

# Sharing Selves: Developing an Ethical Framework for Curating Social Media Data

Sara Mannheimer  
Montana State University

Elizabeth A. Hull  
Dryad Digital Repository

## Abstract

Open sharing of social media data raises new ethical questions that researchers, repositories and data curators must confront, with little existing guidance available. In this paper, the authors draw upon their experiences in their multiple roles as data curators, academic librarians, and researchers to propose the STEP framework for curating and sharing social media data. The framework is intended to be used by data curators facilitating open publication of social media data. Two case studies from the Dryad Digital Repository serve to demonstrate implementation of the STEP framework. The STEP framework can serve as one important 'step' along the path to achieving safe, ethical, and reproducible social media research practice.

*Received 20 October 2016 ~ Revision received 23 January 2017 ~ Accepted 23 January 2017*

Correspondence should be addressed to Elizabeth Hull, PO Box 585, Durham NC 27702. Email: [ehull@datadryad.org](mailto:ehull@datadryad.org)

An earlier version of this paper was presented at the 12<sup>th</sup> International Digital Curation Conference.

The *International Journal of Digital Curation* is an international journal committed to scholarly excellence and dedicated to the advancement of digital curation across a wide range of sectors. The IJDC is published by the University of Edinburgh on behalf of the Digital Curation Centre. ISSN: 1746-8256. URL: <http://www.ijdc.net/>

Copyright rests with the authors. This work is released under a Creative Commons Attribution Licence, version 4.0. For details please see <https://creativecommons.org/licenses/by/4.0/>



## Introduction and Background

In a networked society – and especially in the online communities facilitated by social media – human thoughts and activities take the form of data that can be scraped, downloaded, aggregated, and otherwise collected on a massive scale. Academic researchers have identified this data as a potential source of insight into human behavior, and social media data is increasingly being used for scholarly inquiry (Kietzmann, Silvestre, McCarthy, and Pitt, 2012; Zimmer and Proferes, 2014; Ngai, Tao, and Moon, 2015). At the same time, funding agencies and academic journals are implementing data sharing policies (NSF, 2011; PLOS, 2014; Bill and Melinda Gates Foundation, 2015), and the scientific community is embracing data sharing as a strategy to promote research reproducibility (Collins and Tabak, 2014; Ioannidis, 2005).

Research with social media data doesn't neatly fit into the traditional definition of human subject data outlined decades ago by the Belmont Report (1979) and the Common Rule (1991) (Metcalf and Crawford, 2016; Shilton and Sayles, 2016). When users post to social media, they create data that can be mined by researchers using computational methods rather than more conventional social science research methods like interviews, surveys, ethnographic observation, or close reading of texts (Bruns, 2013). While social media data is often publicly available, social media users may not understand that their posts are being collected and used for research purposes. Moreover, social media users may not intend for their posts to reach beyond their online community.

Some researchers have experienced negative reactions when publishing social media data without proper protections to subjects. The 'Tastes, Ties, and Time' dataset (Lewis et al., 2008), comprising of Facebook user data and published on Harvard's Dataverse, was ultimately taken down due to privacy concerns (Zimmer, 2010). In 2016, when an Aarhus University graduate student scraped the online dating website OkCupid and released the data using the Open Science Framework, the public response was swift and critical (Markham, 2016); the dataset was subsequently taken down. To avoid such backlash and to protect human subjects, the data curation community needs better documentation and guidelines surrounding what Anatoliy Gruzd calls "social media data stewardship" (2016).

The Society of American Archivists (2016) and the Council on Library and Information Resources (Besek, 2003) have both released resources to guide ethical practice for digital archives in general, and the Social Media Archiving Toolkit from North Carolina State University provides ethical and legal guidelines for social media archives in particular (2014). Mannheimer, Young, and Rossmann (2016) propose an ethical framework for researchers using social media data, structured around three points: (1) context, including social media platform and disciplinary norms in the researchers' fields; (2) expectation of social media users; and (3) a value analysis that weighs the benefits of the research against the potential privacy risks to users. Weller and Kinder-Kurlanda's (2016) framework for sharing social media data is an excellent resource aimed at social media researchers. However, the literature does not yet include ethical guidelines tailored specifically to data curators.

