

2013

Big Data for All: Privacy and User Control in the Age of Analytics

Omer Tene

Haim Striks School of Law

Jules Polonetsky

Future of Privacy Forum

Recommended Citation

Omer Tene and Jules Polonetsky, *Big Data for All: Privacy and User Control in the Age of Analytics*, 11 NW. J. TECH. & INTEL. PROP. 239 (2013).

<http://scholarlycommons.law.northwestern.edu/njtip/vol11/iss5/1>

This Article is brought to you for free and open access by Northwestern University School of Law Scholarly Commons. It has been accepted for inclusion in Northwestern Journal of Technology and Intellectual Property by an authorized administrator of Northwestern University School of Law Scholarly Commons.

N O R T H W E S T E R N
JOURNAL OF TECHNOLOGY
AND
INTELLECTUAL PROPERTY

**Big Data for All: Privacy and User Control in the
Age of Analytics**

Omer Tene and Jules Polonetsky



Big Data for All: Privacy and User Control in the Age of Analytics

By Omer Tene¹ and Jules Polonetsky²

We live in an age of “big data.” Data have become the raw material of production, a new source for immense economic and social value. Advances in data mining and analytics and the massive increase in computing power and data storage capacity have expanded by orders of magnitude the scope of information available for businesses and government. Data are now available for analysis in raw form, escaping the confines of structured databases and enhancing researchers’ abilities to identify correlations and conceive of new, unanticipated uses for existing information. In addition, the increasing number of people, devices, and sensors that are now connected by digital networks has revolutionized the ability to generate, communicate, share, and access data. Data creates enormous value for the world economy, driving innovation, productivity, efficiency, and growth. At the same time, the “data deluge” presents privacy concerns which could stir a regulatory backlash dampening the data economy and stifling innovation. In order to craft a balance between beneficial uses of data and individual privacy, policymakers must address some of the most fundamental concepts of privacy law, including the definition of “personally identifiable information,” the role of individual control, and the principles of data minimization and purpose limitation. This article emphasizes the importance of providing individuals with access to their data in usable format. This will let individuals share the wealth created by their information and incentivize developers to offer user-side features and applications harnessing the value of big data. Where individual access to data is impracticable, data are likely to be de-identified to an extent sufficient to diminish privacy concerns. In addition, since in a big data world it is often not the data but rather the inferences drawn from them that give cause for concern, organizations should be required to disclose their decisional criteria.

INTRODUCTION	240
I. BIG DATA: BIG BENEFITS	243
A. Healthcare	245
B. Mobile	247
C. Smart Grid.....	248
D. Traffic Management.....	248
E. Retail	249

¹ Associate Professor, College of Management Haim Striks School of Law, Israel; Senior Fellow, Future of Privacy Forum; Visiting Researcher, Berkeley Center for Law and Technology; Affiliate Scholar, Stanford Center for Internet and Society. I would like to thank the College of Management Haim Striks School of Law research fund and the College of Management Academic Studies research grant for supporting research for this article.

² Co-chair and Director, Future of Privacy Forum.

F. Payments	249
G. Online	250
II. BIG DATA: BIG CONCERNS	251
A. Incremental Effect.....	251
B. Automated Decision-Making.....	252
C. Predictive Analysis	253
D. Lack of Access and Exclusion	254
E. The Ethics of Analytics: Drawing the Line	256
F. Chilling Effect.....	256
III. THE LEGAL FRAMEWORK: CHALLENGES	256
A. Definition of PII.....	257
B. Data Minimization	259
C. Individual Control and Context	260
IV. THE LEGAL FRAMEWORK: SOLUTIONS	263
A. Access, Portability, and Sharing the Wealth.....	263
B. Enhanced Transparency: Shining the Light.....	270
V. CONCLUSION.....	272

INTRODUCTION

¶1 Big data is upon us.³ Over the past few years, the volume of data collected and stored by business and government organizations has exploded.⁴ The trend is driven by reduced costs of storing information and moving it around in conjunction with increased capacity to instantly analyze heaps of unstructured data using modern experimental methods, observational and longitudinal studies, and large scale simulations.⁵ Data are generated from online transactions, email, video, images, clickstream, logs, search queries, health records, and social networking interactions; gleaned from increasingly pervasive sensors deployed in infrastructure such as communications networks, electric grids, global positioning satellites, roads and bridges,⁶ as well as in homes, clothing, and mobile phones.⁷

³ See, e.g., Steve Lohr, *The Age of Big Data*, N.Y. TIMES, Feb. 11, 2012, <http://www.nytimes.com/2012/02/12/sunday-review/big-datas-impact-in-the-world.html?pagewanted=all>; Steve Lohr, *How Big Data Became So Big*, N.Y. TIMES, Aug. 11, 2012, <http://www.nytimes.com/2012/08/12/business/how-big-data-became-so-big-unboxed.html>; Janna Anderson & Lee Rainie, *The Future of Big Data*, PEW INTERNET & AM. LIFE PROJECT (July 20, 2012), http://pewinternet.org/~media/Files/Reports/2012/PIP_Future_of_Internet_2012_Big_Data.pdf.

⁴ Kenneth Cukier, *Data, Data Everywhere*, THE ECONOMIST, Feb. 25, 2010, <http://www.economist.com/node/15557443>; see, e.g., World Economic Forum, *Big Data, Big Impact: New Possibilities for International Development* (2012), available at http://www3.weforum.org/docs/WEF_TC_MFS_BigDataBigImpact_Briefing_2012.pdf.

⁵ See, e.g., TREVOR HASTIE, ROBERT TIBSHIRANI & JEROME FRIEDMAN, THE ELEMENTS OF STATISTICAL LEARNING: DATA MINING, INFERENCE, AND PREDICTION (2009).

⁶ For the erosion of privacy in the public sphere, see *United States v. Jones*, 565 U.S. ___, 132 S. Ct. 945 (2012).

⁷ Omer Tene, *Privacy: The New Generations*, 1 INT'L DATA PRIVACY LAW 15 (2011), available at

¶2 The Obama Administration has recently announced a new, multi-agency big data research and development initiative aimed at advancing the core scientific and technological means of managing, analyzing, visualizing, and extracting information from large, diverse, distributed, and heterogeneous data sets.⁸ This initiative is based on recognition of the immense social and economic value captured in information and the intention to unleash it in order to progress from data to knowledge to action.⁹ Big data boosts the economy, transforming traditional business models and creating new opportunities through the use of business intelligence, sentiment analysis, and analytics. It advances scientific research, transforming scientific methods from hypothesis-driven to data-driven discovery.¹⁰ Big data furthers national goals such as optimization of natural resources, response to national disasters, and enhancement of critical information infrastructure.¹¹

¶3 The extraordinary societal benefits of big data—including breakthroughs in medicine, data security, and energy use—must be reconciled with increased risks to individuals' privacy.¹² As is often the case, technological and business developments in big data analysis have far outpaced the existing legal frameworks, which date back from an era of mainframe computers, predating the Internet, mobile, and cloud computing.¹³

<http://idpl.oxfordjournals.org/content/1/1/15.full>.

⁸ News Release, Office of Science and Technology Policy, Executive Office of the President, Obama Administration Unveils “Big Data” Initiative: Announces \$200 Million in New R&D Investments (Mar. 29, 2012), http://www.whitehouse.gov/sites/default/files/microsites/ostp/big_data_press_release.pdf.

⁹ *Personal Data: The Emergence of a New Asset Class*, WORLD ECONOMIC FORUM (2011), http://www3.weforum.org/docs/WEF_ITTC_PersonalDataNewAsset_Report_2011.pdf; Steve Lohr, *New U.S. Research Will Aim at Flood of Digital Data*, N.Y. TIMES, Mar. 29, 2012, http://www.nytimes.com/2012/03/29/technology/new-us-research-will-aim-at-flood-of-digital-data.html?_r=2.

¹⁰ See Chris Anderson, *The End of Theory: The Data Deluge Makes the Scientific Method Obsolete*, WIRED, June 23, 2008, available at http://www.wired.com/science/discoveries/magazine/16-07/pb_theory; see also Peter Norvig, UBC Department of Computer Science's Distinguished Lecture Series: The Unreasonable Effectiveness of Data, (Sept. 23, 2010), available at <http://www.youtube.com/watch?v=yvDCzhbjYWs>.

¹¹ Farnam Jahanian, Assistant Director, National Science Foundation, NSF Keynote at TechAmerica's Big Data Congressional Briefing, (May 2, 2012), available at http://www.youtube.com/watch?v=Do_IPA6-E9M.

¹² Omer Tene & Jules Polonetsky, *Privacy in the Age of Big Data: A Time for Big Decisions*, 64 STAN. L. REV. ONLINE 63 (2012).

¹³ See OECD Guidelines on the Protection of Privacy and Transborder Flows of Personal Data, ORG. FOR ECON. CO-OPERATION & DEV. (Sept. 23, 1980), http://www.oecd.org/document/18/0,3343,en_2649_34255_1815186_1_1_1_1,00.html [hereinafter: OECD Guidelines]; Council of Europe Convention 108 for the Protection of Individuals with Regard to Automatic Processing of Personal Data, Strasbourg, (Jan. 28, 1982), <http://conventions.coe.int/treaty/en/treaties/html/108.htm>; *Directive 95/46/EC of the European Parliament and of the Council of 24 October 1995 on the Protection of Individuals with Regard to the Processing of Personal Data and on the Free Movement of Such Data*, 1995 O.J. (L 281) 31 (Nov. 23, 1995), available at <http://eur-lex.europa.eu/LexUriServ/LexUriServ.do?uri=OJ:L:1995:281:0031:0050:EN:PDF> [hereinafter: European Data Protection Directive]; and in the United States: The Privacy Act of 1974, Pub. L. No. 93-579, 88 Stat. 1897 (Dec. 31, 1974). All of the major frameworks are being reviewed this year. See The White House, *Consumer Data Privacy in a Networked World: A Framework for Protecting Privacy and Promoting Innovation in the Global Digital Economy*, (Feb. 2012), <http://www.whitehouse.gov/sites/default/files/privacy-final.pdf> [hereinafter: White House Blueprint]; Federal Trade Commission Report, *Protecting Consumer Privacy in an Era of Rapid Change: Recommendations for Businesses and Policymakers* (Mar. 2012), <http://ftc.gov/os/2012/03/120326privacyreport.pdf> [hereinafter: FTC Final Report]; *Proposal for a*

For the past four decades, the tension between data innovation and informational privacy has been moderated by a set of principles broadly referred to as the Fair Information Practice Principles (FIPPs), based on a framework set in the 1980 OECD Guidelines.¹⁴ In the latest version presented by the White House this year, the FIPPs include the principles of individual control, transparency, respect for context, security, access and accuracy, focused collection, and accountability.¹⁵ The big data paradigm challenges some of these fundamental principles, including the scope of the framework (often addressed by framing the term “personally identifiable information” (PII)), the concepts of data minimization (“focused collection”) and consent (“individual control” and “respect for context”), and the right of individual access (“access and accuracy”).¹⁶

¶4 This article addresses the legal issues arising from the big data debate. It suggests that the FIPPs should be viewed as a set of levers that must be adjusted to adapt to varying business and technological conditions. Indeed, the ingenuity of the FIPPs is manifest in their flexibility, which has made them resilient to momentous change—some principles retract while others expand depending on the circumstances. In the context of big data, this means relaxing data minimization and consent requirements while emphasizing transparency, access, and accuracy. The shift is from empowering individuals at the point of information collection, which traditionally revolved around opting into or out of seldom read, much less understood corporate privacy policies, to allowing them to engage with and benefit from information already collected, thereby harnessing big data for their own personal usage. Further, such exposure will prevent the existence of “secret” databases and leverage societal pressure to constrain any unacceptable uses.

¶5 This article assesses the definition of PII in a world where de-identification is often reversible and sometimes detrimental to the integrity of the very data it aims to protect. It seeks to reconcile the current technological and business realities with the data minimization and purpose limitation principles. These principles are antithetical to big data, which is premised on data maximization—a theory that posits that the more data processed, the finer the conclusions—and seeks to uncover surprising, unanticipated correlations.

¶6 This article suggests that to solve the big data privacy quandary, individuals must be offered meaningful rights to access their data in a usable, machine-readable format. This, in turn, will unleash a wave of innovation for user-side applications and services based on access to PII, a process we refer to as the “featurization” of big data.¹⁷

Regulation of the European Parliament and of the Council on the protection of individuals with regard to the processing of personal data and on the free movement of such data (General Data Protection Regulation), COM(2012) 11 final (Jan. 25, 2012), available at http://ec.europa.eu/justice/data-protection/document/review2012/com_2012_11_en.pdf [hereinafter: EU General Data Protection Regulation].

¹⁴ OECD, *supra* note 13, at 4. The OECD (Organization for Economic Cooperation and Development) Guidelines include the principles of collection limitation, data quality, purpose specification, use limitation, security safeguards, openness, individual participation, and accountability.

¹⁵ White House, *supra* note 13, at 4.

¹⁶ Julie Brill, Commissioner, Fed. Trade Comm’n, Remarks at Fordham University School of Law: Big Data, Big Issues (Mar. 2, 2012) (transcript available at <http://ftc.gov/speeches/brill/120228fordhamlawschool.pdf>). Federal Trade Commission Commissioner Julie Brill said: “Big Data’s impact on privacy is requiring some new and hard thinking by all of us.”

¹⁷ See discussion *infra* notes 138 to 162 and accompanying text.

Featurization will allow individuals to declare their own policies, preferences and terms of engagement, and do it in ways that can be automated both for them and for the companies they engage.¹⁸ Where individual access to data is impracticable, data are likely to be de-identified to an extent sufficient to diminish privacy concerns.¹⁹ Where access is possible, organizations must provide it with robust mechanisms for user authentication and through secure channels to prevent leakage. This implies the development of user-centric or federated identity management schemes, which include single sign-on capability and at the same time do not become vehicles for universal surveillance.²⁰

¶7 To minimize concerns of untoward data usage, organizations should disclose the logic underlying their decision-making processes to the extent possible without compromising their trade secrets or intellectual property rights. As danah boyd and Kate Crawford recently noted: “In reality, working with Big Data is still subjective, and what it quantifies does not necessarily have a closer claim on objective truth”²¹ It is imperative that individuals have insight into the decisional criteria of organizations lest they face a Kafkaesque machinery that manipulates lives based on opaque justifications. While we recognize the practical difficulties of mandating disclosure without compromising organizations’ “secret sauce,” we trust that a distinction can be drawn between proprietary algorithms, which would remain secret, and decisional criteria, which would be disclosed.

¶8 Part One will describe some of the benefits of big data to individuals and society at large, including medical research, smart grid information, and traffic management. Some instances of big data use are so compelling that few would argue they should be forgone in light of the incremental risk to individuals’ privacy. Part Two discusses some of the risks of big data, including the unidirectional, incremental chipping away at informational privacy; the social stratification exacerbated by predictive analysis; and the exclusion of individuals from the value generated by their own information. Part Three addresses the challenges big data poses to existing privacy rules, including the definition of PII, the principle of data minimization, and the concept of meaningful, informed consent. Part Four argues the benefits of providing individuals with useful access to their data, allowing them to share the gains generated by the combination of their information with resources invested by businesses and government. Part Four also introduces arguments for requiring organizations to be transparent with respect to the decisional criteria underlying their big data choices.

