

# Formal Privacy Models and Title 13

Research Project Proposal Submitted to the United States Bureau of the Census  
Funding Opportunity: Administrative Data and Data Linkages Research Program  
(CENSUS-ADR-CARRA-2016-2004966)

**Lead Organization Principal Investigator:** Kobbi Nissim  
Department of Computer Science  
Georgetown University  
kobbini@hotmail.com

**Sub-awardee Principal Investigators:** Urs Gasser  
Harvard Law School  
Berkman Klein Center for Internet & Society  
Harvard University  
ugasser@cyber.law.harvard.edu

Adam Smith  
Computer Science and Engineering Department  
Pennsylvania State University  
asmith@psu.edu

Salil Vadhan  
John A. Paulson School of Engineering and Applied Science  
Harvard University  
salil@seas.harvard.edu

**Amount Requested:** TBD

**Requested Project Period:** April 1, 2017 - March 31, 2020

---

Authorized Representative

---

Date

## Formal Privacy Models and Title 13: Project Synopsis

**Overview.** The US Bureau of the Census (BOC) collects large quantities of data that can be useful for research and decision making by policymakers, businesses, and academics. The BOC is responsible for analyzing and publishing useful statistical data. As much of the collected data pertain to individuals, households, and establishments, the BOC also has a legal obligation to protect their privacy. The BOC has long employed statistical disclosure limitation (SDL) techniques in order to guarantee the confidentiality of the data it maintains and releases. However, with the SDL techniques traditionally used, it is difficult and labor-intensive to determine with confidence that these two requirements—privacy and utility—have been satisfied.

The proposed research project seeks to further the use of formal approaches such as differential privacy that have the potential to provide rigorous guarantees that legal requirements for privacy and utility are met. Applying such approaches requires (a) bridging legal privacy requirements with mathematical privacy requirements, and (b) designing analysis methods that provide statistical utility while satisfying the privacy requirements from (a). Our team has developed tools satisfying these two goals in the context of Harvard’s *Privacy Tools* project, and with this proposal we seek to collaborate with BOC staff to develop similar solutions that are tailored to the bureau’s specific requirements.

We address two major challenges confronting wider adoption of formal privacy models: (a) There is a wide conceptual and practical gap between the approaches found in formal privacy models and the heuristic approaches in current use and contemplated by existing regulatory and policy frameworks. (b) There is a gap between theoretical developments showing that formal privacy models like differential privacy permit, in principle, a wide collection of analyses and the actual use of analysis and publication techniques by the BOC.

**Furtherance of BOC Priorities.** The proposed project would result in methods for publishing data in ways that satisfy both formal mathematical privacy requirements and legal standards for privacy protection, thereby furthering “improvements to existing methods that protect privacy, avoiding the release of any information that would identify an individual or business in public statistics.”

**Principal Investigators.** Multidisciplinary in nature, this research project brings together experts in privacy, cryptography, law, and policy. Kobbi Nissim is a Prof. of Computer Science at Georgetown University and a Sr. Research Fellow at the Center for Research on Computation and Society at Harvard. His research interests lie in privacy, cryptography, and their connections to machine learning, statistics, game theory, and policy. Adam Smith is a Prof. of Computer Science and Engineering at Pennsylvania State University. His research interests lie in cryptography, privacy and their connections to information theory, quantum computing, and statistics. Nissim’s work from 2003 initiated rigorous foundational research of privacy and in 2006 he and Smith presented with Dwork and McSherry *differential privacy*, a strong definition of privacy in computation. Salil Vadhan is a Prof. of Computer Science at Harvard. His research interests lie in complexity theory, privacy, and cryptography. Urs Gasser is a Prof. of Practice at Harvard Law School and the Executive Director of the Berkman Klein Center for Internet & Society at Harvard. His research focuses on information law, policy, and society issues. The research team has established a successful interdisciplinary collaboration through an ongoing NSF Frontier project *Privacy Tools for sharing Research Data* led by Prof. Vadhan.

## Formal Privacy Models and Title 13: Project Description

### 1 Introduction

The US Bureau of the Census (BOC) collects large quantities of data that can be useful for research and decision making by policymakers, businesses, and academics. The BOC is responsible for analyzing and publishing useful statistical data. As much of the collected data pertain to individuals, households, and establishments, the BOC also has a legal obligation to protect their privacy. The BOC has long employed statistical disclosure limitation (SDL) techniques in order to guarantee the confidentiality of the data its maintains and releases. However, with the SDL techniques traditionally used, it is difficult and labor-intensive to determine with confidence that these two requirements—privacy and utility—have been satisfied. Traditional SDL techniques are heuristic in nature, and it is increasingly harder to demonstrate their efficacy as the information environment rapidly changes, and inaccurate statistical results have been reported from public-use microdata samples (PUMS files) constructed by the BOC using traditional SDL [10].

The proposed research project seeks to further the use of formal approaches that have the potential to provide rigorous guarantees that legal requirements for privacy and utility are met. Applying such approaches requires (a) bridging legal privacy requirements with mathematical privacy requirements, and (b) designing analysis methods that provide statistical utility while satisfying the privacy requirements from (a). Our team has developed tools satisfying these two goals in other contexts [29], and with this proposal we seek to collaborate with BOC staff to develop similar solutions that are tailored to the bureau’s specific requirements.

**Background.** BOC collects and analyzes data about the nation, its people, and economy, as codified in Title 13 of the United States Code [6], and its publications serve as the basis for research and decision making by policymakers, researchers, and businesses. Because much of these data pertain to individuals, households, and establishments, BOC has an obligation to employ disclosure limitation practices that strike a balance between enabling socially beneficial uses of data and protecting individual and establishment privacy interests. To this end, BOC applies a suite of statistical disclosure limitation (SDL) techniques to ensure that the risk of individual information being identified in BOC publications is small.

