

DIFFERENTIAL PRIVACY: RAISING THE BAR

Anna Myers* & Grant Nelson•

CITE AS: 1 GEO. L. TECH. REV. 135 (2016)

<https://perma.cc/8QDH-43V8>

INTRODUCTION.....	135
DIFFERENTIAL PRIVACY: NOT DE-IDENTIFICATION.....	136
DIFFERENTIAL PRIVACY ENCODES PRIVACY LAW & POLICY IN ITS SYSTEMS	138
Differential Privacy in Research.....	140
Differential Privacy in Commerce.....	140

INTRODUCTION

Uncle Ben’s sage advice in *Spiderman* that “with great power comes great responsibility,” no doubt applies to today’s great power: big data. Like Peter Parker, privacy advocates and technologists are racing to harness the power of big data’s web of connections, but are sorely lagging in handling the power responsibly. Existing privacy protecting strategies, including de-identification, anonymization, pseudonymization, and encryption, have encountered bumps in the road. Data thought to be sufficiently de-identified has been re-identified;¹ anonymization and pseudonymization are considered privacy failures;² and encrypted email services have shut down in response to

* Attorney in Washington, D.C. where her practice focuses on privacy & technology legal matters. Ms. Myers’ privacy-related experience includes the International Association of Privacy Professionals, Harvard University’s Berkman-Klein Center for Internet and Society, the Network Advertising Initiative, and the U.S. Department of the Treasury’s Office of Privacy, Transparency, and Records. She holds a J.D. from The George Washington University Law School and her B.A. in Rhetoric and Media Studies from Willamette University. © 2016, Anna Myers & Grant Nelson.

• Passed the D.C. bar and works for the Network Advertising Initiative in Washington, D.C. His experience includes launching several successful webapps, working for the Privacy Tools project at Harvard’s Berkman-Klein Center for Internet and Society, and building predictive models. Mr. Nelson holds a J.D. from The George Washington University Law School. He thanks his parents for their loving support. © 2016, Anna Myers & Grant Nelson.

¹ See generally Arvind Narayanan & Edward Felton, *No Silver Bullet: De-Identification Still Doesn’t Work*, <http://randomwalker.info/publications/no-silver-bullet-de-identification.pdf> (2014); Latanya Sweeney, *Foundations of Privacy Protection from a Computer Science Perspective*, DATA PRIVACY LAB (2000), <http://dataprivacylab.org/projects/disclosurecontrol/index.html>.

² Paul Ohm, *Broken Promises of Privacy: Responding to the Surprising Failure of Anonymization*, 57 UCLA L. REV. 1701 (2010); see also Arvind Narayanan & Vitaly

government subpoenas to protect their users' information.³ The landscape is not, however, without hope: with every failure or data breach, technologists and advocates are evolving and building better privacy protections.

One such new protection is differential privacy. Differential privacy has been used in academic and research settings for nearly a decade but is just starting to break into the commercial space. Differential privacy describes a system that provides a protective layer between data and a user of the data in which the protective layer mathematically distorts the data with minor falsities so that it masks sensitive aspects of the data while retaining the statistically significant characteristics. Differential privacy is raising the bar for effective data responsibility by redefining the balance and reducing the trade-off between privacy and data utility.

DIFFERENTIAL PRIVACY: NOT DE-IDENTIFICATION⁴

Differential privacy is de-identification's cynical sibling. Differential privacy gained momentum in the wake of several high-profile failures of de-identification strategies, and its strengths reflect the frustration with the failure of de-identification. Whereas de-identified datasets are subject to re-identification attacks using other available datasets, differential privacy's threat model often assumes that bad actors or researchers "accessing any differentially private dataset are omniscient, omnipotent and constantly co-conspiring data snoops."⁵ Differential privacy reduces the ambiguity of determining when data is sufficiently de-identified, and goes a level further than de-identification

Shmatikov, *Robust De-anonymization of Large Sparse Datasets*, SP '08 PROC. OF THE 2008 IEEE SYMP. ON SEC AND PRIVACY 111 (2008).

³ James Ball, Julian Borger, & Glenn Greenwald, *Revealed: how US and UK spy agencies defeat Internet privacy and security*, THE GUARDIAN (Sept. 6, 2013), <http://www.theguardian.com/world/2013/sep/05/nsa-gchq-encryption-codes-security>.

⁴ Differential privacy has primarily been identified as unique from de-identification. See Narayanan & Felton, *supra* note 1. The National Institute of Standards and Technology ("NIST") categorizes differential privacy as "privacy preserving data mining" and de-identification as "privacy preserving data publishing." U.S. DEP'T OF COMMERCE, NAT'L INST. OF STANDARDS AND TECH., NISTIR 8053, DE-IDENTIFICATION OF PERSONAL INFORMATION (2015).

