

Response to the NTIA Request for Comments on Pervasive Data

From the National Science Foundation-funded PERVADE Project: Katie Shilton, Matthew Bietz, Casey Fiesler, Jacob Metcalf, Jessica Vitak, and Michael Zimmer

Questions

1. What are the potential benefits of developing national-level ethical guidelines for researchers collecting, analyzing, and sharing pervasive data?

Ethical concerns about appropriate data use have been cited as a key obstacle to the progress of computational social science (Lazer, 2020). Developing national-level ethical guidelines for pervasive data research will help level the playing field for researchers who, at the moment, must navigate departmental, institutional, and field-specific best practices. This increases the costs of doing computational research for less experienced researchers and researchers who are at institutions without experienced senior mentors, or researchers in emerging fields that have not yet established norms for research with pervasive data.

Lazer, D. M. J., Pentland, A., Watts, D. J., Aral, S., Athey, S., Contractor, N., Freelon, D., Gonzalez-Bailon, S., King, G., Margetts, H., Nelson, A., Salganik, M. J., Strohmaier, M., Vespignani, A., & Wagner, C. (2020). Computational social science: Obstacles and opportunities. *Science*, 369(6507), 1060–1062. <https://doi.org/10.1126/science.aaz8170>

2. What are the potential drawbacks of developing national-level ethical guidelines for researchers collecting, analyzing, and sharing pervasive data?

A major challenge for national-level ethical guidelines is that our work has consistently found that the answer to whether data collection or use is ethical is, *it depends*. The ethics of pervasive data collection and use depend on a variety of contextual factors including platform, data type, consent conditions and data norms, research goals, and researcher positionality (Gilbert et al, 2023). Creating ethics guidelines that respect these variations is a challenge.

To help address this challenge, the PERVADE team has conducted research to outline these factors and created decision-support tools to lead researchers through the factors that matter to the ethics of pervasive data use. The factors include *regulation* (including both IRB and Terms of Service considerations); *user expectations* (including what users know and understand, what users feel is appropriate on various platforms; and how users interact with pervasive data); *community norms* (including accepted and controversial practices in diverse computing communities); *risks* (including current and future risks to both individuals and groups); and finally *justice* (including fairness and equity).

Gilbert, S., Shilton, K. and Vitak, J. 2023. When research is the context: Cross-platform user expectations for social media data reuse. *Big Data & Society*. 10, 1 (Jan. 2023), 32 pages. DOI:<https://doi.org/10.1177/20539517231164108>.

3. To what extent does the definition of *pervasive data* in this Request for Comments capture the appropriate scope for national ethical guidelines? Are there particular types of data or other digital artifacts^[42] that should be carefully considered or included/excluded in the definition?

Examples of pervasive data research include social media research, passive sensing research, personal sensing research, digital phenotyping, and computational social science using search histories, geolocation data, or wearables that record and/or sense personal behavior.

6. Consent and autonomy are key principles in human subjects research ethics. However, users of online services may be required to divulge certain personal information and/or have no ability to freely make decisions about its use.^[44] How should researchers working with pervasive data consider consent and autonomy?

The PERVADE project has published a discussion of problems of consent and autonomy for pervasive data ([Shilton et al, 2021](#)). We excerpt the argument here:

Trustworthy practice for pervasive data research—ensuring that researchers meet commitments like consent and autonomy—is problematized by the ecosystem where digital research takes place. Research participants routinely deny knowledge of widespread research conducted with digital data and express that, while they might be willing to participate in digital data research, they expect to be asked for consent ([Fiesler and Proferes, 2018](#); [Gilbert et al., 2021](#); [Hudson and Bruckman, 2004](#)). However, informed consent is not always logistically or philosophically appropriate for research in the big data age. Logistically, there is now a large amount of data about people available online. Securing individual consent to use this data would be incredibly challenging—if not impossible—where individual identity was knowable, and arguably unethical where doing so would require collecting even more personal data ([Ioannidis, 2013](#)).