I. BIG DATA: BIG BENEFITS

¶9 Big data is a big industry. Research conducted at the Massachusetts Institute of Technology shows that companies that use “data-directed decisionmaking” enjoy a 5%–6% increase in productivity.²² There is a strong link between effective data management

¹⁸ See Doc Searls, *The Customer as a God*, WALL ST. J., July 20, 2012, available at <http://online.wsj.com/article/SB10000872396390444873204577535352521092154.html>.

¹⁹ See discussion *infra* note 169 and accompanying text.

²⁰ See, e.g., Ann Cavoukian, *7 Laws Of Identity: The Case for Privacy-Embedded Laws Of Identity in the Digital Age* (2006), http://www.identityblog.com/wp-content/resources/7_laws_whitepaper.pdf.

²¹ danah boyd & Kate Crawford, *Critical Questions for Big Data*, 15 INFO. COMM. & SOC’Y 662, 667 (June 2012).

²² Erik Brynjolfsson, Lorin Hitt & Heekyung Kim, *Strength in Numbers: How Does Data-Driven*

strategy and financial performance. Companies that use data most effectively stand out from the rest. A report by the McKinsey Global Institute (MGI) demonstrates the transformative effect that big data has had on entire sectors ranging from health care to retail to manufacturing to political campaigns.²³ Just as it helps businesses increase productivity, big data allows governments to improve public sector administration and assists global organizations in analyzing information to devise strategic planning. Demand for big data is accelerating. MGI projected that the United States already needs 140,000 to 190,000 more workers with “deep analytical” expertise and 1.5 million more data-literate managers.²⁴

¶10 This chapter presents some anecdotal examples of the benefits of big data. When considering the risks that big data poses to individual privacy, policymakers should be mindful of its sizable benefits. Privacy impact assessments (PIA), systematic processes undertaken by government and business organizations to evaluate the potential risks to privacy of products, projects or schemes, often fail to bring these benefits into account.²⁵ Concluding that a project raises privacy risks is not sufficient to discredit it. Privacy risks must be weighed against non-privacy rewards. And while numerous mechanisms exist to assess privacy risks,²⁶ we still lack a formula to work out the balance.²⁷

¶11 At the same time, under existing market conditions, the benefits of big data do not always (some say, ever) accrue to the individuals whose personal data are collected and harvested.²⁸ This creates a twofold problem: on the one hand, individuals should not be required to volunteer their information with little benefit beyond feeding voracious corporate appetites; on the other hand, self interest should not frustrate societal values and benefits such as economic development or improved capabilities for law enforcement

Decision-Making Affect Firm Performance? A51 (Apr. 2011), http://www.a51.nl/storage/pdf/SSRN_id1819486.pdf; see *supra* note 9, at 3 (commenting on recent WEF report referring to personal data as “the new oil,” a new asset class emerging as the most valuable resource of the 21st century).

²³ James Manyika et al., *Big Data: The Next Frontier for Innovation, Competition, and Productivity*, MCKINSEY GLOBAL INSTITUTE (May 2011), http://www.mckinsey.com/Insights/MGI/Research/Technology_and_Innovation/Big_data_The_next_frontier_for_innovation [hereinafter MGI Report]; see also Thomas B. Edsall, *Let the Nanotargeting Begin*, NY TIMES CAMPAIGN STOPS BLOG (Apr. 15, 2012, 10:39 PM), <http://campaignstops.blogs.nytimes.com/2012/04/15/let-the-nanotargeting-begin>.

²⁴ Ben Rooney, *Big Data's Big Problem: Little Talent*, WALL ST. J. TECH EUROPE (Apr. 26, 2012), http://blogs.wsj.com/tech-europe/2012/04/26/big-datas-big-problem-little-talent/?mod=google_news_blog.

²⁵ E-Government Act of 2002, Pub. L. No. 107-347, § 208, 44 U.S.C. § 101 (2003); EU General Data Protection Regulation, art. 33-34.

²⁶ See, e.g., Roger Clarke, *An Evaluation of Privacy Impact Assessment Guidance Documents*, 1 INT'L DATA PRIVACY LAW 111 (2011); U.S. Securities and Exchange Commission, PRIVACY IMPACT ASSESSMENT (PIA) GUIDE (Jan. 2007), <http://www.sec.gov/about/privacy/piaguide.pdf>; U.S. Department of Homeland Security, PRIVACY IMPACT ASSESSMENTS, THE PRIVACY OFFICE OFFICIAL GUIDANCE (June 2010), http://www.dhs.gov/xlibrary/assets/privacy/privacy_pia_guidance_june2010.pdf; U.S. Department of Justice, Office of Privacy and Civil Liberties, PRIVACY IMPACT ASSESSMENTS, OFFICIAL GUIDANCE (Aug. 2010), http://www.justice.gov/opcl/pia_manual.pdf.

²⁷ For example, if analysis of de-identified online search engine logs enabled identification of a life threatening epidemic in $x\%$ of cases thus saving y lives, should such analysis be permitted assuming a $z\%$ chance of re-identification of a certain subset of search engine users? This is a meta-privacy question, which must be answered by policymakers implementing more than just a PIA; the PIA only solves one side of the equation.

²⁸ See, e.g., Natasha Singer, *Consumer Data, but Not for Consumers*, N.Y. TIMES, July 21, 2012, available at <http://www.nytimes.com/2012/07/22/business/acxiom-consumer-data-often-unavailable-to-consumers.html>.

and public health authorities. If individuals could reap some of the gains of big data, they would be incentivized to actively participate in the data economy, aligning their own self-interest with broader societal goals.

A. Healthcare

¶12 Dr. Russ Altman, a professor of medicine and bioengineering at Stanford University, and his colleagues made a groundbreaking discovery last year. They found that when taken together, Paxil®—the blockbuster antidepressant prescribed to millions of Americans—and Pravachol®—a highly popular cholesterol-reducing drug—have a dreadful side effect: they increase patients’ blood glucose to diabetic levels. Each drug taken alone does not have the diabetic side effects; hence, the Food and Drug Administration (FDA) approved the drugs for use. The FDA, which has limited resources, cannot afford to test each and every drug for every conceivable interaction.

¶13 Altman and his team made their discovery by pursuing statistical analysis and data mining techniques to identify patterns in large datasets. They analyzed information in the Adverse Event Reporting System (AERS), a database maintained by the FDA to collect adverse drug event reports from clinicians, patients, and drug companies for more than thirty years.²⁹ Using the AERS, they created a “symptomatic footprint” for diabetes-inducing drugs (i.e., the side effects a patient might report if she had undiagnosed diabetes), then searched for that footprint in interactions between pairs of drugs not known to induce such effects when taken alone. Four pairs of drugs were found to leave the footprint; of those, Paxil and Pravachol were the most commonly prescribed.

¶14 Next, the scientists approached Microsoft Research to examine de-identified Bing search engine logs,³⁰ querying whether a higher proportion of users who searched for *both* “Paxil” and “Pravachol” also typed in words related to the “symptomatic footprint” (such as “headache” or “fatigue”) than those who searched for just Paxil or Pravachol separately. Sure enough, their research hypothesis found support in that big data set as well. Users who searched Bing for the name of both drugs together were much likelier to search for diabetes-related side effects than users who searched for only one of the drugs.³¹

²⁹ *Adverse Event Reporting System (AERS)*, FOOD AND DRUG ADMIN., <http://www.fda.gov/drugs/guidancecomplianceregulatoryinformation/surveillance/adversedrugs/effects/default.htm> (last updated Sept. 10, 2012).

³⁰ Through de-identification, organizations can reduce privacy risks associated with data while still salvaging such data for beneficial use. De-identification could be achieved through various techniques such as data masking (stripping out obvious personal identifiers such as names from a piece of information, to create a data set in which no person identifiers are present); pseudonymization (de-identifying data so that a coded reference or pseudonym is attached to a record to allow the data to be associated with a particular individual without the individual being identified); aggregation (data is displayed as totals, so no data relating to or identifying any individual is shown; small numbers in totals are often suppressed through ‘blurring’ or by being omitted altogether); and more. See Information Commissioner’s Office, ANONYMISATION: MANAGING DATA PROTECTION RISK CODE OF PRACTICE, Nov. 2012, http://www.ico.gov.uk/news/latest_news/2012/~/media/documents/library/Data_Protection/Practical_application/anonymisation_code.ashx.

³¹ See 2012 Stanford Law Review Symposium, *The Privacy Paradox: Health and Medical Privacy*, YOUTUBE (Feb. 27, 2012), http://www.youtube.com/watch?v=ntL4WVGkiXo&feature=player_embedded#! (Altman describing the research process, including the search engine logs analysis, from minute 32 of the video).

¶15 By implementing a novel signal detection algorithm that identifies statistically significant correlations, the researchers were thus able to parse out latent adverse effect signals from spontaneous reporting systems.³² In 2009, for example, “there were an estimated 15 million prescriptions for paroxetine [Paxil] and 18 million prescriptions for pravastatin [Pravachol] in the United States”; there were an estimated one million individuals who used both drugs in combination.³³ For these users, the work of Altman and his colleagues was potentially life-saving.³⁴

¶16 In addition to the findings of Altman and his team, there are numerous other examples of significant healthcare breakthroughs based on big data analysis. The discovery of Vioxx’s adverse drug effects, which led to its withdrawal from the market, was made possible by analysis of clinical and cost data collected by Kaiser Permanente, the California-based managed-care consortium.³⁵ Had Kaiser Permanente not aggregated clinical and cost data, researchers might not have been able to attribute 27,000 cardiac arrest deaths occurring between 1999 and 2003 to use of the drug.

¶17 In another example, researchers in South Africa discovered a positive relationship between therapeutic vitamin B use and delay of progression to AIDS and death in HIV-positive patients.³⁶ This was a critical finding at a time and in a region where therapies for people living with HIV are well beyond the financial means of most patients. The researchers noted that “[n]onlinear statistical analysis . . . can help elucidate clinically-relevant relationships within a large patient population such as observational databases.”³⁷ Another oft-cited example is Google Flu Trends, which predicts and locates outbreaks of the flu making use of information—aggregate search queries—not originally collected with this innovative application in mind. Of course, “[e]arly detection of disease activity, when followed by rapid response, can reduce the impact of both seasonal and pandemic influenza.”³⁸ Yet another example is the National Retail Data Monitor (NRDM), which keeps tabs on sales of over-the-counter healthcare items from 21,000 outlets across the United States. By analyzing the remedies people purchase, health officials can anticipate short-term trends in illness transmission. “Data from the NRDM show that sales of over-the-counter products like cough medicine and electrolytes . . . spike before visits to the emergency room do,” and that the lead-time can be significant—two and a half weeks in the case of respiratory and gastrointestinal illnesses.³⁹ According

³² See also David Reshef et al., *Detecting Novel Associations in Large Data Sets*, 334 SCIENCE 1518, 1520 (2011).

³³ Nicholas Tatonetti et al., *Detecting Drug Interactions From Adverse-Event Reports: Interaction Between Paroxetine and Pravastatin Increases Blood Glucose Levels*, 90 CLINICAL PHARMACOLOGY & THERAPEUTICS 133, 133, 139 (2011).

³⁴ See Nicholas Tatonetti et al., *A Novel Signal Detection Algorithm for Identifying Hidden Drug-Drug Interactions in Adverse Event Reports*, 12 J. AM. MED. INFORMATICS ASS’N 79, 79–80 (2011).

³⁵ See, e.g., Rita Rubin, *How Did Vioxx Debacle Happen?*, USA TODAY (Oct. 12, 2004, 12:00 AM), available at http://www.usatoday.com/news/health/2004-10-12-vioxx-cover_x.htm.

³⁶ Andrew Kanter et al., *Supplemental Vitamin B and Progression to AIDS and Death in Black South African Patients Infected With HIV*, 21 JOURNAL OF ACQUIRED IMMUNE DEFICIENCY SYNDROMES 252, 253 (1999).

³⁷ *Id.*

³⁸ Jeremy Ginsberg et al., *Detecting Influenza Epidemics Using Search Engine Query Data*, 457 NATURE 1012, 1012 (2009).

³⁹ Brian Fung, *Using Data Mining to Predict Epidemics Before They Spread*, THE ATLANTIC, May 2, 2012, <http://www.theatlantic.com/health/archive/2012/05/using-data-mining-to-predict-epidemics-before-they-spread/256605>.

to a study published in a medical journal, it took weeks for official sources in Haiti to report details of a cholera epidemic in 2010, resulting in more than 7,000 casualties and 500,000 infections, whereas on Twitter, news of the disease traveled far more quickly.⁴⁰

¶18 The potential advantages of big data analytics within the medical field have resulted in public policy initiatives to mine and leverage such data. David Cameron, Prime Minister of the United Kingdom, recently announced that every NHS patient would henceforth be a “research patient” whose medical record would be “opened up” for research by private healthcare firms.⁴¹ The Prime Minister emphasized that privacy-conscious patients would be given opt out rights. He added that “this does not threaten privacy, it doesn't mean anyone can look at your health records, but it does mean using anonymous data to make new medical breakthroughs.” While a significant driver for research and innovation, the health sector is not the only arena for groundbreaking big data use.

B. Mobile

¶19 Mobile devices—always on, location aware, and with multiple sensors including cameras, microphones, movement sensors, GPS, and Wi-Fi capabilities—have revolutionized the collection of data in the public sphere and enabled innovative data harvesting and use. A group of scientists working on a collaborative project at MIT, Harvard, and additional research universities is currently analyzing mobile phone communications to better understand the needs of the one billion people who live in settlements or slums in developing countries.⁴² They explore ways to predict food shortages using variables such as market prices, drought, migrations, previous regional production, and seasonal variations;⁴³ to quantify crime waves by tracking the time, place, and nature of criminal activity in locations across a city;⁴⁴ and to decide which intervention is the most effective means for improving learning outcomes in developing country schools.⁴⁵

⁴⁰ See Rumi Chunara et al., *Social and News Media Enable Estimation of Epidemiological Patterns Early in the 2010 Haitian Cholera Outbreak*, 86 AM. J. TROP. MED. HYG. 39 (2012); see also Alessio Signorini et al., *The Use of Twitter to Track Levels of Disease Activity and Public Concern in the U.S. During the Influenza H1N1 Pandemic*, PLOS ONE (May 2011), <http://www.divms.uiowa.edu/~asignori/papers/use-twitter-track-level-disease-activity-and-concern-in-us-during-h1n1.pdf>.

⁴¹ See, e.g., *Everyone 'to be research patient', says David Cameron*, BBC NEWS, Dec. 5, 2011, <http://www.bbc.co.uk/news/uk-16026827>.

⁴² See *Big Data for Social Good Initiative*, HARVARD SCHOOL OF PUBLIC HEALTH: ENGINEERING SOCIAL SYSTEMS, <http://www.hsph.harvard.edu/ess/bigdata.html> (last visited April 2, 2013); see also Amy Wesolowski & Nathan Eagle, *Parameterizing the Dynamics of Slums*, PROCEEDINGS OF AAAI ARTIFICIAL INTELLIGENCE FOR DEVELOPMENT (AI-D'10), <http://ai-d.org/pdfs/Wesolowski.pdf> (last visited December 2, 2012).

