Research Knowledge Graphs and scholarly information extraction @ BERD@NFDI & GESIS

Focused Tutorial on Capturing, Enriching, Disseminating Research Data Objects

*Stefan Dietze, 25.11.2022*

# GESIS @ National Research Data Infrastructure (NFDI)

Relevant consortia with GESIS in leading roles

- BERD@NFDI
  https://www.berd-nfdi.de/

- NFD4DataScience – National Research Data Infrastructure for Data Science & AI
  https://www.nfdi4datascience.de/

- KonsortSWD
  https://www.konsortswd.de/en/

- Base4NFDI
  https://base4nfdi.de/

# Provenance & Dependencies of Research Data, Resources, Knowledge



Relations between scientific resources, data, knowledge



Research Data Cycle

# Provenance & Dependencies of Research Data, Resources, Knowledge

- **Research Data**
- **Publications**
- **Code/Scripts**
- **ML Models**
- **Methods**
- **Claims**
- **Metrics**

Relations between scientific resources, data, knowledge

Common questions for researchers

- Which top-tier <u>publications</u> cite which <u>data/method</u>? („dataset authority")

- Which <u>data</u> was used to train/evaluate which <u>method</u>? Which <u>method</u> to produce what <u>data</u>?

- Which <u>claims</u> are supported/cited/rejected by what <u>dataset</u> or <u>publication</u>?

# Provenance & Dependencies of Research Data, Resources, Knowledge

- Research Data
- Publications
- Code/Scripts
- ML Models
- Methods
- Claims
- Metrics

Relations between scientific resources, data, knowledge
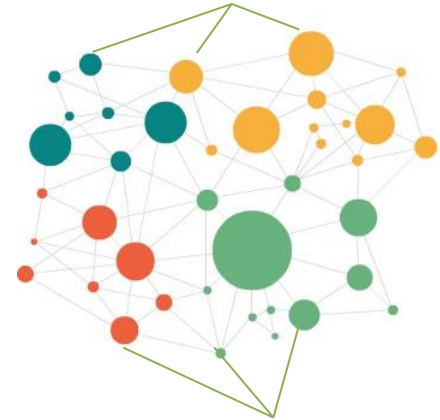
Challenges

- Data & metadata about resources and concepts not represented in **structured, machine-interpretable, integrated manner** (hidden in publications, web pages etc)

- **Persistent identifiers** (e.g. DOIs) used inconsistently (e.g. on publications/datasets, to small degree on ML models)

- **Relations and semantics** not explicit

- **Reproducibility crisis** in CS/DS/AI

# Knowledge Graphs for FAIR Research Data

- **Data interoperability and reuse** through established W3C standards for data sharing (on the Web), e.g. RDF, JSON, shared vocabularies (e.g. schema.org, DCAT, DDI), APIs for data reuse and linking

- Making **links** between resources and concepts explicit & **machine-interpretable**
  (e.g. which publications cite what dataset?)

- Consistent **use of persisent IDs** (e.g. URIs, DOIs) across all data, e.g. concepts, resources etc („*DOIs for all*")

**Resources**
- Datasets
- Publications
- Code
- Software

**Concepts**
- Terms & Definitions
- Claims
- Methods
- Topics
- Entities

# Research KGs in Practice: integrated search @ GESIS

Dataset

Rel. Publications

# From publications to machine-interpretable metadata KGs
## Disambiguation of dataset & software/script citations

https://data.gesis.org/softwarekg
https://data.gesis.org/somesci

- Manual annotation ("SomeSci")

- Training deep learning-based model for extraction software & data references in large-scale data (3.5 M publications)

- Data lifting into KG ("SoftwareKG")

- 300+ M triples / statements

- Search across data/software/publications (GESIS Search)



Schindler, D., Bensmann, F., Dietze, S., Krüger, F., SoMeSci—A 5 Star Open Data Gold Standard Knowledge Graph of Software Mentions in Scientific Articles, (CIKM2021), ACM 2021

# From publications to machine-interpretable <u>metadata KGs</u>
## Understanding scientific software/data usage

(Schindler et al., CIKM2021)