Philosophically, informed consent for pervasive data research suffers from a number of problems. [Metcalf and Crawford \(2016\)](#) point out that codes of informed consent were established specifically to govern physician-researchers, who balance the broad social interest in research results with their individual duty of care for a patient. The procedures and norms that IRBs use to generate trust operate with an unstated assumption that these social conditions hold for all types of research. However, the trust relationships between computational social scientists, data scientists, and the public seldom conform to the social conditions that hold between physician-researchers and research subjects. The norms of pervasive data researchers (unlike, say, the norms of ethnographers) do not currently require preexisting, personal, or even explicitly declared relationships with the communities they study to collect data. The typical scale of data science manifests in numerous ethical and epistemic challenges for understanding the ethical interests of data subjects that extend beyond matters of logistics ([Hanna and Park, 2020](#)). Finally, [Richards and Hartzog \(2017\)](#) identify “pathologies” of consent resulting from overuse in the digital age. They argue that consent works best when it is given infrequently, when the harms are visceral and easily imagined, and when the stakes of a decision are significant. Pervasive data research meets some, but not all, of these standards. Explicit, informed consent to research participation happens infrequently for many participants. However, the harms of data research are rarely visceral or easily imagined. And it is

unclear to what degree individuals consider the stakes of participating in research. Therefore, while it is unclear whether informed consent is philosophically the right mechanism to navigate the relationship between data scientists and data subjects, it is the case that many of the norms and mechanisms that other forms of research use to achieve informed consent do not translate well to pervasive data research.

Pervasive data research is not the first field confronted with the shortcomings of traditional research ethics' conceptualization of informed consent. Ethnographic researchers have grappled with research ethics both within and beyond the framework of institutional review (Davies, 2012). Questions of whether and how to make the ethnographer's presence known to research subjects have been long debated in the literature (Bernard, 2006). Pervasive data research has more in common with ethnography than is immediately obvious. The instruments are different—human senses instead of digital sensors, individual sensemaking instead of algorithmic pattern-matching—but both forms of research rely on integration and interpretation of multiple data streams, and both require judgment about what features of a context are relevant for making meaning. The ethical challenges of research with pervasive data—the richly personal nature, the emphasis on observation, integration of multiple data types, and the drawing of inferences and conclusions based on patterns—are the same challenges that can be found in the world of ethnography and participant observation. Ethnographers have deep experience in building trust with research subjects and respecting autonomy with practices beyond informed consent. Data scientists can use this experience, and the practices of ethnographic intervention, to help define trustworthy practice for pervasive data research.

We outline how data scientists can draw on ethnographic intervention to navigate consent and autonomy in more depth in (Shilton et al, 2021).

Shilton, K., Moss, E., Gilbert, S.A., Bietz, M.J., Fiesler, C., Metcalf, J., Vitak, J. and Zimmer, M. 2021. Excavating awareness and power in data science: A manifesto for trustworthy pervasive data research. *Big Data & Society*. 8, 2 (Jul. 2021), 1–12.
DOI:<https://doi.org/10.1177/20539517211040759>.

a. What, if any, would be an appropriate consent model for research with pervasive data? How and how often should consent occur?

For particularly sensitive forms of pervasive data, researchers might consider donation models, where participants proactively opt to share data with researchers.

b. Are there alternative models to traditional consent that either support autonomy or provide protections for data subjects in cases where autonomy is limited?

Meaningful informed consent is one standard for raising data subjects' awareness of research. But for pervasive data researchers who can't secure meaningful consent because of scale, pervasiveness, or other issues, adaptations of entrée and participant checking should evolve with the field. As a starting point, they might include website pop-ups, such as those used as part of GDPR notification requirements, that explain the research, risk to participants, and allow

potential subjects to opt out. Another way to build awareness is to increase research subjects' participation in the research. Increasing participation will require researchers to think about what participation looks like for the community or population they study ([Sloane et al., 2020](#)).

7. What ethical issues and risks to privacy and other rights, and mitigation strategies, should be considered during the research design phase?

a. Users' concerns about researcher data access vary based on contextual factors.^[46] What contextual factors increase or alter the risks to data subjects in research using pervasive data? ^[47]

Contextual factors that alter the risk to data subjects in research include research domain, purpose of research and sensitivity of the research questions, and the sensitivity of the data type collected (Gilbert et al, 2023). In addition, researchers should consider whether data could be used to make predictions about a person which could cause future harms.

Gilbert, S., Shilton, K. and Vitak, J. 2023. When research is the context: Cross-platform user expectations for social media data reuse. *Big Data & Society*. 10, 1 (Jan. 2023), 32 pages. DOI:<https://doi.org/10.1177/20539517231164108>.

c. What power differences exist between researchers and data subjects, or between online service providers and data subjects, that could create unique risks and potential for harm.^[49] How should these differences be considered and mitigated during the research design phase?