The open data movement operates under the belief that open data is a common good, and data sharing is becoming more widespread, encouraged in large part by funding agency and journal policies. However, there remains a lack of clarity about

human subject privacy for data that lies outside the traditional realm of Institutional Review Boards. In particular, sharing social media data presents unique challenges regarding sensitive topics, transparency of documentation, user privacy expectations, and social media platform policies. This paper introduces the STEP (Sensitivity, Transparency, Expectation of privacy, Platform) Framework, designed to help data curators in open access repositories operate within these gray areas, balancing the benefits of open data with the potential risks to social media users. Two case studies from the Dryad Digital Repository serve to demonstrate implementation of the STEP framework.

### Social Media Data in the Dryad Digital Repository

The Dryad Digital Repository is a useful point of reference for exploring the ethics of data sharing. Dryad is a general purpose repository that provides unrestricted access to data. Dryad content includes openly published datasets associated with social media research, including data collected from Twitter, Facebook, Instagram, YouTube and Flickr. Dryad submitters are responsible for aligning the content of their data publications with Dryad's policies (Dryad, 2016), which state that "human subject data must be properly anonymized and prepared under applicable legal and ethical guidelines" (Dryad, 2016). In addition, Dryad's curation team reviews datasets prior to publication and assists researchers in achieving a level of subject anonymity that can be considered 'safe.' An increasing number and diversity of submissions of this type have highlighted the need for a framework to help structure curator inquiry around ethical publishing of social media data.

## Guiding Principles and STEP Framework

The STEP Framework helps guide curators through ethical inquiry when assessing social media data for the purpose of open archiving. While some repositories (e.g. ICPSR<sup>1</sup>, Qualitative Data Repository<sup>2</sup>, and UK Data Service<sup>3</sup>) can provide restricted access for sensitive data, this framework focuses on curating fully open access data. The STEP Framework aims to help curators think through ethical challenges regarding social media data, with the ultimate goal of encouraging open data sharing for social media researchers.

### Guiding Principles

The framework operates under three high-level principles:

- **Value analysis:** When sharing social media data, researchers and data curators must measure the benefits of sharing data against the potential risks to human subjects.
- **Responsibility:** Data curators can help educate researchers about ethical data sharing, but researchers themselves are ultimately responsible for the data they share.

---

<sup>1</sup> ICPSR: <http://icpsr.umich.edu/>

<sup>2</sup> Quantitative Data Repository: <https://qdr.syr.edu/>

<sup>3</sup> UK Data Service: <https://www.ukdataservice.ac.uk/>

- **Continual inquiry:** Ethical practice requires ongoing dialogue and examination.

### **Principle 1: Value analysis**

Open data and user privacy are both ethical imperatives. But data sharing and research reproducibility may stand at odds with ethical and legal concerns regarding social media data. In many cases, as more privacy measures are implemented, social media data becomes less fit for confirming reproducibility (Weller and Kinder-Kurlanda, 2016). As the UKAN anonymisation decision making framework suggests, “zero risk is not a realistic possibility if you are to produce useful data” (Elliot, Mackey, O’Hara, and Tudor, 2016). When sharing social media data, researchers and curators must therefore strike a balance between data openness and user privacy.