Since the 1990s, computer scientists and legal scholars have been calling attention to ways in which traditional approaches to statistical disclosure limitation often fail to address emerging privacy risks in the sharing of research data [100, 36, 92]. The urgency of this problem has increased in recent years in response to a number of high-profile privacy breaches enabled by significant advances in analytical capabilities and the wider availability of personal data from different sources [101, 82, 83, 84, 62, 14, 41]. There is a growing recognition that current approaches to regulating privacy in data releases rely on concepts that reflected an information regime that is very different from the environment today. In particular, many approaches in current use address only a limited scope of potential privacy threats, as they were developed based on models of specific strategies and auxiliary information likely to be used by potential attackers. In addition, the patchwork of requirements for data releases in different contexts, organizations, and jurisdictions, and the mismatch between these requirements and scientific understandings of privacy, creates uncertainty for practitioners responsible for data releases. For instance, BOC is governed by Title 13, and in certain contexts may handle data governed by agency- and sector-specific laws such as the Confidential Information Protection and Statistical Efficiency Act [1], the Privacy Act of 1974 [5], and the E-Government Act of 2002 [2]. In addition to variations in the legal and policy standards

for privacy protection and utility that may apply in different settings, there is uncertainty regarding how such standards apply to the adoption of new and emerging technologies.

Research results from the computer science literature offer new approaches to privacy that overcome the shortcomings of traditional privacy measures. Formal privacy models such as differential privacy [42] provide strong, quantitative notions of privacy that are meaningful regardless of the auxiliary information or strategy used by a potential attacker. Federal statistical agencies are increasingly recognizing the need to improve their understanding of the science underlying robust solutions to privacy, and they are currently exploring ways to expand their privacy toolsets by interpreting and applying scientific concepts to their disclosure limitation practices. For example, BOC’s *OnTheMap* tool is based on a synthetic dataset carefully generated from confidential data in a way that satisfies a variant of differential privacy. It enables users to interactively study the commuting patterns of workers across the United States through the Longitudinal Employer-Household Dynamics program [76].

In this project, we propose to address two major challenges confronting wider adoption of formal privacy models with BOC and similar organizations.

- (a) There is a wide conceptual and practical gap between the approaches found in formal privacy models and the heuristic approaches in current use and contemplated by existing regulatory and policy frameworks.
- (b) There is a gap between theoretical developments showing that formal privacy models like differential privacy permit, in principle, a wide collection of analyses and the actual use of analysis and publication techniques by the BOC.

Successfully implementing tools and policies that bridge these gaps will require close collaboration between theoretical computer scientists, legal scholars, and BOC staff.

## 1.1 Core Themes and Specific Research Areas

This research has two primary directions, corresponding to the gaps articulated above: (a) assessing through legal-technical methods the use of formal privacy models by BOC; and (b) developing privacy-preserving statistical analyses for use by BOC. These threads and their points of intersection are illustrated in Figure 1 and outlined in turn below.

**Technical-Legal Analysis (PIs Gasser & Nissim).** Applying our combined expertise in law, policy, and computer science, we will conduct a legal-technical analysis of the use of formal privacy models by BOC. This will involve a legal analysis of the statutes and policies specifying the privacy and utility requirements for BOC’s analyses and publications of data, an analysis of historically-accepted standards for privacy in BOC publications, and a legal-technical analysis of formal privacy models in light of these statutory and policy requirements.

**Private Statistical Analysis (PIs Nissim, Smith, & Vadhan).** Building on our expertise in formal privacy models such as differential privacy, we will construct new data-release algorithms that satisfy formal privacy requirements and are optimized for the types of statistical utility that are most important in data released by the BOC.

An essential aspect of this proposal lies in the intersection points between these themes. Our analysis of BOC’s legal and policy requirements for privacy and utility will help us tune the formal mathematical requirements to address in specific technical contributions. In addition, investigating the utility-privacy tradeoffs in specific statistical analyses in conjunction with legal and policy requirements will inform our choice of privacy parameters.

In Section 3, we discuss in more detail the specific areas that we plan to explore.

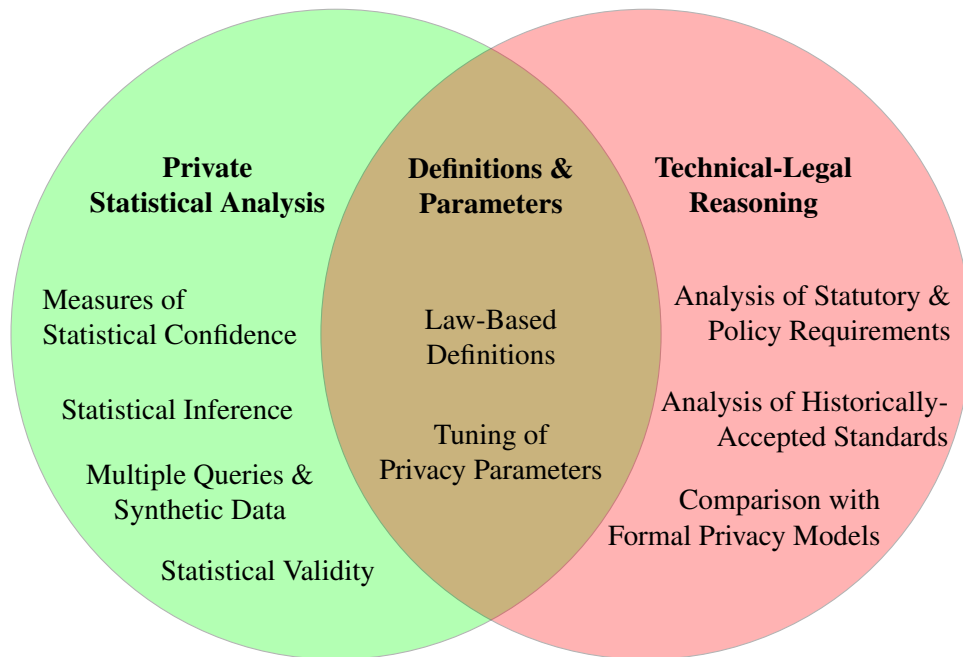


Figure 1: Illustration of the main components of the proposed project.

