



BIG DATA TECHNIQUES APPLIED TO THE **STUDY AND **CHARACTERISATION** OF **SCIENTIFIC ACTIVITY** ON **SOCIAL MEDIA****

by Wenceslao Arroyo Machado
PhD Advisors Enrique Herrera Viedma & Daniel Torres Salinas

CONTENTS

1. ■ INTRODUCTION
2. ■ LITERATURE REVIEW
3. ■ OBJECTIVES
4. ■ METHODOLOGY
5. ■ SUMMARY
6. ■ DISCUSSION OF RESULTS
7. ■ CONCLUDING REMARKS

WENCESLAO ARROYO MACHADO

***Big data techniques applied to
the study and characterisation of
scientific activity on social media***

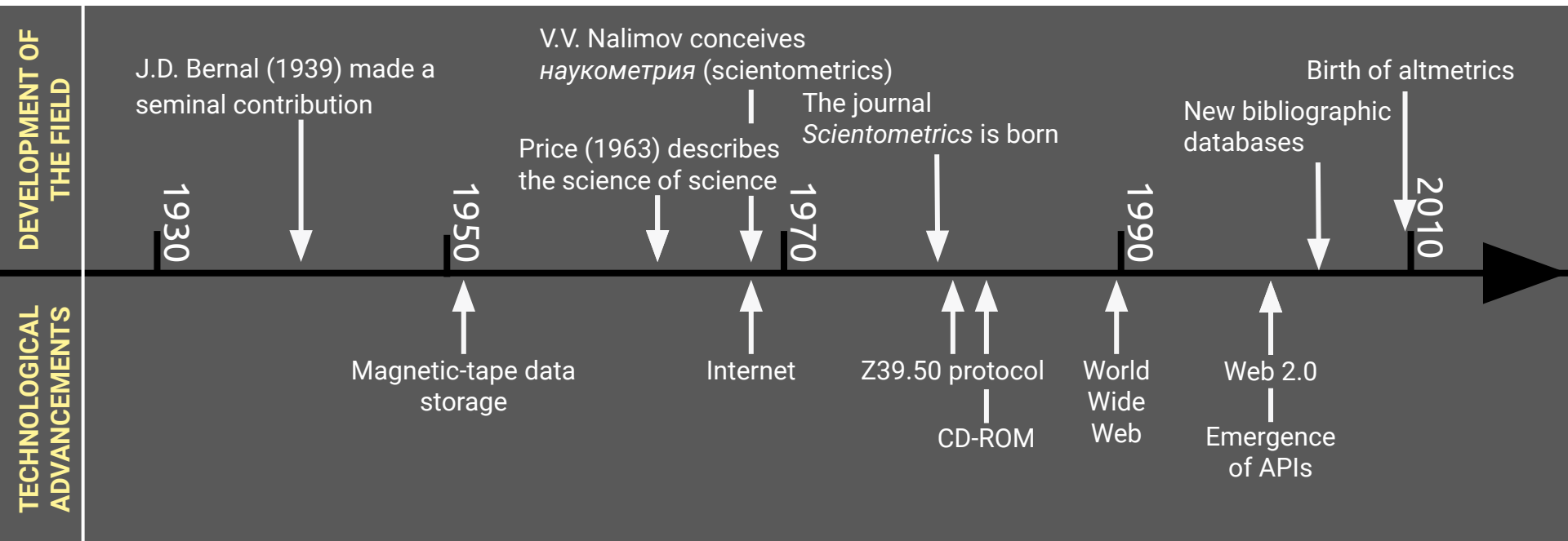


UNIVERSIDAD
DE GRANADA

THESIS BY COMPENDIUM OF PUBLICATIONS



The birth of scientometrics

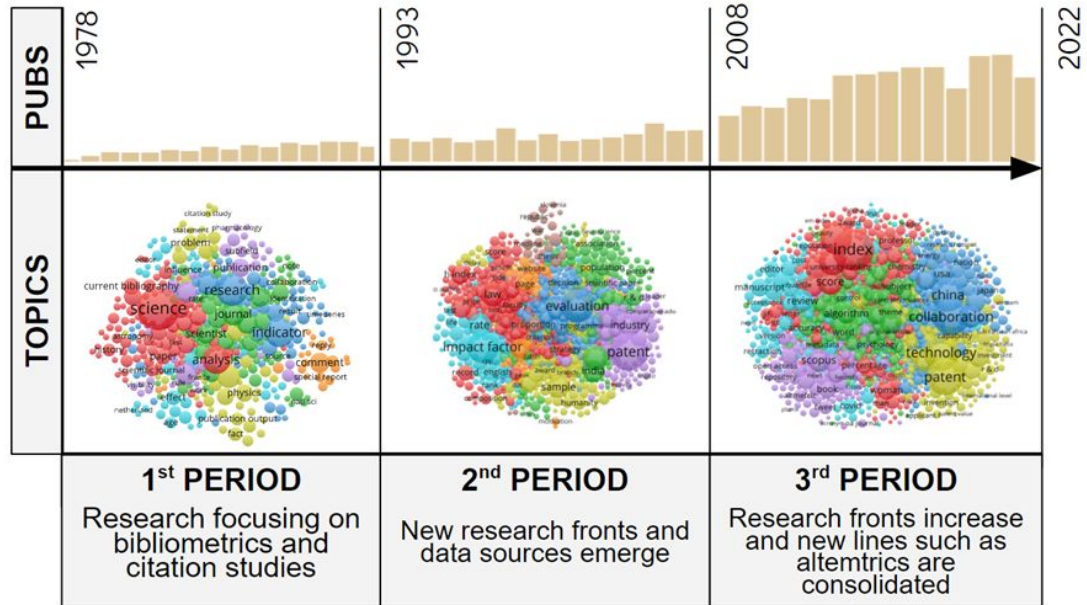
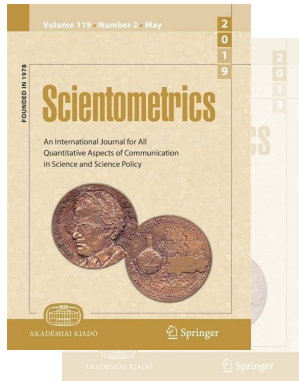


REFERENCES

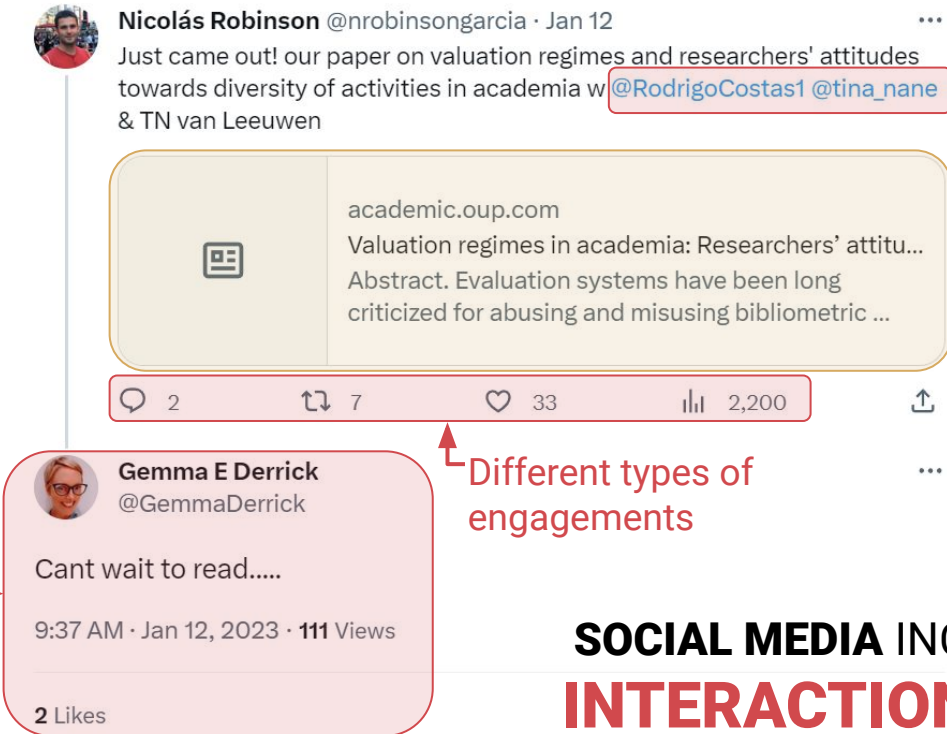
- Bernal, J. D. (1939). The social function of science. *The Social Function of Science*.
Price, D. J. D. S. (1963). *Little Science, Big Science*. Columbia University Press. <https://doi.org/10.7312/pric91844>

The development of scientometrics

The evolution of the journal **Scientometrics** reflects the evolution of the field



The birth of altmetrics



Nicolás Robinson @nrobinsongarcia · Jan 12

Just came out! our paper on valuation regimes and researchers' attitudes towards diversity of activities in academia with @RodrigoCostas1 @tina_nane & TN van Leeuwen

academic.oup.com
Valuation regimes in academia: Researchers' attitudes...
Abstract. Evaluation systems have been long criticized for abusing and misusing bibliometric ...