### **Principle 2: Responsibility**

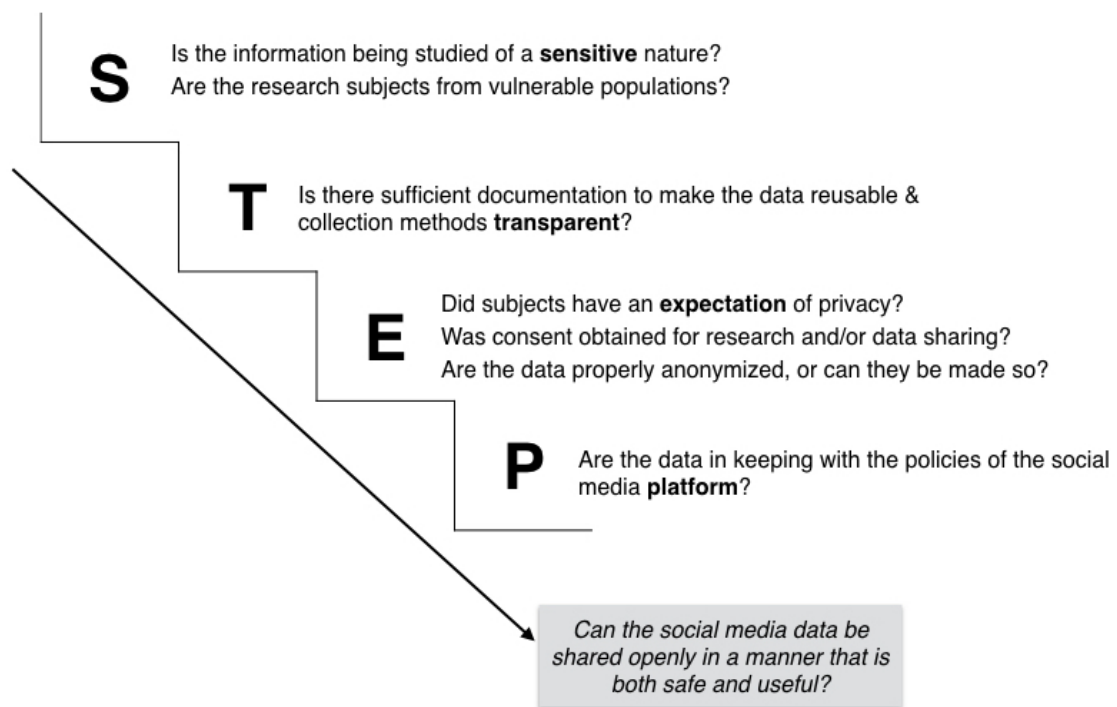
Data curation and data review are key quality-control elements in the data publication process. Curators should err on the side of caution when curating social media data and other ethically-complex data – even if the research was conducted ethically, curators can’t assume that researchers or IRBs have considered the ethical implications specific to sharing data. When necessary, curators should contact researchers to request better data documentation or de-identification of variables. However, data curators cannot be expected to be ethics experts. Mistakes will happen, and even good faith efforts can fall short (Zimmer, 2010). While data curators and data repositories have a role to play in educating researchers and promoting ethical data publishing, it is ultimately the responsibility of researchers to ensure that their data is shared ethically.

### **Principle 3: Continual inquiry**

The Society of American Archivists’ Code of Ethics encourages archivists to “consult with colleagues, relevant professionals, and communities of interest to ensure that diverse perspectives inform their actions and decisions” (2012). Data curators – as archivists of research data – will benefit from similar practice. Data curators should consult within their curation team and discuss details with researchers who submit data. Data curators may also benefit from consulting with data curators at other repositories and reaching out to professionals in data- and ethics-related fields. Policies and practices surrounding social media data are very much in flux, and will likely remain so. Continually discussing and reevaluating ethical standards will help curators stay up-to-date with ethical norms.

## **The STEP Framework for Data Curators**

The STEP framework is structured around four key areas of inquiry for data curators: Sensitivity, Transparency, Expectation of privacy, and Platform (STEP) (see Figure 1). This framework is not meant to provide hard and fast rules, but rather aims to improve practice and manage risk for data repositories, researchers, and social media users.



**Figure 1.** Visualization of the STEP framework for curating social media data.

### **Sensitivity**

Social media data relating to sensitive topics or collected from vulnerable populations requires that data curators examine the data with a particular focus on potential risks to users.

### ***Sensitive topics***

Sensitive topics require increased vigilance regarding privacy and anonymity. Lee and Renzetti suggest four areas in which research is likely to be threatening to subjects: (1) when research intrudes into the private sphere or delves into some deeply personal experience; (2) when the study is concerned with deviance or social control; (3) when the study impinges on the vested interests of powerful persons or the exercise of coercion or domination; (4) when the research deals with things that are sacred to those being studied that they do not wish profaned (1993).

### ***Vulnerable populations***

Research data collected from vulnerable populations who are susceptible to exploitation should also be considered sensitive (Belmont Report, 1979; World Medical Association, 2008). Mechanic and Tanner suggest that subject vulnerability can result from “developmental problems, personal incapacities, disadvantaged social status, inadequacy of interpersonal networks and supports, degraded neighborhoods and environments, and the complex interactions of these factors over the life course” (2007). Vulnerable populations have less power in the research process and less power over what happens to their data. Researchers and data curators therefore take on more responsibility regarding data privacy (Elliot, Mackey, O’Hara, and Tudor, 2016). When dealing with social media data, this aspect arises most often with regard to minors, who

tend to be active users of social networking sites and have different privacy expectations than adults (boyd, 2014).