- **Understanding SW usage, citation habits and their evolution across disciplines**
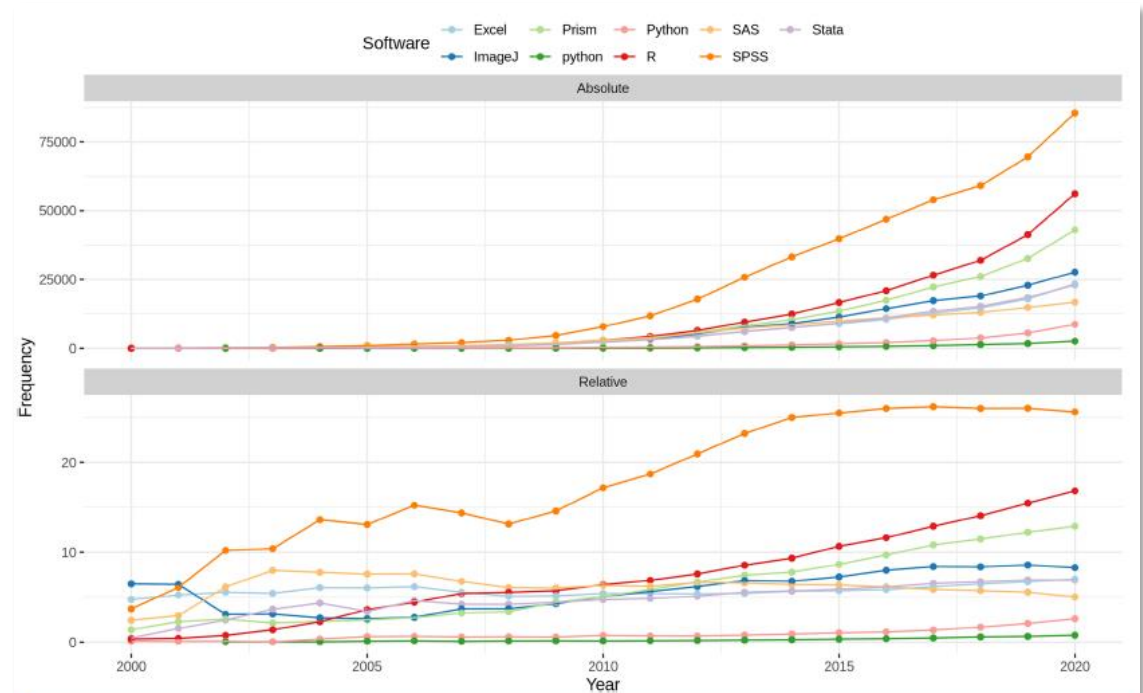
- **Rise of data science = rise of software usage**



Figure 15 Relative and absolute amount of articles per year mentioning the top statistical software.
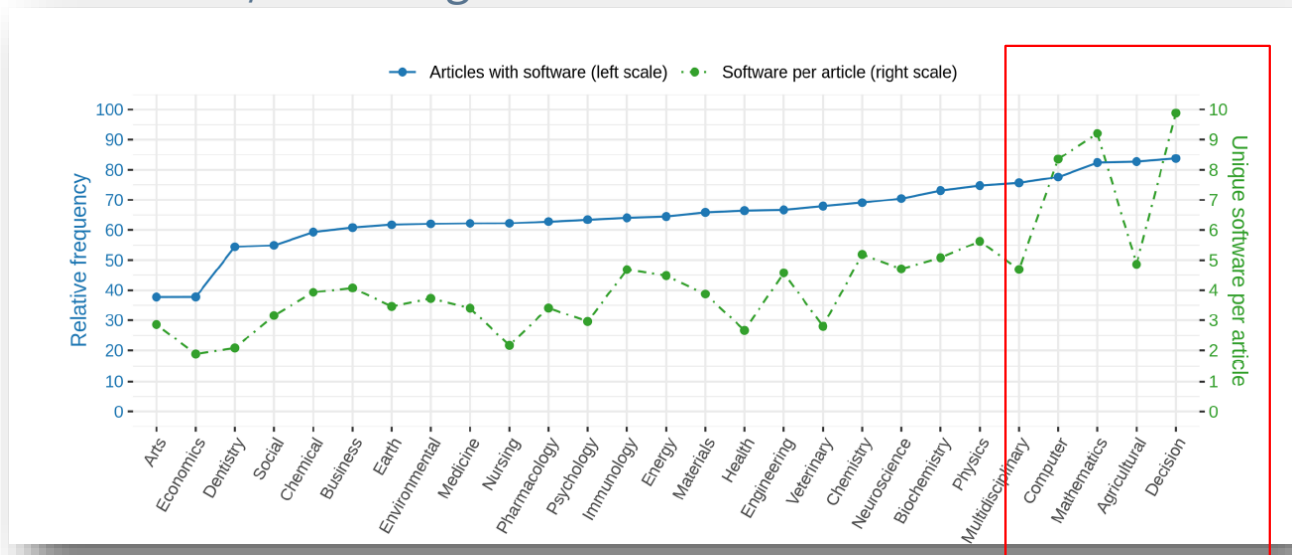Full-size ◪ DOI: 10.7717/peerj-cs.835/fig-15

Schindler D, Bensmann F, Dietze S, Krüger F., The role of software in science: a knowledge graph-based analysis of software mentions in PubMed Central. **PeerJ Computer Science 8:e835**

