# Ethical questions in research with digital trace data

Johannes Breuer

Leibniz Association

CS3 Lab – Computational Survey and Social Science

June 18th , 2025

# GESIS Services
## for
## Digital Behavioral Data

GESIS Leibniz-Institute for the Social Sciences

https://rrr.is/gesisdbd

GESIS Panel.dbd
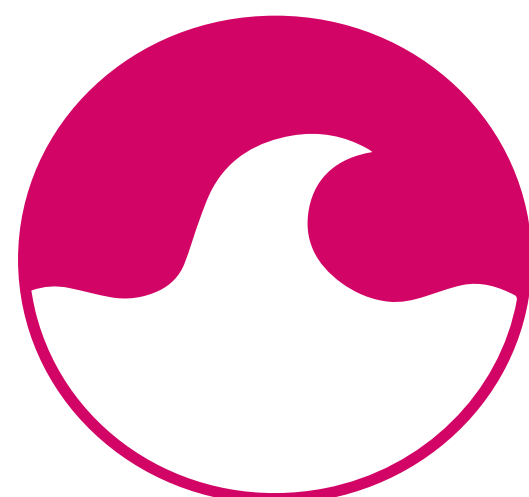
DBD Consulting

GESIS Guides to Digital Behavioral Data

GESIS Web Tracking
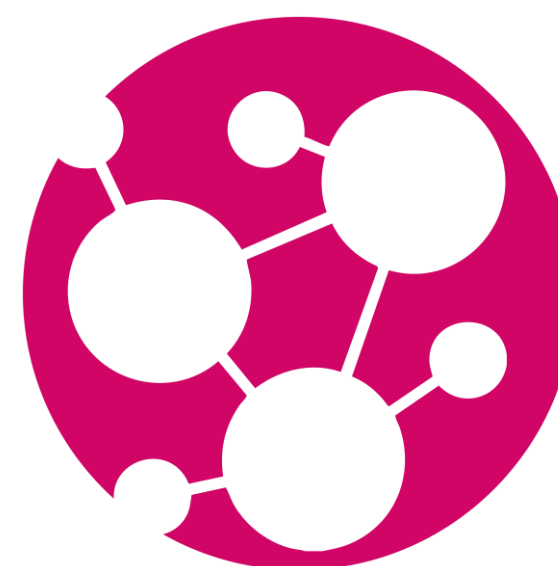
GESIS Methods Hub

GESIS AppKit

GESIS Web Data

# Main basis of this talk

Breuer, J., Stier, S., Lukito, J., Mangold, F., Wieland, M., Radovanović, D., Radovanović, D., Zens, M., Breuer, J., Weller, K., & Wagner, C. (2025). *Overview of Ethical Considerations when Working with Digital Behavioral Data (GESIS Guides to Digital Behavioral Data, 14)* (Version 1.0). GESIS - Leibniz-Institute for the Social Sciences. https://doi.org/10.60762/GGDBD25014.1.0

# A Case For Ethical and Transparent Research Experiments in the Public Interest

**April 30, 2025**

**By Sarah Gilbert, Michael Zimmer, Nathan Matias, Ethan Zuckerman**

On April 26, moderators of r/ChangeMyView, a community on Reddit dedicated to understanding the perspectives of others, revealed that academic researchers from the University of Zürich conducted a large-scale, unauthorized AI experiment on their community. The researchers had used AI bots to secretly impersonate people for experiments in persuasion.

Crucially, this experiment was carried out without the Reddit community's knowledge or consent, and the bots were not labeled as AI. Reddit users were unknowingly exposed to sometimes deceptive AI-generated content designed to shape their opinions—raising significant ethical and transparency concerns.

When the researchers shared preliminary results with the community, r/ChangeMyView moderators contacted the researchers' ethics board out of concern. When told the project met institutional ethics standards, moderators disclosed the study to their shocked and outraged community. Reddit banned the associated accounts and issued a statement condemning the activity.

This case has drawn sharp criticism from the tech research and ethics community. It highlights a mistaken belief by some researchers that public-interest intentions justify ignoring the ethical responsibilities of researchers toward participants and the public.

SEARCH...

## Support Our Work

We are currently seeking funding to support the work of the Coalition. If you'd like to help, please get in touch.

You can also donate to the Coalition online.