## 1.2 Research Team

Our research team includes experts in privacy and cryptography (PIs Nissim, Smith, & Vadhan), and experts in information privacy law, policy, and society issues (PI Gasser, Sr. Researcher O’Brien, & Research Fellow Wood). The collaboration of the group emerged from our successful collaboration through Harvard University’s *Privacy Tools for Sharing Research Data* project, for which Vadhan is the lead PI, Gasser is a co-PI, and Nissim, O’Brien, Smith, and Wood are senior research collaborators. In addition, Nissim, Smith, and Vadhan have collaborated on several research projects related to theoretical aspects of differential privacy. Gasser, Nissim, O’Brien, Vadhan, and Wood have also collaborated on related research bridging legal and mathematical definitions of privacy.

### Principal Investigators

*Kobbi Nissim (Georgetown University)* is one of the inventors of differential privacy [42] and was a coauthor on all of the work leading up to its definition [35, 38, 25]. His paper initiating formal research on privacy [35] received the ACM PODS Alberto O. Mendelzon Test-of-Time Award in 2013, and his paper presenting differential privacy [42] received the IACR TCC Test-of-Time award in 2016 for providing a solid mathematical foundation for a vast body of subsequent work on private data analysis.

*Salil Vadhan (Harvard University)* is the Vicky Joseph Professor of Computer Science and Applied Mathematics in the Harvard Paulson School of Engineering and Applied Sciences. He is the lead PI on *Privacy Tools for Sharing Research Data*, a five-year NSF Secure and Trustworthy Cyberspace Frontier Project (2012–17), which is a collaboration between computer science, law, social science,

and statistics. In recent years, he and his collaborators have obtained numerous results delineating the border between what is possible and impossible with differential privacy, through numerous computational and statistical lower bounds. He is a recipient of a Simons Investigator Award, a Gödel Prize, a Guggenheim Fellowship, a Phi Beta Kappa Award for Excellence in Teaching, and the ACM Doctoral Dissertation Award.

*Adam Smith (Penn State University)* is one of the inventors of differential privacy and a Professor of Computer Science and Engineering at Penn State. He received Presidential Early Career Award for Scientists and Engineers (PECASE) in 2009. His paper presenting differential privacy [42] received the IACR TCC Test-of-Time award in 2016. His contributions include an array of techniques for designing differentially private algorithms as well as foundational results on private machine learning and statistical inference.

*Urs Gasser (Harvard University)* is the Executive Director of the Berkman Klein Center for Internet & Society and a Professor of Practice at Harvard Law School. He is also a visiting professor at the University of St. Gallen (Switzerland) and at KEIO University (Japan), and he teaches at Fudan University School of Management (China). Professor Gasser's research and teaching activities focus on technology and information law, policy, and societal issues. Throughout his research, Professor Gasser has closely examined privacy both in the US and internationally in a number of contexts. He has an established history of collaboratively working with leading organizations and computer scientists, technologists, and sociologists and routinely participates in interdisciplinary research and problem solving.

### **Other Key Personnel**

*David O'Brien (Harvard University)* is a senior researcher at the Berkman Klein Center. He has contributed legal research to and led the Berkman Klein Center's efforts across a variety of projects, publications, and initiatives, spanning the topics of privacy, intellectual property, cloud computing, cybersecurity, digital publishing, and internet governance. Under the direction of Executive Director Urs Gasser, David O'Brien currently leads the Berkman Klein Center's research contributions to the *Privacy Tools for Sharing Research Data* project.

*Alexandra Wood (Harvard University)* is a research fellow at the Berkman Klein Center for Internet & Society and a collaborator with the *Privacy Tools for Sharing Research Data* project legal team. Her research explores new and existing legal, regulatory, and contractual approaches to data privacy and contributes to the development of new legal instruments, frameworks, and policy recommendations through the *Privacy Tools* project collaboration. Before joining the Berkman Klein Center, she served as a legal fellow with US Senator Barbara Boxer, and as a law clerk with the Center for Democracy & Technology and the Electronic Privacy Information Center.

**Project Coordinator.** We will recruit a coordinator at Georgetown University to help us collaborate across multiple sites and disciplines. The coordinator will deal with all logistical aspects of the project, including scheduling meetings between the researchers (setting up online meetings when needed), travel arrangements, collection and organization of material for project reports and project specific financial summaries. The project coordinator will record notes at project meetings, collect presentations and reports, and share them to keep all participants updated on project progress.

**Established collaborations.** Success in the research directions will rely both on the deep expertise of the PIs in their specific areas and on substantive interdisciplinary collaboration between them. This proposal emerged from the team's successful collaboration through the *Privacy Tools for*

*Sharing Research Data* project (National Science Foundation Frontier project, CNS-1237235), which is a joint effort between computer science, law, social science, and statistics. In this project, Vadhan is the lead PI, Gasser is a co-PI, and Nissim, O’Brien, Smith, and Wood participate as senior research collaborators.

The unifying focus of the *Privacy Tools* project is the design of tools for privacy-protective sharing of research data in social science data repositories, and we expect a number of the tools to be deployed in Harvard’s Dataverse repository by the time the project ends in fall 2017. A significant amount of effort in the Project has been devoted towards overcoming the challenges in deploying differential privacy in practice, and we will leverage this experience in the proposed work. In particular, we are constructing a differentially private data analysis system, PSI, that enables social science users (without specialized expertise in differential privacy, law, social science, and statistics) to release and interpret differentially private statistical information about privacy-sensitive datasets for which it is not safe to share the raw data. The vision and design of PSI is described in a working paper coauthored by Nissim, Vadhan, and others [53]. As part of our work on PSI, we have designed new differentially private algorithms that are tailored to the kinds of analyses most common in social science, providing both descriptive statistics as well as inference procedures with measures of statistical confidence (e.g.,  $p$  values and confidence intervals) that are missing from much of the previous literature on differential privacy.