### **Transparency**

Transparent data documentation facilitates ethical data sharing and ethical data reuse. For researchers, transparency includes clearly documenting the data collection methodology, anonymization processes, and ethical considerations, as well as providing ReadMe files or codebooks that help others understand the data being shared. Curators should encourage researchers to include documentation as part of their data publication. When researchers are transparent about their process, they support a culture of openness, facilitate data reuse, and help educate other researchers about methods for ethical data sharing. Further, Rivers and Lewis suggest that transparency regarding social media research can help foster ‘privacy literacy’ so that the users can make informed decisions about participating (2014). Ideally, curators should also clearly document their own decisions and activities over the course of reviewing and publishing the data.

### **Expectation of privacy**

In the context of online social networks, the public and the private become intertwined (Zimmer, 2010; Rivers and Lewis, 2014). While social media posts are available in public forums, social media users may not expect that their posts are being seen beyond their perceived online community (boyd, 2014). As Zimmer writes, “just because personal information is made available in some fashion on a social network, does not mean it is fair game for capture and release to all” (2010).

Each social media platform functions differently with regard to privacy. Some social media platforms – such as Facebook – are ‘closed networks,’ with customizable privacy settings. Other platforms – such as Twitter – are publicly-visible by default. (Twitter users may opt to protect their accounts, limiting access to a select group of followers, but few do so.<sup>4</sup>) Many social media sites support hashtags, which reach a broader audience, and @-mentions, which address specific users. Some social media sites allow pseudonyms<sup>5</sup>, while others require real names<sup>6</sup>. Some social media platforms provide easy access to user data, which encourages data collection and research<sup>7</sup>. All of these platform-specific usage norms can affect a user’s expectation of privacy.

Regardless of platform, user expectations are key to determining the sensitivity of the data – lower user expectation of privacy makes the data less sensitive. Politicians, celebrities, or organizations likely expect that their social media posts will be read by a wide audience. Private citizens, on the other hand, may not expect that their posts will be viewed by audiences beyond her immediate social network. For example, when Freelon, McIlwain, and Clark collected Twitter data documenting the Black Lives Matter movement, they attempted to honor users’ expectations of privacy by publishing only Tweets that had been widely shared, from Twitter users with a large number of followers (2016). While strategies like these are helpful, there will always be ambiguity in determining user intention, and user expectations may change over time. The most

4 In 2009, Tech Crunch concluded that the percentage of protected accounts on Twitter was about 10% (<https://techcrunch.com/2009/10/05/twitter-data-analysis-an-investors-perspective-2/>). A 2013 Pew survey found that 24% of teens had protected Twitter accounts. (<http://www.pewinternet.org/2013/05/21/teens-social-media-and-privacy/>).

5 For example, Google+ (<https://plus.google.com/+googleplus/posts/V5XkYQYYJqy>) and Twitter (<https://twitter.com/en/privacy>).

6 For example, Facebook (<https://www.facebook.com/help/112146705538576>).

7 Twitter’s API is a notable example (<https://dev.twitter.com/rest/public>).

unambiguous method for aligning research with user expectations is to obtain informed consent.

### ***Informed consent***

Curators should consider whether and how consent was obtained for the research before archiving social media data. The literature is split regarding the level of consent necessary for social media research. Rivers and Lewis (2014) assert that informed consent must be granted by each social media user whose posts are used for research purposes, suggesting that researchers “avoid qualitatively analyzing [social media] communications as if they are offered for research consumption without consent, because it does not align with the context in which the tweets were created.” Elliot, Mackey, O’Hara, and Tudor (2016) are more lenient, writing that, “given the current state of the information society, [obtaining informed consent] is both impractical and undesirable”; they suggest that a lack of informed consent for data-driven research does not necessarily preclude sharing, but merely makes data more sensitive. Hutton and Henderson (2015) suggest a model that applies Nissenbaum’s theory of contextual integrity, which states that people have “a right to live in a world in which [their] expectations about the flow of personal information are, for the most part, met” (2009). In their study of Facebook users, Hutton and Henderson used pop-up messages to evaluate participants’ willingness to share certain types of data, thus tailoring informed consent to each user’s expectations of privacy on Facebook (2015). The conversation surrounding informed consent will likely continue to evolve; data curators should stay abreast of the latest developments to inform dataset review.