# From publications to machine-interpretable metadata KGs
## Understanding scientific software/data usage

- Top adopters of data science/AI/software…



Schindler D, Bensmann F, Dietze S, Krüger F., The role of software in science: a knowledge graph-based analysis of software mentions in PubMed Central.
**PeerJ Computer Science 8:e835**

# From publications to machine-interpretable metadata KGs
## Understanding scientific software/data usage

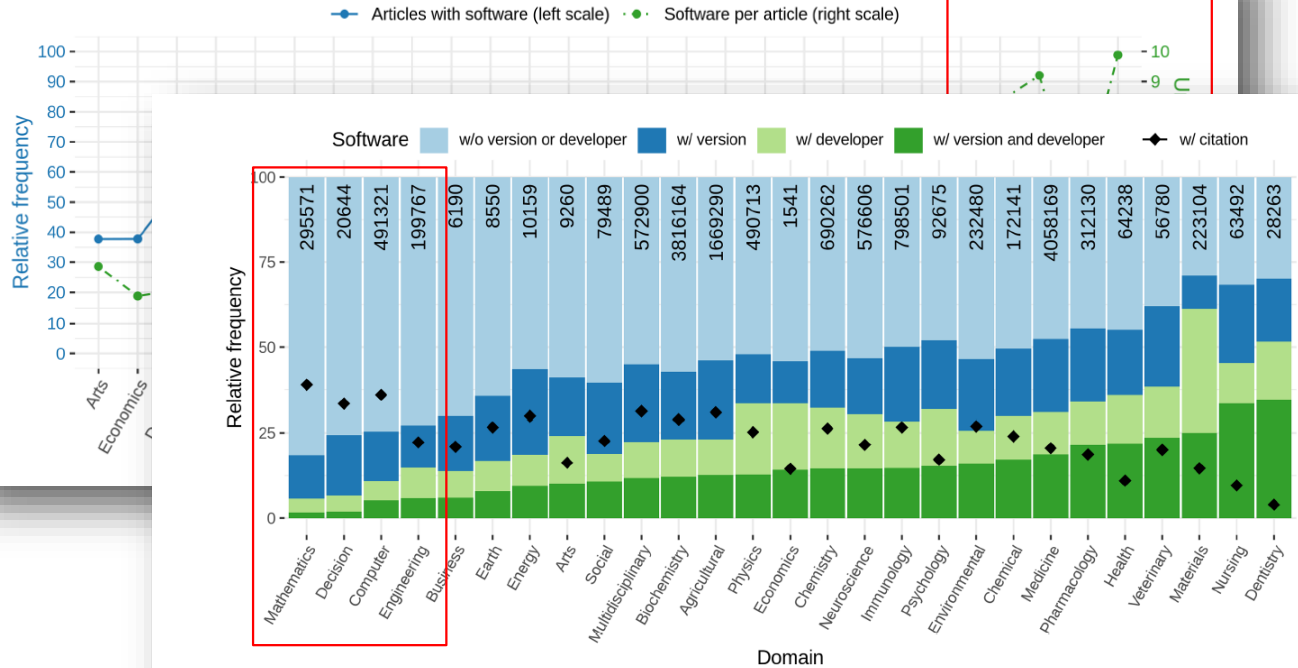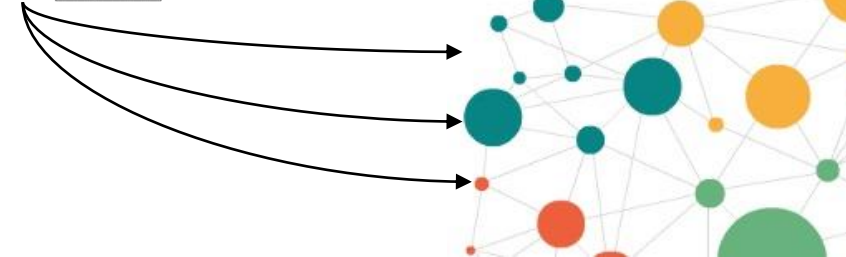- Top adopters of data science/AI/software…
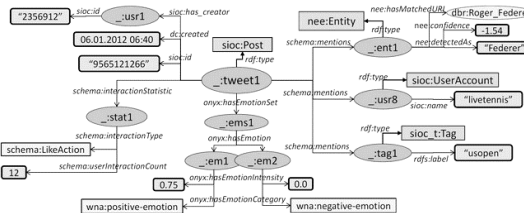- …follow the worst citation habits