In addition, a working group led by PI Nissim has developed a methodology for formally proving that a technological method for privacy protection satisfies the requirements of a particular law involving (a) translating a legal standard into a formal mathematical requirement of privacy and (b) constructing a rigorous proof for establishing that a technical approach satisfies the mathematical requirement derived from the law. We first applied this methodology to argue that differential privacy satisfies the privacy requirements of the Family Educational Rights and Privacy Act (FERPA), as described in a working paper authored by Nissim, Wood, Gasser, O’Brien, Vadhan and others [86] and presented at a variety of venues, including the Second Census-MIT Workshop on Big Data in 2015 and the Privacy Law Scholars Conference in 2016.

Through the *Privacy Tools* project and other collaborations, the current proposal’s team has produced a number of joint publications in top computer science venues [42, 87, 69, 70, 89, 33, 28, 19, 90], as well as papers directed at law, policy, and social science audiences that integrate law and computer science perspectives on privacy [91, 110, 13, 106]. These collaborations have also generated a number of joint legal-computer science policy commentaries [11, 12, 111] and facilitated mutual feedback on interdisciplinary research efforts. In addition, our team collectively has co-supervised many undergraduate and graduate students from computer science, law, social science, and statistics, providing junior scholars with opportunities to engage in interdisciplinary privacy research.

## 2 Background

Successful attacks have demonstrated that some traditional SDL techniques can fail to provide adequate privacy. In particular, de-identification techniques for sanitizing data prior to release, which are common in information privacy laws such as FERPA [3] and HIPAA [4], are open to re-identification [101, 83]. Moreover, PIs Nissim, Nissim, Vadhan, and others have shown that even publishing statistical estimates, without releasing the underlying data, can lead to a massive leakage of individual information [35, 43, 40, 65, 48]. Motivated by these revelations, a formal theory of privacy in computation has been pursued in computer science for the last twelve years, yielding

the concept of differential privacy, as introduced by PIs Nissim, Smith, and others [35, 38, 25, 42].

## 2.1 Differential privacy

*Differential privacy* is a property of a computation that takes a sensitive dataset  $D$  as input and releases statistical information about  $M$  (in a probabilistic manner, injecting random noise in order to protect privacy). Differential privacy requires that the information released “looks essentially the same” regardless of any one individual’s data. Formally:

**Definition 1** ([42]). *A probabilistic computation  $M$  satisfies  $\epsilon$ -differential privacy if, for every two datasets  $D, D'$  that differ only on the data of one individual, the output distributions  $M(D)$  and  $M(D')$  are similar in the sense that for every set  $T$  of potential outcomes we have  $\Pr[M(D) \in T] \leq e^\epsilon \cdot \Pr[M(D') \in T] \approx (1 + \epsilon) \cdot \Pr[M(D') \in T]$ .*

So, when  $\epsilon$  is small, an adversary  $A$  viewing the output of  $M$  would draw essentially the same conclusions if we removed any individual from the dataset (or replaced her data with that of an arbitrary other individual), and in this sense, the privacy of each individual is protected. This interpretation holds regardless of what computational strategy is used by  $A$  (since if  $M$  is differentially private, then so is  $A \circ M$  for any function  $A$ ), or the auxiliary information that  $A$  has about the dataset or the individual being targeted (as we quantify over all datasets  $D, D'$  and strategies  $A$ ). See [68] for a Bayesian formalization of this interpretation of differential privacy.

The parameter  $\epsilon$  is often referred to as the *privacy loss parameter*. For small values of  $\epsilon$  (e.g.,  $\epsilon = 0.1$ ), differential privacy guarantees that anything an observer of the outcome  $M(D)$  could learn about an individual in  $D$ , he could also learn from  $M(D')$  where  $D'$  does not contain that individual’s information.

An important property of differential privacy is that it *composes*. In other words, an individual participating in two or more differentially private analyses experiences a degradation in privacy that is bounded by (a function of) the privacy parameters of the individual analyses.<sup>1</sup> The simplest composition theorem states that the accumulation of the privacy loss parameters is at most linear:

**Theorem 1** (Basic Composition). *A mechanism that permits interactions with programs  $M_1, M_2$  with parameters  $\epsilon_1$  and  $\epsilon_2$  respectively satisfies  $(\epsilon_1 + \epsilon_2)$ -differential privacy.*

PI Vadhan’s research has resulted in more efficient composition theorems showing that the increase in parameters is often slower than linear. These include the *advanced composition theorem* [45] and the *optimal composition theorem* [64, 81], where the latter result gives a tight characterization of composition for differential privacy. For a more detailed exposition of differential privacy, we refer the reader to [39].

**Background on tools for differential privacy.** Research on differential privacy has yielded a plethora of techniques that can be used to design computations that satisfy this definition. Several of these techniques have been implemented in software applications allowing users who are not experts in privacy to design data analysis that satisfy the differentially private definitions. This research has mainly followed three paths:

- Implementation of specific algorithms tailored to specific problems or domains, trying to get the best privacy-utility trade-off from the specialization. These include OnTheMap [76], MWEM [60], DualQuery [52], PrivBayes [115], DP-WHERE [79], and RAPPOR [49].

---

<sup>1</sup>To date, privacy concepts other than differential privacy do not compose, i.e., the degradation in privacy can be unbounded.

- Implementation of general techniques for combining basic mechanisms as building blocks for a more complex analysis. These tools include PINQ [77], Airavat [94], GUPT [80], and our work within the Harvard Privacy Tools project [53].
- Utilization of techniques from programming language research to design methods for users to program and prove their own analyses are differentially private. Examples are Fuzz [57, 51] and Certipriv [16, 17].