### ***Anonymization***

Most open data repositories require that data be de-identified prior to submission (Dryad, 2016; ICPSR, 2012). The Anonymisation Decision-Making Framework from UKAN (Elliot, Mackey, O’Hara, and Tudor, 2016) provides detailed guidance that – although targeted at researchers – can also be helpful to data curators as they review data for publication. Social media data can be very difficult to anonymize (Zimmer, 2010). However, anonymization may not be strictly necessary with social media data, depending on social media users’ expectation of privacy. As noted in Principle 1: Value Analysis, curators should consult with data submitters to weigh the benefits of publishing the data against the risk that data that could be re-identified. And as noted in Principle 2: Responsibility, while curators review social media data to the best of their knowledge for de-identification issues, the ultimate responsibility falls on the data submitter.

### ***Platform***

Social media data is hosted by social media sites, each of which has unique privacy policies, terms of service, and developer agreements (Thomson, 2016). Some social media platforms’ terms of service limit what content can be published. For example, Twitter’s Developer Agreement and Policy states that developers who use their API “will only distribute or allow download of Tweet IDs and/or User IDs” (Twitter, 2016). Some researchers (Summers, 2014; Freelon, McIlwain, and Clark, 2016) have published only Tweet IDs, not only to align with Twitter’s policy, but also as a strategy to honor the intent of Twitter users. Since Twitter allows users to adjust their privacy settings at any time, users may delete posts or adjust privacy settings to limit the accessibility of posts. Published Twitter data should reflect the privacy choices of Twitter users.

Each social media platform's policies include different rules for data sharing. These terms and policies change over time, and some academic researchers choose to sidestep platform policies if they consider the benefit of their research to be worth the risk of violating terms of service (Kelley, Sleeper, and Cranshaw, 2013). Weller and Kinder-Kurlanda (2016) suggest establishing a dialogue between social media companies, researchers, and data repositories – ultimately aiming to “establish feasible interpretations [of terms of service] that allow researchers to at least share data for the sake of quality control and reproducibility.”

Ultimately, curators should aim to be aware of platform policies, but should take into consideration Principle 2: Responsibility – researchers are ultimately responsible for the data they collect and publish.

## Case Studies

Two case studies from Dryad – both of which deal with Twitter data – provide examples of using the STEP framework to review research data for publication.

### Case Study 1 – Data from *The Topology of a Discussion: The #occupy Case*

Gargiulo, Bindi, and Apolloni (2015a) used Twitter hashtag data to study the evolution of political discussion during and after the Occupy Wall Street movement. The associated Dryad data package (Gargiulo, Bindi, and Apolloni, 2015b) includes one .csv file containing three variables: ‘user,’ ‘hashtag,’ and ‘time.’

- **Sensitivity:** The research deals with active and public participation in a social movement and does not focus on a particular population.
- **Transparency:** No documentation is provided with the data package. Some information about data collection is included in the associated article, but details are insufficient and the content of the .csv file is not adequately explained. The method of analysis is laid out in detail in the article, which would hypothetically allow others to reproduce the results.
- **Expectation of privacy:** The use of hashtags on Twitter generally indicates one's desire to participate in a larger conversation and/or be identified with a concept or cause. The sample size is large (more than 37,000 users), and the risk of contributors being identified from the contents of the data file is low.
- **Platform policy compliance:** The .csv file contains a user ID which is described in the article as being “anonymous,” but there is no explanation of how this was derived. The file also contains the actual hashtags used, and Twitter policies are unclear on whether this information can be distributed to third parties.

### Conclusion

Given the low sensitivity of the research, the public nature of the discussion (and the platform) and the fact that an attempt was made to anonymize the data, the Dryad team concurred that this data package could be safely shared. Better documentation would have made the methods more reproducible and the data more useful. However, the



article provides a good model for similar network analyses of social movements. It is unclear whether the publication of this data strictly follows Twitter's policies, but any non-compliance is the responsibility of the authors.

