

gesis

Leibniz Institute  
for the Social Sciences



# A Manifesto for Data Sharing in Social Media Research

Katrin Weller

Katharina Kinder-Kurlanda

*ACM Web Science 2016*

# WHY?

**1. Advance reproducibility and comparability**

**2. Avoid 'digital divides' in data accessibility**

**3. Save time and money in data collection**

# HOW?

# Successful data sharing

Depends on three critical perspectives:

- ▶ **Methodological:**  
documentation to ensure reproducibility
- ▶ **Legal:**  
share in accordance with various legal frameworks, e.g. data protection legislation, copyright
- ▶ **Ethical:**  
consider privacy expectations, sharing should „do no harm“

# Examples of current practices

Approach	Example	Retrievability	Documentation + standards	Long-term perspective
„Grey market“	Data shared with colleagues (often upon request)	<b>Low</b> , via personal connections	<b>Low</b> , no standardized documentation or data format	<b>Low</b> , no guaranteed long term availability
Researchers' personal/professional websites		<b>Medium</b> , URLs may be referenced	<b>Low</b> , no standardized documentation or data format	<b>Low</b> , no guaranteed long term availability
Social media providers	e.g. Wikipedia dumps	<b>High</b> , if directly provided from social media platform	<b>Medium/high</b> , depending on the producer	<b>Medium/high</b> , depending on the producer
Project-based or thematic collections	e.g. KONECT, CrisisLex	<b>Medium</b> , URLs may be referenced	<b>Medium</b> , if same principles are applied within entire collection	<b>Low</b> , no guaranteed long term availability
Conferences & journals	e.g. ICWSM datasets with conference papers	<b>High</b> , usually related to accepted publications	<b>Medium/high</b> , depending on the publisher	<b>Medium/high</b> , depending on the publisher
Professional archives	e.g. datasets at GESIS data archive	<b>Medium/high</b> , datasets may be referenceable with DOIs.	<b>Medium/high</b> , depending on the archive's requirements	<b>High</b> , guaranteed availability for different time spans

- For additional examples see our paper: <http://dx.doi.org/10.1145/2908131.2908172>
- See also: Thomson, S.D. 2016. Preserving Social Media. DPC Technology Watch Report. Retrieved from <http://dpconline.org/publications/technology-watch-reports>

# How much should I share?

**Most reproducibility**

**What is being shared?**

- whole dataset plus additional research information (e.g. scripts)
- whole dataset
- whole dataset, but without direct identifiers (pseudonymization)
- parts of the dataset removed (anonymization)
- changed dataset (e.g. only tweet IDs)

**Most privacy**

# WHAT'S NEXT?



# Open challenges & next actions

Stakeholder	Open challenges	Next actions
Researchers (individuals/groups)	<ul style="list-style-type: none"><li>Better documentation practices of how the data was collected *exactly*</li></ul>	<ul style="list-style-type: none"><li>Share internally, i.e. within groups and with collaborators</li><li>Share data (and code), in particular datasets used for accepted papers</li></ul>
Archival institutions	<ul style="list-style-type: none"><li>Practices for anonymization</li></ul>	<ul style="list-style-type: none"><li>Develop standards for documentation and metadata</li><li>Establish different models for data access (e.g. embargo, secure data access)</li></ul>
Publishers / Conference organizers	<ul style="list-style-type: none"><li>Ways to include checks for data quality in the review process</li></ul>	<ul style="list-style-type: none"><li>Publish non-standard works, such as descriptions of datasets, reproduction studies etc.</li></ul>
Research associations	<ul style="list-style-type: none"><li>Dialogue with major social media companies</li></ul>	<ul style="list-style-type: none"><li>Continuous dialogue with members on best practices, guidelines</li></ul>

# Open challenges & next actions

Stakeholder	Open challenges	Next actions
All	<ul style="list-style-type: none"><li>• Raise awareness of legal and ethical issues</li></ul>	<ul style="list-style-type: none"><li>• Discuss data sharing practices</li></ul>

# Questions and Feedback

**Dr. Katrin Weller**

[katrin.weller@gesis.org](mailto:katrin.weller@gesis.org)

[@kwelle](#)

<http://katrinweller.net>



**Dr. Katharina Kinder-Kurlanda**

[Katharina.kinder-kurlanda@gesis.org](mailto:Katharina.kinder-kurlanda@gesis.org)

[@ka\\_kinder](#)

<http://www.gesis.org/sdc>