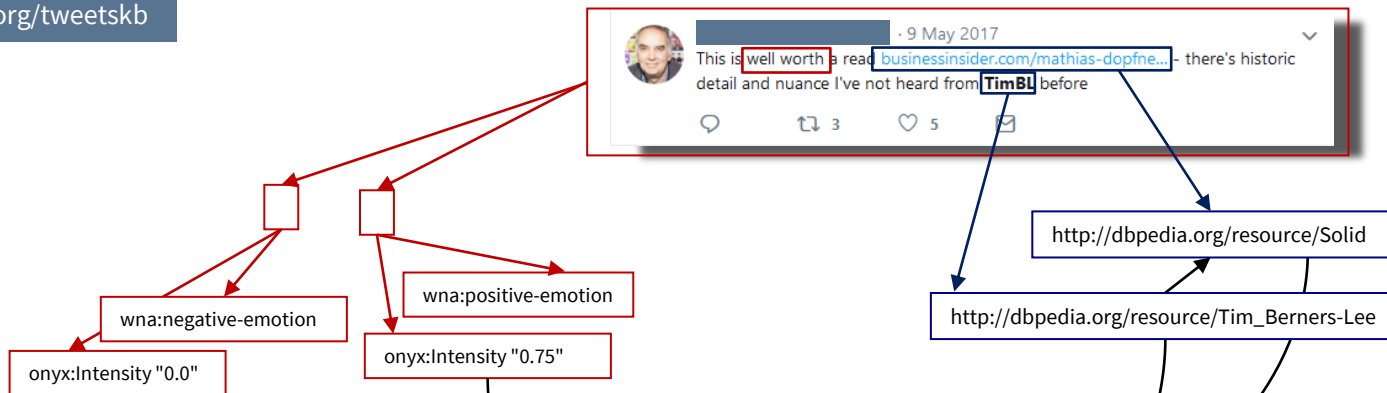**Figure 12** **Distribution of software completeness per research domain.** The numbers at the top of the bars represent the absolute numbers of software considered per domain. Please note that articles may belong to multiple categories. Full-size ■ DOI: 10.7717/peerj-cs.835/fig-12

# From social media to machine-interpretable research data KGs
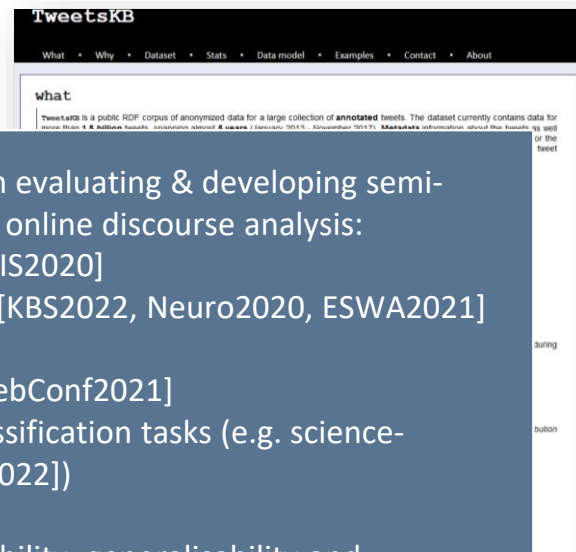## Building a public research knowledge graph from Twitter data

# From social media to machine-interpretable research data KGs
## TweetsKB – a large-scale research KG of societal opinions

https://data.gesis.org/tweetskb



- **Harvesting & archiving of** 10 Billion tweets (permanent collection from Twitter 1% sample since 2013)

- **Information extraction pipeline** to build a KG of entities, interactions & sentiments (distributed Map/Reduce batch processing)
  - Entity linking with knowledge graph/DBpedia (*"president"/"potus"/"trump"* => *dbp:DonaldTrump)*
  - Sentiment analysis/annotation
  - Geotagging
  - Lifting into knowledge graph schema

KTS research focused on evaluating & developing semi-supervised methods for online discourse analysis:
- Stance detection [IJIS2020]
- Sentiment analysis [KBS2022, Neuro2020, ESWA2021]
- Entity linking
- Georeferencing [WebConf2021]
- More fine-grain classification tasks (e.g. science-relatedness [CIKM2022])

But: focus here on scalability, generalisability and robustness towards evolving data/vocabulary => unsupervised approaches
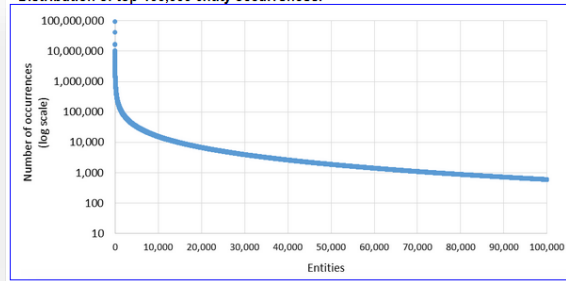
Dimitrov, D., Baran, E., Fafalios, P., Yu, R., Zhu, X., Zloch, M., Dietze, S., TweetsCOV19 – A Knowledge Base of Semantically Annotated Tweets about the COVID-19 Pandemic, **CIKM2020**