Source: https://independenttechresearch.org/a-case-for-ethical-and-transparent-research-experiments-in-the-public-interest/

**Further discussion & reporting:**

- Blog post by Kevin Munger

- Blog posts by *Retraction Watch*[1][2]

- Bluesky thread by Casey Fiesler

# Research ethics & digital trace data in practice

- despite the relevance of the issue, several systematic review studies have shown that **research ethics are often not (explicitly/properly) addressed in publications** based on digital trace data (Fiesler et al., 2024; Knöpfle et al., 2024, Lisker & Mihaljevic, 2025)

# Key questions guiding our research activities

- **What can we do?**

    - legal regulations/frameworks

    - methods, data access, & resources


- **What should we do?**

    - **research ethics**

    - methodological rigor & data quality

# Research ethics

- "Research ethics primarily concerns itself with the **responsible conduct of research**, emphasizing the **protection of human participants**, the **integrity of data**, and the **avoidance of harm**." (Knöpfle et al., 2024, p. 335, emphasis added)

  - **closely related to but not synonymous with research integrity** (Emmerich, 2020) → honesty, transparency, and adherence to professional standards

  - despite the use in everyday language, ethics are **not a binary concept**

# Different perspectives on research ethics

- deontology vs. consequentialism

- in a nutshell...

  - **deontology**: adherence to specific **norms and fundamental values** to guide decision-making

  - **consequentialism**: evaluation of **anticipated outcomes and their ethical implications**

  (see, e.g., Knöpfle et al., 2024;  Salganik, 2019)

- in practice typically a combination of both perspectives (Schlütz & Möhring, 2018)

# Belmont Principles

National Commission for the Protection of Human Subjects of Biomedical and Behavioral Research (1979)

- **Respect for Persons**

- **Beneficence**

- **Justice**

→ **minimizing risks and harms** while **prioritizing the value of research**

- **different values and goals** that can be **evaluated differently** and **may even be conflicting** in specific situations (Iphofen, 2020; Israel, 2015)

# Digital behavioral data / digital trace data

"DBD encompass **digital observations of human or algorithmic behavior**. DBD are generated (1) through **interactions and content production online** (e.g., on platforms such as Google, Facebook or websites on the World Wide Web) or (2) by **software or sensors for recording specific processes** (e.g., smartphones, RFID sensors, satellites, street view cameras, or web tracking)." (Wagner et al., 2025, p. 2; emphasis added)
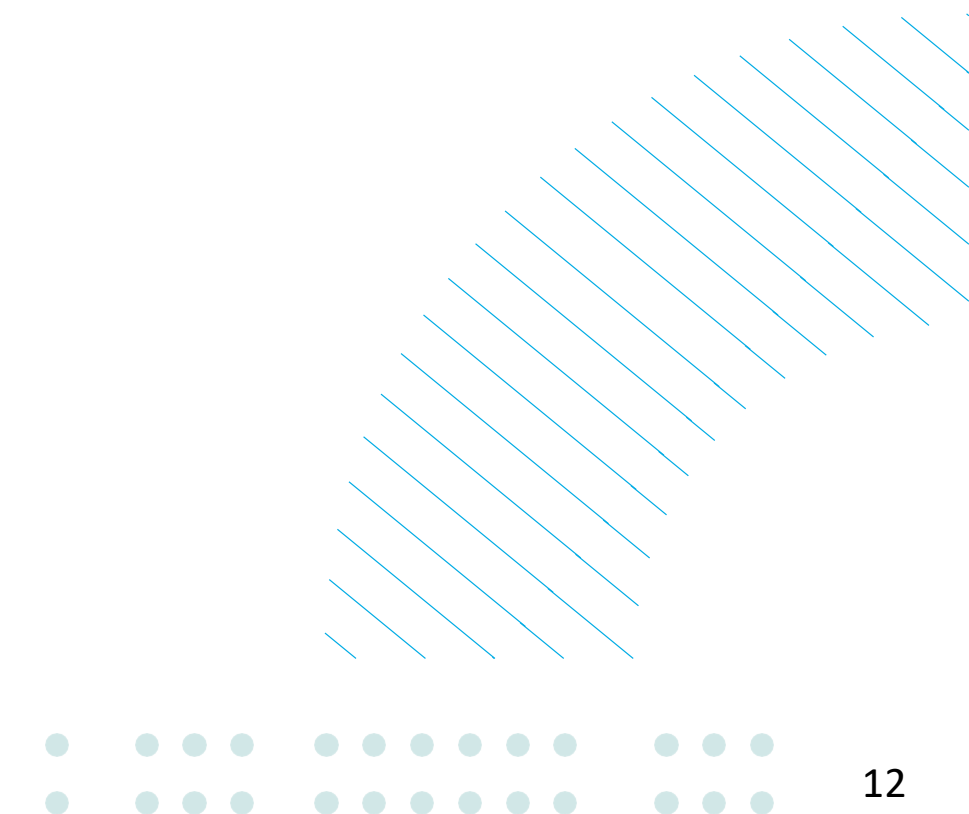
# Types of digital trace data