**Note on our use of differential privacy.** We focus on differential privacy because of its developed theory but this should not be understood as advocating that differential privacy be the sole privacy concept in use by BOC—a significant part of our research would be devoted to understanding and formalizing privacy requirements in the context of BOC work and to developing analyses that satisfy these requirements.

## 2.2 Private statistical analysis

There is a growing body of work on the design and limitations of differentially private algorithms for tasks related to statistical inference and machine learning. It includes significant lines of work on PAC (probably approximately correct) learning [25, 69, 20, 23, 21, 22, 23, 24, 28] and point estimation for parametric models (e.g., [30, 31, 95, 37, 97, 18]). To a lesser extent, the literature also explores topics such as nonparametric estimation [108], model selection [71, 98, 99, 74], goodness of fit measures, network models, Bayesian sampling-based methods (e.g., [34, 107]), and hypothesis testing [103].

However, the extant literature has several important limitations. To the extent that the utility of these methods is analyzed rigorously, the literature has largely focused on understanding asymptotic behavior as the data set size,  $n$ , grows. It has also focused on measures of accuracy tied to the data at hand, asking how well an algorithm controls *empirical* risk (such as misclassification rate on the input training data). In statistical applications, however, the data is often a sample from an underlying population, and *population* risk measures are more relevant. This focus on empirical error arose in part because, for very large data sets, the empirical error is small enough (typically,  $O(1/\epsilon n)$ ) that it is dominated by the sampling error—the extent to which the data set’s characteristics differ from those of the population—which generally scales as  $O(1/\sqrt{n})$ . In such a regime, differentially private statistics may be used as direct substitutes for their nonprivate analogues without loss of statistical power (see, e.g., [97] for general statements of this form).

The usefulness of this type of analysis is limited by several factors. First, asymptotics say little about finite sample sizes—especially when statistics are broken down by subpopulation—and ignoring the distortion due to differential privacy leads to incorrect inferences—especially inasmuch as noise can dramatically change the coverage of confidence intervals and  $p$ -values [103, 54]. This effect is exacerbated in the typical case where the data are high-dimensional; the error introduced to ensure privacy must generally grow quickly with the dimension [27, 47]. Even if noise is negligible, existing methods may assume that the input has particular characteristics (integrality, for example), leading to unpredictable behavior when outputs violate these assumptions. Second, little research is currently known about to provide interpretative information—confidence intervals, goodness of fit measures, and hypothesis tests—that must accompany statistical estimates for them to be useful. Third, much of the existing literature has focused on the privacy-utility tradeoff for carrying out a *single* statistical inference procedure. There are several approaches to allowing analysts to carry out multiple inference procedures, but they carry a significant cost in utility. One objective of this proposal is to address these limitations—see Research Goal 2 below.

### **2.3 Regulation of privacy in the use and release of statistical information**

Numerous laws and policies govern the confidentiality of information held by federal statistical agencies. For instance, BOC is bound by Title 13 of the U.S. Code [6]. Section 9 of Title 13 [8] guarantees the confidentiality of information maintained by BOC, prohibits the use of collected data by the government for any purpose other than the statistical purposes for which it is supplied, and prohibits anyone other than the sworn officers and employees of the Department of Commerce or bureau or agency thereof to examine the individual reports. The parties to whom Section 9 applies are prohibited from making any publication whereby the data furnished by any particular establishment or individual under this title can be identified [8]. Other sections of the Title also provide guidance regarding the types of information disclosure risks BOC must protect against. For example, Section 8 provides that “[s]ubject to the limitations contained in sections 6(c) and 9 of this title, the Secretary may furnish copies of tabulations and other statistical materials which *do not disclose the information reported by, or on behalf of, any particular respondent*” to particular recipients “upon payment of the actual or estimated cost of such work” [7] (emphasis added). Other laws such as the Privacy Act of 1974 [5], the Confidential Information Protection and Statistical Efficiency Act [1], and the E-Government Act of 2002 [2] protect personal information maintained in government records, and allow the release of information in a form that is not individually identifiable to support scientific research and public policy decisions.

Legal scholars have noted that privacy laws in these areas rely on ambiguous concepts such as personally identifiable information, draw seemingly arbitrary distinctions between different types of information, and encourage the use of heuristic de-identification techniques [92, 96, 112]. Arguably, these laws create uncertainty regarding the scope of information covered and the legal standard that governs in a given context (especially in cases where an organization may fall under multiple legal regimes). In many cases, laws and related guidance often provide inadequate guidance on appropriate privacy measures and their implementation in specific settings, and lead to limited privacy protection in practice [13]. In particular, these characteristics of law and policy pose potential challenges for the implementation of techniques that provide formal privacy guarantees, as discussed in more detail below.

### **2.4 Bridging legal and computer science approaches to privacy**

Privacy laws rely on narrower and less technically formal conceptions of risk than those described by the computer science literature. The boundaries of the law are sometimes unclear and the common practices endorsed by some regulations and policies fail to protect against the full range of data privacy risks. In addition, given that the law is generally context-specific, organizations are often tasked with implementing, on a case-by-case basis privacy measures that are tailored to the requirements of the applicable statute or regulation (or, in some cases, multiple laws).

In a recent collaboration our team presented a new methodology [86] for bridging between the very different approaches to defining privacy in regulations such as FERPA [3] and HIPAA [4] and differential privacy. In particular, this new methodology leads to a formal argument that differential privacy satisfies the requirements of these two regulations. The methodology involves two steps: The legal standard is first translated into a formal mathematical requirement of privacy following game-based cryptographic definitions. Ambiguity is overcome by taking a worst-case approach to interpretation with respect to the adversarial computational resources, access to the privacy mechanism, and goal. Given the resulting formulation, a rigorous proof is given that differential privacy

satisfies the mathematical requirement derived from the law.<sup>2</sup> The analysis takes a conservative approach and extracts a mathematical requirement that is robust to some of the inherent ambiguities in interpreting the law.