# From social media to machine-interpretable <u>research data KGs</u>

## TweetsKB – a large-scale research KG of societal opinions

https://data.gesis.org/tweetskb

- **Harvesting & archiving of** 10 Billion tweets
  (permanent collection from Twitter 1% sample since 2013)

- **Information extraction pipeline** to build a KG of entities, interactions & sentiments
  (distributed Map/Reduce batch processing)
  - Entity linking with knowledge graph/DBpedia
    (*"president"/"potus"/"trump"* =>
    *dbp:DonaldTrump*)
  - Sentiment analysis/annotation
  - Geotagging
  - Lifting into knowledge graph schema

- **Public, privacy-aware, large-scale research corpus of public opinions and their evolution
  => interdisciplinary research**



**Distribution of top-100,000 *entity* occurrences:**

| Number of tweets | 1,560,096,518 |
| Number of distinct users | 125,104,569 |
| Number of distinct hashtags | 40,815,854 |
| Number of distinct user mentions | 81,238,852 |
| Number of distinct entities | 1,428,236 |
| Number of tweets with sentiment | 772,044,599 |
| Number of RDF triples | 48,207,277,042 |

Dimitrov, D., Baran, E., Fafalios, P., Yu, R., Zhu, X., Zloch, M., Dietze, S., TweetsCOV19 – A Knowledge Base of Semantically Annotated Tweets about the COVID-19 Pandemic, **CIKM2020**

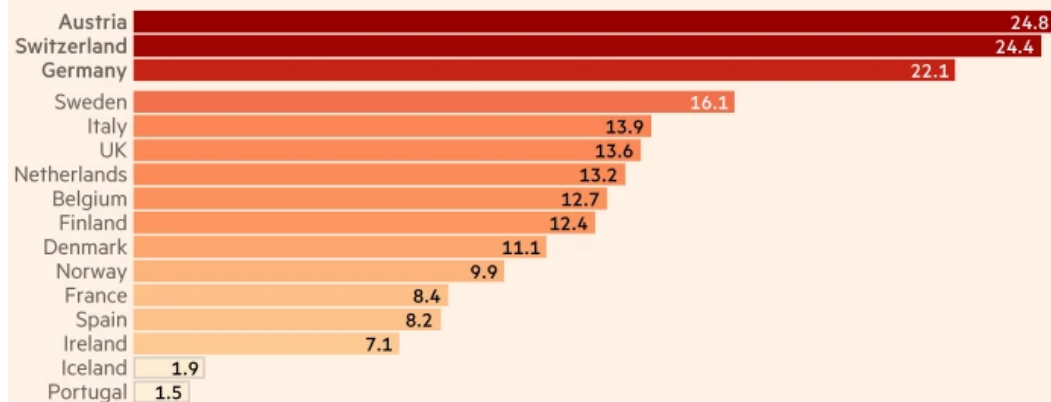# RKG-based social science research using TweetsKB
## Investigating Vaccine Hesitancy in DACH countries

### German-speaking countries have the highest shares of unvaccinated people in western Europe

Share of population aged 12+ that has not had any Covid vaccine dose (%)

| Country | % |
| --- | --- |
| Austria | 24.8 |
| Switzerland | 24.4 |
| Germany | 22.1 |
| Sweden | 16.1 |
| Italy | 13.9 |
| UK | 13.6 |
| Netherlands | 13.2 |
| Belgium | 12.7 |
| Finland | 12.4 |
| Denmark | 11.1 |
| Norway | 9.9 |
| France | 8.4 |
| Spain | 8.2 |
| Ireland | 7.1 |
| Iceland | 1.9 |
| Portugal | 1.5 |