- **found data** (= not genuinely produced for research) vs. **designed data** (= genuinely produced for research)

- Hox (2017): **unintentional** vs. **intentional traces**

- Menchen-Trevino (2013)

  - **participation traces** (e.g., comments or posts) and **transactional data** (e.g., login data or logs more generally)

  - **horizontal** (e.g., all posts from a social media platform containing a specific hashtag) vs. **vertical trace data** (comprehensive usage data - potentially from different sources - for a limited group of users)

→ different implications for questions related to research ethics

# Digital trace data & research ethics

- many ethical questions are the same for digital trace data and other types of data, but **some issues are unique or at least particularly pronounced/important**

- ethical questions arise in **all phases of the research process** when working with digital trace data

  - **Study design & data collection**
  - **Data processing & analysis**
  - **Publication & data sharing**

# Study design & data collection

- **different data collection methods** raise specific ethical questions: e.g., APIs, web scraping, data donation (Breuer et al., 2020; Ohme et al., 2023)

- important distinction for ethical considerations:

  - **platform-centric** approaches:
    - sampling from one or more platforms based on relevant entities, such as users, topics, hashtags, search queries or time
    - data collection directly from platforms (typically via APIs or scraping)

  - **user-centric** approaches:
    - users recruited for the study (e.g., via existing panels)
    - data collected via dedicated research software (e.g., browser plugin) or donations of so-called data download packages (Carrière et al., 2024; van Driel et al., 2022)
    - often involves data linking: e.g., digital traces (possibly from multiple platforms) + surveys (Stier et al., 2020)

# Study design & data collection

- aspects to consider:

  - **topics** (e.g., regarding sensitivity)

  - **data types** (e.g., text, image, video)

  - **sample** (e.g., vulnerable groups included)

  - **platform attributes & user expectations**
    - researchers are typically not part of the "imagined audience" of social media users (Marwick & boyd, 2011; Fiesler & Proferes, 2018)
    - "Concerns over consent, privacy and anonymity do not disappear simply because subjects participate in online social networks; rather, they become even more important" (Zimmer, 2010, S. 324)

  - **contact** with the people whose data are being collected
    - in case of platform-centric collection, the term participant may not be appropriate (Breuer et al., 2023)
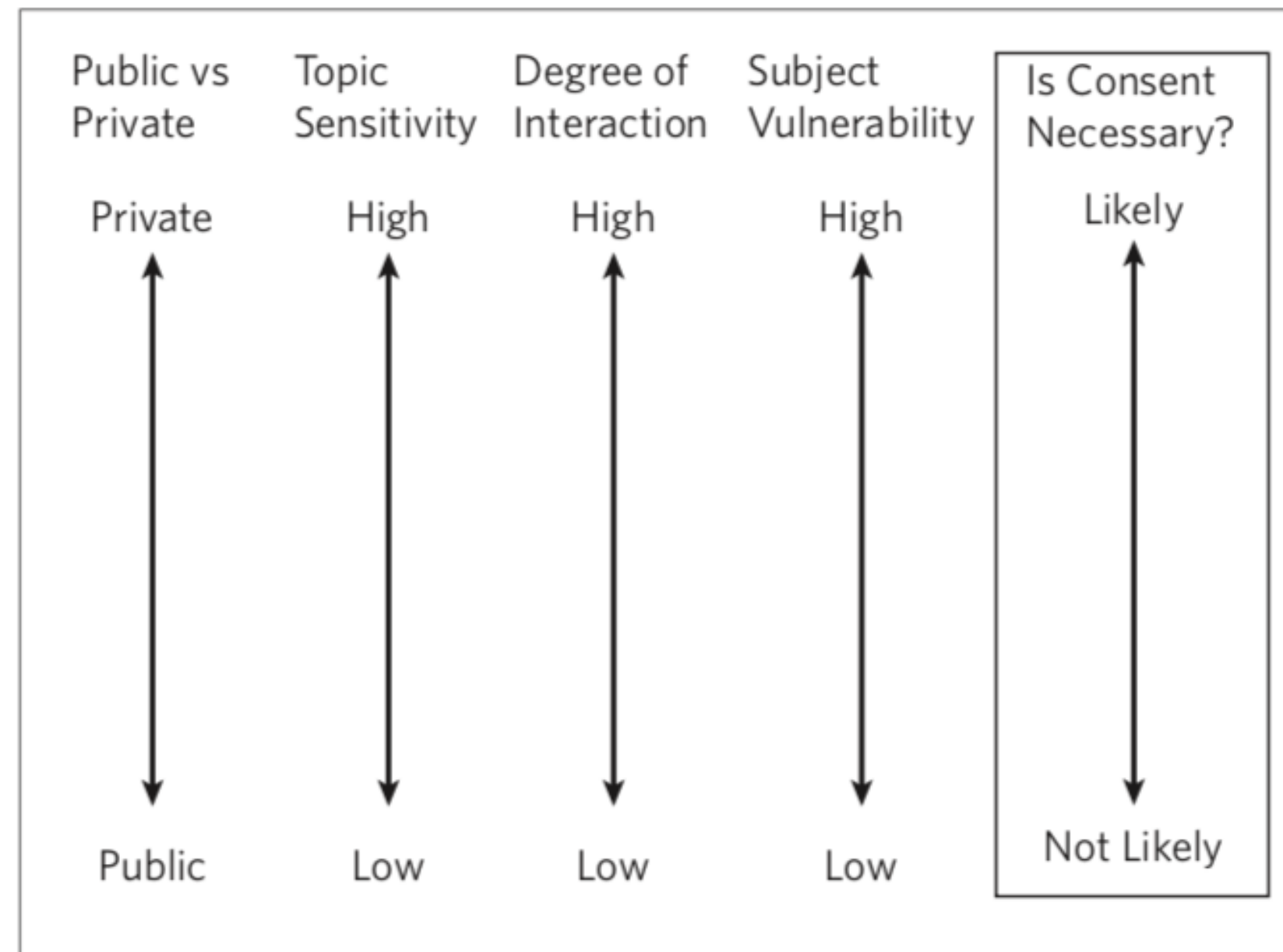    - important for the possibility of obtaining **informed consent**

