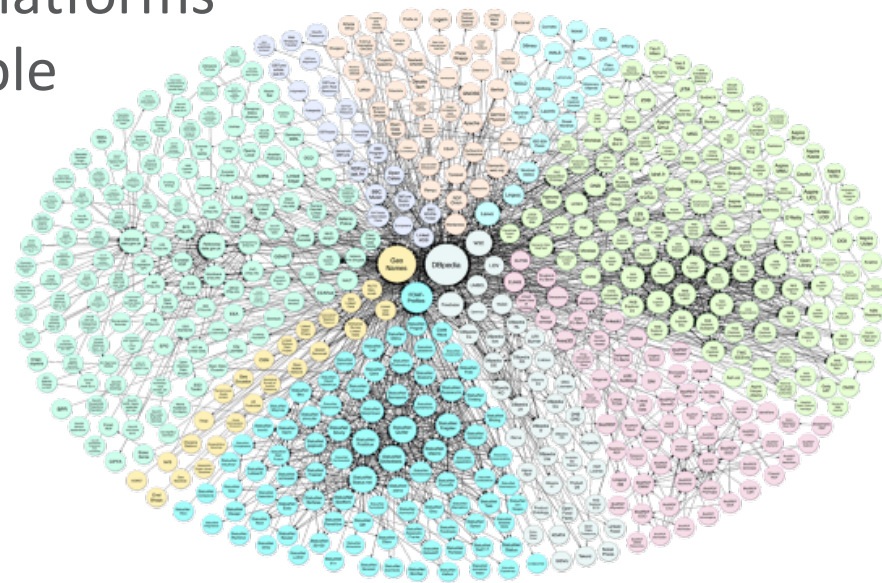
# #1.2



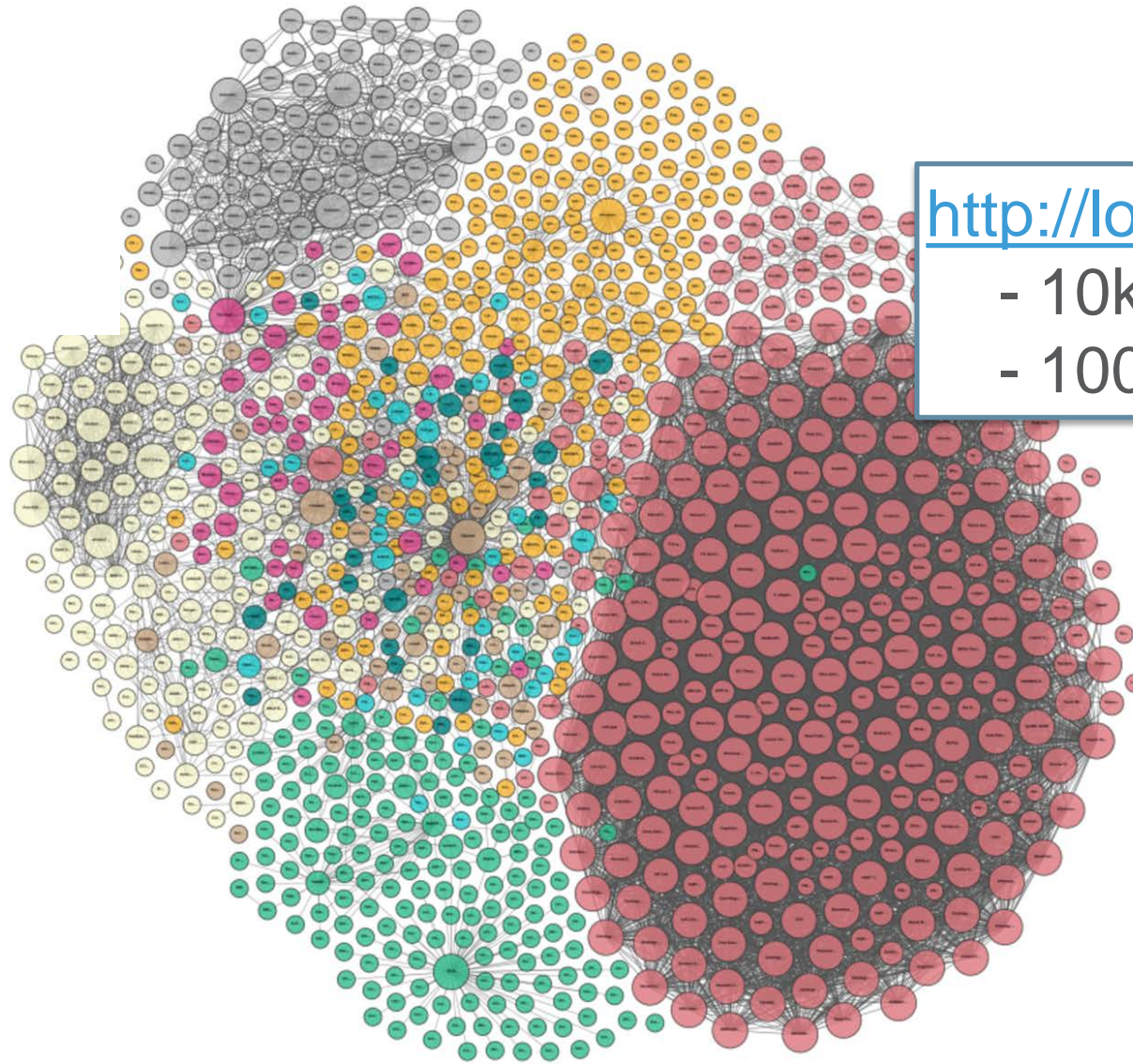
# The Web of Data: Benefits of Linked Data at a glance

## Why publish Linked Open Data?

- Increase discoverability and ranking through linking
- “The future library catalogue is the Internet!”
- Drive traffic and usage to our platforms
- Make our data machine-readable
- Be part of the LOD cloud



# The Global Linked Data Cloud (2017)



# Libraries using Linked Data

The Library of Congress > Linked Data Service

**LIBRARY OF CONGRESS LINKED DATA SERVICE**

**LC Linked Data Service**  
Authorities and Vocabularies

Search

Enter Keyword or Phrase

All  
LC Subject Headings  
LC Name Authority File  
LC Classification  
LC Children's Subject Headings

Search

**Available Datasets**

The Linked Data Service provides access to commonly found standards and vocabularies promulgated by the Library of Congress. This includes data values and the controlled vocabularies that house them. The following are currently available:

- LC Subject Headings
- LC Name Authority File
- LC Classification
- LC Children's Subject Headings
- LC Genre/Form Terms
- LC Medium of Performance Thesaurus for Music
- MARC Relations
- MARC Countries
- MARC Geographic Areas
- MARC Languages
- MARC Genre Terms
- ISO639-1 Languages
- ISO639-2 Languages
- Schemes
- Identifiers
- Carriers
- Content Types
- Media Types
- Resource Types
- Description Conventions

## Library of Congress Linked Data Service (2009)

- A library catalog “must be designed by considering its context of the Web”
- Access to data at no cost.
- Ability to link to Library of Congress data values within your metadata via Linked Data.