2 7 33 2,200

Gemma E Derrick @GemmaDerrick

Cant wait to read.....

9:37 AM · Jan 12, 2023 · 111 Views

2 Likes

Mentions to accounts

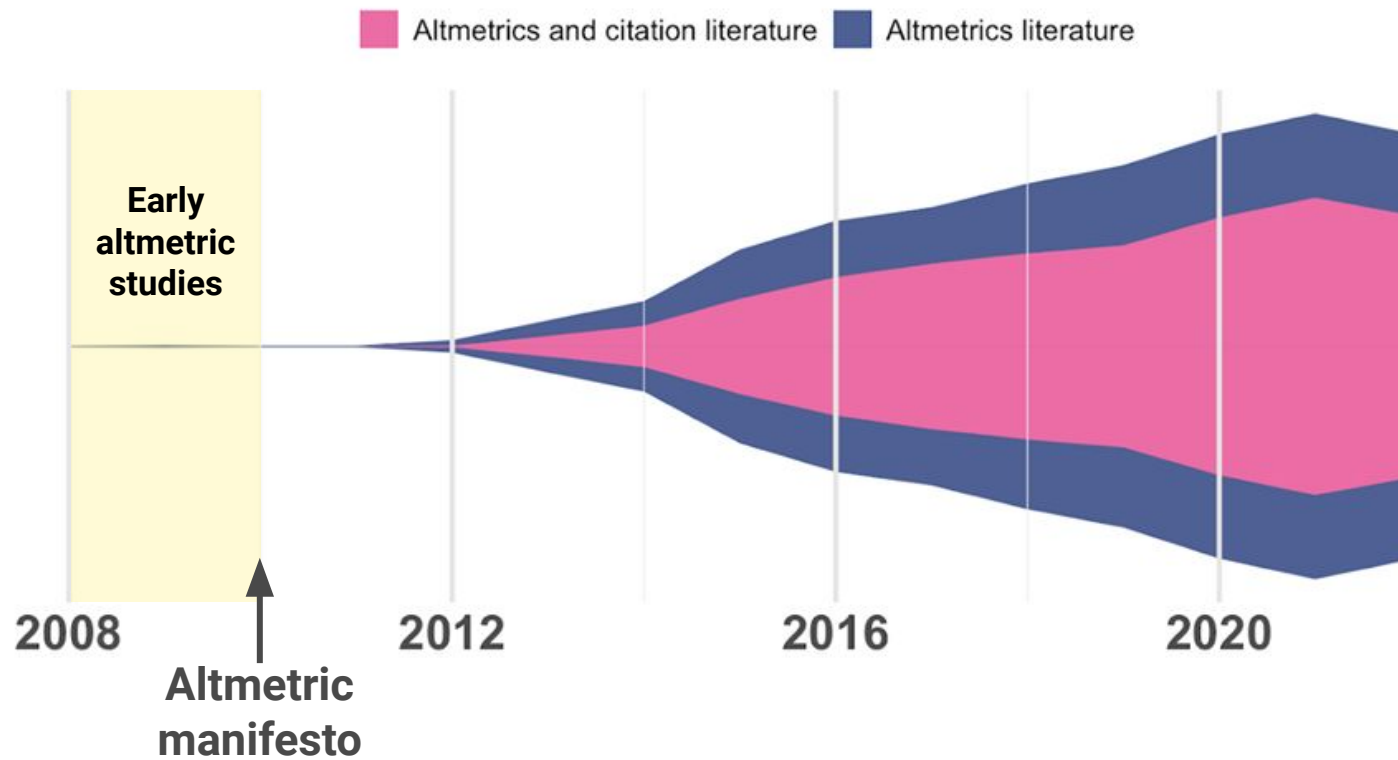
Mentions to scholarly outputs

Different types of engagements

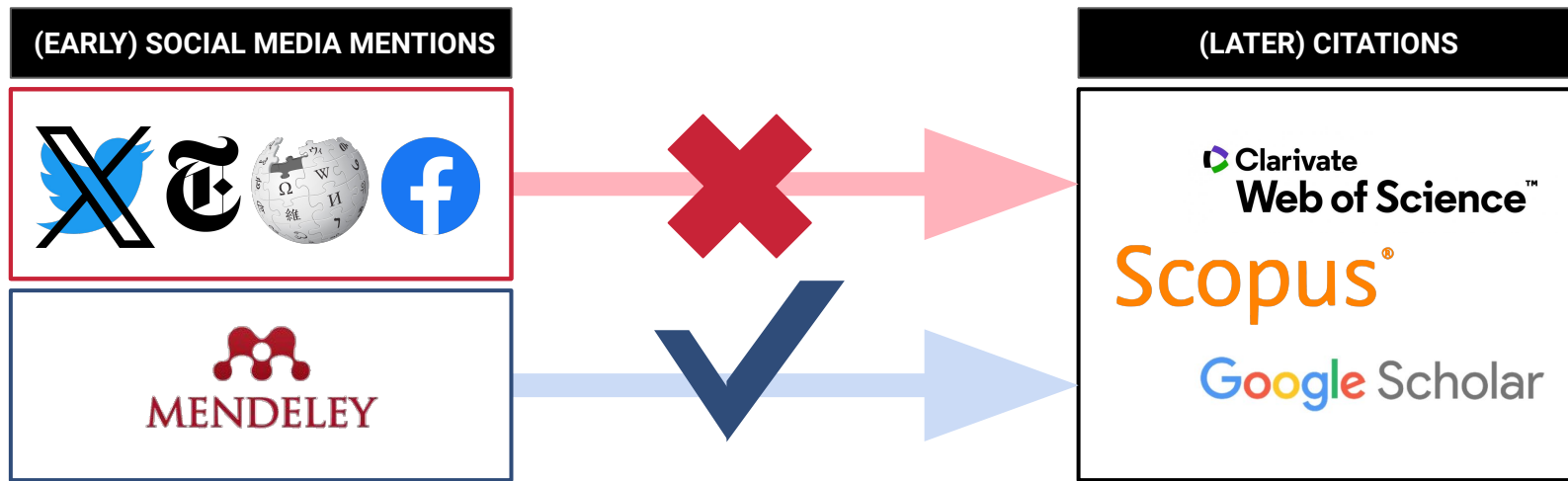
Responses

SOCIAL MEDIA INCLUDES A LARGE NUMBER OF INTERACTIONS AROUND SCIENCE

The birth of altmetrics



Towards a new generation of altmetrics



The **lack of correlation** between (almost all) altmetrics and citations gives rise to research focused on understanding this phenomenon and the context of mentions

Major
challenges
identified
in
altmetric
research

1

A large data
landscape of
social interactions
around science
that remain
unexplored