⁵ Daniel C. Barth-Jones, *Ethical Concerns, Conduct and Public Policy for Re-Identification and Deidentification Practice: Part 3 (Re-Identification Symposium)*, HARV. L. BILL OF HEALTH (Oct. 2, 2013), <http://blogs.harvard.edu/billofhealth/2013/10/02/ethical-concerns-conduct-and-public-policy-for-re-identification-and-de-identification-practice-part-3-re-identification-symposium/>.

because it “seeks to mathematically prove that a certain form of data analysis can’t reveal anything about an individual”⁶

Differential privacy does not prescribe the use of a specific algorithm or encryption technique. Unlike de-identification, which typically relies on omission or mutation of data, differential privacy can be conceptualized as a gatekeeping mechanism that serves as a privacy-protecting layer between raw data and a user of the data. The differential privacy layer can be applied to data at the point of collection or at the point of querying the data.⁷ Applying the differential privacy layer at the point of collection provides additional protection while the data is stored and in transit. Applying the noise at the point of query allows the flexibility to later repurpose the data.

Protected datasets require all potential users to submit queries through the differential privacy-providing mechanism to access the dataset. When a user queries the data, the system evaluates that request against all previous queries and determines the sensitivity of the query. The system then applies noise or small-falsities to the data to protect the individual data subjects and returns an answer to the user. The noise injecting algorithm can be mathematically tuned to guarantee minimum levels of protection against reverse-engineering the underlying data. The key input to the algorithm is the *privacy budget*.

Every differential privacy system operates on a privacy budget—how much time, resources, and potentially traded utility the data controller is willing to trade in exchange for added privacy protection. The privacy budget of a differential privacy mechanism is a measurement of how much noise the algorithm injects to differentiate the data passed along from the true raw data. Determining the privacy budget is a social decision more than a mathematical one: the dataset’s owner can increase the privacy budget (injecting more noise) on a dataset that contains sensitive information and decrease the privacy budget (resulting in more accurate responses) for a dataset that contains more innocuous data. If a query requires the system to exceed the privacy budget the system will not provide the answer to the user. A differential privacy layer can be tuned to prevent leakage of data even in a situation where every query of the data is from bad actors with an infinite timeline or query budget, collaborating with each other. If a privacy budget is depleted or exceeded that dataset may no

⁶ Andy Greenburg, *Apple’s ‘Differential Privacy’ is about Collecting your Data – But not Your Data*, WIRED (June 13, 2016), <https://www.wired.com/2016/06/apples-differential-privacy-collecting-data/>.

⁷ This allows for entities to be strategic about their data vulnerabilities and use differential privacy to adapt to their different environments and privacy needs. See e.g. Anthony Tockar, *Differential Privacy: The Basics*, NEUSTAR (Sept. 8, 2014), <https://research.neustar.biz/2014/09/08/differential-privacy-the-basics/>.

longer be usable. In a production database, however, the chances of a budget being depleted are slim given the rate at which new data can be added to datasets.

Despite the strong protections offered by differential privacy, it requires users to put their faith in the dataset owner's algorithms, typically without strong means to validate the integrity of the algorithm's noise injection. This is especially true when the data collector aggregates unencrypted data in a database and applies the differential privacy layer at the point of database query, rather than applying the differential privacy filter at the point of collection or contribution to the database. The need for a consumer to entrust a company with at least some data is all but unavoidable, and a shift towards using differential privacy provides more manageable and robust protection than its alternatives.

DIFFERENTIAL PRIVACY ENCODES PRIVACY LAW & POLICY IN ITS SYSTEMS

One of the main challenges of the privacy industry has been transforming complex concepts into technological tools. Privacy concepts are more challenging to implement technologically because they are not as straightforward as security concepts, such as user authentication. Security protections are objective and mechanical in nature with a united goal of keeping the data in and the bad actors out. Basic privacy concepts used by both the private and public sectors, such as the Fair Information Practice Principles ("FIPPs"),⁸ are more subjective and therefore more challenging to translate into code or technological tools.

The FIPPs framework originated from a 1973 report issued by the precursor to the U.S. Department of Health and Human Services,⁹ and was later codified in the Organisation for Economic Co-operation and Development ("OECD") privacy guidelines.¹⁰ The FIPPs are the core of the Privacy Act of

⁸ Robert Gellman, *Fair Information Practices: A Basic History*, (June 17, 2016), <http://bobgellman.com/rg-docs/rg-FIPShistory.pdf>.