### Other libraries:

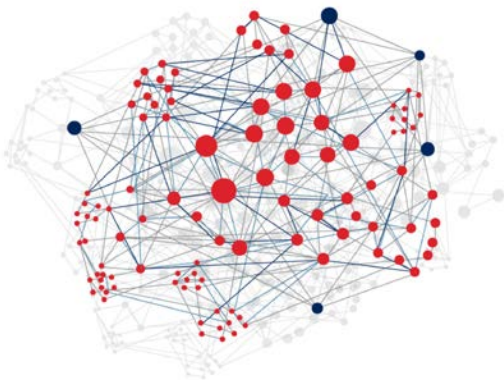
- British Library (BL)
- French National Library (BNF)
- German National Library (DNB)
- National Library of Spain (BNE)
- National Library of Sweden (LIBRIS)
- Hungarian National Library (NSL)

**just  
released****SPRINGER NATURE**

# Springer Nature SciGraph

A Linked Open Data platform for the scholarly domain

**SN** SciGraph



- > Collaborative effort between Springer Nature and Digital Science
- > Supporting internal use cases, but also contributing to an emerging web of **linked science data**
- > Integrating data from a **variety of sources** using Linked Data technology
- > Not just publications but a **wealth** of other related data

# Linked Open Data Publishing: making science more accessible

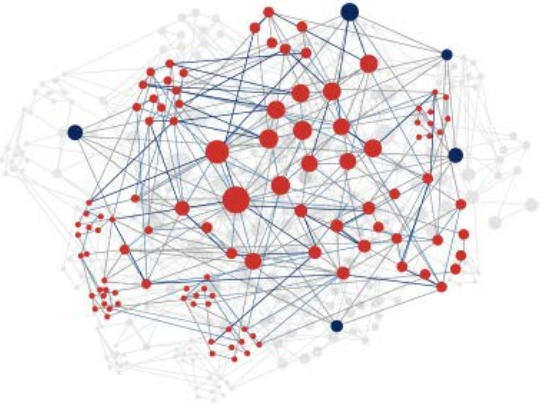
**SPRINGER NATURE**

## Springer Nature SciGraph

A Linked Open Data platform for the scholarly domain

We are pleased to introduce Springer Nature SciGraph, the new Linked Open Data platform aggregating data sources from Springer Nature and key partners from the scholarly domain. The Linked Open Data platform will initially collate information from across the research landscape, such as funders, research projects, conferences, affiliations and publications. Additional data, such as citations, patents, clinical trials and usage numbers will follow over time. This high quality data from trusted and reliable sources provides a rich semantic description of how information is related, as well as enabling innovative visualizations of the scholarly domain.

By doing so, Springer Nature SciGraph overcomes former boundaries by relating comprehensive information about the research landscape. It represents a further step in data integration and it will continue to grow organically. This platform will increase the discoverability of high quality data as larger parts of our datasets will be made freely available under a CC BY-NC 4.0 license.



The data in Springer Nature SciGraph is projected to contain 1.5 to 2 billion triples. It will comprise metadata from journals and articles, books and chapters, organizations, institutions, funders, research grants, patents, clinical trials, substances, conference series, events, citations and reference networks, Altmetrics, links to research datasets and much more.

Any questions?  
Please contact us.

Dataset Download

Licensing Information

Further Info

Conference Presentation 2016 (PDF, 11.56 MB)

## Feb. 2017 Data Release

### At a glance:

- 150 M triples / 32 GB download size
- CC-BY-NC License

### Metadata about:

- Articles 2012-2016 (5M) + Abstracts
- Grants (200k)
- Journals (3k)
- Subjects (3k)
- Core Ontology

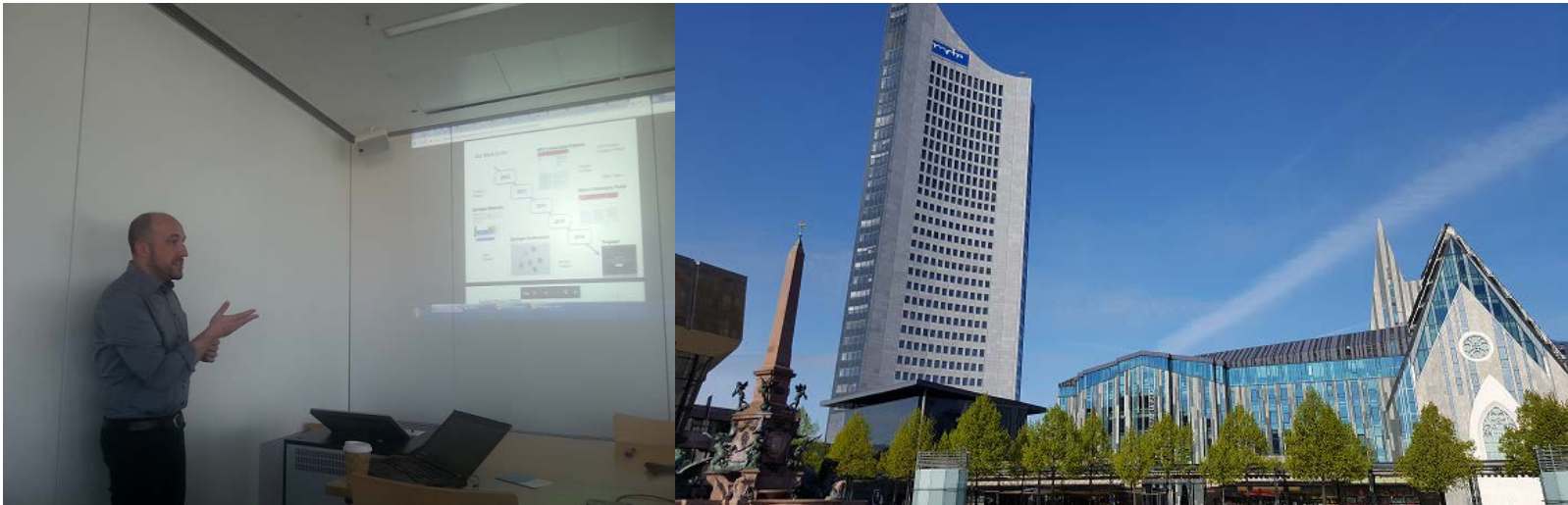
# Springer Nature SciGraph: Feedback on Twitter after PR

- **Selected tweets:**
  - “@SpringerNature Thank you for SciGraph integrated research data for helping the knowledge revolution #Semantic #Reason #ThinkingWeb”
  - “@SpringerNature @NatureCellBio SciGraph excellent new tool”
  - “Interesting, but only has info from one publisher. Is it extensible to other publishers? Single-publisher solutions are no solution!”
  - “Springer Nature SciGraph Dataset ネイチャー誌5年分の記事メタデータ1.55億トリプルをCC-BY-NCで提供。年内にさらに追加と。モデル説明あ”
  - “I hope this will go far beyond the mere understanding of ‘which conferences are taking place in subject areas with rising funding’”
  - „Springer Nature SciGraph: Supporting open science and the wider understanding of research”