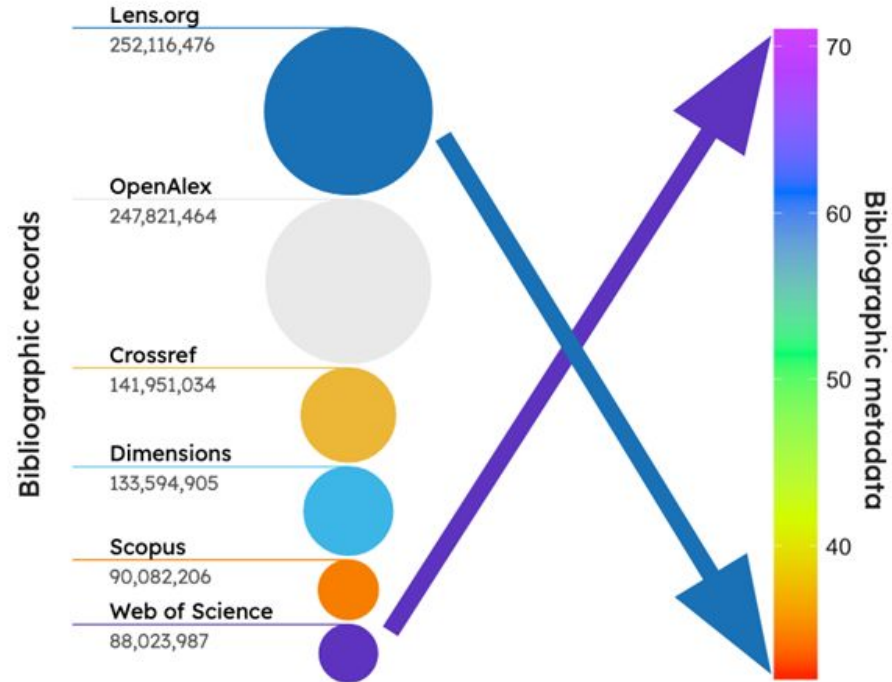
2

Lack of
scientometric
methods that
exploit social
media data

Major challenges in bibliographic data

The bibliographic universe has undergone an **avalanche** of databases, metadata and records

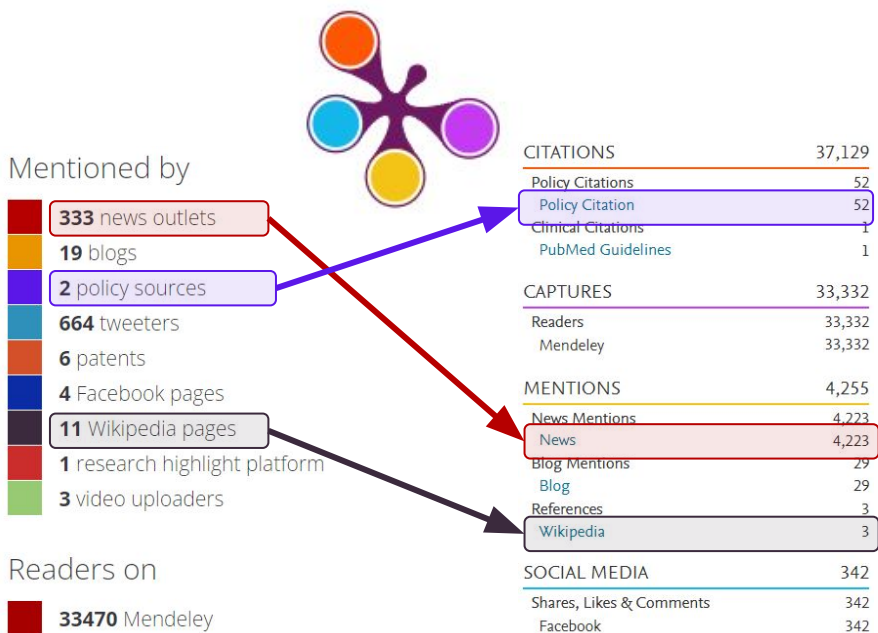
There is a noticeable gap between the **quantity** and **quality** of data



Major challenges in social media data

Global Cancer Statistics 2020

GLOBOCAN Estimates of Incidence and Mortality Worldwide for 36 Cancers in 185 Countries



In altmetrics, **data aggregators** show important differences that can determine the direction of the research

Aggregators only include the main and most **popular metrics**



MAIN CHALLENGES

Social
media data
collecting
and
processing

- 1. Dependence** on data provided by aggregators
- 2. Lack of source exploration**

Applying classic scientometrics

Scientometrics offers a wide range of methods that have demonstrated their **usefulness** in **quantitative analysis** of science through bibliographic records

Many are based on the analysis of **citation relationships** and **patterns**

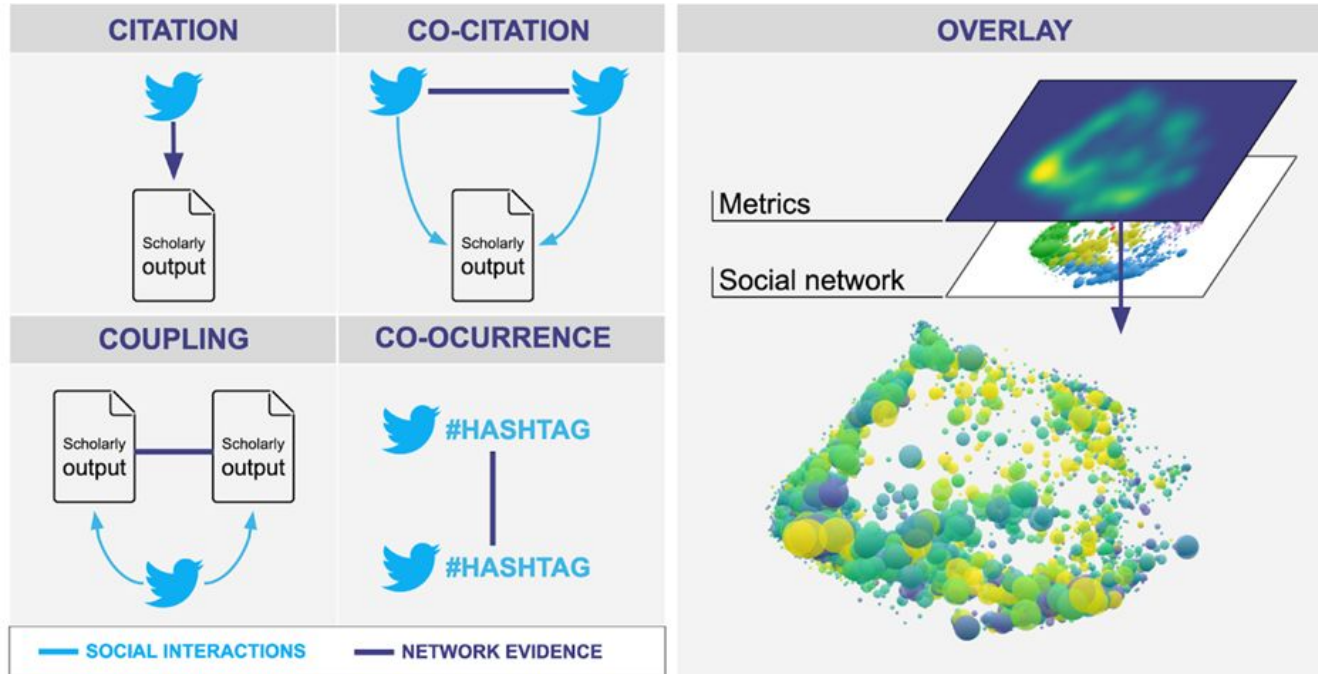
TABLE I

	Total	1921-1925	1916-1920	1911-1915	1906-1910	1901-1905	1896-1900	1891-1895	1886-1890	1881-1885	1876-1880	1871-1875
<i>Ber.</i>	686	78	30	67	115	79	64	60	56	53	44	33
<i>J. Chem. Soc.</i>	390	122	37	60	45	47	21	20	5	2	1	...
<i>Ann.</i>	278	26	8	37	23	23	22	21	19	18	13	...
<i>Z. physik. Chem.</i>	191	53	6	21	39	29	19	28	16	6
<i>Compt. rend.</i>	126	26	3	23	15	23	15	21	7	9	8	...
<i>J. Phys. Chem.</i>	93	42	13	13	5	1	1
<i>Ann. Physik</i>	93	18	4	28	13	6	0	0	6	5	2	...
<i>J. Biol. Chem.</i>	80	41	16	14	7
<i>Am. Chem. J.</i>	70	9	21	20	14	8	4	2	1	...
<i>Z. anorg. Chem.</i>	68	21	11	5	8	11	6	2
<i>Ann. Chim.</i>	68	5	0	6	9	7	3	5	1	8	4	2
<i>Bull. Soc. Chim.</i>	60	16	3	4	7	10	4	4	3	4	2	1
<i>Proc. Roy. Soc.</i>	55	30	5	4	8	5	1	0	1
<i>J. Ind. Eng. Chem.</i>	53	33	10	5	1
<i>Z. Phys.</i>	51	41	5
<i>Monatsch.</i>	51	2	1	21	5	9	3	2	5	3
<i>J. prakt. Chem.</i>	50	6	1	2	2	6	3	12	6	6	2	2
<i>Phil. Mag.</i>	49	17	14	4	2	3	3	1	1	0	0	1
<i>Gazz. chim. ital.</i>	44	10	6	2	6	4	8	4	3	0	1	...
<i>Phys. Rev.</i>	44	23	8	3	5	4
<i>Physik. Zeit.</i>	41	26	0	7	3
<i>Z. Elektrochem.</i>	37	11	13	4	4	4	1
<i>Biochem. Z.</i>	37	18	2	9	10
<i>Rec. trav. chim.</i>	36	14	5	2	2	2	5	4	1	1
SCIENCE	27	22	3
<i>Trans. Far. Soc.</i>	24	18	0	1	0	1
<i>Proc. Nat'l Acad.</i>	22	19	0
<i>Nature</i>	21	13	5	1

The abbreviations used above and in the tables to follow are those accepted by *Chemical Abstracts* and may be found in their list of periodicals abstracted, issued October 20, 1926.