While FERPA, HIPAA, and differential privacy are used to illustrate application of this methodology, we believe it is a general approach that can be extended to bridge between technologies beyond differential privacy and privacy laws beyond FERPA and HIPAA. The degree of rigor employed enables us to make strong arguments about the privacy requirements of statutes and regulations. Furthermore, with the level of generalization afforded by this conservative approach to modeling, differences between sector- and institution-specific standards are blurred, making the methodology broadly-applicable. We believe this methodology can be used to argue that differential privacy (and potentially other privacy technologies) can be used to satisfy various applicable legal standards and overcome, for example, differences in definitions of *personally identifiable information* across different statutes and regulations.

### 3 Research Methods and Objectives

Our research takes the view that addressing privacy rigorously from both a technical and a legal perspective is essential for BOC’s activities. Hence, our research program involves employing analysis tools from the very different disciplines of computer science and law, and a substantial part of our work would be devoted to developing novel methodologies for bridging the conceptual gaps between technical and legal approaches to privacy and demonstrating adherence to both scientific principles and legal standards for privacy protection.

Our technical work is grounded in theoretical computer science, and hence follows formal, rigorous approaches to privacy, which have proved to be successful in providing deep insights into privacy questions (and, in particular, drawing boundaries between what can be computed with robust privacy protection and what is inherently non-private). These formal approaches have also been shown to provide a framework for producing a huge variety of analyses with provable privacy guarantees, as well as providing connections between privacy and well-grounded statistical analysis practices.

Our work is a collaboration across disciplines and institutions, and we will employ mechanisms for establishing and maintaining active collaborations among the key personnel involved in this project, building on connections our team has established through related joint research projects. We are also looking forward to a substantial involvement by BOC in this project, to better tune our research to BOC’s needs, and to bring knowledge from academic research to BOC. Although much of our work is theoretical in nature, we will be implementing and experimenting with some of our algorithms. We will consult BOC regarding representative datasets to use in this process.

In the remainder of this section, we describe our research methods and objectives, following the project structure as described in Section 1 and illustrated in Figure 1.

---

<sup>2</sup>This approach contrasts with research by others to close this gap by modifying the definition of differential privacy to directly address the concept of identifiability [72].

## **Research Objective 1: Technical-legal reasoning on the use of formal privacy models by the US Bureau of the Census**

The first research objective involves analyzing the BOC’s legal and policy requirements for privacy protection using methods that combine legal and technical reasoning. There are two key components to this research objective: (a) analyzing the privacy requirements of Title 13 and related laws and policies applicable to BOC’s analysis and publication of data, and (b) analyzing the use of formal privacy models with respect to the privacy requirements obtained in (a).

The main product of this research objective will be an integrated legal-technical argument for the use of formal privacy models by BOC that is rigorous both legally and mathematically. This line of research will enhance the understanding of the relationship between legal and mathematical definitions of privacy both in general and in the context of BOC work. It will also highlight issues that could be addressed in future guidance documents and policies, so as to better facilitate the adoption of novel disclosure limitation techniques in addition to (or, in some cases, replacing) traditional disclosure limitation techniques used by BOC.

**Subgoal 1.1: Legal analysis of the privacy requirements of Title 13.** We will analyze the nature of the privacy protection required by Title 13 and related laws and policies that are applicable to BOC’s analysis and publication of data. The expected outcome is a legal memorandum analyzing the privacy-related definitions and requirements of Title 13, focusing on their potential applicability to tools that satisfy formal privacy requirements. This memorandum will be used by the project team as a working document that serves as a reference point for the analysis bridging between Title 13 and formal privacy models described under Subgoal 1.2.

The confidentiality of information maintained by BOC is mandated by Section 9 of Title 13 of the U.S. Code [8]. However, these requirements that “[no] particular establishment or individual under this title can be identified” provide a privacy goal that is succinct and open to interpretation, which will present challenges for our analysis in Subgoal 1.2 below. Hence we plan to inform our analysis by examining other relevant standards including the Confidential Information Protection and Statistical Efficiency Act (CIPSEA) [1], as well as federal guidance documents such as Federal Committee on Statistical Methodology’s Statistical Policy Working Paper 22 [50], and internal policies and reports produced by BOC. Analysis of this combination of sources will allow us to gain a fuller picture of the nature of privacy protection intended, and also highlight points of ambiguity and inconsistency in applicable laws and policies. We plan to analyze these sources with respect to the types of information that are protected (or not protected), the capabilities and resources of potential adversaries acknowledged implicitly or explicitly, and the types of privacy breaches recognized, among other factors relevant to the modeling described in the discussion of the next component of the research.

This subgoal relies on our past experience in carrying out similar legal analyses for the Family Educational Rights and Privacy Act (FERPA) [3] and the Health Insurance Portability and Accountability Act (HIPAA) [4], which were supported by detailed memoranda systematically analyzing the language of the regulations, the rulemaking history, and relevant agency guidance.

**Subgoal 1.2: Analysis of formal privacy models in relation to the privacy requirements of Title 13.** The privacy definition set forth by a legal standard such as Title 13 is very different from formal privacy definitions such as differential privacy. For this reason, it is important to demonstrate that the use of techniques satisfying a specific formal definition of privacy also satisfies the requirements of the law. Some of the differences that need to be bridged include: (1) the level of formality and precision—formal privacy definitions such as differential privacy are mathematically

precise whereas legal definitions are not mathematically precise and often inherently or intentionally ambiguous; (2) what constitutes an attack—legal definitions of privacy consider an attacker that attempts to *identify* individuals and establishments, a notion that is not directly applicable, e.g., to statistical computation, whereas formal models describe attackers that attempt to perform a more general distinguishing task; and (3) legal definitions do not fully define the attacker’s capability and resources, and instead refer to standards of reasonableness, whereas formal models need to define the attack model precisely.