# Informed consent

- informed consent is relevant from a legal as well as an ethical perspective
  - sidenote: informed consent is arguably the "safest" option but only one of the possible legal bases for collecting and processing personal data according to the GDPR (another relevant one for academic research is legitimate interest)

- informed consent should **appropriately inform participants about the nature and purpose of the data collection in a comprehensible manner**
  - possibility to provide additional "technical details" via optionally available supplementary information (Breuer, Al Baghal et al., 2021)

# Informed consent

- informed consent may be very difficult or even impossible to obtain in platform-centric approaches

  - there are/can be solutions, however, such as the *Bartleby* tool by Zong & Matias (2022) or working together with moderators or instance admins on decentral platforms like *Mastodon* (Wähner et al., 2024)

- **different options of implementing informed consent** (esp. for user-centric approaches): e.g., granular consent for multiple data types or dynamic consent for longitudinal studies (Breuer et al., 2025)

# Need for informed consent



Factors affecting consent

Source: McKee & Porter (2009), p. 88

# Data processing & analysis

- anonymization/pseudonymization
    - digital trace data can contain many types of direct or indirect identifiers
        - these may also include data from others
        - combinations of metadata can be used for identifying individuals (Sloan et al., 2020)
        - example of URLs from web tracking: profile names/IDs in the path or parameters, search strings, geo coordinates, location names...

- data linking
    - can increase sensitivity of the data + (re-)identification risk

- inferred attributes (using ML, NLP, computer vision)
    - may be wrong
    - can increase sensitivity of the data
    - have not been provided by the individuals (→ consent)

# Publication & data sharing

- **protecting participants**/individuals whose data were collected vs. **increasing transparency** & **maximizing the value of data** by making it reusable

- **full raw data are often difficult or impossible to share** due to legal and/or ethical considerations
  - digital trace data are often **personal** and can be **proprietary** and/or **sensitive**
  - legal & ethical considerations are key reasons why researchers do not share social media data (Akdeniz et al., 2023)

- several publications provide some **guidance on sharing digital trace data** (esp. social media data), including ethical considerations (Breuer, Borschewski et al., 2021; Bishop & Gray, 2017; Williams et al., 2017)