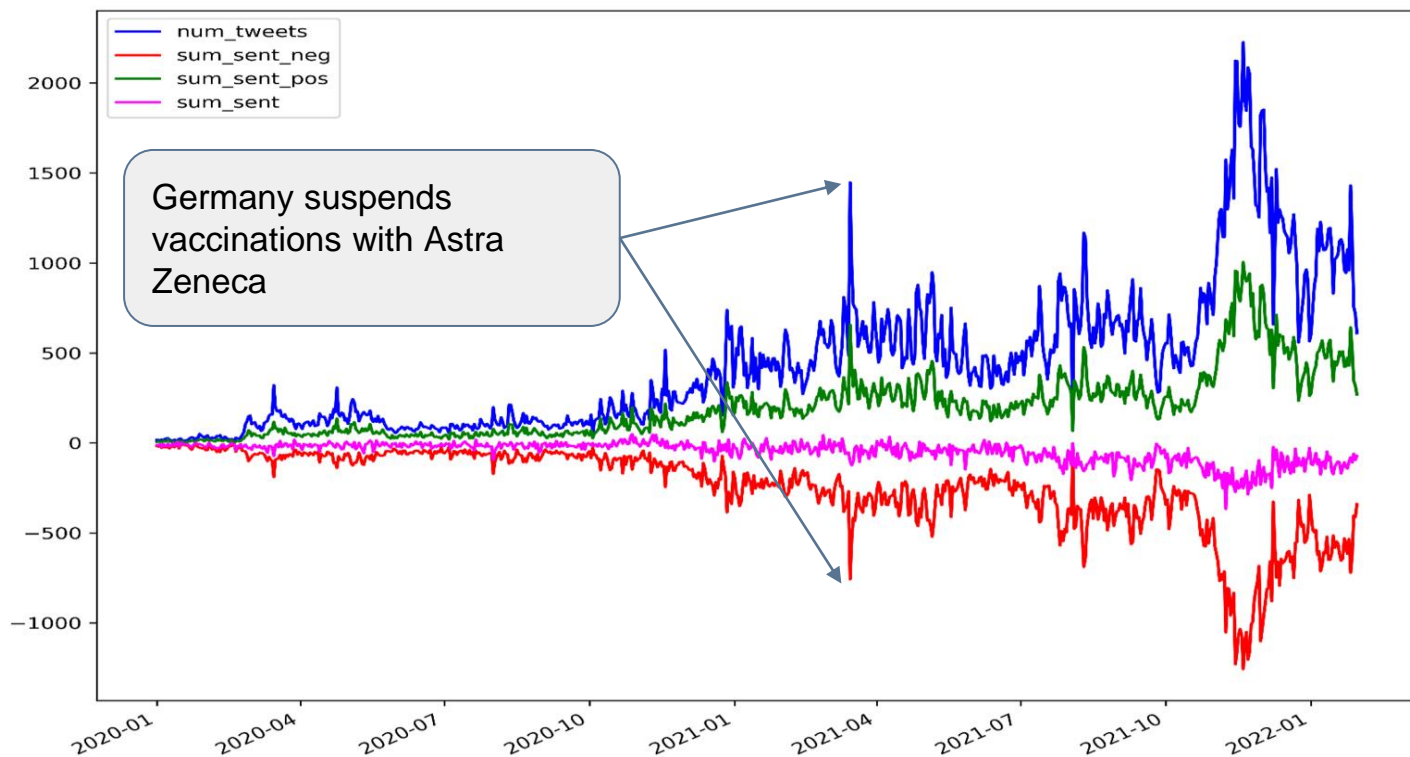Source: FT analysis of figures from national sources and Our World in Data. Rates shown are as of November 9

# RKG-based social science research using TweetsKB
## Investigating Vaccine Hesitancy in DACH countries

Twitter discourse zu "Impfbereitschaft"



Germany suspends vaccinations with Astra Zeneca

# RKG-based discourse analysis using TweetsKB

## Vaccine Hesitancy– key topics in "safety" category

**DD4P**
DiscourseData4Policy

# How about mentions of science resources on the Web?

## Example: Twitter

Table 1: Examples (tweets 1 to 4) and Counterexamples (tweet 5) of scientific online discourse tweets

| | |
|---|---|
| Science claim | (1) Donating blood not only helps others, but reduces the rate of cancer and heart disease in the donor. |
| Science reference | (2) via @medical_xpress A new in vitro (test tube) study, ""Dietary functional benefits of Bartlet http://t.co/Qv1C1GjQin #UFO4UBlogHealth |
| Science relevance | (3) How is @UChicagoIME shaping the future of science? Find out on April 6! |
| Science reference | (4) Study: Shifts in electricity generation spur net job growth, but coal jobs decline - via @DukeU http://t.co/AXGmKUPata |
| No science | (5) My father got COVID-19. |

Hafid, S., Schellhammer, S., Bringay, S., Todorov, K., Dietze, S., "SciTweets - A Dataset and Annotation Framework for Detecting Scientific Online Discourse", **CIKM2022**
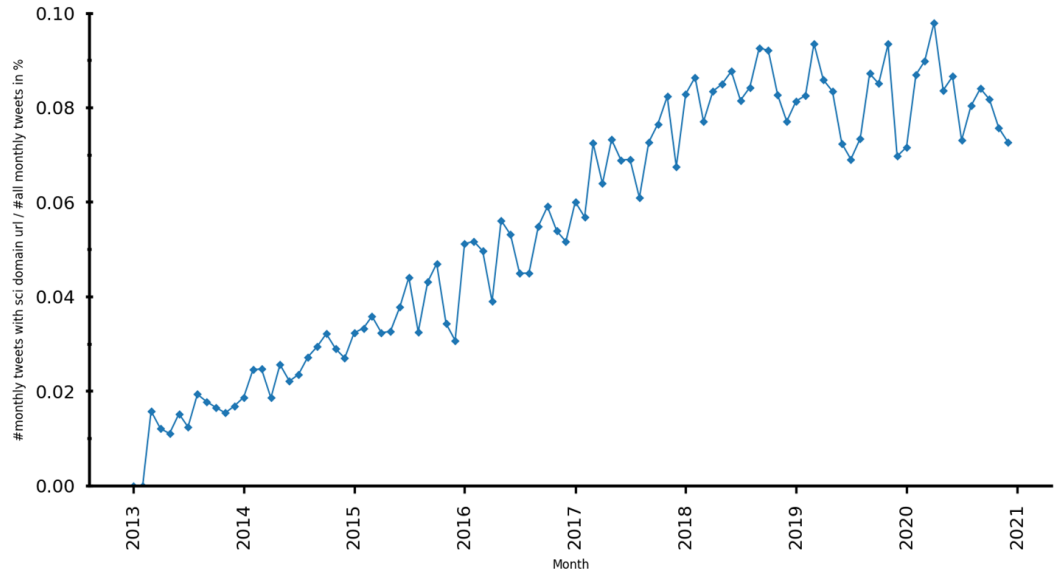
# How about mentions of science resources on the Web?
## Example: Twitter