We will examine how a methodology for bridging between legal and mathematical standards of privacy that we developed in the context of Harvard’s *Privacy Tools for Sharing Research Data* project [86] can be applied to make such formal arguments with respect to Title 13 and other regulations, guidance, and internal documentation governing the BOC’s activities. We expect to make adjustments to this methodology if needed and then apply it with respect to differential privacy (and other formal privacy models that would emerge as relevant during the lifetime of the project) to develop a legally and mathematically rigorous argument that the formal definitions we use (with an appropriate setting of the privacy parameters) satisfy applicable legal requirements.

The products of Subgoal 1.1 will be used as a basis for establishing our argument. We expect our analysis to be based, in particular, on a modeling of factors such as the attacker’s computational power, the attacker’s goal, the external information available to the attacker, and the definition of what constitutes a privacy breach. We expect the products of Subgoal 1.2 to inform our work in research objectives 2 and 3.

## **Research Objective 2: Developing privacy-preserving statistical analysis for use by the Bureau of the Census**

Our second research objective is to develop statistical methodology for tasks important to the BOC that satisfies rigorous privacy guarantees, such as differential privacy, while providing accurate and interpretable outputs. Below, we spell out specific directions related to differential privacy. The challenges we address apply equally well to other rigorous notions of privacy, however, and the proposed research will also investigate their resolution in the context of the formal models developed in Research Objective 3.

**Subgoal 2.1: Statistical inference.** Methodology for differentially private statistical inference still faces a number of crucial challenges—laid out in Section 2.2—that require a combination of theoretical and empirical research. Our main goal is a better understanding of optimal methodology for finite sample sizes. We propose to pursue several avenues along these lines. For each one, we intend a combination of theoretical work and empirical evaluation, with a strong feedback loop of empirical results guiding further theoretical work. We also hope to collaborate closely with the BOC as the project progresses to select statistical problems and data sets that are representative of Census workloads and relevant to the BOC’s needs.

An initial thrust will be understanding how further processing of differentially private outputs can improve their statistical usefulness. Recent research [15, 61, 109, 75, 73, 67] has shown that post-processing of simple noise addition techniques leads to significant accuracy gains in practice. The proposed project will seek a deeper theoretical rationale for this phenomenon (which can be analyzed theoretically in some cases [85]). We will also investigate how the phenomenon can best be used with more sophisticated noise addition techniques such as those that arise in MWEM [59, 60], iterative optimization algorithms [109, 32, 18, 102] and objective perturbation [31]. The PIs already have work along these lines in the context of hypothesis testing (see Subgoal 2.2).

More broadly, our aim is to develop new methods for accurate and private statistical inference. Most work on private inference has considered estimation problems with either succinct sufficient statistics nor convex formulations, and most techniques are limited to these settings. One promising avenue is to apply recently developed noisy versions of iterative optimization methods (such as the PIs’ work on gradient descent [18]) to nonconvex problems such as estimating mixture and tensor models. Recent work on deep learning [9] provides a large-scale and computationally intensive example of this idea; we conjecture that simpler procedures suffice for more concise models. Another important avenue is the investigation of model selection problems, building on the PIs’ work [71, 98, 74] on model selection for linear regression. In addition to the difficulties above, model selection problems exhibit a multiple hypothesis testing problem that makes statistically valid inferences delicate.

A related goal is to better understand the limits of differentially private inference and characterize exactly optimal algorithms for specific problems. This task is challenging. For example, a basic exact optimality result in statistical theory is the Neyman-Pearson lemma, which characterizes the optimal hypothesis tests distinguishing between two given distributions. There is no differentially-private analogue of this lemma. Such a result, or results along similar lines on the minimum variance private estimators, would guide the design of private inference algorithms as their nonprivate analogues guided the development of classical statistics. The work of PIs Vadhan, Nissim and Smith has established many of the known limits on private data analysis, showing asymptotic regimes in which accurate statistics leak individual information (see the recent survey [105]). The proposed project will seek to tighten these bounds, showing concrete attacks that leak information, for example when many two-way contingency tables are released. Geometric and information-theoretic techniques (as in [56, 64]) also provide a promising foothold.

**Subgoal 2.2: Measures of statistical confidence.** As shown in [10], the SDL techniques used by the BOC have several times resulted in public-use microdata samples (PUMS files) that give inaccurate statistical information (in particular about people around or above retirement age). Reducing accuracy is necessary to protect privacy; indeed, it has been demonstrated that *any* release of “too many” aggregate statistics with “too much” accuracy necessarily compromises privacy (e.g. allowing an adversary to reconstruct almost all of the entries in a sensitive dataset) [36, 62, 47].

However, formal privacy models such as differential privacy offer the possibility of mitigating these inaccuracy concerns. A key point is that the methods used to construct the releases (namely, where noise is introduced and according to what distribution) can be made completely public and transparent; the privacy protections do not rely on the methods or noise parameters being secret. Thus, in principle, one can fully take the additional noise into account in any statistical inference procedure, and thereby avoid drawing false conclusions even when the sample size  $n$  is too small for the asymptotics to kick in (e.g. by appropriately adjusting  $p$  values and the sizes of confidence intervals).

PI Vadhan and collaborators have recently shown how to account for differentially private noise in basic hypothesis testing procedures ( $\chi^2$  tests for goodness of fit and independence testing) [54] and in constructing confidence intervals for a normal mean [66]. In the latter case, lower bounds are also proven, showing that the lengths of the obtained confidence intervals are asymptotically nearly optimal among all differentially private algorithms that achieve a desired coverage probability.