Gross, P. L. K., & Gross, E. M. (1927). College libraries and chemical education. *Science*, 66(1713), 385-389.
<https://doi.org/10.1126/science.66.1713.385>

Applying classic scientometrics



Many of these classic methods have been successfully **adapted** to Twitter



MAIN CHALLENGES

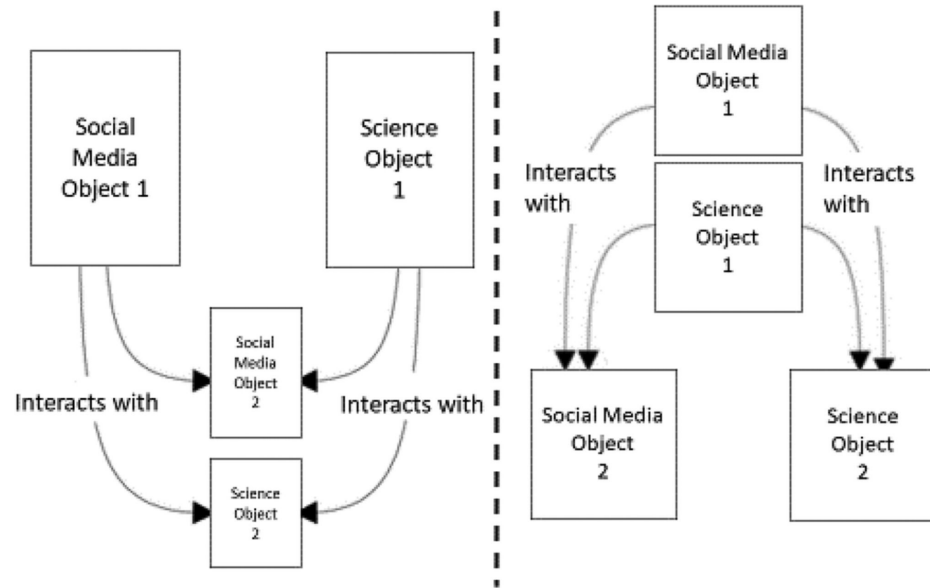
Adapting
**standard
scientometrics**
methods in
altmetrics
research

- 1.** Absence of approaches applying these methods to sources other than **Twitter**
- 2.** Lack of **comparisons** between the social and scientific perspective of scientific knowledge

Towards New Horizons

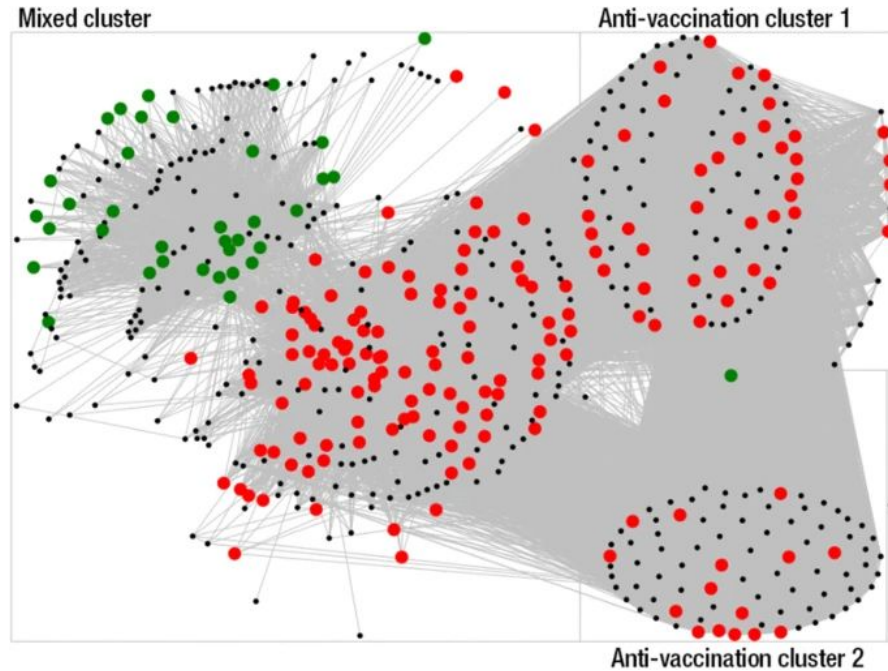
Beyond classic methods, **new methods** are needed that truly leverage the rich context and environment in which interactions around science occur

Evidence from different sources are also not combined



Costas, R., De Rijcke, S., & Marres, N. (2021). "Heterogeneous couplings": Operationalizing network perspectives to study science-society interactions through social media metrics. *Journal of the Association for Information Science and Technology*, 72(5), 595-610. <https://doi.org/10.1002/asi.24427>

Towards New Horizons



- *pro-science*
- *anti-vaccination*

Van Schalkwyk, F., Dudek, J., & Costas, R. (2020). Communities of shared interests and cognitive bridges: The case of the anti-vaccination movement on Twitter. *Scientometrics*, 125(2), 1499-1516. <https://doi.org/10.1007/s11192-020-03551-0>

There is great potential in studying engagement around academic objects to better understand **science-society interactions**



MAIN CHALLENGES

Many
opportunities
and
possibilities

- 1.** Many unexplored opportunities in science maps and social network analysis
- 2.** Updated research required due to the ever-changing social media landscape

Main objectives

OBJECTIVE

1

To explore challenges in processing large bibliographic and social media **data**, with a focus on combining them for altmetric studies

MAIN EXPECTED RESULT

Providing **tools** and **curated datasets** for altmetric research

OBJECTIVE

2

To adapt **scientometric methods** for social media, aiming to create science maps from Wikipedia that reflect social attention

MAIN EXPECTED RESULT

Map the science structure through the lens of Wikipedia

OBJECTIVE

3

To develop novel methodologies for scientific mapping by combining **social and semantic data** from diverse sources

MAIN EXPECTED RESULT

Implementation of **innovative methodologies** that integrate interactions and interests

Traditional scientific method

OBSERVATION 1

Exploration and mapping of the interaction between science and society on social media

HYPOTHESIS FORMULATION 2

Adaptation of traditional scientometric methods to fit social media environments

OBSERVATION GATHERING 3

Using results from applied methods to social media and validating with indicators from social network analysis

CONTRASTING THE HYPOTHESIS 4

Comparison of the acquired results with those published by other novel related proposals