- Percentage of tweets containing links to scientific articles (journals, publishers, science blogs etc)

- Uses list of > 30 K science web domains

- Data source: TweetsKB (https://data.gesis.org/tweetskb/), > 10 bn tweets archived since 2013
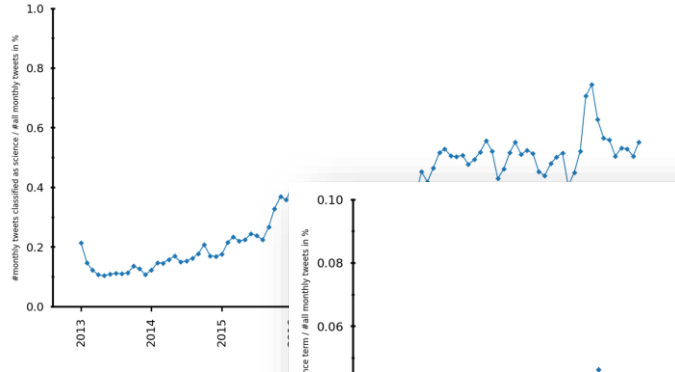


- AI4Sci project: understanding and classification of science discourse online (news, social Web)

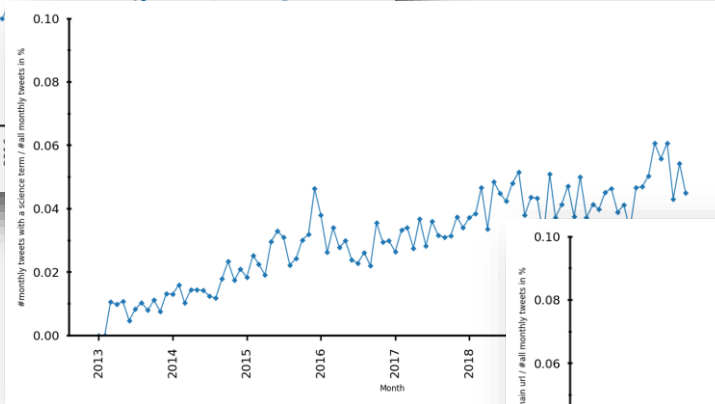# How about mentions of science resources on the Web?
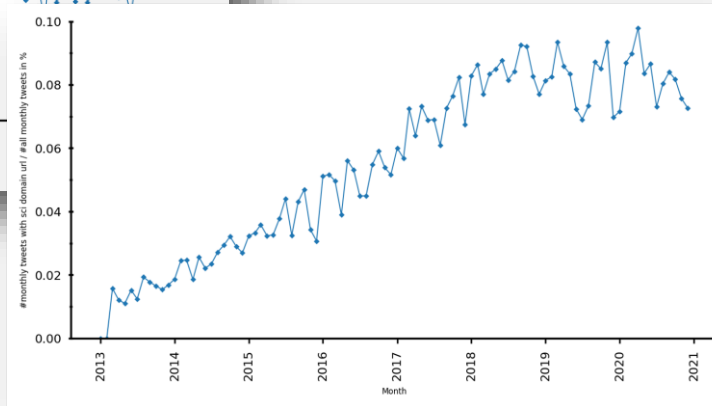
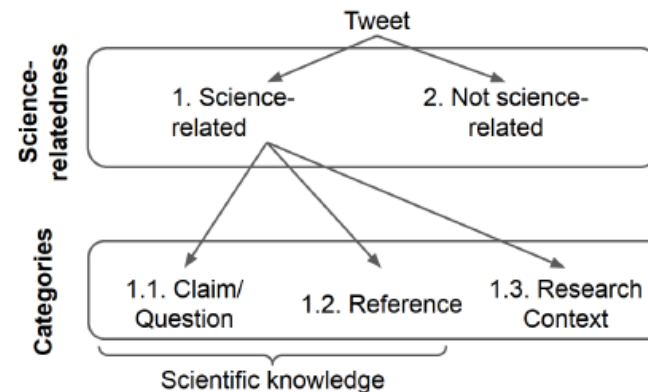## Example: Twitter