# Workshop with **DBpedia** (April 2017)

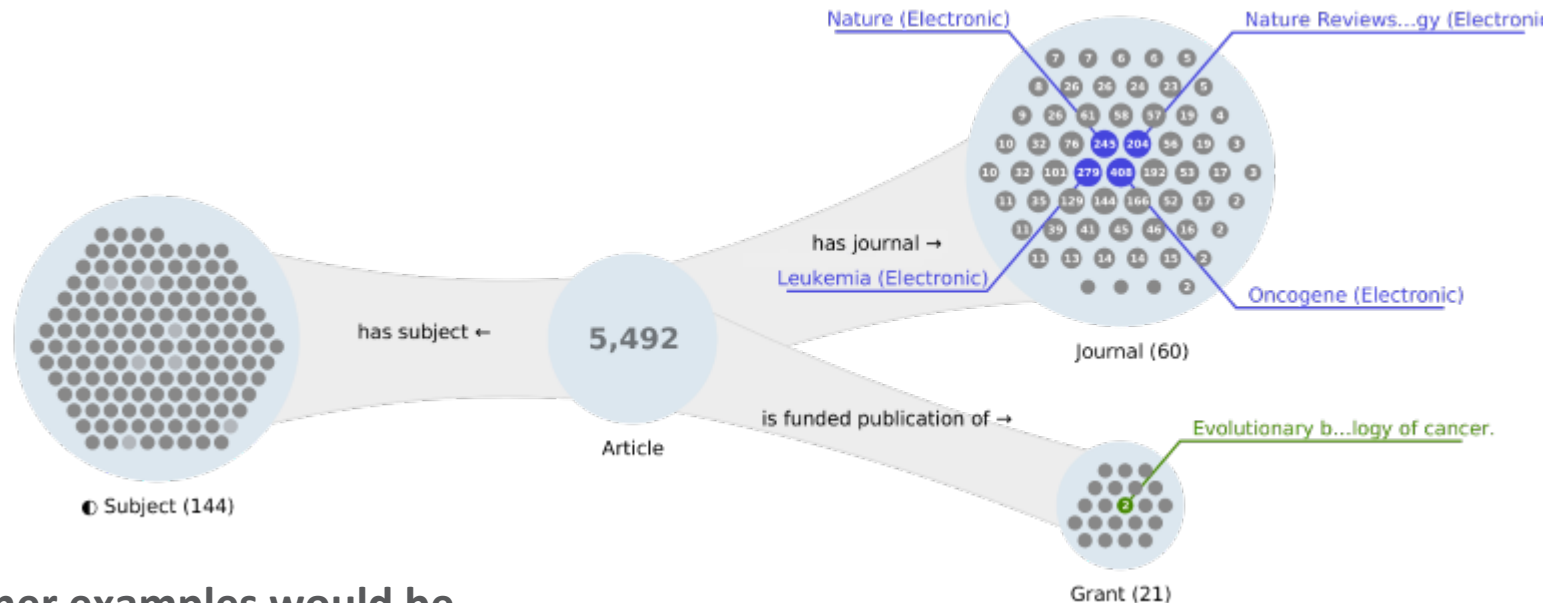
- DBpedia is the Linked Open Data transformation of facts extracted from Wikipedia
- We are uploading our data sets with links to Wikipedia on the DBpedia GitHub repository in order to generate backlinks from DBpedia and increase traffic to our content platforms and usage of our SN SciGraph LOD.
- Currently in the process of hiring an intern that works both for DBpedia and SciGraph





## Examples of users working with SN SciGraph data

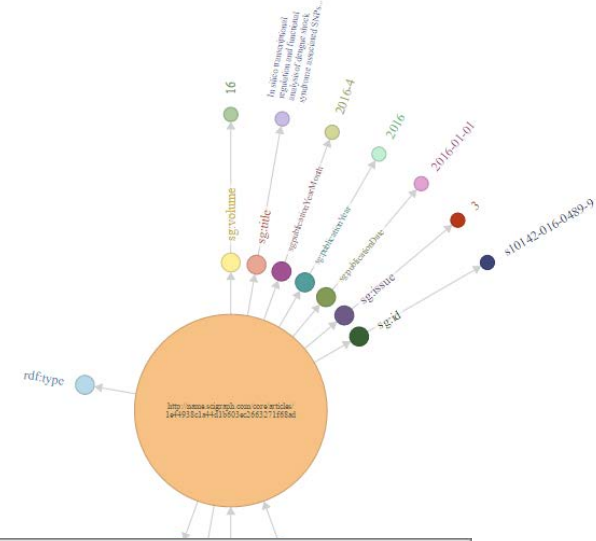
- [SemSpect](#): Uncovering the Hidden in Springer Nature's SciGraph
  - The essential problem is to get an idea of the queries that deliver real insight.
  - [This video](#) shows a sample exploration of SN SciGraph data with SemSpect.



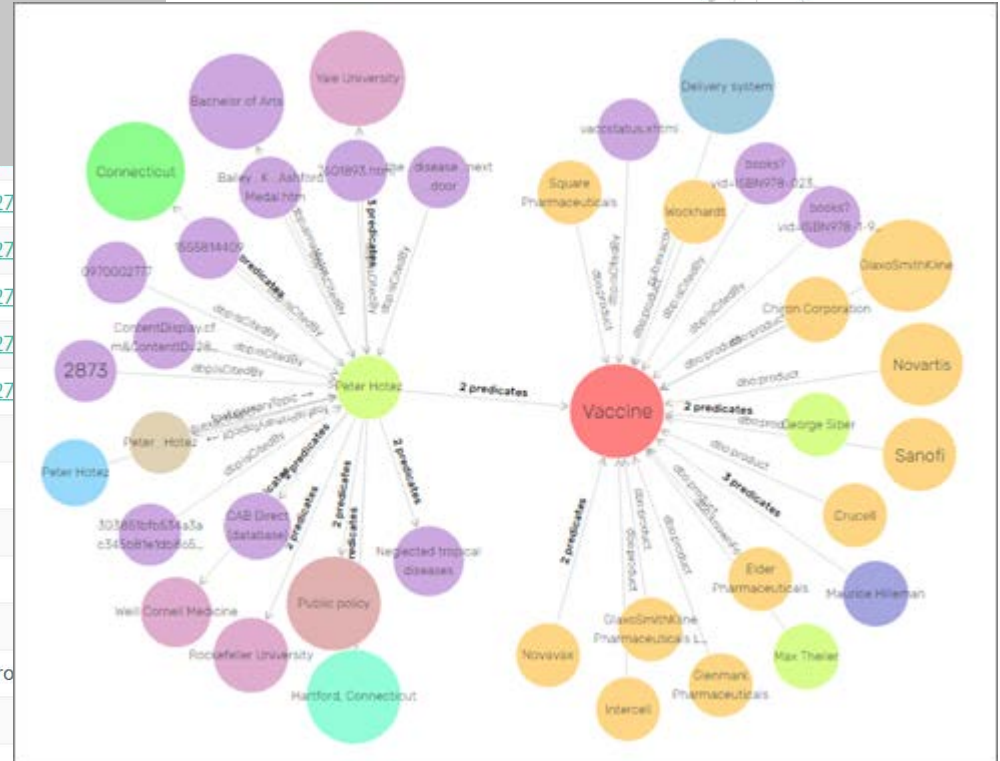
- Other examples would be
  - ResearchGraph (Australia)
  - Science Media Center (Cologne)



# Springer Nature SciGraph: LOD Browsers



Property	Value
<a href="#">sg:doi</a>	10.1007/s101...
<a href="#">sg:doiLink</a>	<http://dx.doi...
<a href="#">sg:hasAbstract</a>	<http://name.s...
<a href="#">sg:hasContribution</a>	<http://name.scigraph.com/core/contributions/1e44938c1a44d1b603ec266327
<a href="#">sg:hasContribution</a>	<http://name.scigraph.com/core/contributions/1e44938c1a44d1b603ec266327
<a href="#">sg:hasContribution</a>	<http://name.scigraph.com/core/contributions/1e44938c1a44d1b603ec266327
<a href="#">sg:hasContribution</a>	<http://name.scigraph.com/core/contributions/1e44938c1a44d1b603ec266327
<a href="#">sg:hasContribution</a>	<http://name.scigraph.com/core/contributions/1e44938c1a44d1b603ec266327
<a href="#">sg:hasContribution</a>	<http://name.scigraph.com/core/contributions/1e44938c1a44d1b603ec266327
<a href="#">sg:id</a>	s10142-016-0489-9
<a href="#">sg:issue</a>	3
<a href="#">sg:publicationDate</a>	2016-01-01
<a href="#">sg:publicationYear</a>	2016
<a href="#">sg:publicationYearMonth</a>	2016-4
<a href="#">sg:title</a>	In silico transcriptional regulation and functional analysis of dengue shock syndro
<a href="#">sg:volume</a>	16
<a href="#">rdf:type</a>	<a href="#">sg:Article</a>



# Status

- SN SciGraph Hack Day
- Analytics Dashboards

# #2

# Status

- SN SciGraph Hack Day
- Analytics Dashboards

# #2.1

## Key Partners



Digital Science is a primary data provider and development partner on the project - also helping to kick-start the project with senior staff joining our team.



derivo is supporting the project with their expertise in knowledge modeling. They also provide an instance of SemSpect for visually exploring SN SciGraph data: <http://scigraph.semspect.de>



We are using Ontotext's triple store GraphDB as scalable semantic graph database. Their support during the initiation of the project and afterwards has been excellent.