# Specific vs. general(izable) guidance

- no "one-size-fits-all" solutions
- decisions in research ethics depend on the specific case
- trade-off between concreteness and applicability
  across contexts
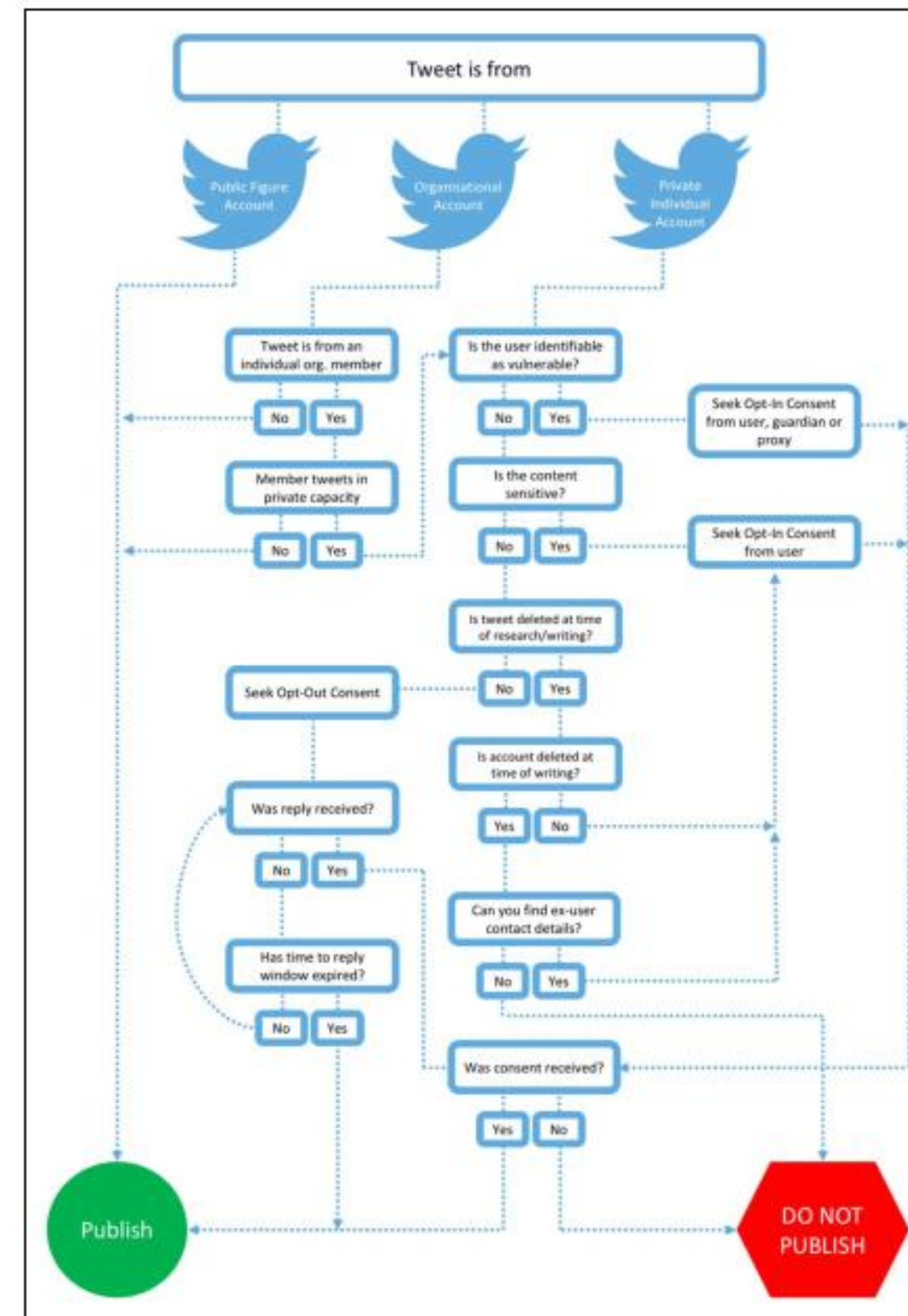
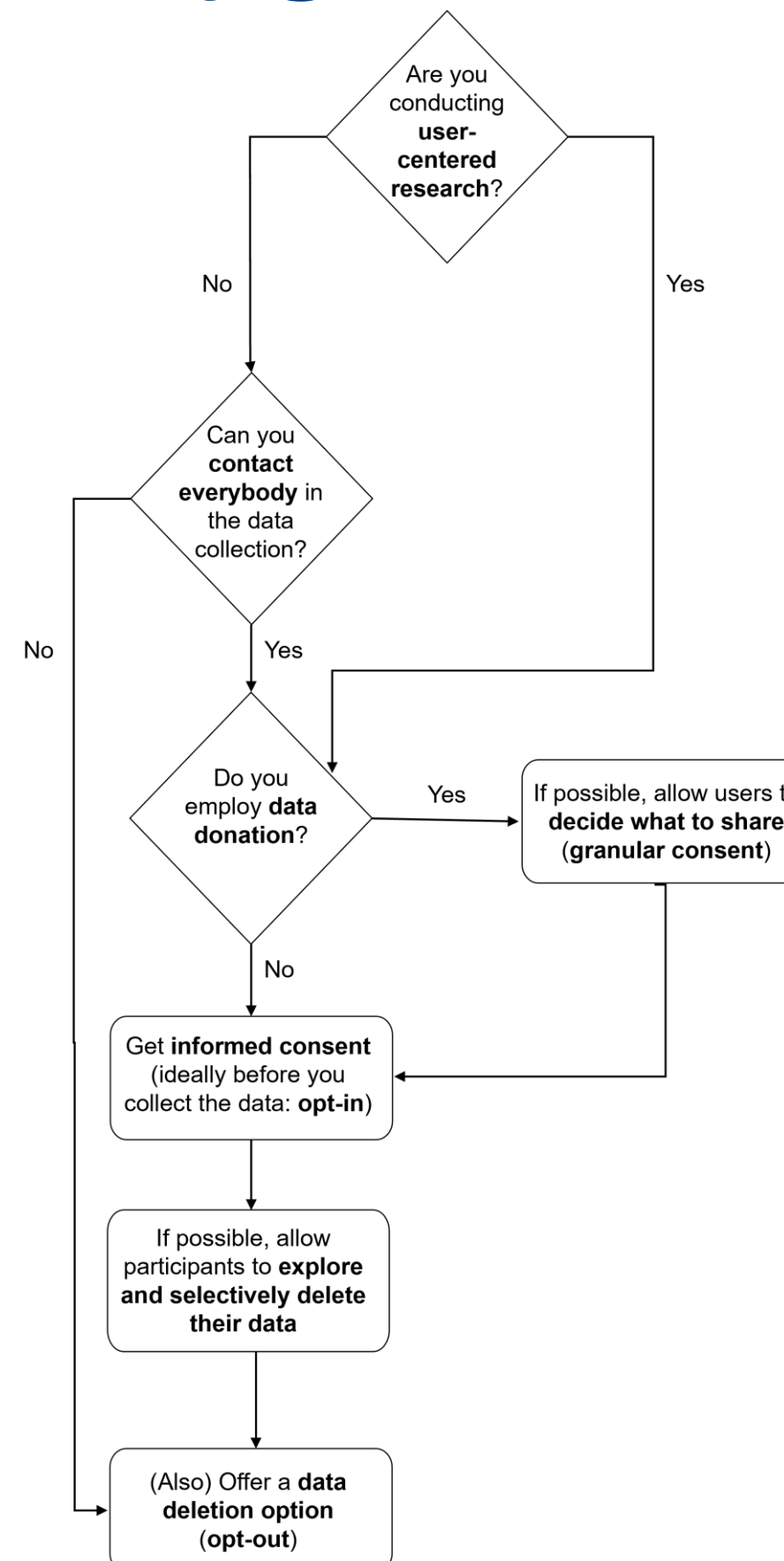# Specific vs. general(izable) guidance