### **Case Study 2 – Data from *In the Mood: The Dynamics of Collective Sentiments on Twitter***

This Dryad data package (Charlton, Singleton, and Greetham, 2016b) and its associated article (Charlton, Singleton, and Greetham, 2016a) present a study of the relationship between UK Twitter users' 'sentiment levels' and the network structure created by @-mentions. Based on statistical analysis of Twitter data, the researchers selected 18 'communities' to monitor and used these to formulate a model for "reproducing measures of emotive response." The data package contains several dynamic mention networks split over six tables; variables include an anonymised tweet ID, anonymised user IDs, and timestamps of tweets.

- **Sensitivity:** Topics being discussed by the selected communities are wide-ranging – from 'friends chatting' and 'dogs' to 'Islam versus atheism,' 'Gamergate' and 'smoking/e-cigarettes.'
- **Transparency:** The authors provided a ReadMe with the data package that explains the content of each file, and a section of the article describes in detail how the data were obtained.
- **Expectation of privacy:** The authors assert that tweets with @-mentions are public and may be read and commented on by any other user. They also argue that ethical approval was unnecessary for their research because "the human data ... analysed is in the public domain." However, the use of @-mentioning indicates communications intended for specific people, and implies an expectation of discussion within the user's specific network. While the tweet IDs and user IDs provided in the data package were anonymized, exact timestamps present a potential (though low) risk for re-identification.
- **Platform policy compliance:** Twitter policies are unclear on whether timestamps may be distributed to third parties.

### **Conclusion**

This case is an interesting one in terms of user expectations when engaged in what some might consider a 'private' conversation on a public platform. Some of the topics being discussed are sensitive, and many of those participating probably did not consider that their comments would be 1) broadcast to an audience beyond their immediate network and 2) collected and analyzed by researchers. Taking a hard line on informed consent, this study would likely not pass muster. However, given the fact that IDs were anonymized and that the research was presented in a transparent and reproducible way, the benefits of data publication were deemed greater than the risks.

## Future Work

The STEP framework should evolve over time. Future work could expand upon the current discussion of the theory and concepts involved in evaluating social media data for publication. In addition, the STEP framework would be strengthened by additional case studies examining social media data from a wide variety of repositories and social media platforms; additional case studies will help demonstrate expanded applicability for the framework. The authors also see a need for additional guidance that can complement the STEP framework's focus on social media data. The data sharing community will benefit from expanded frameworks that apply to general big data research, including social science data journalism.

## Conclusion

Sharing social media data helps ensure research reproducibility, advances science, and encourages research efficiency (Weller and Kinder-Kurlanda, 2016). Data sharing also facilitates equity of data access, narrowing the divide between the 'big data rich' and the 'big data poor' (boyd and Crawford, 2012; Metzler, Kim, Allum, and Denman, 2016). The STEP framework encourages open data for the public good by providing curators with guidelines for assessing data submissions according to Sensitivity, Transparency, Expectation of privacy, and Platform. Curators using the framework are encouraged to think critically and carefully when reviewing social media data for publication, taking into consideration the three guiding principles of the framework: Value Analysis, Responsibility, and Continual Inquiry. Curators must continue to stay informed about social media research practice, and should keep an active dialogue with researchers, other data curators, archivists and librarians, and ethicists. Due to the quickly-evolving nature of the field, the authors envision the STEP framework as just one important 'step' along the path to achieving safe, ethical, and reproducible social media research practice.

## Acknowledgements

The authors would like thank the Dryad curation team for their input and assistance with case studies. We would also like to thank Todd Vision and the Skylight writing group (Doralyn Rossmann, Justin D. Shanks, and Scott W. H. Young) for their feedback and support. This work was supported by NSF grants DBI-1612608 and DBI-1564925 to Dryad.

## References

- Besek, J. (2003). Copyright issues relevant to the creation of a digital archive: A preliminary assessment. Council on Library and Information Resources and Library of Congress. Retrieved from <https://www.clir.org/pubs/reports/pub112/body.html>