InfoChem provides named entity recognition in the chemical domain, annotating relevant substances, chemical compounds & molecules.



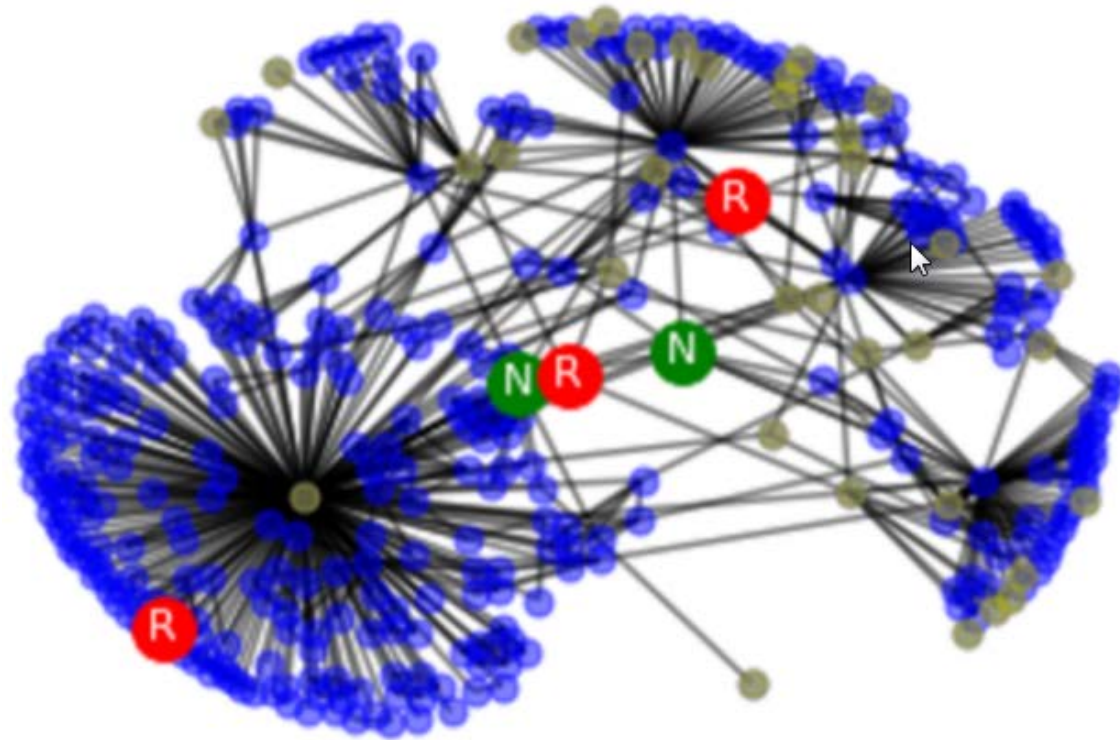
We will be storing the main concepts of our publications as extracted by Unsilo in our graph to cluster documents around certain topic areas.



# Springer Nature SciGraph: Hack Day June 23<sup>rd</sup> in London

## What have participants worked on?

- “Helping you get into the right journal”
- “Connecting Articles and Data Repositories”
- “Peer Reviewer assignment system”



More project details of this productive day can be found [in these slides](#).



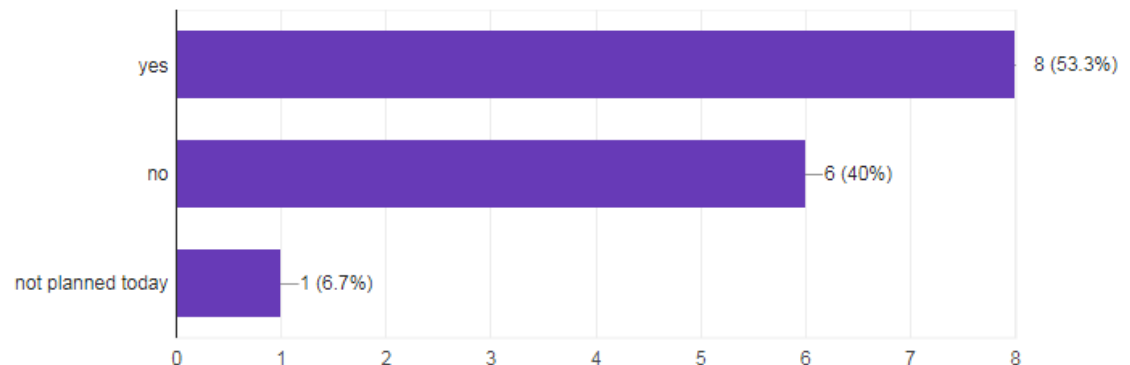
# Springer Nature SciGraph: Hack Day June 23<sup>rd</sup> in London

## What were the key insights?

- One day was enough time; the event was well structured (but some attendees of course wanted more time for hacking )
- Our data model is very much understandable; there is room for improvement around data quality, data coverage and data APIs – all of which were considered adequate
- A clear mandate that citations would be most useful to users – this is aligned with our goal to include citations in the graph next
- Finally, more than 50% would want to use the data commercially – (un)fortunately, our current CC-BY-NC license doesn't allow that
- It is really affirming that 87% would like to remain in contact with us and be informed about upcoming events.

License: are you interested in using the data commercially?

15 responses



You can read more details about the day in [my blog post on Hive](#).