**Figure 1.** Decision flow chart for publication of Twitter communications.

Source: Williams et al. (2017), p. 1163

# Specific vs. general(izable) guidance



Adapted from Breuer, Weller, & Kinder-Kurlanda (2023)

# Some general recommendations

- Prioritize **informed consent** when possible.

- Implement **data minimization**: Collect only the data necessary for addressing a specific research question.

- **Only link** data **when** it is **necessary for an analysis**.

- When sharing data, be **"As open as possible, as closed as necessary"**
  - if possible/applicable, use anonymization/pseudonymization procedures, access control, trusted repositories, and consider novel/alternative sharing approaches, such as synthetic data or "non-consumptive use"/remote code execution (see, e.g., van Atteveldt et al., 2020)

- Engage in **ongoing ethical reflection** throughout the research process.

- Consult **ethics guidelines** and **IRBs/review boards.**

# Further resources

- franzke, a. s., Bechmann, A., Zimmer, M., Ess, C. & the Association of Internet Researchers (2020). *Internet research: Ethical guidelines 3.0*. https://aoir.org/reports/ethics3.pdf

- Leslie, D. (2023). The ethics of computational social science. In E. Bertoni, M. Fontana, L. Gabrielli, S. Signorelli, & M. Vespe (Eds.), *Handbook of computational social science for policy* (pp. 57–104). Springer International Publishing. https://doi.org/10.1007/978-3-031-16624-2_4

- Moreno, M. A., Goniu, N., Moreno, P. S., & Diekema, D. (2013). Ethics of social media research: Common concerns and practical considerations. *Cyberpsychology, Behavior, and Social Networking*, *16*(9), 708–713. https://doi.org/10.1089/cyber.2012.0334

- Rau, J., Münch, F., & Asli, M. (2021). SOCRATES: Social Media Research Assessment Template for Ethical Scholarship. https://leibniz-hbi.github.io/socrates/