- Bill and Melinda Gates Foundation. (2015). Open access policy. Retrieved from <http://www.gatesfoundation.org/How-We-Work/General-Information/Open-Access-Policy>
- boyd, d. (2014). *It's complicated: The social lives of networked teens*. New Haven: Yale University Press.
- boyd, d., & Crawford, K. (2012). Critical questions for big data: Provocations for a cultural, technological, and scholarly phenomenon. *Information, Communication & Society*, 15(5): 662-679. doi:10.1080/1369118X.2012.678878
- Bruns, A. (2013). Faster than the speed of print: Reconciling 'big data' social media analysis and academic scholarship. *First Monday*, 18(10). doi:10.5210/fm.v18i10.4879
- Charlton, N., Singleton, C., & Greetham, D.V. (2016a). In the mood: The dynamics of collective sentiments on Twitter. *Royal Society Open Science*, 3(6): 160162. doi:10.1098/rsos.160162
- Charlton, N., Singleton, C., & Greetham, D.V. (2016b). *In the mood: The dynamics of collective sentiments on Twitter* [data set]. Dryad Digital Repository. doi:10.5061/dryad.5302r
- Collins, F.S., & Tabak, L.A. (2014). NIH plans to enhance reproducibility. *Nature*, 505(7485): 612. Retrieved from <http://www.nature.com/news/policy-nih-plans-to-enhance-reproducibility-1.14586>
- Dryad Frequently Asked Questions. (2016). Submitting data: What kinds of data does Dryad accept? Retrieved from <http://datadryad.org/pages/faq>
- Dryad Terms of Service. (2016). Submitter obligations, representations and warranties. Retrieved from <http://datadryad.org/pages/policies#submitter-obligations>
- Elliot, M., Mackey, E., O'Hara, K., & Tudor, C. (2016). *The anonymisation decision-making framework*. UKAN Publications. Retrieved from <http://ukanon.net/wp-content/uploads/2015/05/The-Anonymisation-Decision-making-Framework.pdf>
- Freelon, D.G., McIlwain, C.D., & Clark, M.D. (2016). *Beyond the hashtags: #Ferguson, #Blacklivesmatter, and the online struggle for offline justice*. Center for Media & Social Impact. Retrieved from <http://cmsimpact.org/resource/beyond-hashtags-ferguson-blacklivesmatter-online-struggle-offline-justice>
- Gargiulo, F., Bindi, J., & Apolloni, A. (2015a). The topology of a discussion: The #occupy case. *PLOS ONE*, 10(9): e0137191. doi:10.1371/journal.pone.0137191
- Gargiulo, F., Bindi, J., & Apolloni, A. (2015b). *The topology of a discussion: The #occupy case* [data set]. Dryad Digital Repository. doi:10.5061/dryad.q1h04

- Gruzd, A. (2016). Defining social media data stewardship (SMDS). *Social Media Lab*. Retrieved from <http://socialmedialab.ca/2016/defining-social-media-data-stewardship-smds>
- Hutton, L., & Henderson, T. (2015). “I didn’t sign up for this!”: Informed consent in social network research. In Proceedings of the 9th International AAAI Conference on Web and Social Media. Retrieved from <http://www.aaai.org/ocs/index.php/ICWSM/ICWSM15/paper/view/10493>
- Inter-university Consortium for Political and Social Research (ICPSR). (2012). *Guide to social science data preparation and archiving phase 5: Preparing data for sharing*. Retrieved from <https://www.icpsr.umich.edu/icpsrweb/content/deposit/guide/chapter5.html>
- Ioannidis, J.P. (2005). Why most published research findings are false. *PLoS Medicine*, 2(8): e124. doi:10.1371/journal.pmed.0020124
- Kelley, P.G., Sleeper, M., & Cranshaw, J. (2013). *Conducting research on Twitter: A call for guidelines and metrics*. CSCW Measuring Networked Social Privacy Workshop. Retrieved from <http://patrickgagekelley.com/papers/twitter-pmj.pdf>
- Kietzmann, J.H., Silvestre, B.S., McCarthy, I.P., & Pitt, L.F. (2012). Unpacking the social media phenomenon: Towards a research agenda. *Journal of Public Affairs*, 12(2): 109–119. doi:10.1002/pa.1412
- Lee, R.M., & Renzetti, C.M. (1993). *Researching sensitive topics*. Newbury Park, CA: SAGE Publications, Inc.
- Lewis, K., Kaufman, J., Gonzalez, M., Wimmer, A., & Christakis, N. (2008). Tastes, ties, and time: A new social network dataset using Facebook.com. *Social networks*, 30(4): 330-342. doi:10.1016/j.socnet.2008.07.002
- Mannheimer, S., Young, S.W.H., & Rossmann, D. (2016). On the ethics of social network research in libraries. *Journal of Information, Communication, and Ethics in Society*, 14(2): 139–151. doi:10.1108/JICES-05-2015-0013
- Markham, A. (2016). OkCupid data release fiasco: It’s time to rethink ethics education. *Points: Data & Society*. Retrieved from <https://points.datasociety.net/okcupid-data-release-fiasco-ba0388348cd#>
- Mechanic, D., & Tanner, J. (2007). Vulnerable people, groups, and populations: Societal view. *Health Affairs*, 26(5): 1220-1230. doi:10.1377/hlthaff.26.5.1220
- Metcalf, J., & Crawford, K. (2016). Where are human subjects in big data research? The emerging ethics divide. *Big Data & Society*, 3(1): 1–14. Retrieved from <http://bds.sagepub.com/content/3/1/2053951716650211>