SciBERT classifier

Heuristic: Sci term

Sci subdomain

# *SciTweets* dataset & classifier

- Ground truth dataset, heuristics-based sampling strategy and annotation framework for testing classification models

- 1261 expert-labeled tweets across all classes/labels

- Baseline classifiers based on SciBERT transformer model (fine-tuned/tested on SciTweets)

- Ongoing: analysis of large-scale science discourse and its evolution

| Task | Category | Precision | Recall | F1 |
|------|----------|-----------|--------|-----|
| binary | 1 - Science-related | 84.70 | 83.99 | 84.34 |
| | 2 - Not Science-related | 92.67 | 93.03 | 92.85 |
| multi | 1.1 - Scientific Claim | 75.00 | 81.18 | 77.97 |
| | 1.2 - Reference | 76.19 | 77.01 | 76.60 |
| | 1.3 - Research Context | 81.06 | 79.65 | 80.35 |

Hafid, S., Schellhammer, S., Bringay, S., Todorov, K., Dietze, S., *SciTweets - A Dataset and Annotation Framework for Detecting Scientific Online Discourse*, **CIKM2022**

# Summary: Research KGs @ GESIS

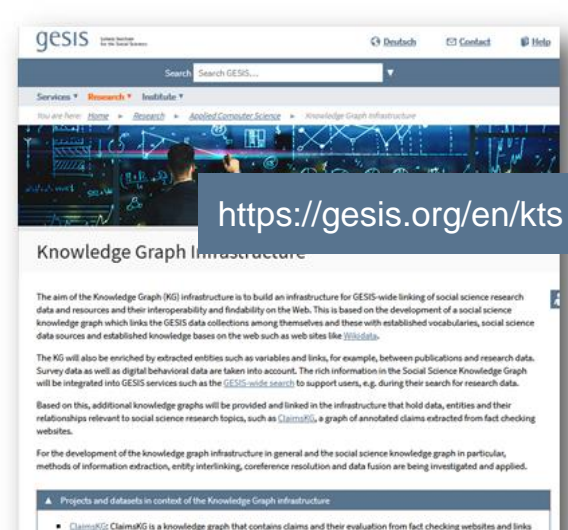**Tools for constructing scholarly knowledge graphs**

- NLP and deep learning-powered methods for extracting large-scale KGs about methods, claims, data, software involved in the scientific process
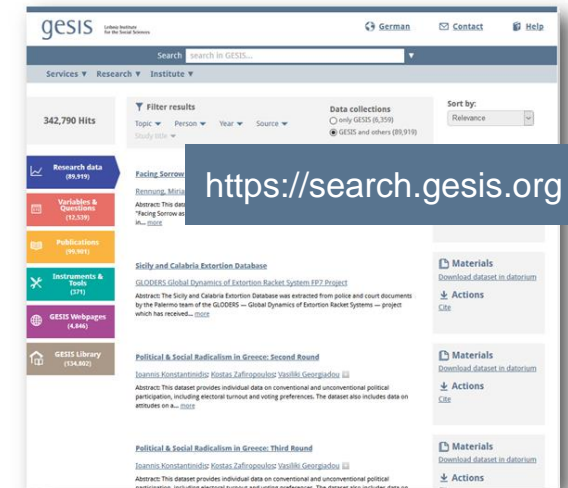
**Large-scale scholarly KGs, e.g.**

- KGs about scholarly use of software & research data
  (e.g. **SoftwareKG**: 1.8 M disambiguated software mentions extracted from 3 M publications, https://data.gesis.org/softwarekg/)
- Web mined KGs of social science research data, e.g. public opinions, claims and attitudes expressed on social media
  (e.g. **TweetsKB**: > 10 Bn semantically annotated tweets, sentiments, https://data.gesis.org/tweetskb)

**Semantic Search powered by KGs and related tools**

- RKG-powered search across scholarly publications, datasets, methods and their relations (e.g. **GESIS Search**, https://search.gesis.org)



https://gesis.org/en/kts

https://search.gesis.org

# Outlook: shared tasks on scholarly information extraction

Enganging with the community to advance progress in RKGs & scholarly IE

Creating large training/testing corpora and run shared tasks for
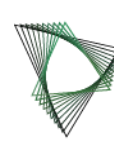
- Software / code detection and disambiguation

- Leaderboard extraction / task-dataset-metric detection (TDM)

- Dataset mention detection & disambiguation

- Machine learning model detection & disambiguation

- Research field classification

More to be announced soon.

@stefandietze
http://stefandietze.net