# Status

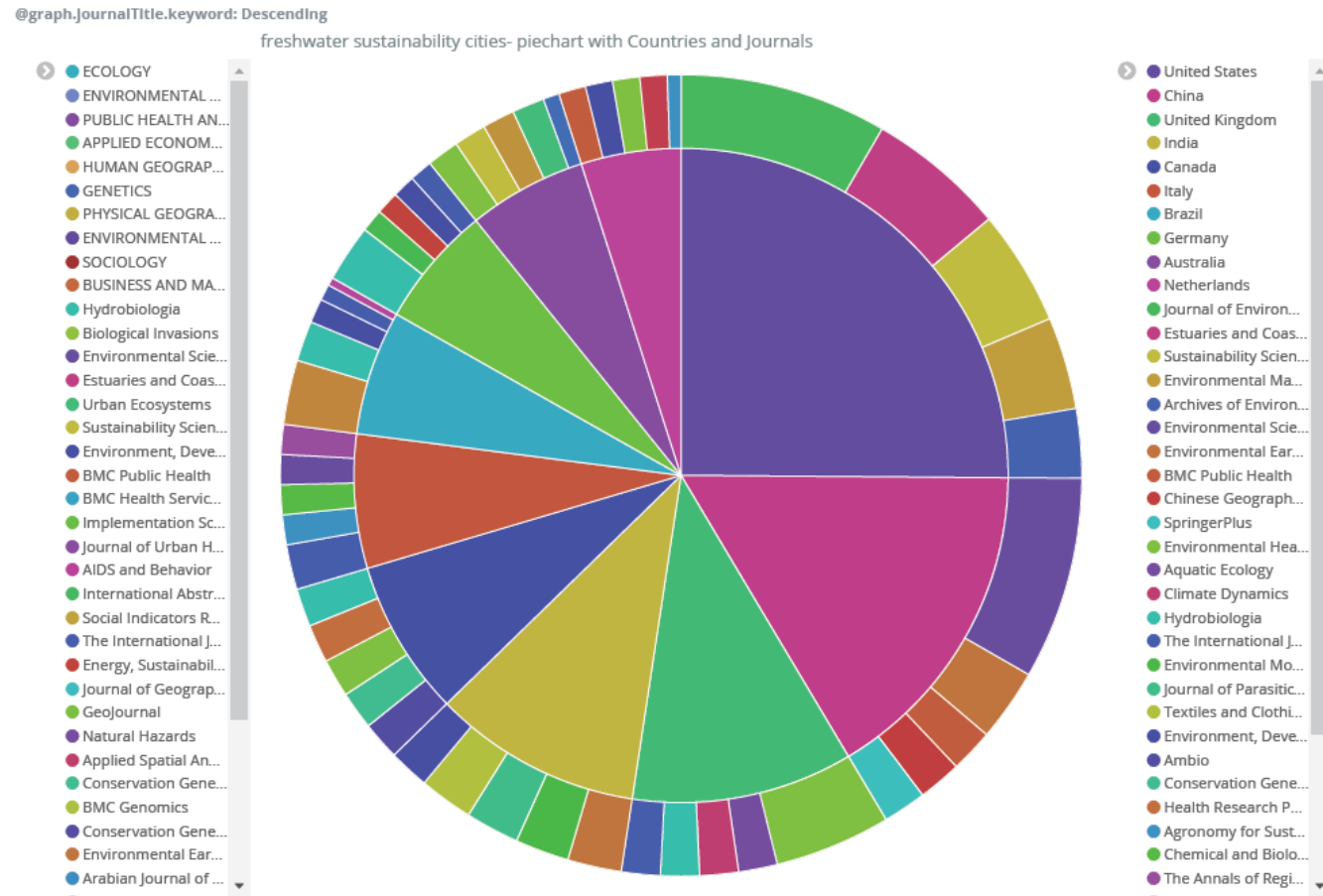
- SN SciGraph Hack Day
- **Analytics Dashboards**

# #2.2



# Springer Nature SciGraph: What are we currently working on?

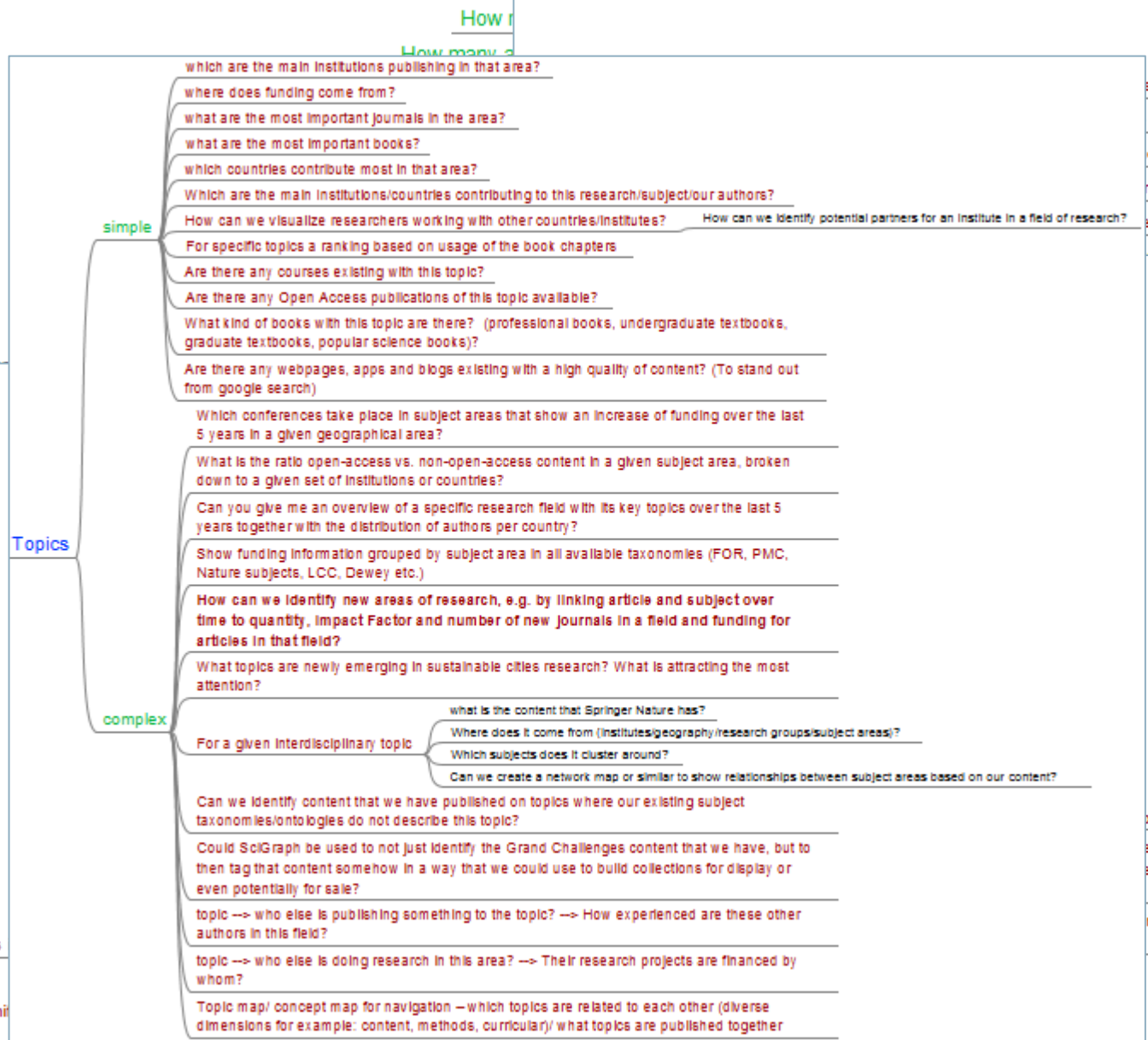
- **Analytics dashboard**
  - Make data available and queryable within Elasticsearch
  - Use Kibana to create visualizations on top of the data



## SciGraph: questions from sales, marketing and editorial

- What's the **Open Access adoption ratio**, i.e. how many articles are open access in this journal and what was the trend over the last 5 years?
- How many articles of this journal were **indexed in Web of Science / Scopus** in the last 3 years?
- What is the **funding pattern and amount** for the top 10 funded institutes in a given country?
- Which institution shows the **most growth in article output** over the past 3 years?
- Who has **interacted with Springer Nature** at a certain institution?
- Which **conferences** take place in subject areas that show an increase of funding over the last 5 years in a given geographical area?
- Can you give me an **overview of a specific research field** with its key topics over the last 5 years together with the distribution of authors per country?

i.e. "where you can see an institution's relationships with Springer Nature"



Can you make this

Can you also extract significant titles)

What was the funding pattern taxonomies (e. g. Field of Res

Is it possible to show the trends are main areas of research with development of research in fields

- supplementary books, courses)
- How large is the target group? (If possible)
- What other books were bought by customers, who bought this book
- Is the topic of the book at the cutting edge of science? And are there any new books in this area?

course and recommends which books

specific topics and the book that are recommended in ranked on the basis of usage and ideally on the basis

ing SpringerLink-Chapter-DOIs, which can be than used



Authors with ORCID  
Authors without ORCID

Article - FieldC

PUBLIC HEALTH A...  
PLANT BIOLOGY  
PHARMACOLOGY ...  
MEDICAL AND HE...  
COMPLEMENTARY...  
CLINICAL SCIENCES  
BIOLOGICAL SCIE...  
BIOCHEMISTRY A...  
SPECIALIST STUDI...  
PSYCHOLOGY AN...  
PSYCHOLOGY  
PHYSICAL CHEMIS...  
PEDIATRICS AND...  
OTHER MEDICAL A...  
ONCOLOGY AND ...  
NUTRITION AND D...  
NEUROSCIENCES  
MICROBIOLOGY  
MEDICAL PHYSIOL...  
MEDICAL MICROBI...