⁴³ See, e.g., Washington Okori & Joseph Obua, *Machine Learning Classification Technique for Famine Prediction*, PROCEEDINGS OF THE WORLD CONGRESS ON ENGINEERING 2011 (July 6–8, 2011), http://www.iaeng.org/publication/WCE2011/WCE2011_pp991-996.pdf.

⁴⁴ See, e.g., Jameson Toole et al., *Quantifying Crime Waves*, PROCEEDINGS OF AAAI ARTIFICIAL INTELLIGENCE FOR DEVELOPMENT (AI-D'10), <http://ai-d.org/pdfs/Toole.pdf> (last visited December 2, 2012).

⁴⁵ See Massoud Moussavi & Noel McGinn, *A Model for Quality of Schooling*, PROCEEDINGS OF AAAI ARTIFICIAL INTELLIGENCE FOR DEVELOPMENT (AI-D'10), <http://www.aaai.org/ocs/index.php/SSS/SSS10/paper/view/1126/1351> (last visited December 2, 2012).

C. Smart Grid

¶20 Big data use within the “smart grid”⁴⁶ context also illustrates the benefits of sophisticated data analysis. The smart grid is designed to allow electricity service providers, users, and other third parties to monitor and control electricity use. Utilities view the smart grid as a way to precisely locate power outages or other problems, including cyber-attacks or natural disasters, so that technicians can be dispatched to mitigate problems.

¶21 Consumers benefit from more choices on means, timing, and quantity of electricity they use.⁴⁷ Pro-environment policymakers view the smart grid as key to providing better power quality and more efficient delivery of electricity to facilitate the move towards renewable energy. Other benefits, such as accurately predicting energy demands to optimize renewable sources, may be reaped by society at large. Not only will future renewable sources benefit from the use of smart grid data, but also the current energy infrastructure will as well, for example, by utility companies accurately determining when to use peak versus baseload power plants.

D. Traffic Management

¶22 An additional area for data-driven environmental innovation is traffic management and control. Governments around the world are establishing electronic toll pricing systems, which determine differentiated payments based on mobility and congestion charges.⁴⁸ These systems apply varying prices to drivers based on their differing use of vehicles and roads.

¶23 Urban planners benefit from the analysis of personal location data for decisions involving road and mass transit construction, mitigation of traffic congestion, and planning for high-density development.⁴⁹ Such decisions can not only cut congestion but also control the emission of pollutants.⁵⁰ At the same time, individual drivers benefit from smart routing based on real-time traffic information, including accident reports and information about scheduled roadwork and congested areas.

¶24 Automotive telematics is another area of innovation. Vehicles equipped with navigation systems with embedded communication modules propose a range of telematics services to improve fuel-efficient driving and allow drivers to plan trips taking into account the location of charging stations or activate their air conditioner remotely.⁵¹

⁴⁶ The “smart grid” refers to the modernization of the current electrical grid to introduce a bi-directional flow of information and electricity.” *E.g.*, *Information and Privacy Commissioner of Ontario & Future of Privacy Forum*, SMART PRIVACY FOR THE SMART GRID: EMBEDDING PRIVACY INTO THE DESIGN OF ELECTRICITY CONSERVATION (Nov. 2009), <http://www.ipc.on.ca/images/Resources/pbd-smartpriv-smartgrid.pdf>.

⁴⁷ *See, e.g.*, Katie Fehrenbacher, *Introducing the Facebook Social Energy App*, GIGAOM (Oct. 17, 2011, 7:57 AM), <http://gigaom.com/cleantech/introducing-the-facebook-social-energy-app>.

⁴⁸ *See, e.g.*, Directive 2004/52/EC of the European Parliament and of the Council of 29 April 2004 on the Interoperability of Electronic Road Toll Systems in the Community, 2004 O.J. (L 166) 124, 125–27; *see also* Commission Decision 2009/750/EC of 6 October 2009 on the Definition of the European Electronic Toll Service and Its Technical Element, 2009 O.J. (L 268) 11, 11–14.

⁴⁹ *See, e.g.*, Carlo Ratti et al., *Mobile Landscapes: Using Location Data from Cell-Phones for Urban Analysis*, 33 ENV'T. AND PLAN. B: PLAN. AND DESIGN 727, 745 (2006).

⁵⁰ MGI Report, *supra* note 23, at 92.

⁵¹ For various examples, *see* special issue *Automotive Pervasive Computing*, IEEE PERVASIVE COMP.,

E. Retail

¶25 Big data is also transforming the retail market. It was Wal-Mart's inventory-management system ("Retail Link") which pioneered the age of big data by enabling suppliers to see the exact number of their products on every shelf of every store at each precise moment in time.⁵² Many shoppers use Amazon's "Customers Who Bought This Also Bought" feature, prompting users to consider buying additional items selected by a collaborative filtering tool. The most prevalent business model for the Internet is based on financing products and services with targeted ads whose value correlates directly with the amount of information collected from users.⁵³ Businesses care not so much about the identity of each individual user but rather on the attributes of her profile, which determine the nature of ads she is shown.⁵⁴

¶26 Analytics can also be used in the offline environment to study customers' in-store behavior to improve store layout, product mix, and shelf positioning. A 2011 report by McKinsey & Company explains that "[r]ecent innovations have enabled retailers to track customers' shopping patterns (e.g., foot traffic and time spent in different parts of a store), drawing real-time location data from smartphone applications (e.g., Shopkick), shopping cart transponders, or passively monitoring the location of mobile phones within a retail environment."⁵⁵ Increasingly, organizations are seeking to link online activity to offline behavior, both in order to assess the effectiveness of online ad campaigns, as judged by conversion to in-store purchases, and to re-target in-store customers with ads when they go online.

F. Payments

¶27 Another major arena for valuable big data use is fraud detection in the payment card industry. With electronic commerce capturing an increasingly large portion of the retail market, the merchants that bear ultimate responsibility for fraudulent card payments⁵⁶ must implement robust mechanisms to identify suspect transactions often

July-Sept. 2011, at 12. *See also* Sastry Duri et al., *Data Protection and Data Sharing in Telematics*, 9 MOBILE NETWORKS & APPLICATIONS, 693, 695 (2004).

⁵² *See, e.g., A Different Game: Information is Transforming Traditional Businesses*, THE ECONOMIST, Feb. 25, 2010, <http://www.economist.com/node/15557465>.

⁵³ *See* Article 29 Working Party, *Opinion 2/2010 on Online Behavioral Advertising*, at 5, WP 171 (June 22, 2010), *available at* http://ec.europa.eu/justice_home/fsj/privacy/docs/wpdocs/2010/wp171_en.pdf; *see also*, FTC STAFF REPORT, SELF-REGULATORY PRINCIPLES FOR ONLINE BEHAVIORAL ADVERTISING (2009), *available at* <http://www.ftc.gov/os/2009/02/P085400behavadreport.pdf>.

⁵⁴ *See* Omer Tene, *For Privacy, European Commission Must Be Innovative*, CDT BLOG (Feb. 28, 2011), <https://www.cdt.org/blogs/privacy-european-commission-must-be-innovative> ("[I]t is the singling out of an individual for unique treatment (e.g., the pricing of a loan or targeting of an ad) based on his or her profile, even without the ability to unmask his or her name, which has significant privacy implications."); *see generally*, Omer Tene & Jules Polonetsky, *To Track or 'Do Not Track': Advancing Transparency and Individual Control in Online Behavioral Advertising*, 13 MINN. J.L. SCI. & TECH. 282 (2012).

⁵⁵ MGI Report, *supra* note 23, at 68.

⁵⁶ A set of laws and regulations serve to protect consumer users of credit and debit cards from bearing the consequences of fraud losses associated with lost or stolen cards. *See* Truth in Lending Act, Title I of the Consumer Credit Protection Act, 15 U.S.C. § 1601; *see also* Regulation Z, 12 C.F.R. § 226, promulgated by the Federal Reserve Board pursuant to authority granted under 15 U.S.C. § 1607. The Electronic Fund Transfer Act and Federal Reserve Board Regulation E place a floating cap on a consumer cardholder's liability for unauthorized debit card use under which the maximum liability amount is determined when the cardholder notifies the card issuer of the loss or theft of the card used to perpetrate the

performed by first-time customers. To this end, some companies have developed solutions to provide merchants with predictive fraud scores for “Card-Not-Present transactions” in order to measure in real time the likelihood that a transaction is fraudulent.⁵⁷ To do that, the services analyze buyer histories and provide evaluations, much like a summarized list of references but in the form of a single score. As fraudsters become more sophisticated in their approach, online merchants must remain ever more vigilant in their efforts to protect the integrity of the online shopping experience.

G. Online

¶28

Finally, perhaps the most oft-cited example of the potential of big data analytics lies within the massive data silos maintained by the online tech giants: Google, Facebook, Microsoft, Apple, and Amazon. These companies have amassed previously unimaginable amounts of personal data. Facebook, for example, has more than 900 million users who upload more than 250 millions photos and click the “Like” button more than 2.5 billion times per day.⁵⁸ Google offers a plethora of data-intensive products and services, including its ubiquitous search engine, mobile operating system (Android), web browser (Chrome), email service (Gmail), video streaming site (YouTube), mapping service (Google Maps), social networking service (Google Plus), website analytics tool (Google Analytics), cloud platform service (Google Apps), and many others.⁵⁹ In addition, Google owns the largest online advertising serving company, DoubleClick, which it purchased in 2007, much to the consternation of privacy advocates,⁶⁰ as well as AdMob, the leading mobile advertising company. As a result, Google now has a presence on well over 70 percent of third party websites.⁶¹ Amazon and Yahoo are seeking new ways to leverage and monetize their treasure trove of customer data.⁶² Apple and Microsoft make

fraud. If the cardholder notifies the card issuer within two business days of learning of the loss or theft of the debit card, the cardholder’s maximum liability is limited to the lesser of the actual amount of unauthorized transfers or \$50.00. *See* Electronic Fund Transfer Act, 15 U.S.C. §§ 1693–1700; *see also* Federal Reserve Board Regulation E, 12 C.F.R. § 205.6(b)(1). Liability is further allocated between card issuers and merchants, generally shifting the risk away from the card issuers and onto the merchants, based on a complicated set of rules that vary based on the type of transaction at issue. *See* Duncan Douglass, *An Examination of the Fraud Liability Shift in Consumer Card-Based Payment Systems*, 33 ECONOMIC PERSPECTIVES 43, 45 (2009).

⁵⁷ For various solution providers operating in this space, *see* *Solution Provider Search*, MERCHANT RISK COUNCIL, <https://www.merchantriskcouncil.org/Resources/Pages/Solution-Provider-Search.aspx>.

⁵⁸ Melissa Fach, *Stats on Facebook 2012*, SEARCH ENGINE J. (Feb. 17, 2012), <http://www.searchenginejournal.com/stats-on-facebook-2012-infographic/40301>; *see also* Margot Bonner, *10 Key Statistics About Facebook*, EXPERIAN HITWISE BLOG (Feb. 2, 2012), <http://www.experian.com/blogs/hitwise/2012/2/2/10-key-statistics-about-facebook>.

⁵⁹ *Google Products*, GOOGLE, <http://www.google.com/intl/en/about/products/index.html> (last visited Dec. 2, 2012).

⁶⁰ *See* FEDERAL TRADE COMMISSION, STATEMENT OF FEDERAL TRADE COMMISSION CONCERNING GOOGLE/DOUBLECLICK, F.T.C. FILE NO. 071-0170, (Dec. 20, 2007), *available at* <http://www.ftc.gov/os/caselist/0710170/071220statement.pdf>; *see also* FEDERAL TRADE COMMISSION, DISSENTING STATEMENT OF COMM. PAMELA JONES HARBOUR, IN THE MATTER OF Google/DoubleClick, F.T.C. FILE NO. 071-0170, (Dec. 20, 2007), *available at* <http://www.ftc.gov/os/caselist/0710170/071220harbour.pdf>.

⁶¹ *See* Public Comment from Balachander Krishnamurthy & Craig Wills, *Privacy Diffusion on the Web: A Longitudinal Perspective*, (Oct. 20, 2009) (in response to Federal Trade Commission Project No. P095416), <http://www.ftc.gov/os/comments/privacyroundtable/544506-00009.pdf>.

⁶² *See* Nicole Perlroth, *BITS; Revamping at Yahoo to Focus on Its Media Properties and Customer*

operating systems as well as browsers, both of which are important focal points for collecting online and mobile user information.

II. BIG DATA: BIG CONCERNS

¶29 Big data poses big privacy risks. The harvesting of large sets of personal data and the use of state of the art analytics implicate growing privacy concerns. Protecting privacy will become harder as information is multiplied and shared ever more widely among multiple parties around the world. As more information regarding individuals' health, financials, location, electricity use, and online activity percolates, concerns arise regarding profiling, tracking, discrimination, exclusion, government surveillance, and loss of control.⁶³ This Part lays out some of the unique privacy risks presented by big data.

A. Incremental Effect

¶30 The accumulation of personal data has an incremental adverse effect on privacy.⁶⁴ A researcher will draw entirely different conclusions from a string of online search queries consisting of the words “paris,” “hilton” and “louvre” as compared to one featuring “paris,” “hilton” and “nicky.” Add thousands and thousands of search queries, and you can immediately sense how the data become ever more revealing.⁶⁵ Moreover, once data—such as a clickstream or a cookie number—are linked to an identified individual, they become difficult to disentangle.⁶⁶ This was demonstrated by University of Texas researchers Arvind Narayanan and Vitaly Shmatikov, who re-associated de-identified Netflix movie recommendations with identified individuals by crossing a de-identified database with publicly available resources accessible online.⁶⁷ Narayanan and Shmatikov explained, “Once any piece of data has been linked to a person’s *real* identity, any association between this data and a *virtual* identity breaks anonymity of the latter.”⁶⁸ Paul Ohm warned that this incremental effect will lead to a “database of ruin,” chewing

Data, N.Y. TIMES BITS BLOG, Apr. 11, 2012,

<http://query.nytimes.com/gst/fullpage.html?res=9B0DE2D91631F932A25757C0A9649D8B63&partner=rsnyt&emc=rss>.

⁶³ See Daniel Solove, *A Taxonomy of Privacy*, 154 U. PA. L. REV. 477 (2006) (for a taxonomy of privacy harms).

⁶⁴ Solove in his “taxonomy” calls this “aggregation.” *Id.* at 505–09

⁶⁵ See, e.g., Michael Barbaro & Tom Zeller, *A Face Is Exposed for AOL Searcher No. 4417749*, N.Y. TIMES, Aug. 9, 2006, <http://www.nytimes.com/2006/08/09/technology/09aol.html?pagewanted=all>.