PERVADE found that pervasive data research has entwined ethical problems: participant unawareness of research, and the relationship of pervasive data research to social, political and economic power. Below, we excerpt and summarize PERVADE's approach to addressing these ethical and social challenges (Shilton et al, 2021).

Researchers using pervasive data must consider 1) participant awareness and 2) power beyond traditional research ethics concerns. Awareness of pervasive data can be mapped on two spectra based on how digital traces are created: traces created in private to public settings, and traces created by intentional to automatic means. Private, intentional data trails are "secrets": for example, texts to a spouse or family photos. People are aware they are creating communications or documentation, and also make efforts not to share them widely. Public, intentional data trails are "broadcasts": for example, posts to social platforms. People are aware they are creating communication or documentation and purposefully share them widely, even if they may not know their reach (Proferes & Fiesler, 2018). Private, automatic data trails are "espionage": communications or documentation created without human intervention or awareness and not widely shared. Examples include the data collected by a smart fridge or thermostat. "Espionage" data can also include geolocation or telemetry data collected as part of the functioning of devices like smartphones. Though users may be aware of functions that require such documentation, they may not know the extent to which it is being collected (Hannay and Baatard, 2011). Finally, public automatic data can be thought of as "exhaust": documentation captured in public that individuals are not aware they are putting out into

the world. Examples include CCTV camera recordings, satellite images, or automated license plate readers. We recommend that researchers account for where their data fall on the private/public and automatic/intentional spectra and use this reflection as a guide for considering both awareness and power implications of their research.

Pervasive data researchers and the institutions that support them should also explicitly consider power relations and representational justice: they must consider the appropriateness of converting digital traces into research data. Pervasive data researchers should consider whether it is appropriate to make a given community, stakeholder group, or population more vulnerable either by creating new forms of data (which may be used by other parties to increase their vulnerability) or through secondary uses of data: by making research data out of traces or communications created for other purposes. This consideration might involve spending time in a community to understand their norms, collaborating with a community to serve their needs, or speaking to community gatekeepers to understand specific harms—for example, lack of research reflexivity for the culture and historical context of data (Klassen & Fiesler, 2022).

Klassen, S., & Fiesler, C. (2022). “This isn’t your data, friend”: Black Twitter as a case study on research ethics for public data. *Social Media+ Society*, 8(4), 20563051221144317.

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

Hannay, P., & Baatard, G. (2011). GeoIntelligence: Data Mining Locational Social Media Content for Profiling and Information Gathering. International Cyber Resilience Conference. <https://ro.ecu.edu.au/icr/20>

Shilton, K., Moss, E., Gilbert, S. A., Bietz, M. J., Fiesler, C., Metcalf, J., Vitak, J., & Zimmer, M. (2021). Excavating awareness and power in data science: A manifesto for trustworthy pervasive data research. *Big Data & Society*, 8(2), 1–12.

<https://doi.org/10.1177/20539517211040759>

e. What other vulnerable communities or vulnerability risk factors warrant additional consideration when conducting research with pervasive data? Please explain.

We excerpt background research on vulnerability drawn from the PERVADE Project’s publication [Shilton et al, 2021](#) here:

Research indicates that unwillingness to participate in pervasive data research is a larger concern among marginalized communities, where issues range from fear of surveillance and deportation (Nebeker et al., 2017b) to concerns that deployed technologies will fail to represent the needs and realities of user communities ([Winchester, 2018](#)) and to unwanted amplification of content or communities ([Dym and Fiesler, 2020](#)). Moreover, as [Hoffman and Jonas \(2017\)](#) point out, the costs of online participation are unequally borne by women and people of color, obligating researchers to consider the differential needs of vulnerable data subjects. Both distrust in platforms and concern for the uneven risks of surveillant research methods signal challenges of social power—who bears risk, and how partnerships with platforms shape that risk—that researchers must navigate.

Pervasive data researchers should consider whether it is appropriate to make a given community, stakeholder group, or population more vulnerable either by creating new forms of data (which may be used by other parties to increase their vulnerability) or through secondary uses of data: by making research data out of traces or communications created for other purposes. This consideration might involve spending time (virtually or physically) in a community to understand their norms, collaborating with a community to serve their needs, or speaking to community gatekeepers to understand specific harms—for example, amplifying content beyond its intended audience ([Dym and Fiesler, 2020](#)).