⁹ U.S. DEP'T OF HEALTH, EDUC. & WELFARE, NO. (OS)73-94, REPORT OF THE SEC'Y'S ADVISORY COMM. ON AUTOMATED PERSONAL DATA SYSTEMS (1973).

¹⁰ *OECD Guidelines on the Protection of Privacy and Transborder Flows of Personal Data*, ORG. FOR ECON. CO-OPERATION & DEV. (1980), <https://www.oecd.org/internet/ieconomy/oecdguidelinesontheProtectionofPrivacyandTransborderflowsofpersonaldata.htm>

1974,¹¹ and form the basis of other policy frameworks, such as the Department of Homeland Security privacy guidelines.¹² The principles are as follows:¹³

1. **Transparency:** information collectors should be transparent in their collection, use, dissemination, and maintenance practices.
2. **Individual Participation:** consent of the individual for the collection of the data should be obtained.
3. **Purpose Specification:** the specific purpose(s) the information is being collected for should be articulated.
4. **Data Minimization:** only the information necessary to accomplish the specified purpose should be collected.
5. **Use Limitation:** the information should only be used for the specific purpose(s) for which it is being collected.
6. **Data Quality and Integrity:** To the extent practicable collected information should be accurate, relevant, timely, and complete.
7. **Security:** Collected information should be protected from loss, unauthorized access or use, destruction, modification, or unintended or inappropriate disclosure.
8. **Accountability and Auditing:** Collecting organizations should be accountable for compliance with the FIPPs and the use of information should be audited to demonstrate compliance with the FIPPs and all applicable data protection requirements.

Privacy Enhancing Technologies (“PETs”) integrate concepts like the FIPPs, other privacy best practices, and applicable legal regimes in their design. For example, in the United States, the faster a video is uploaded, the better; however, in areas where governments suppress information, slower upload speeds may be desired so that a video upload does not appear different from other internet traffic. A PET for that scenario could be designed to protect the content of the video by masking it as other internet traffic, and thereby avoid raising any red flags. An implementation of differential privacy is a privacy enhancing technology (PET) because developers utilize the FIPPs and take into consideration the types of data in a database and applicable laws & policies in designing a system.

¹¹ 5 U.S.C. § 552a (2014).

¹² HUGO TEUFEL III, U.S. DEP’T OF HOMELAND SEC., MEMO. NO. 2008-01, PRIVACY POLICY GUIDANCE MEMORANDUM at 3-4. (Dec. 29, 2008), http://www.dhs.gov/xlibrary/assets/privacy/privacy_policyguide_2008-01.pdf.

¹³ Descriptions of FIPPs adapted from *id.*

Differential Privacy in Research

Differential privacy was formalized by and is most strongly associated with Cynthia Dwork's work while at Microsoft Research. In 2006, Dr. Dwork published "Differential Privacy," a 12-page paper presented at the 33rd International Colloquium on Automata, Languages and Programming, part II.¹⁴ Since then cryptologists, mathematicians, and computer scientists have pursued academic research on differential privacy resulting in a multi-disciplinary research effort.

Harvard's Berkman-Klein Center and MIT have pushed the multidisciplinary approach by bringing together computer scientists and attorneys from the Berkman-Klein Center, social scientists from the Institute for Quantitative Social Science, and mathematicians and cryptologists from MIT in the PrivacyTools Project.¹⁵ Their research is a multi-faceted approach to protecting privacy while preserving the value of data, with the goal of including promising techniques in the open-source database software, Dataverse. Because of the imperative to maintain data's value while also maximizing user privacy, differential privacy has proven to be a large focus of their attention. Aaron Roth, an associate professor of computer and information science at the University of Pennsylvania, co-authored the essential textbook *The Algorithmic Foundations of Differential Privacy* with Dr. Dwork.¹⁶ Roth's expertise in the mathematical foundations of differential privacy was affirmed when Apple sought his review of its algorithms prior to announcing publicly that it will deeply integrate differential privacy into its devices.¹⁷

Differential Privacy in Commerce

Several companies have started to implement differential privacy into their data acquisition and storage systems. Most notably, Apple recently announced that it will integrate differential privacy mechanisms into its iOS

¹⁴ Cynthia Dwork, *Differential Privacy*, MICROSOFT (Feb. 2016), <https://www.microsoft.com/en-us/research/wp-content/uploads/2016/02/dwork.pdf>

¹⁵ *Harvard University Privacy Tools Project*, HARV. UNIV., <http://privacytools.seas.harvard.edu/> (last visited Nov. 21, 2016).