- Samuel, G., & Buchanan, E. (2020). Guest editorial: Ethical issues in social media research. *Journal of Empirical Research on Human Research Ethics*, *15*(1–2), 3–11. https://doi.org/10.1177/1556264619901215

# References

Akdeniz, E., Borschewski, K. E., Breuer, J., & Voronin, Y. (2023). Sharing social media data: The role of past experiences, attitudes, norms, and perceived behavioral control. *Frontiers in Big Data*, *5*, 971974. https://doi.org/10.3389/fdata.2022.971974

Breuer, J., Al Baghal, T., Sloan, L., Bishop, L., Kondyli, D., & Linardis, A. (2021). Informed consent for linking survey and social media data—Differences between platforms and data types. *IASSIST Quarterly*, *45*(1), 1–27. https://doi.org/10.29173/iq988

Breuer, J., Bishop, L., & Kinder-Kurlanda, K. (2020). The practical and ethical challenges in acquiring and sharing digital trace data: Negotiating public-private partnerships. *New Media & Society*, *22*(11), 2058–2080. https://doi.org/10.1177/1461444820924622

Breuer, J., Borschewski, K., Bishop, L., Vávra, M., Štebe, J., Strapcova, K., & Hegedűs, P. (2021). *Archiving Social Media Data: A guide for archivists and researchers*. https://doi.org/10.5281/zenodo.6517880

Breuer, J., Stier, S., Lukito, J., Mangold, F., Wieland, M., Radovanović, D., Radovanović, D., Zens, M., Breuer, J., Weller, K., & Wagner, C. (2025). *Overview of Ethical Considerations when Working with Digital Behavioral Data (GESIS Guides to Digital Behavioral Data, 14)* (Version 1.0). GESIS - Leibniz-Institute for the Social Sciences. https://doi.org/10.60762/GGDBD25014.1.0

Breuer, J., Weller, K., & Kinder-Kurlanda, K. (2023). The Role of Participants in Online Privacy Research: Ethical and Practical Consideration. In S. Trepte & P. K. Masur (Eds.), *The Routledge Handbook of Privacy and Social Media* (pp. 314–323). Routledge. https://www.taylorfrancis.com/chapters/oa-edit/10.4324/9781003244677-35/role-participants-online-privacy-research-johannes-breuer-katrin-weller-katharina-kinder-kurlanda

Bishop, L., & Gray, D. (2017). Chapter 7: Ethical Challenges of Publishing and Sharing Social Media Research Data. In K. Woodfield (Ed.), *Advances in Research Ethics and Integrity* (Vol. 2, pp. 159–187). Emerald Publishing Limited. https://doi.org/10.1108/S2398-601820180000002007

Carrière, T. C., Boeschoten, L., Struminskaya, B., Janssen, H. L., De Schipper, N. C., & Araujo, T. (2024). Best practices for studies using digital data donation. *Quality & Quantity*. https://doi.org/10.1007/s11135-024-01983-x

Emmerich, N. (2020). A Professional Ethics for Researchers? In R. Iphofen (Ed.), *Handbook of Research Ethics and Scientific Integrity* (pp. 751–767). Springer International Publishing. https://doi.org/10.1007/978-3-030-16759-2_34

# References

Fiesler, C., & Proferes, N. (2018). "Participant" Perceptions of Twitter Research Ethics. *Social Media + Society*, *4*(1), 205630511876336. https://doi.org/10.1177/2056305118763366

Fiesler, C., Zimmer, M., Proferes, N., Gilbert, S., & Jones, N. (2024). Remember the Human: A Systematic Review of Ethical Considerations in Reddit Research. *Proceedings of the ACM on Human-Computer Interaction*, *8*(GROUP), 1–33. https://doi.org/10.1145/3633070

Hox, J. J. (2017). Computational Social Science Methodology, Anyone? *Methodology*, *13*(Supplement 1), 3–12. https://doi.org/10.1027/1614-2241/a000127

Iphofen, R. (2020). An Introduction to Research Ethics and Scientific Integrity. In R. Iphofen (Ed.), *Handbook of Research Ethics and Scientific Integrity* (pp. 3–13). Springer International Publishing. https://doi.org/10.1007/978-3-030-16759-2_62

Israel, M. (2015). *Research ethics and integrity for social scientists: beyond regulatory compliance* (2nd edn.). SAGE. https://doi.org/10.4135/9781473910096.