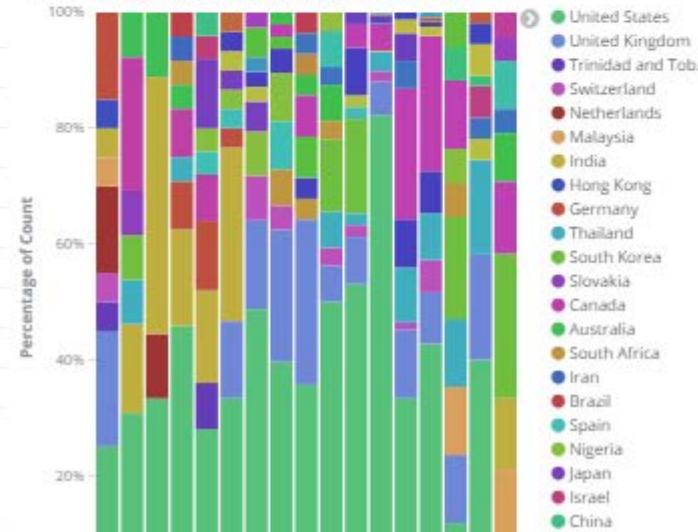


Article - by country

contribution.affiliation.countryName: Descending ⚙ Q

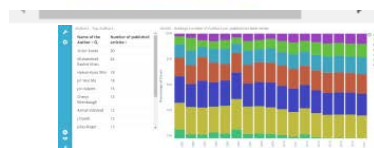
Country	Count
United States	562
China	281
South Korea	237
India	134
Malaysia	128
Germany	108
United Kingdom	95
Australia	94
Canada	92
Taiwan	90

Article - top ten countries distribution as a timeline



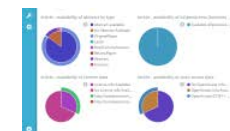
Article - table

Export: [Raw](#) [Formatted](#)



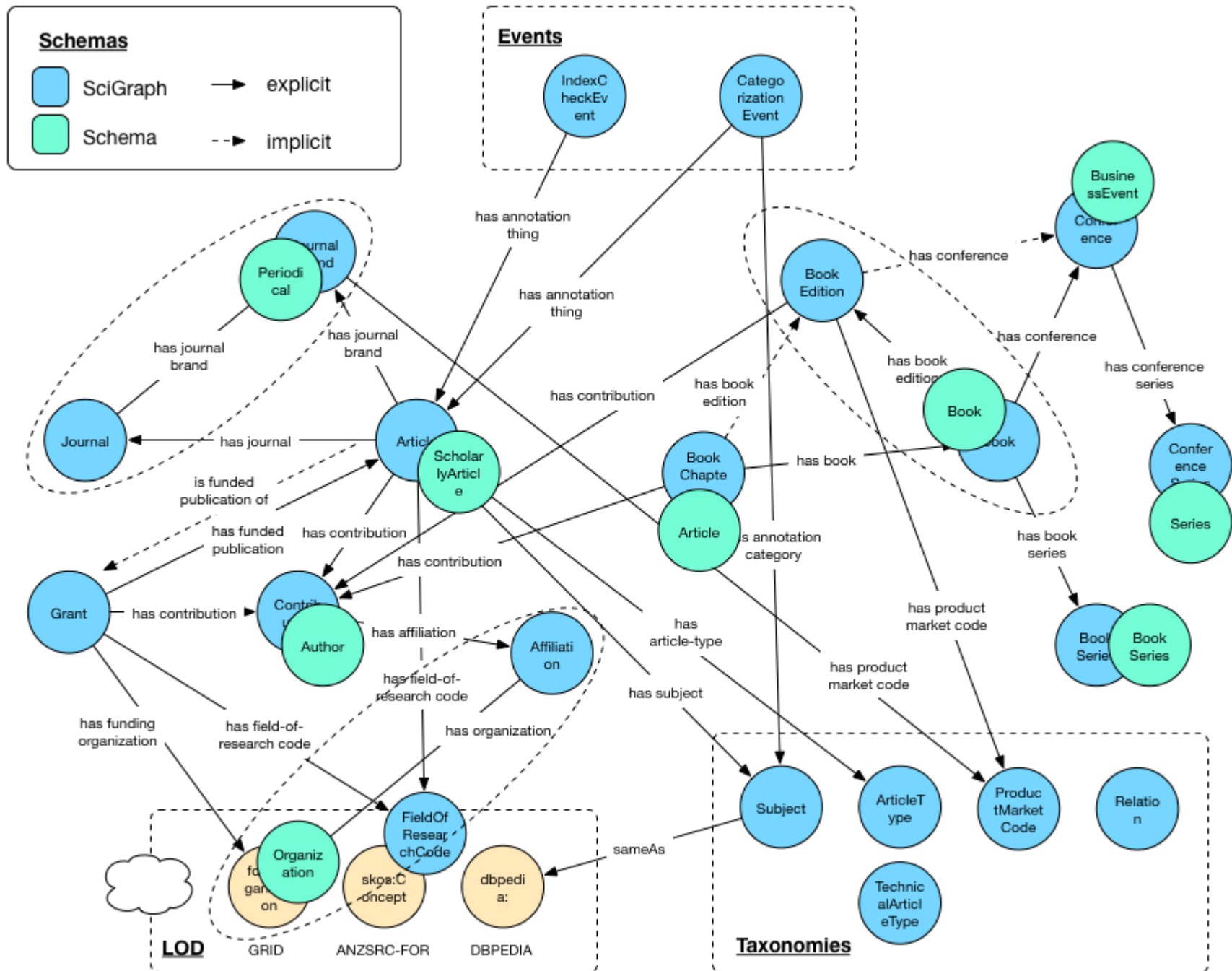
Research Funding

This article has been funded by the National Natural Science Foundation of China (grant number 81573001). The data has been deposited in the Gene Expression Omnibus (GEO) database (accession number GSE100000).