⁶⁶ See Myspace LLC, F.T.C. File No. 102-3058, Agreement Containing Consent Order (May 8, 2012), <http://www.ftc.gov/os/caselist/1023058/120508myspaceorder.pdf> (charging that Myspace “constructively shared” personally identifiable information with third party advertisers by sharing with such advertisers a unique identifier assigned to the profile of each Myspace user (a “Friend ID”), which could then be used to access such user’s profile information – a practice referred to in the industry as “cookie syncing.”). See also Myspace, LLC: Analysis of Proposed Consent Order To Aid Public Comment, 77 Fed. Reg. 28,388 (Federal Trade Commission May 14, 2012), available at <http://www.gpo.gov/fdsys/pkg/FR-2012-05-14/pdf/2012-11613.pdf>. For an analysis of “cookie syncing,” see Ed Felten, *Syncing and the FTC’s Myspace Settlement*, TECH@FTC (May 8, 2012), <http://techatftc.wordpress.com/2012/05/08/syncing-and-the-ftcs-myspace-settlement>.

⁶⁷ Arvind Narayanan & Vitaly Shmatikov, *Robust De-anonymization of Large Sparse Datasets*, 2008 IEEE SYMP. ON SECURITY & PRIVACY 111.

⁶⁸ *Id.* at 119.

away, byte by byte, on an individual's privacy until his or her profile is completely exposed.⁶⁹

¶31 More generally, the ephemeral nature of personal data makes it difficult to recapture after it is exposed in the public or semi-public sphere.⁷⁰ For this reason, the European Commission's proposal of a "right to be forgotten," which would allow individuals to demand organizations to wipe their data slate clean,⁷¹ has been met with fierce resistance from online platforms⁷² and free speech advocates,⁷³ who are concerned about the effect of the proposal on the delicate balance between privacy and regulation of the Internet.

B. Automated Decision-Making

¶32 The relegation of decisions about an individual's life to automated processes based on algorithms and artificial intelligence raises concerns about discrimination, self-determination, and the narrowing of choice.⁷⁴ This is true not only for decisions relating to an individual's credit, insurance, or job prospects,⁷⁵ but also for highly customized choices regarding which advertisements or content a user will see.⁷⁶ In his book *The Daily You: How the New Advertising Industry Is Defining Your Identity and Your Worth*, Joseph Turow argues that increased personalization based on opaque corporate profiling algorithms poses a risk to open society and democratic speech.⁷⁷ He explains that by "pigeonholing" individuals into pre-determined categories, automated decision-making compartmentalizes society into pockets (or "echo chambers") of like-minded individuals.⁷⁸ Turow argues government should regulate information intermediaries to ensure that users have full control over their data and content consumption.

⁶⁹ Paul Ohm, *Broken Promises of Privacy: Responding to the Surprising Failure of Anonymization*, 57 UCLA L. REV. 1701, 1748 (2010).

⁷⁰ A social networking service (SNS) is a semi-public sphere. While an individual user's postings are made according to the SNS privacy settings, other users are not subject to a legal obligation to comply with such user's individual settings. Consequently, a posting made by a user and restricted to her "friends" may later be disseminated broadly by those friends so as to become public or semi-public. See Omer Tene, *Me, Myself and I: Aggregated and Disaggregated Identities on Social Networking Services*, J. INT'L COMM. L. & TECH. (forthcoming 2012).

⁷¹ Council Directive 95/46, art. 17, 1995 O.J. (L 281) 35 (EC); see also Viviane Reding, Vice President, Eur. Comm'n, EU Data Protection Reform 2012: Making Europe the Standard Setter for Modern Data Protection Rules in the Digital Age 5, Speech at the Digital-Life-Design Conference (Jan. 22, 2012), available at Press Release Europa, http://europa.eu/rapid/press-release_SPEECH-12-26_en.htm, (last visited April 2, 2013).

⁷² See Peter Fleischer, *Foggy Thinking About the Right to Oblivion*, PRIVACY...? BLOG (Mar. 9, 2011), <http://peterfleischer.blogspot.com/2011/03/foggy-thinking-about-right-to-oblivion.html>.

⁷³ See Jeffrey Rosen, *The Right to Be Forgotten*, 64 STAN. L. REV. ONLINE 88 (2012).

⁷⁴ See Ruth Gavison, *Privacy*, 89 YALE L.J. 421 (1980); Council Directive 95/46, art. 15, 1995 O.J. (L 281) 43 (EC).

⁷⁵ Such decisions have for many years been regulated by laws such as the Fair Credit Reporting Act, 15 U.S.C. § 1681(a)–(b).

⁷⁶ Kashmir Hill, *Resisting The Algorithms*, FORBES, May 5, 2011, <http://www.forbes.com/sites/kashmirhill/2011/05/05/resisting-the-algorithms>.

⁷⁷ JOSEPH TUROW, *THE DAILY YOU: HOW THE NEW ADVERTISING INDUSTRY IS DEFINING YOUR IDENTITY AND YOUR WORTH* (2011). For similar arguments, see ELI PARISER, *THE FILTER BUBBLE: WHAT THE INTERNET IS HIDING FROM YOU* (2011).

⁷⁸ This phenomenon is sometimes referred to as "cyberbalkanization." See *Cyberbalkanization*, WIKIPEDIA, <http://en.wikipedia.org/wiki/Cyberbalkanization>; see also ANDREW SHAPIRO, *THE CONTROL*

C. Predictive Analysis

¶33 Big data may facilitate predictive analysis with stark implications for individuals susceptible to disease, crime, or other socially stigmatizing characteristics or behaviors. To be sure, predictive analysis can be used for societally beneficial goals, such as planning disaster recovery in an earthquake prone area based on individuals' evacuation paths and purchase needs. Yet it can easily cross the "creepiness" threshold.⁷⁹

¶34 Consider a recent story in the *New York Times*, which uncovered that the retailing giant, Target Inc., assigns a "pregnancy prediction score" to customers based on their purchase habits.⁸⁰ According to the *Times*, Target employed statisticians to sift back through historical buying records of women who had signed up for baby registries. The statisticians discovered latent patterns, such as women's preference for unscented lotion around the beginning of their second trimester or a tendency to buy supplements like calcium, magnesium and zinc within the first 20 weeks of a pregnancy. They were able to determine a set of products that, when grouped together, allowed Target to accurately predict a customer's pregnancy and due date. In one case, the *Times* reported that a father of a teenage girl stormed into a Target store to complain that his daughter received coupons and advertisements for baby products. A few days later, he called the store manager to apologize, admitting that, "There's been some activities in my house I haven't been completely aware of. She's due in August."⁸¹

¶35 Predictive analysis is useful for law enforcement, national security, credit screening, insurance, and employment. It raises ethical dilemmas illustrated, for example, in the film *Minority Report*, where a "PreCrime" police department apprehends "criminals" based on foreknowledge of their future misdeeds. It could facilitate unlawful activity such as "redlining."⁸² Although these practices are illegal under current laws, critics expressed concerns that data are surreptitiously being used in such a manner.⁸³

¶36 Predictive analysis is particularly problematic when based on sensitive categories of data, such as health, race, or sexuality. It is one thing to recommend for a customer books, music or movies she might be interested in based on her previous purchases;⁸⁴ it is

REVOLUTION: HOW THE INTERNET IS PUTTING INDIVIDUALS IN CHARGE AND CHANGING THE WORLD WE KNOW (PublicAffairs, 2000).

⁷⁹ See, e.g., danah boyd, Senior Researcher, Microsoft Research, Speech at the DataEDGE Conference 2012 (cited in Quentin Hardy, *Rethinking Privacy in an Era of Big Data*, N.Y. TIMES, June 4, 2012, <http://bits.blogs.nytimes.com/2012/06/04/rethinking-privacy-in-an-era-of-big-data>) (stating that "privacy is a source of tremendous tension and anxiety in Big Data. It's a general anxiety that you can't pinpoint, this odd moment of creepiness.").

⁸⁰ Charles Duhigg, *How Companies Learn Your Secrets*, N.Y. TIMES MAGAZINE, Feb. 16, 2012, <http://www.nytimes.com/2012/02/19/magazine/shopping-habits.html?pagewanted=all>.

⁸¹ *Id.*

⁸² Redlining refers to the act of denying or increasing the cost of services such as loans, insurance, or healthcare to residents of neighborhoods comprised mostly of minorities. The term was coined to reflect the practice of some lenders of drawing red lines on maps to delineate neighborhoods where they would not lend money. See THE URBAN INSTITUTE, *MORTGAGE LENDING DISCRIMINATION: A REVIEW OF EXISTING EVIDENCE* (Margery Turner & Felicity Skidmore, Eds., 1999).

⁸³ See, e.g., Letter from Center for Digital Democracy, U.S. PIRG & World Privacy Forum, to the Federal Trade Commission, *In the Matter of Real-time Targeting and Auctioning, Data Profiling Optimization, and Economic Loss to Consumers and Privacy* (Apr. 8, 2010), available at <http://www.centerfordigitaldemocracy.org/sites/default/files/20100407-FTCfiling.pdf>.

⁸⁴ Consider Amazon, Netflix and Pandora recommendation systems. See Gediminas Adomavicius & Alexander Tuzhilin, *Towards the Next Generation of Recommender Systems: A Survey of the State-of-the-*

quite another thing to identify when she is pregnant before her closest family knows. In the law enforcement arena, predictive analysis raises the specter of surveying or even incarcerating individuals based on thoughts as opposed to deeds.⁸⁵ This type of activity, while clearly unconstitutional under existing U.S. law, is not so far-fetched in other parts of the world,⁸⁶ and could conceivably cross the line from fiction to reality, given the right circumstances in the United States.⁸⁷

¶37 Even with non-sensitive data categories, predictive analysis may have a stifling effect on individuals and society, perpetuating old prejudices. The wealthy and well-educated will get the fast track; the poor and underprivileged will have the deck stacked against them even more so than before.⁸⁸ By ignoring outliers and assuming that “what has been is what will be,”⁸⁹ predictive analysis becomes a self-fulfilling prophecy that accentuates social stratification.⁹⁰ Predictive analysis leads to morally contentious conclusions, such as those drawn by the (in)famous 2001 article of John Donohue and Steven Levitt, *The Impact of Legalized Abortion on Crime*, which argued that the legalization of abortion in the 1970s contributed significantly to reductions in crime rates experienced in the 1990s.⁹¹

D. Lack of Access and Exclusion

¶38 An additional concern raised by big data is that it tilts an already uneven scale in favor of organizations and against individuals. The big benefits of big data, the argument goes, accrue to government and big business, not to individuals—and they often come at

Art and Possible Extensions, 17 IEEE TRANSACTIONS ON KNOWLEDGE & DATA ENG'G 6, 1–23 (June 2005); *but see* Ryan Singel, *Netflix Spilled Your Brokeback Mountain Secret, Lawsuit Claims*, WIRED (Dec. 17, 2009), <http://www.wired.com/threatlevel/2009/12/netflix-privacy-lawsuit> (plaintiff argues that Netflix made it possible for her to be “outed” when it disclosed insufficiently anonymous information about her viewing habits, including films from the “Gay & Lesbian” genre).

⁸⁵ Rosamunde van Brakel & Paul De Hert, *Policing, Surveillance and Law in a Pre-Crime Society: Understanding the Consequences of Technology Based Strategies*, 20 J. POLICE STUD. 163 (2011).

⁸⁶ *See, e.g.*, Clive Thompson, *Google's China Problem (and China's Google Problem)*, N.Y. TIMES MAGAZINE, Apr. 23, 2006, <http://www.nytimes.com/2006/04/23/magazine/23google.html?pagewanted=all>. Google has since withdrawn from the Chinese market. David Drummond, *A New Approach to China*, GOOGLE OFFICIAL BLOG (Jan. 12, 2010), <http://googleblog.blogspot.com/2010/01/new-approach-to-china.html>

⁸⁷ Marc Rotenberg, *Foreword: Privacy and Secrecy after September 11*, 86 MINN. L. REV. 1115 (2002).

⁸⁸ Omer Tene, *Privacy: For the Rich or for the Poor?*, CONCURRING OPINIONS (July 26, 2012), <http://www.concurringopinions.com/archives/2012/07/privacy-for-the-rich-or-for-the-poor.html> (citing examples); *see also* Jennifer Valentino-Devries, Jeremy Singer-Vine & Ashkan Soltani, *Websites Vary Prices, Deals Based on Users' Information*, WALL ST. J., Dec. 24, 2012, <http://online.wsj.com/article/SB10001424127887323777204578189391813881534.html> (reporting examples of price discrimination based on online profiling and stating “the Journal's testing also showed that areas that tended to see the discounted prices had a higher average income than areas that tended to see higher prices.”).

⁸⁹ *Ecclesiastes* 1:9 (New Revised Standard Version).

⁹⁰ Jay Stanley, *Eight Problems with “Big Data,”* ACLU BLOG (Apr. 25, 2012), <https://www.aclu.org/blog/technology-and-liberty/eight-problems-big-data>.

⁹¹ John Donohue & Steven Levitt, *The Impact of Legalized Abortion on Crime*, 66(2) QUARTERLY J. ECON. 379 (2001). For criticism *see, e.g.*, Christopher Foote & Christopher Goetz, *The Impact of Legalized Abortion on Crime: Comment* (Federal Reserve Bank of Boston, Working Paper No. 05-15, 2005), available at <http://www.bos.frb.org/economic/wp/wp2005/wp0515.pdf> (distinguishing the role of abortion from other potential influences on crime, some of which vary year by year or state by state, including for example the “crack” epidemic, which rose and receded at different times in different places).

individuals' expense. In the words of the adage, "if you're not paying for it, you're not the customer; you're the product."⁹²

¶39 The exclusion of individuals from the benefits of the use of their data manifests in two main ways. First, online interactions are barter-like transactions where individuals exchange personal data for free services.⁹³ Yet those transactions appear to take place in an inefficient market hampered by steep information asymmetries, which are further aggravated by big data. Transacting with a big data platform is like a game of poker where one of the players has his hand open and the other keeps his cards close. The online company knows the preferences of the transacting individual inside and out, perhaps better than the individual knows him or herself. It can therefore usurp the entire value surplus available in the transaction by pricing goods or services as close as possible to the individual's reservation price.

¶40 Second, organizations are seldom prepared to share the wealth created by individuals' personal data with those individuals. In the *Guardian*, Sir Tim Berners-Lee recently remarked:

"My computer has a great understanding of my state of fitness, of the things I'm eating, of the places I'm at. My phone understands from being in my pocket how much exercise I've been getting and how many stairs I've been walking up and so on." Exploiting such data could provide hugely useful services to individuals, he said, but only if their computers had access to personal data held about them by web companies. "One of the issues of social networking silos is that they have the data and I don't."⁹⁴

¶41 The right of access granted to individuals under the European Data Protection Directive⁹⁵ and additional fair information principles has been implemented narrowly. Even where they comply with the law, organizations provide individuals with little useful information.

⁹² This phrase, which has become a staple in online culture, is attributed to a discussion on a MetaFilter community in August 2010. See Jonathan Zittrain, *Meme patrol: "When Something Online is Free, You're Not the Customer, You're the Product,"* THE FUTURE OF THE INTERNET (Mar. 21, 2012), <http://futureoftheinternet.org/meme-patrol-when-something-online-is-free-youre-not-the-customer-youre-the-product>.

⁹³ CHRIS ANDERSON, *FREE: THE FUTURE OF A RADICAL PRICE* (2009).