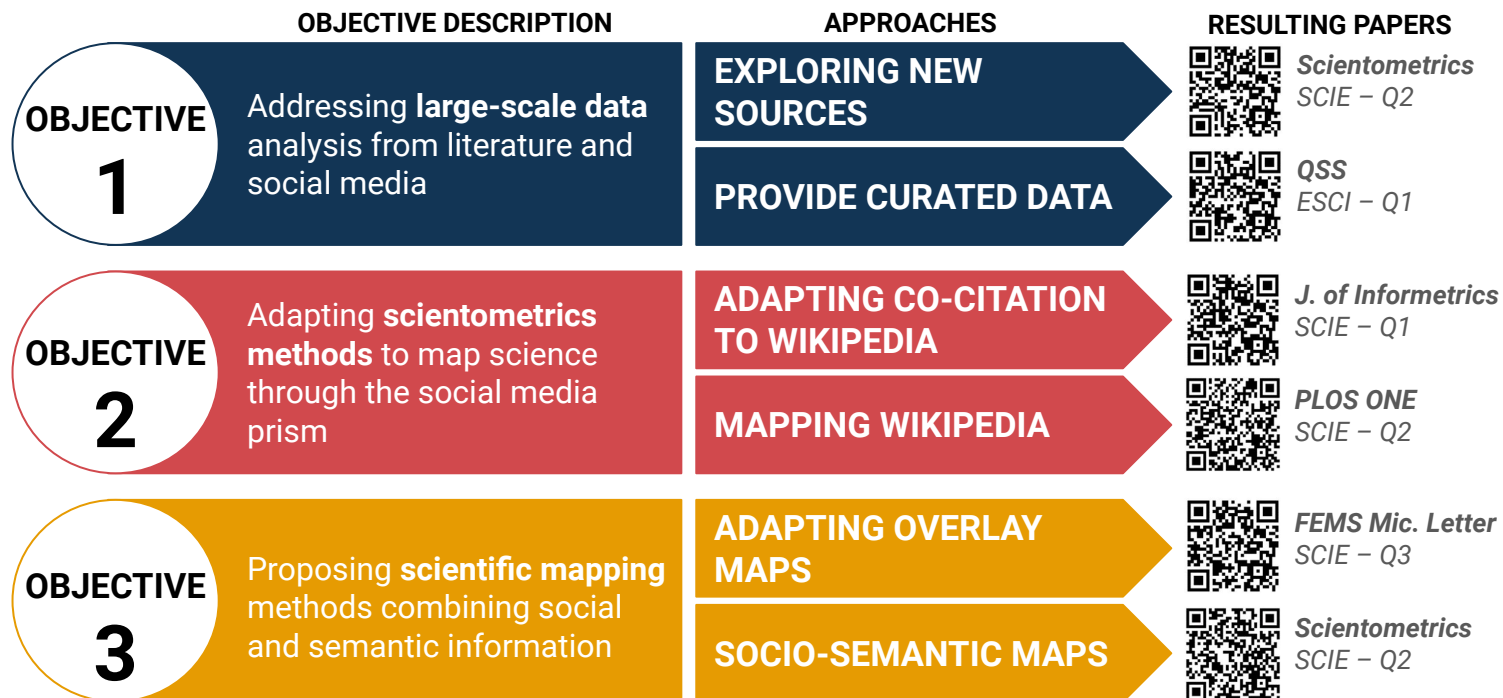
HYPOTHESIS VALIDATION OR REFUSAL 5

Validation of the hypothesis through the conducted experiments and results

SCIENTIFIC THESIS 6

Extraction, redaction and acceptance of the conclusions

Main contributions



Exploring WorldCat Identities

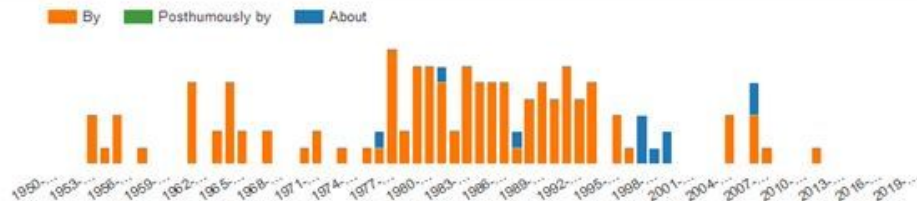


Garfield, Eugene

Overview

Works:	124 works in 383 publications in 4 languages and 2,960 library holdings
Genres:	Dictionaries Reference works Bibliography Periodicals Terminology History
Roles:	Author, Editor, Other
Classifications:	PG2640, 491.7321

Publication Timeline



Academic writing American Chemical Society American Chemical Society--Chemical Abstracts Service Bibliographical citations Bibliometrics
 Canada Chemical engineering Chemistry Chemistry--Abstracting and indexing Chemistry--Notation Citation indexes
 Columbia University Communication in science Cooper, Lewis A. English language English language
 Garfield, Eugene Genetics--Periodicals History--Research--Data processing Humanities--Abstracting and indexing
 Indexing Information science Information scientists Information storage and retrieval systems--Chemistry
 Information storage and retrieval systems--Research Institute for Scientific Information Librarians
 Library science Library science--Periodicals National Federation of Science Abstracting and Indexing Services
 Price, Derek J. de Sola--Derek John de Sola Russian language Science Science--Abstracting and indexing
 Science--Historiography Scientism Scientists Social sciences--Periodicals Technology--Abstracting and indexing
 United States Universities and colleges--Faculty

CONTEXT

OCLC conducted an experimental project, **WorldCat Identities**, in which it generated **author profiles** from the WorldCat catalog with various **indicators**

OBJECTIVES

To **explore** WorldCat Identities as an information source and conduct a **case study** with scientometrics authors

Thesis
objective

11
Contribution



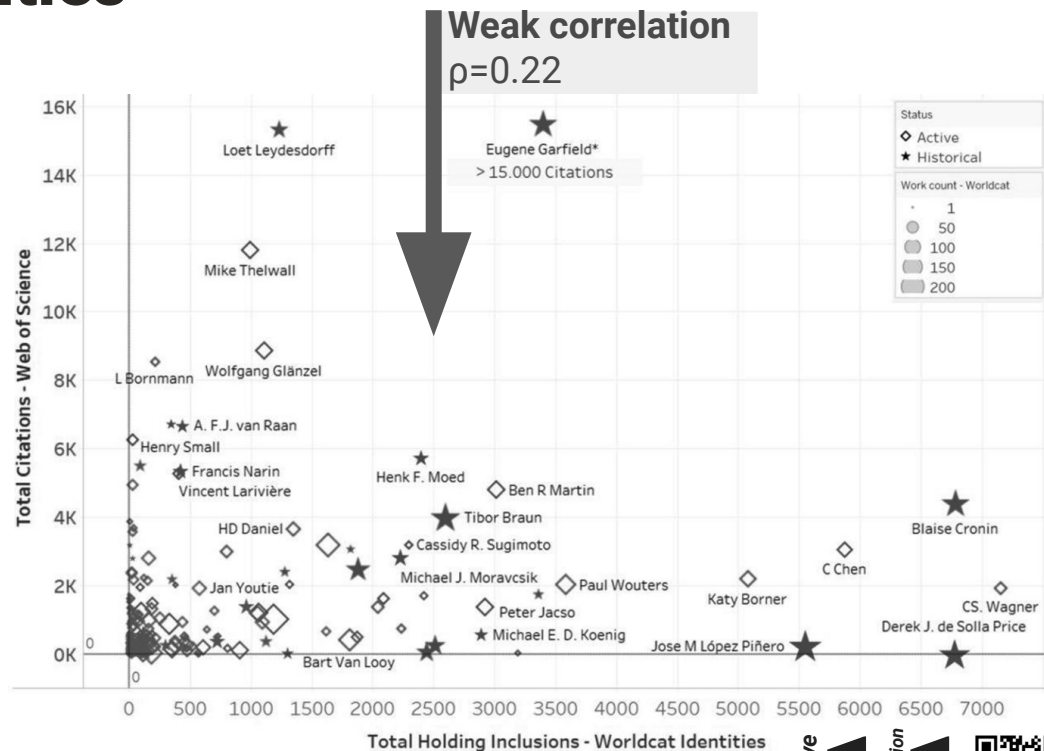
Exploring WorldCat Identities

A python™ **package** was built to retrieve author profiles and their metadata. After retrieving data on scientometricians, we compared the **library holdings** with the citations

METHODS

Despite the **disambiguation issues**, it proved useful as an altmetric tool, offering a **new dimension of influence** distinct from that of citations

FINDINGS



Thesis objective **11** Contribution **11**

The analytical possibilities of Wikipedia

This article is about a health issue (points to the title 'COVID-19')

Open discussions (points to the 'Talk' tab)

Content links (points to the 'SARS-CoV-2' link in the text)

Article length (points to the main text area)

References (points to the 'References' section)

Cited references (points to a reference entry)

Daily views (points to the bar chart showing daily views)

The screenshot shows the Wikipedia article for COVID-19. Red arrows and boxes highlight specific elements: the title 'COVID-19', the 'Talk' tab, a link to 'SARS-CoV-2' in the text, the main text area, the 'References' section, a specific reference, and a bar chart showing daily views from February 2020 to August 2023. The bar chart shows a significant peak in views in early 2020, followed by a decline and then a secondary, smaller peak in late 2022.