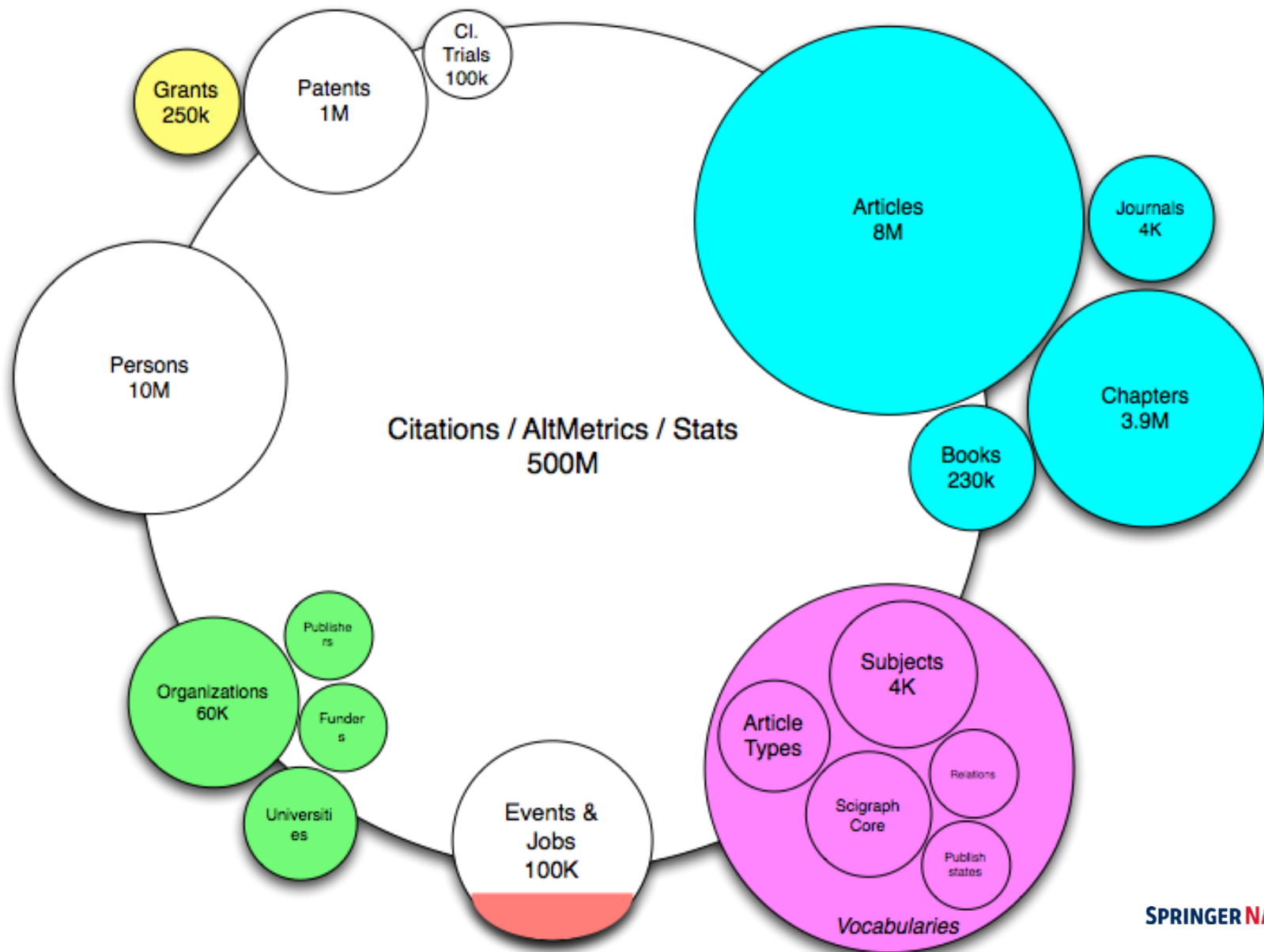


# Data: Roadmap EOY and beyond

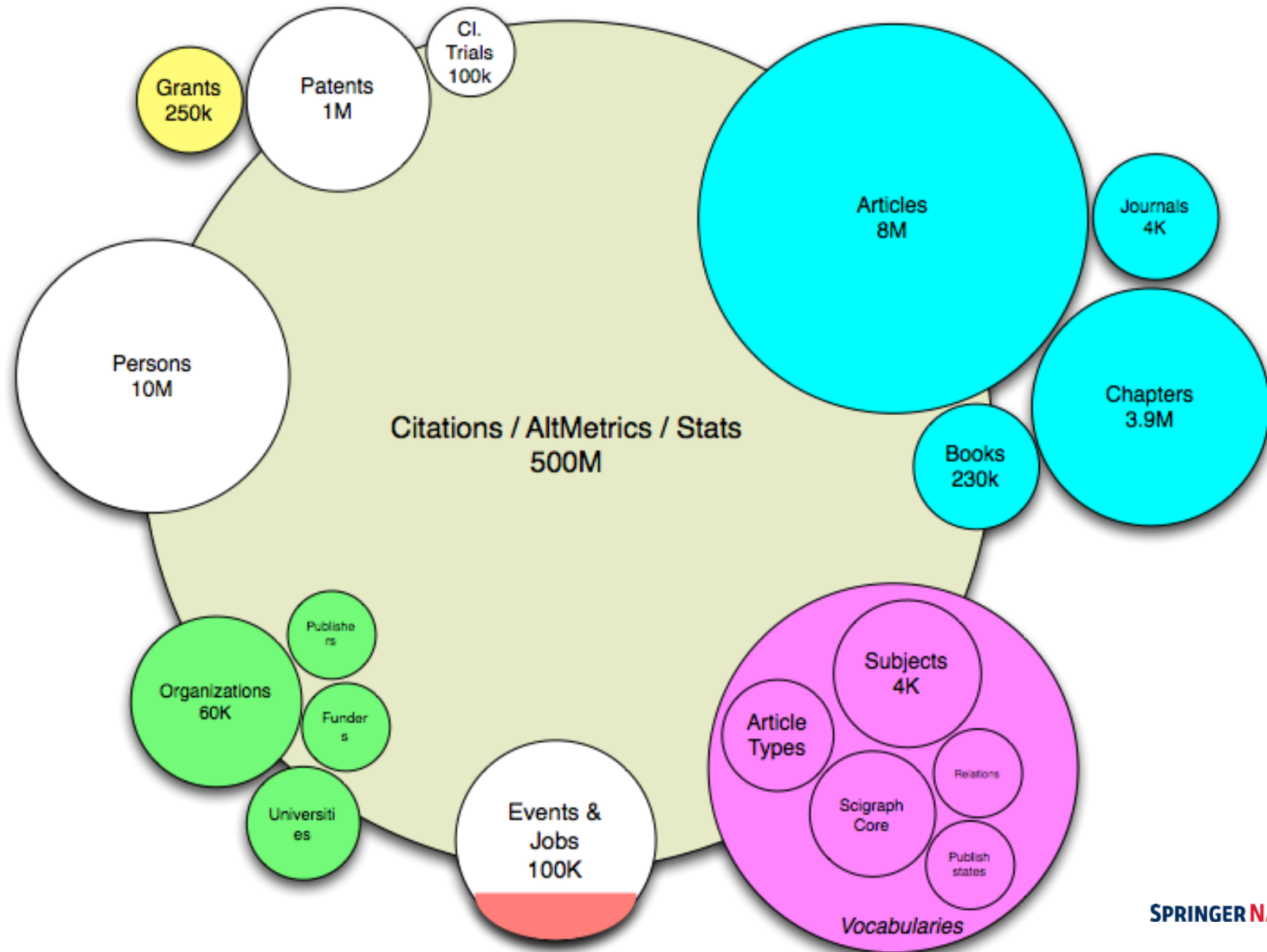
# #3



# Springer Nature SciGraph: Data Landscape - now (Q2)

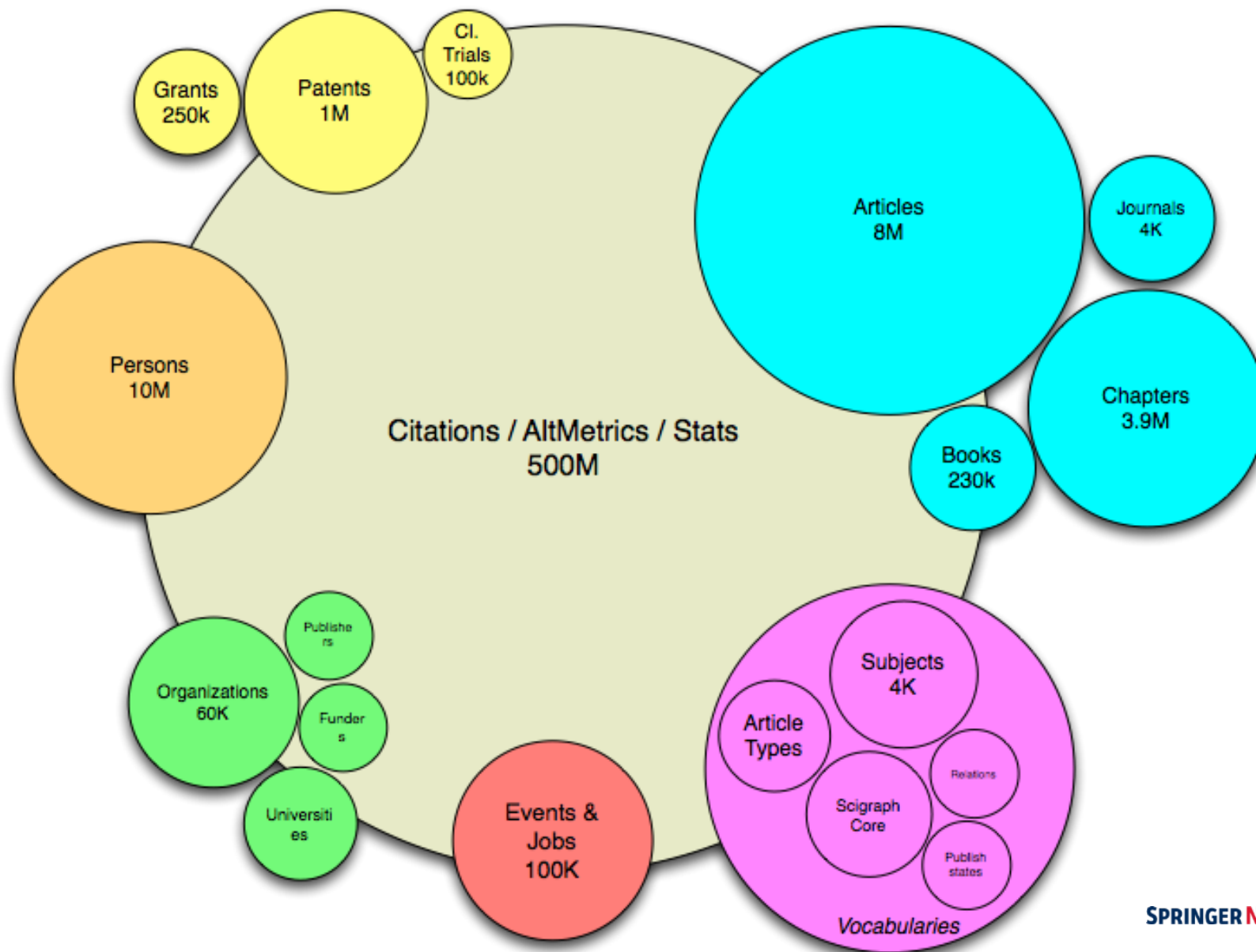


# Springer Nature SciGraph: Data Landscape - Q3





# Springer Nature SciGraph: Data Landscape - Q4



## Data Roadmap 2017

1	<b>Journals + Articles data</b>
2	<b>Institutions (GRID)</b>
3	<b>Books + Chapters data</b>
4	<b>Field of Research categories (FOR)</b>
5	<b>Conferences</b>
6	Disambiguated authors
7	Citations / References
8	<b>Research grants</b> and OA funding information
9	Download + reader numbers
10	Concepts + chemical substances
11	Patents + clinical trials
12	Links to research datasets