⁹⁴ Ian Katz, *Tim Berners-Lee: Demand Your Data from Google and Facebook*, THE GUARDIAN (Apr. 18, 2012), <http://www.guardian.co.uk/technology/2012/apr/18/tim-berners-lee-google-facebook>. See Bruce Upbin, *How Intuit Uses Big Data For the Little Guy*, FORBES (Apr. 26, 2012) <http://www.forbes.com/sites/bruceupbin/2012/04/26/how-intuit-uses-big-data-for-the-little-guy> ("Big Data means big challenges and big opportunities. But, hey, what about me? What do I (meaning the average joe) get out of all this? Companies are flying on the contrails of our spending, hiring and networking behavior, especially at the social/mobile colossi like Facebook, Google and Apple. We ought to see some of that value. Rather than just take take take, why can't more companies give back, reflect our data back on us? Doing this in a real, honest way has to create some business value.").

⁹⁵ Council Directive 95/46, art. 12, 1995 O.J. (L 281) 42 (EC); WHITE HOUSE, *CONSUMER DATA PRIVACY IN A NETWORKED WORLD: A FRAMEWORK FOR PROTECTING PRIVACY AND PROMOTING INNOVATION IN THE GLOBAL DIGITAL ECONOMY* 48 (Feb. 23, 2012) (noting, specifically, the Consumer Privacy Bill of Rights principle of "Access and Accuracy"), available at <http://www.whitehouse.gov/sites/default/files/privacy-final.pdf> [hereinafter WHITE HOUSE, *PRIVACY IN A NETWORKED WORLD*]

E. *The Ethics of Analytics: Drawing the Line*

¶42 Like any other type of research, data analytics can cross the threshold of unethical behavior. Consider the recent research by a Texas University developmental psychology professor, who logged and reviewed every text message, email, photo, and instant message sent by a group of 175 teenagers on Blackberries that she provided to them.⁹⁶ The participants and their parents were required to sign consent forms; yet, regardless of consent form legalese, it is doubtful that the minors could fully assess the implications of the omniscient surveillance.⁹⁷ Like children's data, other categories of sensitive data may be collected and analyzed for ethically dubious research. Consider a service analyzing individuals' preferences on pornography sites for use in behavioral advertising.⁹⁸ More complicated yet, the analysis of apparently innocuous data may create new sensitive facts about an individual, as occurred in the instance of Target's "pregnancy score,"⁹⁹ or which may be possible with a prediction of the onset of Alzheimer's disease. Where should the red line be drawn when it comes to big data analysis? Moreover, who should benefit from access to big data? Could ethical scientific research be conducted without disclosing to the general public the data used to reach the results?

F. *Chilling Effect*

¶43 As recently observed by Jay Stanley of the ACLU, "as the ramifications of big data analytics sink in, people will likely become much more conscious of the ways they're being tracked, and the chilling effects on all sorts of behaviors could become considerable."¹⁰⁰ The result is what the former UK privacy regulator dubbed "a surveillance society," a psychologically oppressive world in which individuals are cowed to conforming behavior by the state's potential panoptic gaze.¹⁰¹

III. THE LEGAL FRAMEWORK: CHALLENGES

¶44 How does the existing privacy framework deal with the big data phenomenon? This part reviews the FIPPs strained by the current technological and business landscape (including the definition of PII), the principles of data minimization and purpose

⁹⁶ Kashmir Hill, *A Texas University's Mind-Boggling Database of Teens' Daily Text Messages, Emails, and IMs over Four Years*, FORBES (Apr. 18, 2012), <http://www.forbes.com/sites/kashmirhill/2012/04/18/a-texas-universitys-mind-boggling-database-of-teens-daily-text-messages-emails-and-ims-over-four-years/>.

⁹⁷ See, e.g., Michael Zimmer, *Research Ethics and the Blackberry Project*, MICHAELZIMMER.ORG (Apr. 25, 2012), <http://michaelzimmer.org/2012/04/25/research-ethics-and-the-blackberry-project>.

⁹⁸ Kashmir Hill, *History Sniffing: How YouPorn Checks What Other Porn Sites You've Visited and Ad Networks Test the Quality of Their Data*, FORBES (Nov. 30, 2010), <http://www.forbes.com/sites/kashmirhill/2010/11/30/history-sniffing-how-youporn-checks-what-other-porn-sites-youve-visited-and-ad-networks-test-the-quality-of-their-data>.

⁹⁹ Duhigg, *supra* note 80 and accompanying text.

¹⁰⁰ Jay Stanley, *The Potential Chilling Effects of Big Data*, ACLU BLOG (Apr. 30, 2012), <http://www.aclu.org/blog/technology-and-liberty/potential-chilling-effects-big-data>.

¹⁰¹ *Watchdog's Big Brother U.K. Warning*, BBC (Aug. 16, 2004), http://news.bbc.co.uk/2/hi/uk_news/politics/3568468.stm (quoting U.K. Information Commissioner Richard Thomas as saying the U.K. could "sleepwalk into a surveillance society"). The paradigmatic example is, of course, GEORGE ORWELL, 1984 (1948).

limitation, and the concept of consent.¹⁰² This part also argues that, inevitably, these elements of the privacy framework should adjust to reflect existing technological and organizational realities, which include ubiquitous data collection and individuals who are ill-placed to meaningfully review privacy policies. Together with the next part, it argues that the FIPPs should be used as a set of levers, which can be modulated to address big data by relaxing the principles of data minimization and individual control while tightening requirements for transparency, access, and accuracy.

A. Definition of PII

¶45 Traditionally, de-identification was viewed as a silver bullet allowing organizations to reap the benefits of analytics while preserving individuals' privacy.¹⁰³ Organizations used various methods of de-identification (anonymization, pseudonymization, encryption, key-coding, data sharing) to distance data from personal identities.¹⁰⁴ Yet, over the past few years, computer scientists have repeatedly shown that even anonymized data can typically be re-identified and associated with specific individuals.¹⁰⁵ De-identified data, in other words, is a temporary state rather than a stable category.¹⁰⁶ In an influential law review article, Paul Ohm observed that “[r]e-identification science disrupts the privacy policy landscape by undermining the faith that we have placed in anonymization.”¹⁰⁷ The implications for government and businesses can be stark, because de-identification has become a key component of numerous business models, most notably in the context of health data (e.g., clinical trials), online behavioral advertising, and cloud computing.

¶46 The first major policy question raised by the big data phenomenon concerns the scope of information subject to privacy law. How robust must de-identification be in

¹⁰² This article deals with the adjustment of the existing privacy framework to accommodate big data realities. Other privacy issues raised by big data, such as government access to or surveillance of private sector databases, are beyond the scope of this paper. See, e.g., James X. Dempsey & Lara M. Flint, *Commercial Data and National Security*, 72 GEO. WASH. L. REV. 1459 (2004).

¹⁰³ See, e.g., Working Party, *Opinion 4/2007 on the Concept of Personal Data*, Article 29 (June 20, 2007), available at http://ec.europa.eu/justice/policies/privacy/docs/wpdocs/2007/wp136_en.pdf; MARKLE FOUNDATION TASK FORCE ON NATIONAL SECURITY IN THE INFORMATION AGE, CREATING A TRUSTED NETWORK FOR HOMELAND SECURITY (2003); Ira S. Rubinstein, Ronald D. Lee & Paul M. Schwartz, *Data Mining and Internet Profiling: Emerging Regulatory and Technological Approaches*, 75 U. CHI. L. REV. 261, 268–29 (2008).

¹⁰⁴ See W. Kuan Hon, Christopher Millard & Ian Walden, *The Problem of 'Personal Data' in Cloud Computing—What Information is Regulated? The Cloud of Unknowing*, 4 INT'L DATA PRIVACY L. 1, 211–228 (Mar. 15, 2011), available at <http://idpl.oxfordjournals.org/content/1/4/211.full.pdf+html>.

¹⁰⁵ This line of research was pioneered by Latanya Sweeney and made accessible to lawyers by Paul Ohm. Ohm, *supra* note 69; Latanya Sweeney, *Uniqueness of Simple Demographics in the U.S. Population*, (Laboratory for International Data Privacy, Working Paper No. 4, 2000). See also Narayanan & Shmatikov, *supra* note 67; Arvind Narayanan & Vitaly Shmatikov, *Myths and Fallacies of “Personally Identifiable Information,”* 53 COMMUNICATIONS OF THE ACM 6, 24 (2010); Arvind Narayanan et al., *On the Feasibility of Internet-Scale Author Identification*, 2012 INSTITUTE OF ELECTRICAL AND ELECTRONICS ENGINEERS SYMPOSIUM ON SECURITY & PRIVACY 300 (2012).

¹⁰⁶ See Ed Felten, *Does Hashing Make Data “Anonymous”?*, TECH@FTC (Apr. 22, 2012, 7:05 AM) <http://techatftc.wordpress.com/2012/04/22/does-hashing-make-data-anonymous>; Ed Felten, *Are pseudonyms “anonymous”?*, TECH@FTC (Apr. 30, 2012, 12:03 PM) <http://techatftc.wordpress.com/2012/04/30/are-pseudonyms-anonymous>; Felten, *supra* note 66 (series of blog posts by Ed Felten, the former Chief Technologist for the FTC and a Professor of Computer Science at Princeton).

¹⁰⁷ Ohm, *supra* note 69, at 1704.

order to “liberate” data from the throes of privacy legislation? One possible conclusion, apparently supported by Ohm himself, is that all data should be treated as PII and subjected to the regulatory framework.¹⁰⁸ Yet, such a result would create perverse incentives for organizations to forgo de-identification altogether and therefore increase, not alleviate, privacy and data security risks.¹⁰⁹ A further pitfall is that with a vastly expanded definition of PII, the privacy framework would become all but unworkable. Difficult enough to comply with and enforce today, the current framework may well be unmanageable if it extends to every piece of information.¹¹⁰ Moreover, while anonymized information always carries some risk of re-identification, many of the most pressing privacy risks exist only if there is reasonable likelihood of re-identification. As uncertainty is introduced into the re-identification equation, we cannot know whether the information truly corresponds to a particular individual, and the dataset becomes more anonymous as larger amounts of uncertainty are introduced.¹¹¹

¶47 More importantly, many beneficial uses of data would be severely curtailed if information, ostensibly not about individuals, comes under full remit of privacy laws based on a remote possibility of being linked to an individual at some point in time through some conceivable method, no matter how unlikely to be used.¹¹² Such an approach presumes a value judgment has been made in favor of individual control over highly beneficial uses of data, such as Dr. Altman’s discovery of the Paxil-Pravachol side effect; yet, it is doubtful that such a value choice has consciously been made.

¶48 PII should instead be defined based on a risk matrix taking into account the risk, intent, and potential consequences of re-identification, as opposed to a dichotomy between “identifiable” and “non-identifiable” data.¹¹³ A bi-polar approach based on labeling information either “personally identifiable” or not, is unhelpful and inevitably leads to an inefficient arms race between de-identifiers and re-identifiers. In this process, the integrity, accuracy, and value of the data may be degraded or lost, together with some of its potential societal benefits.¹¹⁴

¹⁰⁸ See *id.* at 1742.

¹⁰⁹ Ann Cavoukian & Khaled El Emam, *Dispelling the Myths Surrounding De-identification: Anonymization Remains a Strong Tool for Protecting Privacy* 7, INFO. & PRIVACY COMM’R OF ONTARIO (2011) (Can.), available at <http://www.ipc.on.ca/images/Resources/anonymization.pdf>.

¹¹⁰ For example, according to a 2010 report by the EU Agency for Fundamental Rights, even in Europe, data protection authorities lack sufficient independence and funding; impose few sanctions for violations of data protection laws; and “are not endowed with full powers to investigate, intervene, offer legal advice and engage in legal proceedings.” *Data Protection in the European Union: The Role of National Data Protection Authorities*, E.U. AGENCY FOR FUNDAMENTAL RIGHTS, p. 8 (May 7, 2010), available at http://fra.europa.eu/sites/default/files/fra_uploads/815-Data-protection_en.pdf.

¹¹¹ Betsy Masiello & Alma Whitten, *Engineering Privacy in an Age of Information Abundance*, 2010 ASS’N FOR THE ADVANCEMENT OF ARTIFICIAL INTELLIGENCE SPRING SYMPOSIUM SERIES 119, 122 (2010).

¹¹² See, e.g., Kathleen Benitez & Bradley Malin, *Evaluating Re-identification Risks with Respect to the HIPAA Privacy Rule*, 17 J. AM. MED. INFO. ASS’N. 169 (2010) (demonstrating actual risk of re-identification may be low).

¹¹³ Paul M. Schwartz & Daniel J. Solove, *The PII Problem: Privacy and a New Concept of Personally Identifiable Information*, 86 NYU L. REV. 1814 (2011); Omer Tene, *The Complexities of Defining Personal Data: Anonymization*, 8 DATA PROT. L. & POLICY 8, 6 (2011).

¹¹⁴ Daniel C. Barth-Jones, *Balancing Privacy Protection with Scientific Accuracy: Challenges for De-identification Practice*, FIRST ANNUAL COMPARATIVE EFFECTIVENESS RESEARCH SYMPOSIUM 1, 27 (2010).

¶49 A better solution would be, first, to view the identifiability of data as a continuum as opposed to the current dichotomy.¹¹⁵ This means adopting a scaled approach, under which data that are only identifiable at great cost would remain within the legal framework, subject to only a subset of fair information principles.¹¹⁶ Second, the approach that should be adopted is the one proposed by the Federal Trade Commission (FTC) in its recent report *Protecting Consumer Privacy in an Era of Rapid Change*,¹¹⁷ which overlays the statistical probability of re-identifiability with legally enforceable organizational commitments as well as downstream contractual obligations not to re-identify or to attempt to re-identify. According to the FTC, “as long as (1) a given data set is not reasonably identifiable, (2) the company publicly commits not to re-identify it, and (3) the company requires any downstream users of the data to keep it in de-identified form, that data will fall outside the scope of the framework.”¹¹⁸ Recognizing that it is virtually impossible to guarantee privacy by scrutinizing the data alone, without defining and analyzing its intended uses, the FTC shifts the crux of the inquiry from a *factual* test of identifiability to a *legal* examination of an organization’s *intent* and *commitment* to prevent re-identification.

¶50 Finally, we advocate viewing de-identification as an important protective measure to be taken under the data security and accountability principles, rather than a solution to the big data conundrum.¹¹⁹ Organizations collecting and harvesting big data would be wise to de-identify data to the extent possible while not compromising their beneficial use. At the same time, the privacy framework will continue to partially apply to de-identified data because researchers have the ability to re-link almost any piece of data to an individual, if provided appropriate incentive to do so.

B. Data Minimization

¶51 Through various iterations and formulations, data minimization has remained a fundamental principle of privacy law.¹²⁰ Organizations are required to limit the collection of personal data to the minimum extent necessary to obtain their legitimate goals. Moreover, they are required to delete data that is no longer used for the purposes for which they were collected and to implement restrictive policies with respect to the retention of personal data in identifiable form. The big data business model is antithetical to data minimization. It incentivizes collection of more data for longer periods of time. It is aimed precisely at those unanticipated secondary uses, the “crown jewels” of big

¹¹⁵ Schwartz & Solove, *supra* note 113, at 1879.

¹¹⁶ Schwartz & Solove, *supra* note 113.

¹¹⁷ FTC, PROTECTING CONSUMER PRIVACY IN AN ERA OF RAPID CHANGE, RECOMMENDATIONS FOR BUSINESSES AND POLICYMAKERS (2012).