CONTEXT

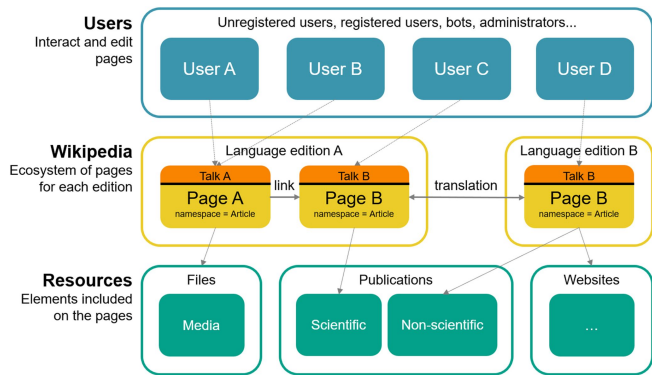
Wikipedia conceals an intricate, **unexplored infometric ecosystem** with vast potential to capture various social interactions

OBJECTIVES

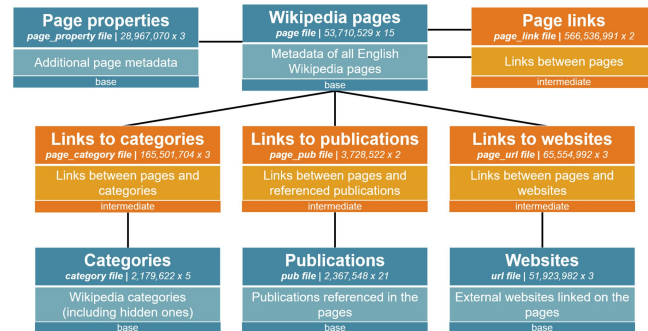
To establish a **framework** for Wikipedia, create a large open knowledge graph based on it, and conduct a **descriptive quantitative study** of Wikipedia

Thesis objective **1** Contribution **2**

The analytical possibilities of Wikipedia



Wikimedia Downloads
+
Wikimedia API
+
zenodo



Wikipedia Knowledge Graph, dataset and description free at: [10.5281/zenodo.6346899](https://zenodo.org/record/105281)

Identification of elements involved in Wikipedia activity

Data processing



Developing of Open Wikipedia Knowledge Graph

Thesis objective 12 Contribution



The analytical possibilities of Wikipedia

Using **data science** methods, heterogeneous data from the English Wikipedia was retrieved and processed to construct a **knowledge graph**

METHODS

Wikipedia has valuable metrics that capture different dimensions of **social attention** that contribute to contextualising how science is consumed by society

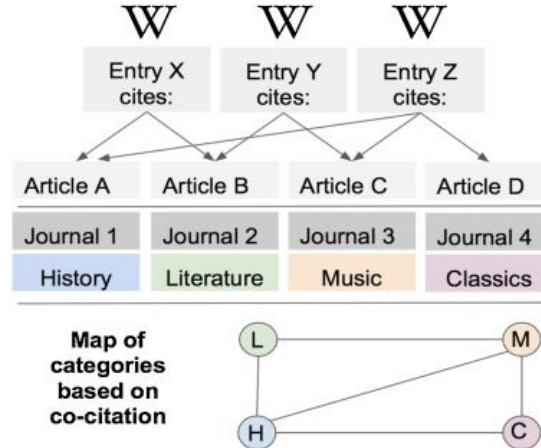
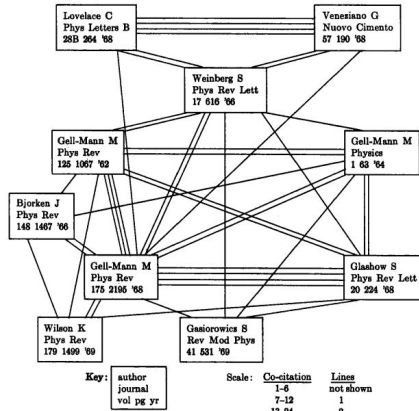
FINDINGS

	All articles	Featured articles	Featured lists	A	Good	B	C	List	Start	Stub
<i>N. of articles</i> → <i>Wiki Metrics</i> ↓	6,328,134	5945	3816	958	34,004	109,019	394,065	253,066	1,818,356	3,079,778
Editors	48.38	516.93	179.13	176.80	275.71	297.62	165.36	56.27	63.13	22.85
Edits	101.92	1491.35	593.61	564.91	724.13	705.41	369.89	159.80	129.52	40.23
Linked	80.53	725.25	175.84	202.01	330.18	417.00	234.08	107.34	93.03	55.70
Links	87.77	329.68	270.16	236.56	224.88	233.87	164.23	174.78	101.28	69.90
Age	9.59	14.33	11.52	12.74	12.06	12.47	10.92	9.13	10.45	9.20
Length	7844.68	61,248	51,549	43,329	39,444	35,009	21,676	18,202	10,033	3748
Talkers	5.38	66.17	16.62	27.90	29.64	28.16	15.03	4.98	6.56	3.64
Talks	9.19	258.40	42.36	92.21	88.56	88.35	35.32	9.07	9.69	4.32
Views	3345.07	64,801	26,685	16,011	29,229	30,359	15,829	3777	4094	710
References	4.6	53.95	55.49	31.76	38.87	26.51	15.40	9.20	5.79	1.84
Pub. Ref.	0.59	14.27	2.34	8.51	5.83	4.77	2.37	0.53	0.69	0.22
URLs	10.33	58.03	67.32	33.32	46.10	40.31	25.95	22.82	12.90	6.09

The **quality** of the articles is linked to metrics such as **page views** or **references**

Adapting co-citations to Wikipedia

Proof of concept



CONTEXT

Wikipedia articles engage with literature in a similar way to scientific publications, providing a good environment for adapting classic scientometric techniques

OBJECTIVES

To transfer co-citation methodology to Wikipedia and test the method by mapping the structure of the Humanities

From paper co-citations

to Wikipedia co-citations

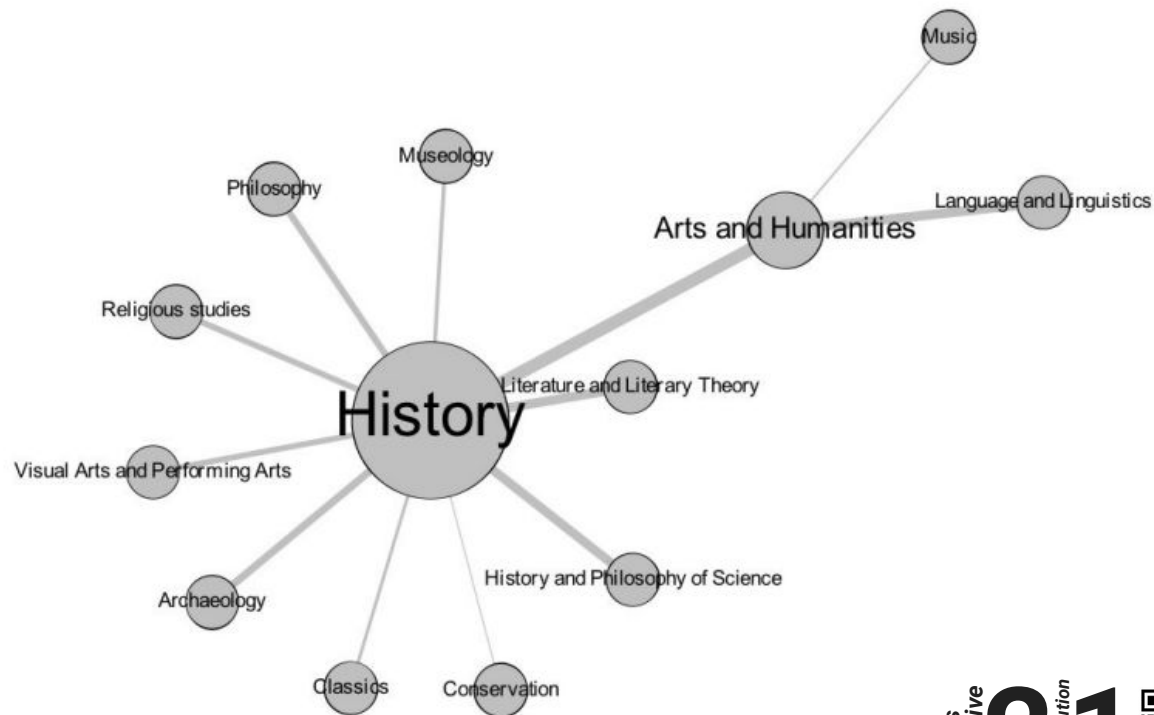
Small, H. (1973). Co-citation in the scientific literature: A new measure of the relationship between two documents. *Journal of the American Society for Information Science*, 24(4), 265-269. <https://doi.org/10.1002/asi.4630240406>