Knöpfle, P., Haim, M., & Breuer, J. (2024). Key topic or bare necessity? How Research Ethics are Addressed and Discussed in Computational Communication Science. *Publizistik*, *69*(3), 333–356. https://doi.org/10.1007/s11616-024-00846-7

Lisker, M., & Mihaljević, H. (2025). *Data Ethics in the Fediverse: Analyzing the Role of Instance Policies in Mastodon Research* (Version 1). arXiv. https://doi.org/10.48550/ARXIV.2505.07606

Marwick, A. E., & boyd, d. (2011). I tweet honestly, I tweet passionately: Twitter users, context collapse, and the imagined audience. *New Media & Society*, *13*(1), 114–133. https://doi.org/10.1177/1461444810365313

McKee, H. A., & Porter, J. E. (2009). *The ethics of internet research: A rhetorical, case-based process* (Vol. 59). Peter Lang.

Menchen-Trevino, E. (2013). Collecting vertical trace data: Big possibilities and big challenges for multi-method research. *Policy & Internet*, *5*(3), 328–339. https://doi.org/10.1002/1944-2866.poi336

Ohme, J., Araujo, T., Boeschoten, L., Freelon, D., Ram, N., Reeves, B. B., & Robinson, T. N. (2023). Digital Trace Data Collection for Social Media Effects Research: APIs, Data Donation, and (Screen) Tracking. *Communication Methods and Measures*, *18*(2), 124–141. https://doi.org/10.1080/19312458.2023.2181319

Rau, J., Münch, F., & Asli, M. (2021): Social Media Research Assessment Template for Ethical Scholarship (SOCRATES): Your politely asking data ethics guide. (Social) Media Observatory. https://leibniz-hbi.github.io/socrates/

# References

Salganik, M. J. (2019). *Bit by bit: social research in the digital age*. Princeton University Press.

Schlütz, D., & Möhring, W. (2018). Between the devil and the deep blue sea: negotiating ethics and method in communication research practice. *SCM Studies in Communication and Media*, *7*(1), 31–58. https://doi.org/10.5771/2192-4007-2018-1-31.

Sloan, L., Jessop, C., Al Baghal, T., & Williams, M. (2020). Linking Survey and Twitter Data: Informed Consent, Disclosure, Security, and Archiving. *Journal of Empirical Research on Human Research Ethics*, *15*(1–2), 63–76. https://doi.org/10.1177/1556264619853447

Stier, S., Breuer, J., Siegers, P., & Thorson, K. (2020). Integrating Survey Data and Digital Trace Data: Key Issues in Developing an Emerging Field. *Social Science Computer Review*, *38*(5), 503–516. https://doi.org/10.1177/0894439319843669

van Atteveldt, W., Althaus, S., & Wessler, H. (2020). The Trouble with Sharing Your Privates: Pursuing Ethical Open Science and Collaborative Research across National Jurisdictions Using Sensitive Data. *Political Communication*, 1–7. https://doi.org/10.1080/10584609.2020.1744780

van Driel, I. I., Giachanou, A., Pouwels, J. L., Boeschoten, L., Beyens, I., & Valkenburg, P. M. (2022). Promises and Pitfalls of Social Media Data Donations. *Communication Methods and Measures*, *16*(4), 266–282. https://doi.org/10.1080/19312458.2022.2109608

Wähner, M., Deubel, A., Breuer, J., & Weller, K. (2024). "Don't research us"—How Mastodon instance rules connect to research ethics. *Publizistik*, *69*(3), 357–380. https://doi.org/10.1007/s11616-024-00855-6

Wagner, C., Stier, S., Zens, M., Radovanović, D., Zens, M., Breuer, J., Weller, K., & Wagner, C. (2025). *What is Digital Behavioral Data? (GESIS Guides to Digital Behavioral Data #1)* (Version 1.0, p. 16 pages) [Application/pdf]. GESIS - Leibniz Institute for the Social Sciences. https://doi.org/10.60762/GGDBD25001.1.0

Williams, M. L., Burnap, P., & Sloan, L. (2017). Towards an Ethical Framework for Publishing Twitter Data in Social Research: Taking into Account Users' Views, Online Context and Algorithmic Estimation. *Sociology*, *51*(6), 1149–1168. https://doi.org/10.1177/0038038517708140

Zimmer, M. (2010). "But the data is already public": On the ethics of research in Facebook. *Ethics and Information Technology*, *12*(4), 313–325. https://doi.org/10.1007/s10676-010-9227-5

Zong, J., & Matias, J. N. (2022). Bartleby: Procedural and Substantive Ethics in the Design of Research Ethics Systems. *Social Media + Society*, *8*(1), 205630512210770. https://doi.org/10.1177/20563051221077021

# Thank you for your attention!

## Looking forward to your comments & questions!

**Contact**

johannes.breuer@gesis.org

https://www.johannesbreuer.com/

Leibniz
Association