¹¹⁸ *Id.* at 22.

¹¹⁹ See ORG. FOR ECON. CO-OPERATION AND DEV., OECD GUIDELINES ON THE PROTECTION OF PRIVACY AND TRANSBORDER FLOWS OF PERSONAL DATA ¶ 14 (1980); See WHITE HOUSE, PRIVACY IN A NETWORKED WORLD, *supra* note 95; Working Party, *Opinion 3/2010 on the Principle of Accountability*, Article 29, (July 13, 2010), available at http://ec.europa.eu/justice/policies/privacy/docs/wpdocs/2010/wp173_en.pdf.

¹²⁰ See ORG. FOR ECON. CO-OPERATION AND DEV., OECD GUIDELINES ON THE PROTECTION OF PRIVACY AND TRANSBORDER FLOWS OF PERSONAL DATA ¶¶ 7–8 (1980); Council Directive 95/46, 1995 O.J. (L 281) 31 (EC); THE WHITE HOUSE, PRIVACY IN A NETWORKED WORLD: A FRAMEWORK FOR PROTECTING PRIVACY AND PROMOTING INNOVATION IN THE GLOBAL DIGITAL ECONOMY 1 (2012).

data. After all, who could have anticipated that Bing search queries would be used to unearth harmful drug interactions?¹²¹

¶152 Here too, legal rules collide with technological and business realities. Organizations today collect and retain personal data through multiple channels including the Internet, mobile, biological and industrial sensors, video, e-mail, and social networking tools. Modern organizations amass data collected directly from individuals or third parties, and they harvest private, semi-public (e.g., Facebook), or public (e.g., the electoral roll) sources. Data minimization is simply no longer the market norm.

¶153 In considering the fate of data minimization, the principles of privacy law must be balanced against additional societal values such as public health, national security and law enforcement, environmental protection, and economic efficiency. A coherent framework should be based on a risk matrix that weighs the value of data against potential privacy risks. Where prospective data uses are highly beneficial and privacy risks minimal, the legitimacy of processing should be assumed even if individuals decline (or are not asked) to consent. For example, web analytics—the measurement, collection, analysis, and reporting of internet data for purposes of understanding and optimizing web usage—creates great value by ensuring that products and services can be improved to better serve consumers. Privacy risks are minimal because analytics, if properly implemented, deals with statistical data, typically in de-identified form.¹²² Yet requiring online users to opt into analytics would no doubt severely limit its application and use.

¶154 This is not to suggest, of course, that data should be collected exclusively in instances where it may become useful or that data collected for one purpose may be repurposed at will. Rather, in a big data world, the principle of data minimization should be interpreted differently, requiring organizations to de-identify data when possible, implement reasonable security measures, and limit uses of data to those that are acceptable from not only an individual but also a societal perspective.

C. Individual Control and Context

¶155 Legal frameworks all over the world continue to emphasize consent, or individual control, as a fundamental principle of privacy law. In the United States, “notice and choice” has been the central axis of privacy regulation for more than a decade.¹²³ In the European Union, consent remains the most commonly used basis to legitimize data processing under Article 7 of the Data Protection Directive.¹²⁴ By emphasizing consent,

¹²¹ *Supra* notes 29–34 and accompanying text.

¹²² *See, e.g.*, Matt McGee, German Government Says Google Analytics Now Verboten, SEARCH ENGINE LAND, January 12, 2011, <http://searchengineland.com/german-govt-says-google-analytics-now-verboten-61109>. Much of the criticism of analytics has been driven by careless practices such as the inadvertent leakage of personal data passed from sites to ad networks, misuse of flash cookies, or concerns that data were being used for behavioral advertising. *See* Paul M. Schwartz, *Data Protection Law and the Ethical Use of Analytics*, BUREAU OF NAT’L AFFAIRS PRIVACY & SEC. LAW REPORT, Jan. 10, 2011.

¹²³ A shift away from notice and choice is underway, as reflected in the Whitehouse Blueprint and FTC Final Report; yet, under both frameworks notice and choice remains a central principle. *See* The Whitehouse, *supra* note 13; FTC, *supra* note 117.

¹²⁴ Council Directive 95/46, 1995 O.J. (L 281) 31 (EC); Working Party, *Opinion 15/2011 on the Definition of Consent*, Article 29 (July 13, 2011), available at http://ec.europa.eu/justice/policies/privacy/docs/wpdocs/2011/wp187_en.pdf. The European Data Protection Regulation would significantly tighten the requirements for consent, effectively permitting only explicit consent and thereby presumably narrowing the scope of consent-based processing. Council

existing privacy frameworks impose significant, sometimes unrealistic, obligations on both organizations and individuals. On the one hand, organizations are expected to explain their data processing activities on increasingly small screens and obtain consent from often-uninterested individuals; on the other hand, individuals are expected to read and understand complicated privacy disclosures and express their “informed” consent.¹²⁵ This takes place against an increasingly complex backdrop in which data flows are handled through intricate arrangements involving dense networks of platforms and applications, including contractors, subcontractors, and service providers operating globally. Moreover, to be meaningful, consent must be specific to the purpose (or context). Yet by its very nature, big data analysis seeks surprising correlations and produces results that resist prediction.

¶156 The consent model is flawed from an economic perspective. Information asymmetries and well-documented cognitive biases cast a shadow on the authenticity of individuals’ privacy choices. For example, Alessandro Acquisti and his colleagues have shown that simply by providing users a *feeling* of control, businesses can encourage the sharing of data regardless of whether or not users actually gained control.¹²⁶ Joseph Turow and others have shown that “[w]hen consumers see the term ‘privacy policy,’ they believe that their personal information will be protected in specific ways; in particular, they assume that a website that advertises a privacy policy will not share their personal information.”¹²⁷ In reality, however, this is not the case. It is common knowledge among practitioners in the field that privacy policies serve more as liability disclaimers for businesses than as assurances of privacy for consumers.

¶157 At the same time, collective action problems threaten to generate a suboptimal equilibrium where individuals fail to opt into societally beneficial data processing in the hope of free-riding on others’ good will. Consider, for example, Internet browser crash reports, which very few users opt into; even when they do make such an election, they are often motivated less by real privacy concerns than by a (misplaced) belief that others will do the job for them. As is often the case in public opinion polling, the precise wording of choice menus presented to individuals has a disproportionate effect on their decisions to opt in or out of such polling. It seems likely that if prompted, most search engine users would decline the search engine permission to analyze their search logs for the detection of harmful drug interactions. Yet, when asked in retrospect about the actions of Dr. Altman and his team, the same users may find them commendable.

¶158 Similar free-riding is common in other contexts where the difference between opt-in and opt-out regimes is stark. This is the case, for example, with organ donation rates.

Directive 95/46, 1995 O.J. (L 281) 31 (EC) (defining “the data subject’s consent” as “any freely given specific informed and explicit indication of his wishes.”).

¹²⁵ Aleecia M. McDonald & Lorrie Faith Cranor, *The Cost of Reading Privacy Policies*, 2008 Privacy Year in Review issue I/S: J. L. & POL’Y FOR INFO. SOC’Y 1, 17 (2008) (finding that to read every privacy policy encountered, an average individual would need to spend approximately 30 working days per year); see also Alexis Madrigal, *Reading the Privacy Policies You Encounter in a Year Would Take 76 Work Days*, THE ATLANTIC, Mar. 1, 2012, 2:25 PM, <http://www.theatlantic.com/technology/archive/2012/03/reading-the-privacy-policies-you-encounter-in-a-year-would-take-76-work-days/253851>.

¹²⁶ Laura Brandimarte et al., *Misplaced Confidences: Privacy and the Control Paradox*, NINTH ANNUAL WORKSHOP ON THE ECONOMICS OF INFORMATION SECURITY (WEIS) (2010).

¹²⁷ Joseph Turow et al., *The Federal Trade Commission and Consumer Privacy in the Coming Decade*, 3(3) I/S: J. L. & POL’Y FOR INFO. SOC’Y 723, 724 (2007).

In countries where citizens must opt in for organ donation, donation rates tend to be very low compared to countries that are culturally similar but have an opt-out regime. This concept is illustrated by the comparative donation rates in Sweden (85.9% under an opt-out regime) versus Denmark (4.25% under an opt-in regime), and in Austria (99.9% under an opt-out regime) versus Germany (12% under an opt-in regime).¹²⁸

¶59 An additional problem is that consent-based processing tends to be regressive because individuals' expectations fall back on existing experiences. For example, if Facebook had not proactively launched its News Feed feature in 2006 and had instead waited for users to opt in,¹²⁹ users might not have enjoyed Facebook as it is known today. It is only when data started flowing that users became accustomed to the change. Similarly, few individuals would have agreed had Google solicited consent (or regulatory approval) for "wardriving"¹³⁰ through cities all over the world to create a comprehensive map of Wi-Fi networks for its geo-location services.¹³¹ Yet in retrospect, after Google provided users an opportunity to opt out their routers, it is doubtful that many users have actually done so.¹³² The decisions by regulators in this case indicate some appreciation for the value of Google's data use, even if this rationale was not clearly expressed.

¶60 This article does not argue that individuals should *never* be asked to expressly consent to the use of their information or offered an option to opt out. Rather, it suggests that the merits of a given data use should be debated as a broader societal issue. Does society believe that direct marketing, behavioral advertising, third-party data brokering, and location-based services are legitimate (or even commendable) models that are worth pursuing or excessive intrusions that should be deterred? When making decisions about the need for individuals' consent and how it should be obtained, policymakers should recognize that default rules often prevail and determine the existence of these data uses. Too often, debates about whether consent should be solicited or opt-out choice provided focus solely on the mechanics of expressing consent.¹³³ But heightened focus on consent and data minimization, with little appreciation for the value of data use, could jeopardize innovation and beneficial societal advances.

¶61 The legitimacy of data use had always intended to take additional values into account beyond privacy. For example, law enforcement has traditionally been allotted a

¹²⁸ Notice that additional factors besides individuals' willingness to participate, such as funding and regional organization, affect the ultimate conversion rate for organ transplants. Hence, Austria, which has an opt-out system, had a deceased organ transplant rate of 20.6 per million people (pmp), whereas the United States, with an opt-in system, had a rate of 26.3 pmp.

¹²⁹ The initial product launch was accompanied by a privacy furor leading Facebook to retract the service, which was rolled out with adjustments. See Mark Zuckerberg, *An Open Letter from Mark Zuckerberg*, THE FACEBOOK BLOG (Sept. 8, 2006 2:48 a.m.), <https://blog.facebook.com/blog.php?post=2208562130>.

¹³⁰ Wardriving, WIKIPEDIA, <http://en.wikipedia.org/wiki/Wardriving> (defining "wardriving" as "the act of searching for Wi-Fi wireless networks by a person in a moving vehicle, using a portable computer, smartphone or personal digital assistant (PDA)").

¹³¹ Google's "wardriving" is featured in a privacy snafu still being investigated by regulators around the globe; yet, it concerns the capture by Google of unencrypted payload (content) data – not the practice of mapping Wi-Fi networks. See Kevin O'Brien, *European Regulators May Reopen Street View Inquiries*, N.Y. TIMES, May 2, 2012, http://www.nytimes.com/2012/05/03/technology/european-regulators-to-reopen-google-street-view-inquiries.html?_r=2.

¹³² Kevin O'Brien, *Google Allows Wi-Fi Owners to Opt Out of Database*, N.Y. TIMES, Nov. 15, 2011, <http://www.nytimes.com/2011/11/16/technology/google-allows-wi-fi-owners-to-opt-out-of-database.html>.

¹³³ For extensive discussion see Tene & Polonetsky, *supra* note 54.

degree of freedom to override privacy restrictions in appropriate cases with the satisfaction of due process requirements.¹³⁴ Consequently, the role of consent should be demarcated according to normative choices made by policymakers with respect to prospective data uses. In some cases, consent should not be required, while in others, consent should be assumed subject to a right of refusal. In specific cases, consent should be required to legitimize data use.

IV. THE LEGAL FRAMEWORK: SOLUTIONS

¶62 This part argues that while relaxing the principles of data minimization and consent, the current privacy framework should stress access and transparency. It explores how individuals can be empowered with enhanced transparency and access rights, thereby rebalancing the framework and creating additional opportunity for efficient value creation and innovation. It argues that if individuals were provided access to their information in machine-readable (heretofore, “usable”) format, the personal information ecosystem would expand; layers upon layers of user-side applications are likely to emerge to harvest information to benefit not only organizations, but also individuals. This part further suggests that, subject to the protection of trade secrets, organizations should be required to reveal the criteria used in their decision-making processes with respect to personal data analysis. Such a requirement will likely discourage unethical, if not illegal, classifications and provide individuals with the due process opportunity to challenge decisions made about them by algorithm-driven machines.

A. Access, Portability, and Sharing the Wealth

¶63 The right to access and rectify one’s individual information—while one of the fundamental principles of information privacy—remains woefully underutilized.¹³⁵ Few individuals are aware of their access rights and even fewer exercise them.¹³⁶ And why should they? Access rights are neither convenient nor particularly useful. Organizations typically provide access to data only in “hardcopy,” after weeks or months of delays arising from correspondence and requests for authentication and payment of fees. Organizations often fail to provide details about sources, uses, and recipients of the information they collect, and seek to rely on a panoply of legal exemptions to mask portions of the data that they do disclose. The increasing complexity of the data ecosystem renders it difficult for individuals to determine to whom an access request should be sent. Furthermore, processors or sub-processors of data are often based in foreign jurisdictions, without a consumer-facing interface to handle individual requests. Indeed, one user’s quest to obtain his personal information from Facebook was so novel that it commanded headlines in newspapers all over the world, including the *New York Times*.¹³⁷

¹³⁴ Law enforcement provisions are also increasingly being limited due to concerns of potential abuse.

¹³⁵ See, e.g., Singer, *supra* note 28.

¹³⁶ A Eurobarometer survey of 2008 found that across the EU, just over a half of the citizens were aware of the right; far fewer had ever exercised it. Eurobarometer, *Data Protection in the European Union Citizens’ perceptions Analytical Report*, 30 (Feb. 2008), http://ec.europa.eu/public_opinion/flash/fl_225_en.pdf.

¹³⁷ Kevin O’Brien, *Austrian Law Student Faces Down Facebook*, N.Y. TIMES, Feb. 5, 2012,

¶164 As a *quid pro quo* for looser data collection and minimization restrictions, organizations should be prepared to share the wealth created by individuals' data with those individuals. This means providing individuals with access to their data in a "usable" format and allowing them to take advantage of third party applications to analyze their own data and draw useful conclusions (e.g., consume less protein, go on a skiing vacation, invest in bonds).

¶165 This "featurization" of big data will unleash innovation and create a market for personal data applications.¹³⁸ The technological groundwork has already been completed with mash-ups and real-time application programming interfaces (APIs),¹³⁹ making it easier for organizations to combine information from different sources and services into a single user experience. Much like open-source software or Creative Commons licenses, free access to personal data is grounded in both efficiency and fairness rationales. Regardless of whether or not you accept a property approach to personal information,¹⁴⁰ fairness dictates that individuals enjoy beneficial use of their data.