Adapting co-citations to Wikipedia

Proof of concept

Pathfinder networks (PFNETs)

Due to the characteristics of these networks, in which there is a high degree of connectivity between all the nodes in the network, it was decided to apply the **Pathfinder algorithm** to eliminate weak links



Adapting co-citations to Wikipedia

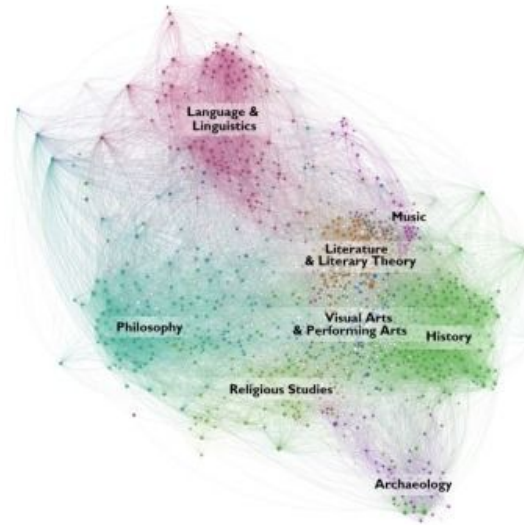
Proof of concept

From Wikipedia references, **relationships between publications** are established that can be **aggregated** by journals and scientific disciplines

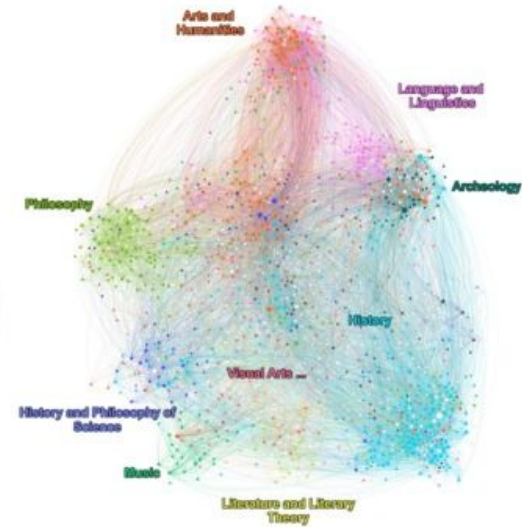
METHODS

The case study **validates** the adaptation of the co-citation method to Wikipedia and demonstrates its relevance by **exposing differences** between the academic and social realms

FINDINGS



Scopus perspective

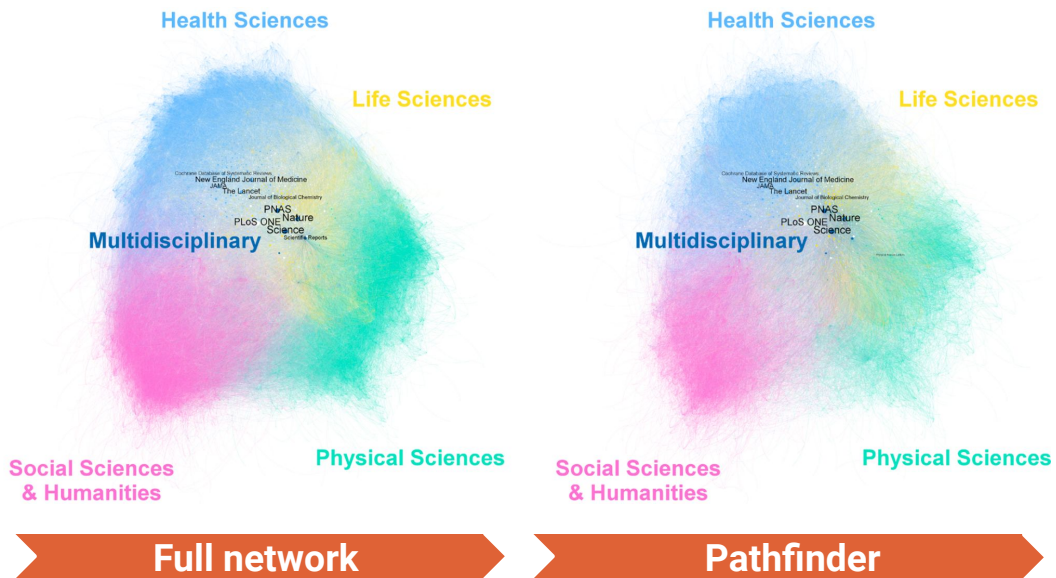


Wikipedia perspective

Richardson, M. (2013). Mapping the multidisciplinary of the Arts & Humanities. *Research Trends*, 1(32), 5. <https://www.researchtrends.com/cgi/viewcontent.cgi?article=1145&context=researchtrends>

Adapting co-citations to Wikipedia

Large-scale mapping



CONTEXT

After validating the co-citation method on Wikipedia, the ambition was to apply it on a **large scale** and to make a comprehensive analysis of the science

OBJECTIVES

Mapping the structure of science and offering a general **portrait** of science through the English Wikipedia

Adapting co-citations to Wikipedia

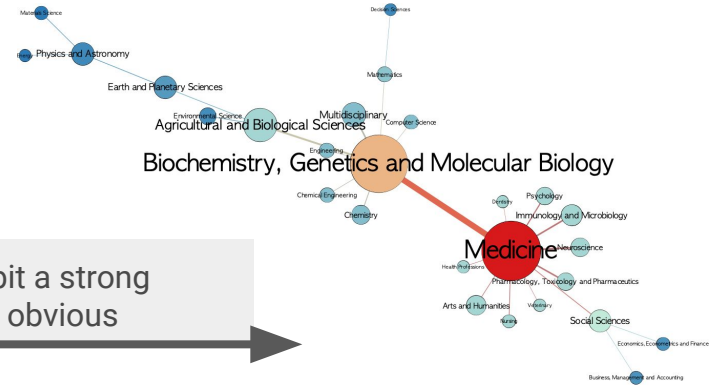
Large-scale mapping

A **co-citation** methodology was adapted to generate **Pathfinder networks (PFNET)** from all references in the English Wikipedia articles

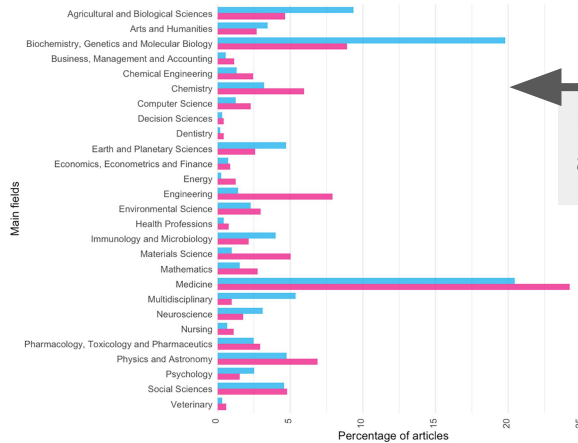
METHODS

Regarding the scientific realm, **discrepancies** were detected in the attention to certain disciplines and a shared interest in **high-impact** publications

FINDINGS



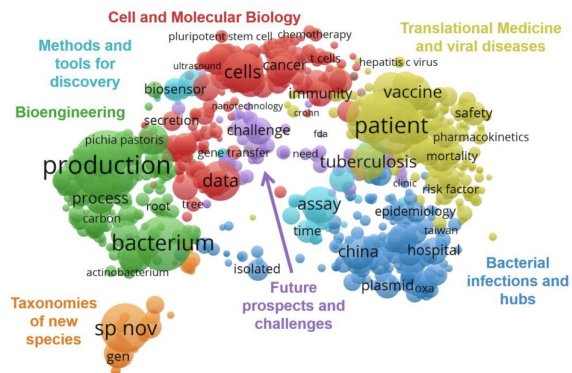
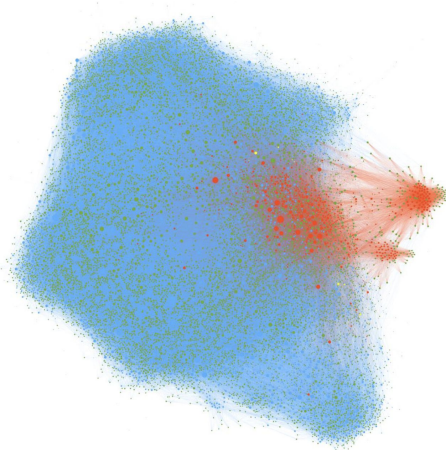
Some disciplines exhibit a strong relationship that is not obvious