# Data Roadmap: General approach

## 1. Scale

1.5 billion facts  
(triples) at the  
end of 2017

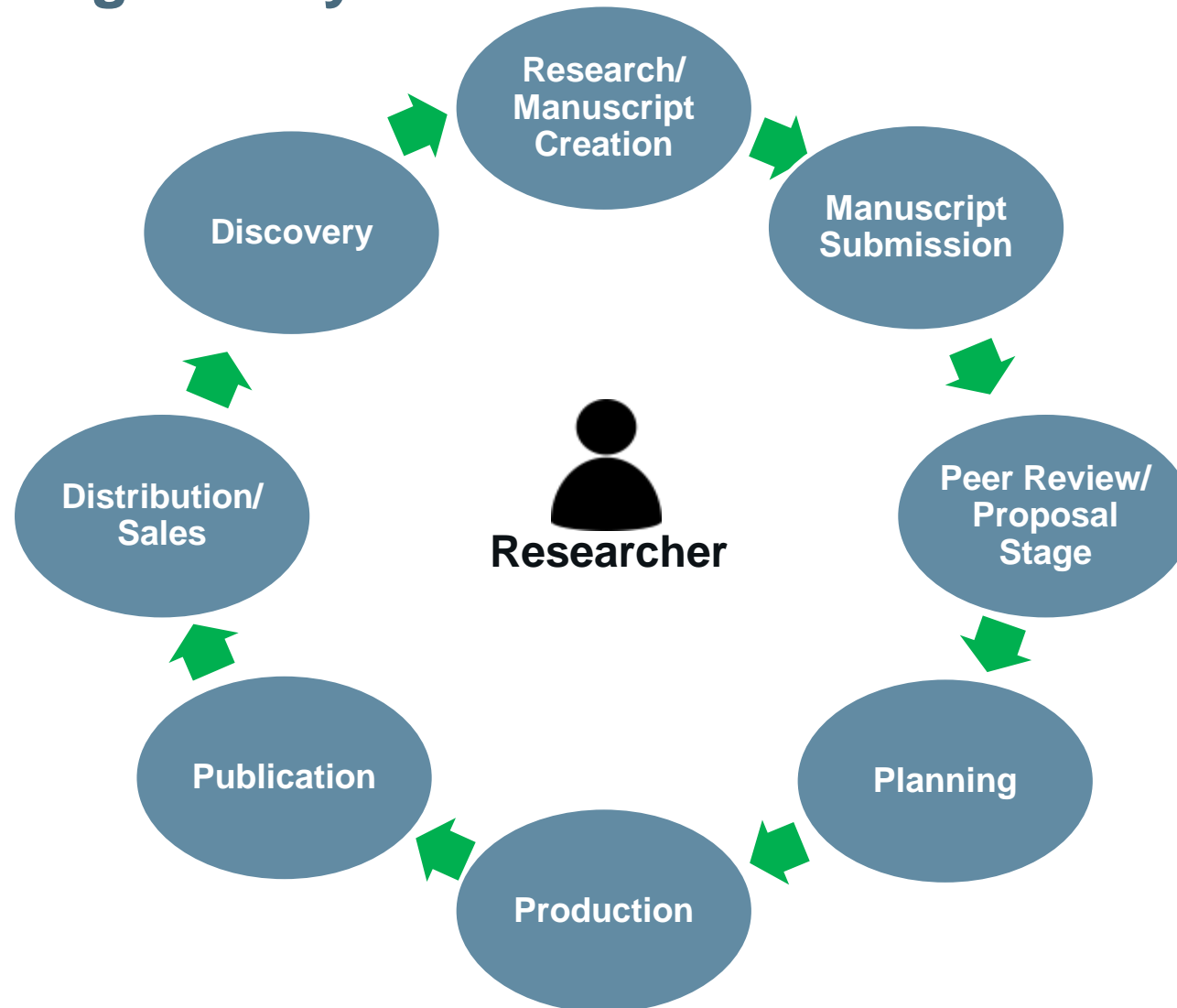
## 2. Size

Extend focus from  
Springer Nature  
publications to the  
entire research  
publishing  
landscape

## 3. Scope

Shift focus from  
published  
documents to  
submission,  
planning and  
production

# Publishing Life Cycle



ANY

QUESTIONS

?

# Thank you



Email : [markus.kaindl@springernature.com](mailto:markus.kaindl@springernature.com)

Senior Manager Semantic Data  
& Product Owner SN SciGraph

## >> How to get in touch:

- Portal  
<http://www.springernature.com/scigraph>
- E-Mail  
[scigraph@springernature.com](mailto:scigraph@springernature.com)
- Twitter  
[#scigraph](https://twitter.com/scigraph)