¶166 The roll out of the smart grid illustrates this point. Electric utilities reap most of the benefits associated with upgrading the electric grid to provide bi-directional communications. This explains why the smart grid was met by pushback from consumers and regulators who are concerned with its implications for privacy, data security, start-up costs, and dynamic pricing. Had consumers felt the beneficial impact of the smart grid themselves, they may have reacted differently. That is precisely the idea behind the Obama Administration's "Green Button" initiative: the initiative establishes that consumers should have access to their own energy usage information in a downloadable, standard, easy-to-use electronic format.¹⁴¹ In a speech on September 15, 2011, Aneesh Chopra, the U.S. Chief Technology Officer, challenged the industry to "publish information online in an open format (machine-readable) without restrictions that would impede re-use."¹⁴² In January 2012, three major California utilities announced their implementation of the Green Button,¹⁴³ and a dozen more utilities followed suit in the first quarter of 2012.¹⁴⁴

¶167 The Administration predicted that making user data available to the public would lead entrepreneurs to develop technologies like energy management systems and smartphone applications that can interpret and use such information.¹⁴⁵ Homeowners, in

<http://www.nytimes.com/2012/02/06/technology/06iht-rawdata06.html?pagewanted=all>.

¹³⁸ Such a market is already picking up. See, e.g., Francesca Robin, *The Emerging Market that Could Kill the iPhone*, FORTUNE (Aug. 1, 2012), <http://tech.fortune.cnn.com/2012/08/01/iphone>.

¹³⁹ An application programming interface allows third party software developers to interface with a given platform or software component.

¹⁴⁰ See discussion and criticism of the property approach in Julie Cohen, *Examined Lives: Informational Privacy and the Subject as Object*, 52 STAN. L. REV. 1373 (2000).

¹⁴¹ National Institute of Standards and Technology, *Green Button Initiative Artifacts Page*, <http://collaborate.nist.gov/twiki-sggrid/bin/view/SmartGrid/GreenButtonInitiative> (Sept. 2011).

¹⁴² Aneesh Chopra, *Remarks to GridWeek* (Sept. 15, 2011),

<http://www.whitehouse.gov/sites/default/files/microsites/ostp/smartgrid09-15-11.pdf>.

¹⁴³ The utilities are Pacific Gas & Electric, Southern California Edison and San Diego Gas & Electric. See Jim Witkin, *Pushing the Green Button for Energy Savings*, N.Y. TIMES (Jan. 20, 2012), <http://green.blogs.nytimes.com/2012/01/20/a-phone-app-for-turning-down-the-thermostat>.

¹⁴⁴ Press Release, SGCC Members Lead Industry in Green Button Initiative: Fifteen members of consumer-focused smart grid nonprofit sign on to Green Button (April 25, 2012), <http://smartgridcc.org/wp-content/uploads/2012/04/Green-Button-Press-Release.pdf>.

¹⁴⁵ Aneesh Chopra, *Modeling a Green Energy Challenge after a Blue Button*, The White House Office of

turn, would seek out applications that enable them to gain greater control over their energy use. Chopra emphasized the importance of providing the data in a *standard format* according to industry-accepted guidelines. A standard, usable format fosters innovation by allowing software developers to create a single version of their product that will work for all utility customers across the country. One developer told the *New York Times* that his company had “created a set of software development tools that had already attracted 150 app developers. His company also plans to set up an online marketplace, similar to Apple’s iPhone App Store or Google’s Android Market, where homeowners could download energy-related applications.”¹⁴⁶

¶168 Accessing information about energy consumption for cost savings and novel usage is not solely the domain of utilities. For example, the Nest Learning Thermostat, developed by Nest Labs, is an energy conserving, self-programming, slickly designed home thermostat. It is also Wi-Fi connected to allow users to adjust their home or office temperature via an iPhone or Android app from anywhere they happen to be.¹⁴⁷ Like the Green Button, the Nest Learning Thermostat lets users tap into their own data trail, which includes their movements about the house and information about their daily routine. Major communications providers such as AT&T, Verizon, and Comcast have also launched innovative home services focused on energy management and home security and control.¹⁴⁸

¶169 The concept of the “Green Button” follows a path charted by a similar initiative in the field of health data. In 2010, the Obama Administration announced the “Blue Button,” a web-based feature through which patients can easily download their health information in usable format and share it with health care providers and trusted third parties. To make the information more useful, the initiative challenged developers to create applications that build on the Blue Button by helping consumers use their data to manage their own health. In turn, applications such as the Blue Button Health Assistant, developed by Adobe, sprung up to facilitate linkage of patient information, including immunizations, allergies, medications, family health history, lab test results, and more.¹⁴⁹

¶170 An additional government program based on a similar mind-set is the “Data.gov” initiative. Government has long been the biggest generator, collector, and user of data (not necessarily PII), keeping records on every birth, marriage, and death, compiling figures on all aspects of the economy, and maintaining statistics on licenses, laws, and the weather. Until recently, all of the data was locked and hard to locate, even if publicly accessible.¹⁵⁰ In many countries, a freedom of information request to obtain information

Science and Technology Policy (Sept. 15, 2011), <http://www.whitehouse.gov/blog/2011/09/15/modeling-green-energy-challenge-after-blue-button>.

¹⁴⁶ *Id.*

¹⁴⁷ David Pogue, *A Thermostat That’s Clever, Not Clunky*, N.Y. TIMES, Nov. 30, 2011, <http://www.nytimes.com/2011/12/01/technology/personaltech/nest-learning-thermostat-sets-a-standard-david-pogue.html>.

¹⁴⁸ See, e.g., Jordan Crook, *AT&T Introduces Digital Life: IP-Based Home Automation and Security System With 24/7 Monitoring Centers*, TECHCRUNCH (May 7, 2012), <http://techcrunch.com/2012/05/06/att-introduces-digital-life-ip-based-home-automation-and-security-system-with-247-monitoring-centers>.

¹⁴⁹ Aneesh Chopra, Todd Park & Peter Levin, ‘Blue Button’ Provides Access to Downloadable Personal Health Data, WHITE HOUSE BLOG (Oct. 7, 2010), <http://www.whitehouse.gov/blog/2010/10/07/blue-button-provides-access-downloadable-personal-health-data>.

¹⁵⁰ See, e.g., Amanda Conley, Anupam Datta, Helen Nissenbaum & Divya Sharma, *Sustaining Privacy and Open Justice in the Transition to Online Court Records: A Multidisciplinary Inquiry*, 71 MD. L. REV.

about the budgetary process, for example, would yield, at best, a voluminous PDF document locked for editing and difficult to explore. The Obama Administration, led by United States Chief Information Officer Vivek Kundra, embraced this innovation by launching “Data.gov.” The stated purpose of the new website was “to increase public access to high value, machine-readable datasets generated by the Executive Branch of the Federal Government.”¹⁵¹ The opening of the government’s data coffers unleashed a wave of innovation and helped create new economic value, as individuals and businesses used raw data to improve existing services and offer new solutions.¹⁵²

¶71 Increased use by individuals of their own data is also evident in the private sector. Various existing business models seek to arbitrate between users and organizations in order to tilt the scale back in favor of individuals. The Harvard Berkman Center’s “ProjectVRM” (“VRM” stands for “vendor relationship management”), which set an admittedly “immodest ambition of turning business on its head,” seeks to “provide customers with tools that provide both independence from vendor lock-in and better ways of engaging with vendors—on terms and by means that work better for both sides.”¹⁵³ In his 2012 book, *The Intention Economy*, ProjectVRM’s leader Doc Searls posits a vision of a world where an individual is in complete control of her digital persona and grants permissions for vendors to access it on her own terms. In this world, individuals would use software applications to signal their needs, which vendors would then compete to fulfill.¹⁵⁴

¶72 Personal.com, for example, is a start up company that enables individuals to own, control access to, and benefit from their personal information.¹⁵⁵ It does so by providing individuals with an online “data vault” divided into compartments called “gems,” where they can store and share information about their shopping habits, travel, log-in credentials on various sites, location information, and more.¹⁵⁶ There are currently more than 100 gems with more than 3,000 fields of data. The food preferences gem, for example, includes allergies, religious and dietary restrictions, and whether a user likes spicy food. Users can share gems with family, friends, employees or colleagues, and more importantly, monetize their own data by selling access to gems to commercial entities. (Personal.com collects a ten percent fee on such sales.) The company’s founders hope that Personal.com will become more than just a data vault, but rather a platform allowing applications to connect to structured user information.¹⁵⁷

772 (2012).

¹⁵¹ *The Open Society: Governments Are Letting in the Light*, THE ECONOMIST, Feb. 25, 2010, <http://www.economist.com/node/15557477>.

¹⁵² *Id.*

¹⁵³ *ProjectVRM*, HARVARD UNIV. BERKMAN CENTER FOR INTERNET AND SOCIETY, <http://cyber.law.harvard.edu/research/projectvr/#>.

¹⁵⁴ Doc Searls, *THE INTENTION ECONOMY: WHEN CUSTOMERS TAKE CHARGE* (Harvard Business Review Press, 2012); see also Joe Andrieu, *Introducing User Driven Services*, JOEANDRIEU.COM, Apr. 26, 2009 (series of ten blog posts), <http://blog.joeandrieu.com/2009/04/26/introducing-user-driven-services>.

¹⁵⁵ Thomas Heath, *Web Site Helps People Profit from Information Collected About Them*, WASH. POST (June 26, 2011), http://www.washingtonpost.com/business/economy/web-site-helps-people-profit-from-information-collected-about-them/2011/06/24/AGPgkRmH_story.html.

¹⁵⁶ For personal data vaults, see also Jerry Kang, Katie Shilton, Deborah Estrin, Jeff Burke & Mark Hansen, *Self-Surveillance Privacy*, 97 IOWA L. REV. 809 (2012) (proposing “personal data guardians” to curate the personal data vaults).

¹⁵⁷ See *Rethinking Personal Data: Strengthening Trust*, World Economic Forum 26–27 (May 2012), https://www.bcgperspectives.com/Images/Rethinking_Personal_Data_1005_light_tcm80-105516.pdf.

¶73 Another example is Intuit’s use of data gleaned from its Quickbooks and TurboTax products, which are used by millions of small businesses and individuals for accounting and tax filings. One new feature added to Quickbooks in 2012 is Easy Saver, which looks for items small business owners purchased frequently and then finds a better price for such items using negotiated high-volume discounts. Users will not see an offer for an item unless they have already bought it and are likely (based on previous purchasing behavior) to need it again soon. “The Trends feature in Quickbooks . . . tells business owners how their key indicators such as sales, operating margin and payroll costs compare with similar small businesses in their area or in the U.S. overall.”¹⁵⁸

¶74 If users fail to exercise their access and rectification rights, why should we expect them to actively engage with their data? The answer is that they are already doing so through a plethora of Apple, Android, and Facebook applications.¹⁵⁹ The entire “app economy” is premised on individuals accessing their own data for novel uses, ranging from GPS programs and restaurant recommendations to self-tailored financial and health services.¹⁶⁰ Applications have become an integral aspect of how users experience social networks and the mobile Internet. They enable individuals to make innovative use of their list of friends on Facebook, address books, Wi-Fi router locations, and many other sources of data. A recent study found that the app economy has created 466,000 jobs in the United States since 2007.¹⁶¹ According to Facebook’s S-1 filing ahead of its IPO, Zynga, an app developer, is responsible for 12% of Facebook’s revenue estimated at more than \$4 billion.¹⁶²

¶75 This article suggests the development of apps for the big data silos of the many companies who have focused on the collection and analysis of personal data for their own use.¹⁶³ Recent market initiatives demonstrate the feasibility of business models based on empowering individual users.¹⁶⁴ What the government seeks to achieve with its Green Button and Blue Button initiatives can and should be replicated in the private sector.

¶76 The call for additional access and transparency echoes one of the fundamental rationales for information privacy law—the prevention of secret databases. From its inception, information privacy law has been modeled to alleviate this concern, which arose in the Watergate period in the United States and the Communist era in Eastern Europe when secret databases were used to curtail individual freedoms.¹⁶⁵ Yet the

¹⁵⁸ Upbin, *supra* note 94.

¹⁵⁹ Michael Liedtke, *Study: App Economy is a Booming Jobs Engine*, USA TODAY, Feb. 7, 2012, <http://www.usatoday.com/tech/news/story/2012-02-07/apps-economy-creates-jobs/52997386/1>.

¹⁶⁰ See, e.g., Steven Overly, *Mobile Health Apps Prompt Questions About Privacy*, WASH. POST, Apr. 29, 2012, http://www.washingtonpost.com/business/capitalbusiness/mobile-health-apps-prompt-questions-about-privacy/2012/04/27/gIQAk17FqT_story.html.

¹⁶¹ Michael Mandel, *Where the Jobs Are: The App Economy*, TECHNET, Feb. 7, 2012, <http://www.technet.org/wp-content/uploads/2012/02/TechNet-App-Economy-Jobs-Study.pdf>.

¹⁶² Anthony Haw, *Zynga Makes Up 12 Percent of Facebook’s Revenue*, TECHCRUNCH, Feb. 1, 2012, <http://techcrunch.com/2012/02/01/zynga-makes-up-12-percent-of-facebooks-revenue>.

¹⁶³ Singer, *supra* note 28.

¹⁶⁴ See, e.g., Joshua Brustein, *Start-Ups Seek to Help Users Put a Price on Their Personal Data*, N.Y. TIMES, Feb. 12, 2012, <http://www.nytimes.com/2012/02/13/technology/start-ups-aim-to-help-users-put-a-price-on-their-personal-data.html>.

¹⁶⁵ While the right to privacy has deep historical roots, information privacy law is generally viewed as a development of the second half of the twentieth century. See Neil Richards, *The Information Privacy Law Project*, 94 GEO. L. J 1087 (2006); Spiros Simitis, *Reviewing Privacy in an Information Society*, 135 U. PA. L. REV. 707 (1987); James Whitman, *The Two Western Cultures of Privacy: Dignity Versus Liberty*, 113

frameworks that emerged in response to such concerns, providing access rights in the United States and requiring database registration in the European Union, failed to engage individuals who remained largely oblivious to their rights.¹⁶⁶

¶77 Big data has reinvigorated the specter of massive data silos accumulating and using information for obscure purposes. Individuals and regulators do not condemn big data as such; rather, they oppose “secret big data,” which raises a Kafkaesque vision of an inhumane bureaucracy.¹⁶⁷ Avoiding potential abuses may require retrofitting transparency obligations and providing more practicable access rights. Any activity performed in the dark raises suspicion of being untoward; what is done in broad daylight must be wholesome and “clean.”

¶78 The call for transparency is not new, of course. Rather the emphasis is on access to data in a *usable format*, which can work to create value to individuals. Transparency and access alone have not emerged as potent tools because individuals do not care for, and cannot afford to indulge in, transparency and access for their own sake without any tangible benefit. For this reason, consumers seldom opt in or opt out of end user license agreements (EULA) or privacy policies, regardless of their merits.¹⁶⁸ The enabler of transparency and access is the ability to *use* the information and *benefit* from it in a tangible way. Such use and benefit may be achieved through “featurization” or “app-ification” of privacy. Useful access to PII will engage individuals, invite scrutiny of organizations’ information practices, and thus expose potential misuses of data. It would be value-minimizing to leave this opportunity untapped. Organizations should build as many dials and levers as needed for individuals to engage with their data.