Disciplines such as Biochemistry are overrepresented in Wikipedia

Developing new methodologies

Shared interests



Social mentions

from Twitter, news outlets, and policy documents



Thematic landscape

from publications' titles

CONTEXT

Although the topics that capture the attention of a discipline in social media have been studied, potential differences in interests between different media have not been delimited

OBJECTIVES

To visualize **key social interest topics** in Microbiology identified via altmetric data



Developing new methodologies

Shared interests

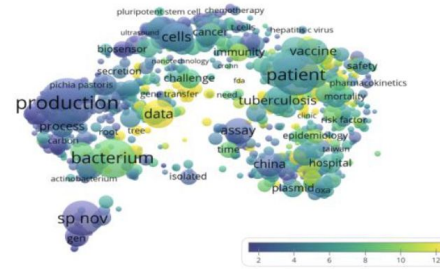
The **overlay maps** have been adapted to the altmetrics to identify the **topics of interest** of each source

METHODS

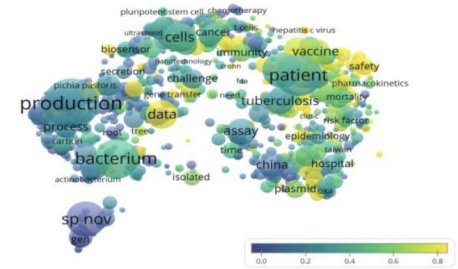
Not only have **differences in interests** among the various studied social media been highlighted, but also **peaks of attention** over time

FINDINGS

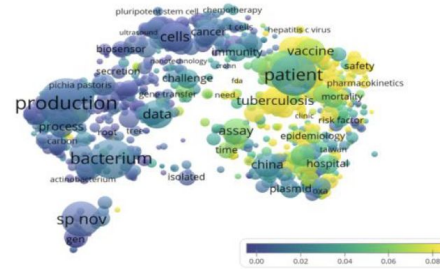
A



C



B



LEGEND

Overlay term maps for **A** Twitter mentions, **B** news media mentions and **C** policy briefs citations to microbial literature based on Figure 4.

 Highly mentioned terms by altmetric source.

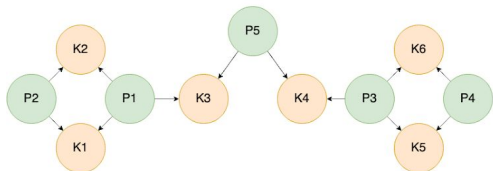
 Terms not mentioned by altmetric source.

Developing new methodologies

Socio-semantic networks

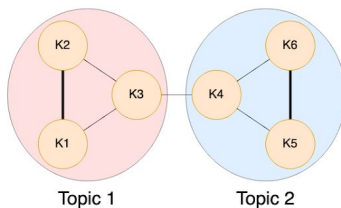
1. Publication-keyword network

Network of scientific publications (Pn) and their Web of Science author keywords (Kn).



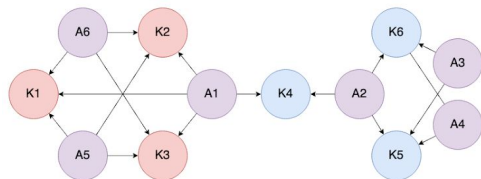
2. Semantic map

Network of author keyword co-occurrence. Topics are identified by community detection.



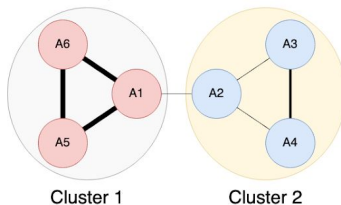
3. Actor-keyword network

Network of social actors (An) that mention keywords, based on the papers mentioned. Keywords belong to one of the topics identified in the semantic map.



4. Socio-semantic network

Network of co-occurrence of actors combined with the semantic map. Each actor belongs to a topic based on its keyword mention. Clusters of actors are identified by community detection.



CONTEXT


Altmetrics allow exploring and profiling social actors who discuss and share scientific literature, but it is a challenge to **identify and characterize communities**

OBJECTIVES

To develop and validate a new method for **profiling social media users** based on their interest on research topics

Developing new methodologies

Socio-semantic networks

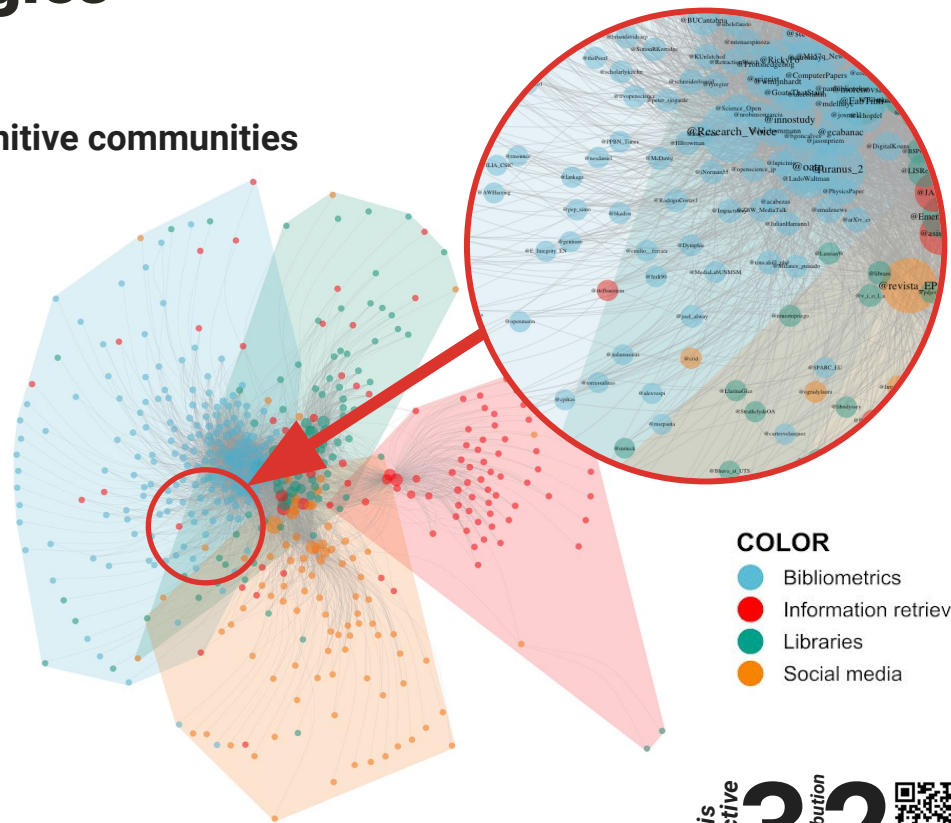
A proposal for **socio-semantic networks** has been developed through a package in  and applied to IS&LS and Microbiology

METHODS

The proposed socio-semantic network can visually simplify social relations and thematic interests, illustrating whether they **align or differ**

FINDINGS

Cognitive communities



Primary findings

OBJECTIVE

1

WorldCat Identities offers unique insights for author analyses but requires careful data validation due to various challenges

Wikipedia has untapped altmetric potential; our framework and metrics highlight its value but emphasize the need for intensive data processing

OBJECTIVE

2

Wikipedia's citation analysis in the Humanities reveals History as a dominant topic, but Humanities citations are just 5% of Wikipedia citations

A broad co-citation analysis mapped Wikipedia's perspectives on science, emphasizing areas like Biochemistry, and showing unique citation patterns, with only 13% to OA journals

OBJECTIVE

3

Overlay maps revealed distinct Microbiology interests across platforms, shedding light on how such attention is generated

Our socio-semantic approach identified genuine scientific interests on social media, hinting at its broader applicability

Social media & science interactions

1. Pioneering exploration of social media to study **science-society relations**
2. Applied **social data mining** for data extraction and processing in altmetrics
3. **Wikipedia**: a significant, yet underutilized, platform for altmetrics
4. **Adapted scientometric methods** to map science through social media's lens
5. Proposed **new methods** merging social and semantic data
6. Emphasized the vast and varied ways science and **society engage on social platforms**

This thesis reveals the transformative role of social data mining for science communication research



BIG DATA TECHNIQUES APPLIED TO THE **STUDY AND **CHARACTERISATION** OF **SCIENTIFIC ACTIVITY** ON **SOCIAL MEDIA****

by Wenceslao Arroyo Machado
PhD Advisors Enrique Herrera Viedma & Daniel Torres Salinas