Using pervasive data for research increases the vulnerability of the people included (and potentially other people like those included), whether by amplifying their behaviors and beliefs, showing new connections or inferences between their activities and habits, or applying categories or labels to their actions. Scholars are increasingly developing approaches to help researchers think through implications and harms of datafication, such as the Omidyar Network's *Ethical Explorer Pack* ([Artefact Group, n.d.](#)) and [Wong et al.'s \(2017\)](#) privacy by design workbooks. And statements of potential vulnerabilities, harms, and biases are increasingly required in technical research communities ([Gibney, 2020](#)). It is important that scientists continue to use pervasive data to study those who need increased representation in knowledge (e.g. people with rare diseases, groups otherwise marginalized in research). Even so, pervasive data research should emphasize the standard drawn from political representation movements and disability activism: nothing about us without us ([Charlton, 2004](#)).

Of course, determining one's power relative to research subjects, platforms, and state actors is a complex process, and we realize that such reflection is a difficult task for many without a background in theories of power. However, there are numerous helpful frameworks for thinking through issues of power adapted specifically for digital data research, including anti-essentialism ([Neyland, 2016](#)), feminism ([D'Ignazio and Klein, 2020](#); [Hoffman and Jonas, 2017](#)), anti-racism ([Benjamin, 2019](#); [Hanna et al., 2020](#)), anti-colonialism ([Dourish and Mainwaring, 2012](#)), and queer theory ([Brim and Ghaziani, 2016](#)). These frameworks can provide concrete guidance to researchers considering how their protocols might unevenly subject participants to increased vulnerability.

h. How can researchers best conduct research with pervasive data in a way that engages the community, users, and data subjects.^[51] What are the best practices for such participatory research that uses pervasive data? What are the challenges and/or barriers to conducting participatory research? What important research questions cannot be answered using participatory mechanisms, and why?

Resources that could support pervasive data researchers in responsible data collection and use include:

- Novel privacy-aware data collection platforms, including sensor networks, wearables, IoT devices, APIs.
- Data archives and infrastructure (such as “clean rooms” for privacy-preserving data analysis) to ensure safe and ethical access to digital data for researchers.

- Data donation platforms for safe and ethical collection and use of mobile phone and device data.
- Open source technical tools and best practice documentation for privacy protection and assessing privacy risks for AI models and applications.

Participatory action research ([Khanlou and Peter, 2005](#)) has grappled extensively with questions of participation, motivation, and accessibility, and can guide data scientists on challenging questions such as how to define a community, how to structure participation, and how to ensure representation of stakeholders across a community.

8. What are the risks and mitigation measures related to pervasive data acquisition and access?

Alarm over digital datafication was not (primarily) created by researchers. Influential works like Zuboff; (2019) *The Age of Surveillance Capitalism: The Fight for a Human Future at the New Frontier of Power* have introduced a broad public to the ways companies and governments use records of online activity to learn about, influence, and predict behavior. But though public alarm is largely a response to corporate data practices, data scientists are directly impacted. [Hallinan et al. \(2020\)](#) found that people's angry reactions to the emotional contagion study were deeply bound up in opinions of the platform, Facebook, as much as or perhaps more so than the research itself. Public distrust in the platform bled into public distrust in the research.

Parallel to the growing public alarm about datafication has been a restriction of research access to some forms of pervasive data by social media platforms. As [Tromble \(2021\)](#) traces, platforms have moved away from open API access for researchers to narrower, more careful data access efforts enabled by tools like differential privacy ([King and Persily, 2020](#)). [Tromble \(2021\)](#) characterizes the “post-API” era as an opportunity for reflection by pervasive data researchers on the rigor and ethics of their data practices.

We must also recognize the ethical ramifications of exclusion from pervasive data research. The medical research field has come to recognize the consequences of excluding various groups (women, children, patients with chronic conditions or disabilities, etc.) from clinical trials (Zuckerman, 2009). While the size of pervasive data sets can give the illusion of a generalized sample, many of the systems and platforms that produce pervasive data have embedded biases that can lead to exclusion both from specific research studies and the larger research ecosystem. Lerman (2013) points out that many groups of people “remain on big data’s periphery. Their information is not regularly collected or analyzed, because they do not routinely engage in activities that big data is designed to capture.” Pervasive data might be collected only if someone has paid to join a specific platform or purchased a particular device. Some communities are over- or under-represented on particular platforms. Some pervasive data collection techniques (like only using public data) may unintentionally end up excluding marginalized communities. This can lead to discrimination in research conducted with these datasets and in algorithms trained on them (Favaretto, et al, 2019; Williams, et al, 2018). Researchers should consider the ways in which pervasive data may intentionally or unintentionally exclude various social groups and communities.