¶79 The extent of transparency and access espoused in this article will no doubt raise serious legal and business complexities. First, organizations (particularly non-consumer facing ones) may argue that in many circumstances providing individual access to massive databases distributed across numerous servers and containing zettabytes of de-identified data is simply not practical. Second, to avoid the creation of a bigger privacy problem than it seeks to solve, direct online accessibility to data requires strong authentication as well as secure channeling, imposing costs and inconveniences on both organizations and individuals. Third, as the ecosystem for personal information expands, building layers upon layers of user-side applications over the existing centralized structure, data security risks of leakage and unauthorized use increase correspondingly. Finally, access to machine-readable data in a usable format appears to promote data portability, a contentious concept which raises further questions regarding intellectual property and antitrust. While further work is required to address these concerns, these issues can be contained.

¶80 First, if data were in fact robustly de-identified, it would be counterproductive to require their re-identification simply in order to provide individuals with access.¹⁶⁹ Yet in

YALE L.J. 1151 (2004); *see also* MICHEL FOUCAULT, DISCIPLINE AND PUNISH (Alan Sheridan trans., Vintage Books 2d ed. 1995) (1977).

¹⁶⁶ Omer Tene, *There is No New Thing Under the Sun*, CONCURRING OPINIONS BLOG, July 30, 2012, <http://www.concurringopinions.com/archives/2012/07/there-is-no-new-thing-under-the-sun.html>.

¹⁶⁷ DANIEL J. SOLOVE, THE DIGITAL PERSON: TECHNOLOGY AND PRIVACY IN THE INFORMATION AGE (2004), chapter 3.

¹⁶⁸ *See* Dan Ariely, *3 Main Lessons of Psychology*, DAN ARIELY BLOG (May 5, 2008), <http://danariely.com/2008/05/05/3-main-lessons-of-psychology>.

¹⁶⁹ Schwartz & Solove, *supra* note 113.

precisely such circumstances the risk to individuals' privacy would be greatly reduced. Access is most needed where de-identification is weak and the data could therefore provide tangible benefits to individuals. Here, too, the flexibility and modularity of the FIPPs' framework proves instrumental: as the degree of data identification increases, so should the level of access rights provided to individuals.

¶81 Second, privacy and data security clearly require that an individual be granted access only to his or her personal data. This means that organizations must authenticate the identity of an individual making a request and that data must be delivered on a secure channel. Implementation may require use of digital signatures and similar identity infrastructures already in existence today, as well as encrypted communication delivery channels.

¶82 Third, the enhancement of big data with interfaces for user interaction increases the number of access points and correspondingly elevates the risk of security breach and data leakage.¹⁷⁰ Yet, that risk is a price worth paying where the goal is data empowerment of individuals. This article disputes the contention that individuals should not be allowed to access their information simply to avoid a potential data leak. To argue otherwise is tantamount to suggesting that a bank should bar customers' access to their accounts to avoid losing their money.

¶83 Finally, although similar to the data portability argument, this article stops short of advocating portability.¹⁷¹ It recognizes that portability is not, strictly speaking, a concept of privacy law but rather one derived from antitrust. It regards personal information as an asset of individuals, which remains under their control unless traded for a fair price. Although the proposed European Data Protection Regulation seeks to weave portability into the fabric of privacy law,¹⁷² this article contends such an approach may go too far. The property metaphor fails to capture the psychological and sociological nuance of the right to privacy. As Julie Cohen wrote a decade ago, "[r]ecognizing property rights in personally-identified data risks enabling more, not less, trade and producing less, not more, privacy."¹⁷³ Moreover, a right to portability could eviscerate the competitive advantage gained by companies that have invested significant skill and resources to collect, organize, and share data in commercially valuable ways, thereby stifling innovation. Companies vying for control of information markets could use it strategically to corner their competitors. Personal information should be regarded as neither an exclusive asset of individuals—treatment which may impinge on business trade secrets and intellectual property rights—nor exclusively the property of businesses, excluding individuals from benefiting. Rather, personal information should be treated as a valuable joint resource and a basis for value creation and innovation.

¶84 Privacy suffers not only when individuals are *unaware* of data practices, but also when they are *uninterested* or *disengaged*. Such an environment, regardless of the regulatory mechanisms in place, provides insufficient checks on data collection and use. Where individuals can access data in a manner that is engaging, useful, or valuable, they

¹⁷⁰ See Tene, *supra* note 70 (in the context of the Facebook ecosystem).

¹⁷¹ A new right to data portability has been introduced by the European Data Protection Regulation, art. 18.

¹⁷² European Data Protection Regulation, art. 18.

¹⁷³ Cohen, *supra* note 140, at 1,391.

will give rise to natural checks on inappropriate behavior, thus serving as a useful compliance mechanism for privacy law.

B. Enhanced Transparency: Shining the Light

¶85 Policymakers have long struggled to draw the line for ethical data use.¹⁷⁴ The discussion has historically revolved around the definition of “sensitive data.” Yet, any attempt to exhaustively define categories of sensitivity typically failed, given the highly contextual nature of personal information. For example, the first data protection case taken by the European Court of Justice, the matter of *Bodil Lindqvist*,¹⁷⁵ dealt with the use of “sensitive” information so benign so as to appear trivial—the fact that the defendant’s fellow churchgoer had a broken leg. A broken leg is clearly a medical condition, which is a category of sensitive data under any legal framework;¹⁷⁶ yet, information about an individual’s broken leg is not generally considered to be sensitive in nature.

¶86 In order to delimit the zone of ethical data analysis we propose that organizations reveal not only the existence of their databases but also the *criteria* used in their decision-making processes, subject to protection of trade secrets and other intellectual property laws.¹⁷⁷ Today, such disclosures are made only when a user is presented with a consumer privacy policy, and even then the logic behind some of the automated processes remains opaque. Louis Brandeis, who together with Samuel Warren introduced the right to privacy into legal discourse in 1890,¹⁷⁸ has also written that “[s]unlight is said to be the best of disinfectants”¹⁷⁹ We trust if the existence and uses of databases were visible to the public, organizations would be more likely to avoid unethical or socially unacceptable uses of data. If organizations were required to disclose their line of reasoning in data processing operations impacting individuals’ lives, they might avoid unethical uses of data pertaining to certain populations, such as children, and certain data such as legally suspect categories—including gender, age, and race—or sensitive data (in the parochial sense), such as sexual preferences or certain medical conditions.

¶87 More broadly, the requirement that organizations reveal their decisional criteria is based on the FIPPs’ transparency and accuracy principles. In a big data world, what calls for scrutiny is often not the accuracy of the *raw data* but rather the accuracy of the *inferences* drawn from the data. Inaccurate, manipulative, or discriminatory conclusions

¹⁷⁴ See boyd & Crawford, *supra* note 21, at 672 (stating, “[V]ery little is understood about the ethical implications underpinning the Big Data phenomenon”).

¹⁷⁵ Case C-101/01, *Bodil Lindqvist*, 2003 ECR I-12971.

¹⁷⁶ See, e.g., European Data Protection Directive, art. 8(1); FTC Final Report, § IV.C.2.e.ii.

¹⁷⁷ See Article 12 of the European Data Protection Directive, which requires organizations to provide an individual with “knowledge of the logic involved in any automatic processing of data concerning him at least in the case of the automated decisions.” Recital 41 of the European Data Protection Directive acknowledges the need to protect organizational assets: “[E]very data subject must also have the right to know the logic involved in the automatic processing of data concerning him, at least in the case of the automated decisions . . . this right must not adversely affect trade secrets or intellectual property and in particular the copyright protecting the software . . . these considerations must not, however, result in the data subject being refused all information.”

¹⁷⁸ Samuel Warren & Louis Brandeis, *The Right to Privacy*, 4 HARV. L. REV. 193 (1890).

¹⁷⁹ Louis D. Brandeis, *What Publicity Can Do*, HARPER’S WEEKLY, Dec. 20, 1913, http://c0403731.cdn.cloudfiles.rackspacecloud.com/collection/papers/1910/1913_12_20_What_Publicity_Ca.pdf.

may be drawn from perfectly innocuous, accurate data. The observer in big data analysis can affect the results of her research by defining the data set, proposing a hypothesis, or writing an algorithm. At the end of the day, big data analysis is an interpretative process, in which one's identity and perspective informs one's results. Like any interpretative process, it is subject to error, inaccuracy, and bias.¹⁸⁰

¶188 The requirement that organizations disclose their decisional criteria (not necessarily the *algorithms*, but rather the *factors* they consider) highlights an important fault line between law and technology. Fairness and due process mandate that individuals are informed of the basis for decisions affecting their lives, particularly those made by machines operating under opaque criteria. In the landmark *Daubert* case, the Supreme Court charged trial judges with the responsibility of acting as gatekeepers to exclude unreliable scientific expert testimony.¹⁸¹ Following *Daubert*, Justice Scalia remarked in *Melendez-Diaz* that “[f]orensic evidence is not uniquely immune from the risk of manipulation.”¹⁸² This was in response to the government's assertion that “there is a difference, for Confrontation Clause purposes, between testimony recounting historical events, which is ‘prone to distortion or manipulation,’ and the testimony at issue here, which is the ‘resul[t] of neutral, scientific testing.’”¹⁸³ We argue that not only the accused, but also any other citizen be afforded a right to confront decisions made about her. *Daubert* and its progeny mandate that, at the end of the day, it is lawyers and judges, not technology, who try individuals.¹⁸⁴

¶189 The rule proposed in this article focuses regulatory attention on the decision-makers who draw conclusions from personal information rather than other parties in the ecosystem. In doing so, it recognizes that some of the risks of big data affect fairness, equality and other values, which may be no less important than—but are theoretically distinct from—core privacy interests. Over the past few years, the debate over privacy has become conflated with broader social values. For example, the increasing tendency of employers to use social networking services to run background checks on prospective job candidates has led critics to condemn the “privacy invasive” nature of such platforms.¹⁸⁵ Yet on closer scrutiny, it is not clear that social networking services should be held accountable for illegal or unethical discrimination by employers. If an employer chooses to screen out job candidates based on race, good looks,¹⁸⁶ or proclivity to drink,¹⁸⁷ then that employer—not the neutral platform used to convey such information—should stand

¹⁸⁰ boyd & Crawford, *supra* note 21, at 668.

¹⁸¹ *Daubert v. Merrell Dow Pharmaceuticals, Inc.*, 509 U.S. 579 (1993); *see also* FED. R. EVID. 702; Paul Giannelli, *Daubert and Forensic Science: The Pitfalls of Law Enforcement Control of Scientific Research*, 2011 U. ILL. L. REV. 53.

¹⁸² *Melendez-Diaz v. Massachusetts*, 557 U.S. 305, 318 (2009).

¹⁸³ *Id.* at 317.

¹⁸⁴ *Cf.* Randall Stross, *The Algorithm Didn't Like My Essay*, N.Y. TIMES, June 9, 2012, <http://www.nytimes.com/2012/06/10/business/essay-grading-software-as-teachers-aide-digital-domain.html>.

¹⁸⁵ *See, e.g.*, Andrew Couts, *Senator Promises Bill to Block Invasive Employer Facebook Checks*, DIGITAL TRENDS (Mar. 23, 2012), <http://www.digitaltrends.com/social-media/senator-promises-bill-to-block-invasive-employer-facebook-checks/#ixzz2IahftPRE>.

¹⁸⁶ *See, e.g.*, *Attractiveness Discrimination: Hiring Hotties*, THE ECONOMIST (July 21, 2012), <http://www.economist.com/node/21559357>.

¹⁸⁷ *See, e.g.*, Jeffrey Rosen, *The Web Means the End of Forgetting*, N.Y. TIMES, July 21, 2010, <http://www.nytimes.com/2010/07/25/magazine/25privacy-t2.html>.

to blame. Accordingly, it is prospective employers—or, in other contexts, insurers, banks and government agencies¹⁸⁸—that need to explain their decisional criteria in reaching personal data driven conclusions.

¶90

Finally, attention must be given to the accessibility of big data sets to the research community at large.¹⁸⁹ Traditionally, when scientists published their research, they also made the underlying data available so that other scientists could verify the results. Yet with big data, it is often only the employees of certain organizations that benefit from access, conducting analysis and publishing results without making the underlying data publicly available.¹⁹⁰ Such scientists may argue, first, that the data are a proprietary asset of their business. Indeed, they may claim that disclosing the data could infringe customers' privacy.¹⁹¹ As boyd and Crawford note, future research must address relevant, fundamental questions, such as who has the right to access big data sets, for what purposes, in what contexts, and with what constraints.¹⁹² Without good answers, we may witness a stratification of the scientific world to the haves and have-nots of big data.¹⁹³

V. CONCLUSION

¶91

Privacy advocates and data regulators increasingly decry the era of big data as they observe the growing ubiquity of data collection and increasingly robust uses of data enabled by powerful processors and unlimited storage holders. Researchers, businesses, and entrepreneurs equally vehemently point to concrete or anticipated innovations that may be dependent on the default collection of large data sets.

¶92

This article has called for the development of a legal model where the benefits of data for organizations and researchers are shared with individuals. If organizations provide individuals with access to their data in usable formats, creative powers will be unleashed to provide users with applications and features building on their data for new innovative uses. In addition, transparency with respect to the logic underlying organizations' data processing will deter unethical, sensitive data use and allay concerns about inaccurate inferences. Traditional transparency and individual access mechanisms have proven to be an ineffective means for motivating individuals to engage their data.

¹⁸⁸ See, e.g., Omer Tene, *What Happens Online Stays Online: Comments on "Do Not Track"*, STANFORD CIS BLOG (Mar. 26, 2011), <http://cyberlaw.stanford.edu/blog/2011/03/what-happens-online-stays-online-comments-do-not-track>.

¹⁸⁹ See John Markoff, *Troves of Personal Data, Forbidden to Researchers*, N.Y. TIMES, May 21, 2012, <http://www.nytimes.com/2012/05/22/science/big-data-troves-stay-forbidden-to-social-scientists.html>.

¹⁹⁰ See Lev Manovich, *Trending: The Promises and the Challenges of Big Social Data*, in DEBATES IN THE DIGITAL HUMANITIES (Matthew Gold ed., 2012), available at http://www.manovich.net/DOCS/Manovich_trending_paper.pdf (claiming that “only social media companies have access to really large social data – especially transactional data. An anthropologist working for Facebook or a sociologist working for Google will have access to data that the rest of the scholarly community will not.”).

¹⁹¹ See Bernardo Huberman, *Big Data Deserve a Bigger Audience*, 482 NATURE 308 (Feb. 16, 2012) (warning that privately held data was threatening the very basis of scientific research, and complaining that “[m]any of the emerging 'big data' come from private sources that are inaccessible to other researchers. The data source may be hidden, compounding problems of verification, as well as concerns about the generality of the results.”).

¹⁹² boyd & Crawford, *supra* note 21, at 673.

¹⁹³ *Id.* at 674 (calling this “the Big Data rich and the Big Data poor”).

The promise of new benefits and value sharing propositions will incentivize individuals to act without compromising organizations' ability to harness big data.