¹⁶ Cynthia Dwork & Aaron Roth, *The Algorithmic Foundations of Differential Privacy*, in FOUNDATIONS AND TRENDS IN THEORETICAL COMPUTER SCIENCE: VOL. 9: NO.3-4 211 (2014).

¹⁷ Kate Conger, *What Apple's differential privacy means you're your data and the future of machine learning*, TECHCRUNCH (June 14, 2016), <https://techcrunch.com/2016/06/14/differential-privacy/>.

devices for some use cases.¹⁸ Apple's implementation aligns with its branding as a privacy-protecting organization: as it will perform the privacy-protecting noise injection at the device-level collection point rather than at the database level, consumer data will remain more secure during transmission and storage. Therefore, protected data leaving any particular iOS device is of minimal use to malicious actors that intercept the transmission, and any database of protected information is of minimal value if breached. Apple is not alone in placing the noise-injection calculations on devices: Google has implemented a differential privacy mechanism, at the device-level for its Chrome browser usage data.¹⁹ Google's Randomized Aggregatable Privacy-Preserving Ordinal Response ("RAPPOR") preserves the predictive power of data in relatively large datasets.²⁰ Some experts believe Google's RAPPOR project is the first commercial deployment of differential privacy.²¹

Additionally, Facebook, no stranger to privacy and big data policy discussions, appears to have implemented a differential privacy mechanism in its advertisement audience estimator tool as early as 2012.²² The tool allows a potential advertiser to estimate how many Facebook users would view an ad based on the ad's target segments, such as location, age, and interests. As shown by Andrew Chin and Anne Klinefelter, Facebook not only rounds estimates to the nearest 20 (and zero if below 40), it appears to apply the rounding to an already-noisy estimate in a pattern that strongly suggests a differential privacy mechanism is at play.²³ Differing from Google and Apple, Facebook does not seem to implement the noise-injection calculation prior to the user sending data to Facebook for retention, but rather keeps all user data in pristine condition and adds noise at the moment of database query.

¹⁸ Andy Greenberg, *supra* note 6.

¹⁹ Ulfar Erlingsson, Vasyl Pihur & Aleksandra Korolova, *RAPPOR: Randomized Aggregatable Privacy-Preserving Ordinal Response*, CCS '14 PROC. OF THE 2014 ACM SIGSAC CONF. ON COMPUT. AND COMMS. SEC. (2011),

<http://static.googleusercontent.com/media/research.google.com/en//pubs/archive/42852.pdf>.

²⁰ *Id.* at 1 ("Randomized Aggregatable Privacy-Preserving Ordinal Response, or RAPPOR, is a technology for crowdsourcing statistics from end-user client software, anonymously, with strong privacy guarantees. In short, RAPPORs allow the forest of client data to be studied, without permitting the possibility of looking at individual trees.").

²¹ Answer by Aaron Roth, QUORA (June 18, 2016), <https://www.quora.com/Does-Google-use-a-differential-privacy-strategy>.

²² Andrew Chin & Anne Klinefelter, *Differential Privacy As A Response To The Reidentification Threat: The Facebook Advertiser Case Study*, 90 N.C. L. REV. 1418, 1456 (May 2012).

²³ *Id.* at 1443.

A popular criticism of differential privacy states that enormous data sets would be required for a system to preserve the utility of a differentiated dataset. By injecting noise using a Laplace distribution,²⁴ as modeled by Dwork, Roth, and others, smaller companies have reported impressive accuracy. For example, Snips, an artificial-intelligence company with an emphasis on privacy showed that a model trained on only 1000 observations filtered through a differential privacy mechanism relying on a Laplace distribution had the same predictive accuracy as a model trained on 1 million observations relying on the RAPPOR distribution.²⁵ In fact, their research showed that the predictive accuracy of a model using data sourced from a differential-privacy system plateaued at as few as 10,000 observations.

Now that the use and development of privacy tools such as differential privacy is growing, the integration of those tools with other technologies provides comprehensive solutions to maximize the potential for privacy by design and user protection. The growing availability of differential privacy mechanisms in academic literature and open source libraries, combined with the fact that even small datasets can be protected using differential privacy and remain valuable makes it likely that more commercial implementations of differential privacy are on the horizon, something that should be encouraged by the legal and regulatory environment.

²⁴ A common probability distribution used in probability theory and statistics, also sometimes known as a double exponential distribution.

²⁵ Morten Dahl & Joseph Dureau, *Differential Privacy for the Rest of Us*, MEDIUM (July 29, 2016), <https://medium.com/snips-ai/differential-privacy-for-the-rest-of-us-665e053cec17>.