- Metzler, K., Kim, D.A., Allum, N., Denman, A. (2016). *Who is doing computational social science? Trends in big data research*. SAGE White Paper. Retrieved from <https://us.sagepub.com/sites/default/files/compsocsci.pdf>
- National Commission for the Protection of Human Subjects of Biomedical and Behavioral Research, Department of Health, Education and Welfare. (1979). *The Belmont report: Ethical principles and guidelines for the protection of human subjects of research*. Washington, DC: US Government Printing Office. Retrieved from <http://www.hhs.gov/ohrp/regulations-and-policy/belmont-report>
- National Science Foundation. (2011). NSF data management plan requirements. Retrieved from <http://www.nsf.gov/eng/general/dmp.jsp>
- Ngai, E.W., Tao, S.S., & Moon, K.K. (2015) Social media research: Theories, constructs, and conceptual frameworks. *International Journal of Information Management*, 35(1): 33-44. doi:10.1016/j.ijinfomgt.2014.09.004
- Nissenbaum, H. (2009). *Privacy in context: Technology, policy, and the integrity of social life*. Palo Alto, CA: Stanford University Press.
- North Carolina State University. (2014). Social media archiving toolkit. Retrieved from <https://www.lib.ncsu.edu/social-media-archives-toolkit>
- Public Library of Science. (2014). Data availability. Retrieved from <http://journals.plos.org/plosone/s/data-availability>
- Rivers, C.M., & Lewis, B.L. (2014). Ethical research standards in a world of big data. *F1000Research*, 3(38). doi:10.12688/f1000research.3-38.v2
- Shilton, K., & Sayles, S. (2016). “We aren’t all going to be on the same page about ethics”: Ethical practices and challenges in research on digital and social media. *2016 49th Hawaii International Conference on System Sciences (HICSS)*: 1909-1918. doi:10.1109/HICSS.2016.242
- Society of American Archivists. (2012). Code of ethics for archivists. Retrieved from <http://archivists.org/statements/saa-core-values-statement-and-code-of-ethics>
- Society of American Archivists. (2016). Case studies in archival ethics. Retrieved from <http://www2.archivists.org/groups/committee-on-ethics-and-professional-conduct/case-studies-in-archival-ethics>
- Summers, E. (2014). On forgetting and hydration. *Archivy*. Retrieved from <https://medium.com/on-archivy/on-forgetting-e01a2b95272>
- Thomson, S.D. (2016). Preserving social media: DPC Technology Watch Report 16-01 February 2016. doi:10.7207/twr16-01
- Twitter. (2016). Developer Agreement & Policy. Retrieved from <https://dev.twitter.com/overview/terms/agreement-and-policy>

Weller, K., & Kinder-Kurlanda, K.E. (2016). A manifesto for data sharing in social media research. *Proceedings of the 8th ACM Conference on Web Science*: 166-172. doi:10.1145/2908131.2908172

World Medical Association. (2008). Declaration of Helsinki. Retrieved from <https://www.wma.net/what-we-do/medical-ethics/declaration-of-helsinki/>

U.S. Department of Health and Human Services. (1991). *Federal policy for the protection of human subjects ('Common Rule')*. Retrieved from <http://www.hhs.gov/ohrp/regulations-and-policy/regulations/common-rule>

Zimmer, M., & Proferes, N.J. (2014). A topology of Twitter research: Disciplines, methods, and ethics. *Aslib Journal of Information Management*, 66(3): 250-261. doi:10.1108/AJIM-09-2013-0083

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