Favaretto, M., De Clercq, E., & Elger, B. S. (2019). Big Data and discrimination: Perils, promises and solutions. A systematic review. *Journal of Big Data*, 6(1), 12. <https://doi.org/10.1186/s40537-019-0177-4>

Lerman, J. (2013). Big Data and Its Exclusions. *Stanford Law Review Online*, 66, 55–64.

Williams, B. A., Brooks, C. F., & Shmargad, Y. (2018). How Algorithms Discriminate Based on Data They Lack: Challenges, Solutions, and Policy Implications. *Journal of Information Policy*, 8, 78–115. <https://doi.org/10.5325/jinfopoli.8.2018.0078>

Zuckerman, D. (2009). The Ethics of Inclusion and Exclusion in Clinical Trials: Race, Sex, and Age. In V. Ravitsky, A. Fiester, & A. L. Caplan (Eds.), *The Penn Center Guide to Bioethics* (pp. 243–258). Springer Publishing Company.

10. What are the risks to privacy and other rights related to the dissemination and archiving of research outputs? What mitigation measures exist?

a. What steps should researchers take to protect data subjects or against societal-level harms prior to the dissemination of research outputs (publications, presentation slides, data visualization, datasets, AI/ML models, etc.)? ^[60]

The recent increase in AI research and production raises new risks to individuals and the public from the use of pervasive data. Existing frameworks (like the NIST AI RMF) have been developed that characterize common risks. Other researchers have identified additional potential harms from pervasive AI sensors (Stewart, et al., 2024). These risks include: increasing data quantity at the expense of data quality; increasing risks to privacy as sensors become smaller and easier to hide; increasing risks to data de-identification arising from AI/ML inference; increased opportunities in a sensor-saturated world for designers to exert power over individuals; increasing sustainability risks associated with both the environmental costs of producing increasingly cheap sensors and the energy consumption associated with their use; and a “chilling effect” on public attitudes and civic participation resulting from over-surveillance and careless deployment of pervasive sensing.

NIST. Artificial Intelligence Risk Management Framework (AI RMF 1.0). Technical Report NIST AI 100-1, National Institute of Standards and Technology (U.S.), Gaithersburg, MD, January 2023.

Matthew Stewart, Emanuel Moss, Pete Warden, Brian Plancher, Susan Kennedy, Mona Sloane, and Vijay Janapa Reddi. 2024. Materiality and Risk in the Age of Pervasive AI Sensors. <https://doi.org/10.48550/arXiv.2402.11183>

12. What are the existing requirements and legal obligations that impact research with pervasive data?

a. What are the risks around research that uses pervasive data, if any, that currently fall beyond the usual considerations of IRBs operating under the Common Rule or FDA regulations?

Existing IRB protections under the Common Rule may not be sufficient ethics protection for research proposals that involve the collection and use of pervasive data. The capacity—and inclination—for IRBs to support ethical data practices beyond the bare minimums designed for a fundamentally different model of data collection and usage is uneven across institutions. Funders or universities might consider adopting additional requirements that PIs affirm they have engaged in an ethical review process specific to pervasive data when appropriate, which might involve considerations (such as data contexts, participant awareness, and power) beyond an IRB review. Recent enhancements to paper review processes by the NeurIPS and ICML communities, or similar ethics requirements established by the SOMAR social media archive, are helpful examples.

13. What structured processes (questionnaires, rubrics, assessment frameworks) could be used to determine which techniques should be used to mitigate risks to data subjects and society in research that relies on pervasive data? ^[70]

The PERVADE Data Ethics Decision Support Tool (<https://pervade.umd.edu/pervade-data-ethics-tool/>) is one method for encouraging research teams to incorporate ethical and social considerations into the design of their research approach. Because a major finding of the PERVADE project was that, when it comes to delineating ethical data collection and use, the answer is almost always context-dependent. Ethical data use depends on the contexts and tools of data collection, the awareness and expectations of data subjects, decisions made about storage, analysis, and sharing, and the positionality of the researcher. The PERVADE tool is an interactive quiz designed to walk researchers through the many factors that matter to digital research ethics, and refer them to resources to support informed decisions about the many complex factors that shape data research ethics.