In the proposed project, we will carry out a similar effort for more complex inference procedures, such as those mentioned in Subgoal 2.1 above. In particular, a natural next step after our work on confidence intervals for a normal mean is providing  $p$  values and confidence intervals for ordinary least-squares (OLS) regression on data drawn from a multivariate Gaussian. Like in [66],

we will seek to find asymptotically optimal differentially private algorithms, and doing so will require the design of new algorithms as well as proving lower bounds (which are specialties of PIs Nissim, Smith, and Vadhan).

One benefit of differentially private procedures over classical ones is that their inferential properties are maintained even when data analysis is *adaptive* (i.e. each query is chosen based on the results of previous queries). Indeed, differentially private analyses automatically prevent overfitting and allow for statistically valid adjustments of  $p$ -values and confidence intervals (as shown in [46] and the PIs’ work [19, 93]). We will use this connection to refine and optimize our measures of statistical confidence (as in PI Smith’s work on population risk measures [18]).

**Subgoal 2.3: Multiple queries and synthetic data generation.** To allow analysts to perform multiple inference procedures with differential privacy, the simplest approach is to apply composition theorems, such as Theorem 1 or its improvements [45, 64, 81], which requires giving potential data analysts a limited “privacy budget,” which implies that analysts can only ask a bounded number of queries and that each query is answered with less accuracy. We have designed such an interactive query system (including tools to help analysts apportion their privacy budget in the most effective way) as part of  $\Psi$ , the differentially private data-sharing interface we have built for social science data repositories. We will work with the BOC to explore the extent to which a similar model could work for them, perhaps as something intermediate between public releases (i.e. PUMS files) and secure research data centers (RDCs), offered, e.g., to a community of researchers that can be trusted to not collude to exceed their individual privacy budgets.

However, it would be preferable for BOC to be able to do a single release of a rich statistical summary or a synthetic dataset that can be safely published in its entirety, and allows analysts to obtain accurate answers to many different statistical queries. In theory (and asymptotically), differential privacy is compatible with producing extremely rich synthetic data releases [26, 59]. However, these general methods are computationally infeasible on high-dimensional data, as was shown to be inherent by PI Vadhan and collaborators [44, 104]. Moreover, it is not clear how to account for the noise introduced by such procedures (which is correlated in sophisticated ways) when performing inference and measuring statistical confidence.

We will explore how to overcome these barriers for the types of data and analyses most relevant to the BOC. In particular, if we restrict attention to parametrized models, a natural approach is to estimate the model parameters in a differentially private way (e.g. estimate the covariance matrix for data drawn from a multivariate Gaussian) and then use those model parameters to generate synthetic data. Works by PI Vadhan and others [116, 34, 107] show that in some cases (e.g. Bernoulli or Gaussian data), differential privacy can be obtained simply by drawing model parameters randomly from a Bayesian posterior, which is very closely related to the multiple imputation approach sometimes used in traditional SDL.

### **Research Objective 3: Definitions and Parameters**

There are two key components to this objective: (1) adjustment of formal privacy models to address a wider set of privacy issues than individual privacy (we will examine potential definitions, their provable privacy consequences, and their limitations on statistical computations, to develop variants that balance these two concerns); and (2) tuning of privacy parameters, taking into account the privacy guarantees of SDL techniques currently in use by BOC.

**Subgoal 3.1: Law-based formal models.** Differential privacy and related formal privacy models focus on protecting the privacy of individuals. With respect to differential privacy, this is captured

by the requirement that an analysis should be insensitive to an arbitrary addition, removal, or change of an individual’s record in the sense that no attacker (even if computationally unlimited and arbitrarily knowledgeable) can distinguish between the outcome of the analysis with and without that individual’s record. Title 13, however, requires that “[no] particular establishment or individual ... can be identified” [8]; thus, establishments also require privacy protection. One approach to deal with this requirement is to modify differential privacy to also require an analysis to be insensitive to the addition or removal of an establishment (as in [58]).

Part of our research will be devoted to an in-depth examination of the types of privacy protection that BOC must provide to individuals, establishments, and, to some extent, groups. This analysis will be based on our findings in Subgoal 1.2 and on BOC’s legal and technical experience. We will translate these findings into new formal models of privacy, examine their relationship to the well-established concept of differential privacy, and determine to what extent these models support the types of statistical analyses we will focus on in Research Goal 2.

**Subgoal 3.2: Tuning of privacy parameters.** One barrier to using differential privacy or other formal privacy models is the need to set privacy parameters (for differential privacy, this is typically the parameter referred to as  $\epsilon$ ). While deciding on permissible privacy parameters is a *normative* question, it can and should be informed by our technical understanding. For example, differential privacy  $\epsilon$  bounds the maximum increase in risk incurred by any individual whose information is used in the computation (as compared with the risk when the computation is performed without using the individual’s information) [78, 55, 113, 88, 114, 63].

Our goal is to search for methodologies that will help initiate a process in which privacy parameters are adjusted over time to achieve an a priori unknown balance between the privacy level provided and the utility obtained from data. In particular, this research thread seeks to examine standards and best practices historically accepted by BOC, and compare the level of privacy protection these techniques provided (or were believed to provide) to the protection obtained by formal privacy models. We will explore how to set formal privacy parameters to ensure (at least) a similar level of privacy protection.

Documentation of the analyses justifying BOC’s adoption of particular traditional SDL techniques is expected to serve as an input for our analyses. This research will also involve developing an understanding of the privacy-accuracy trade-off in light of the requirements that apply to BOC and its historical use of SDL techniques.

## 4 Expected Outcomes

To summarize, the contributions of this project will include:

- An understanding of the relationship between mathematical and legal notions of privacy in the context of BOC work, and proposals to bridge the gap between them in ways that support BOC usage of formal privacy models.
- Formal privacy models for use in contexts where differential privacy is inappropriate or insufficient.
- New theoretical results on the performance of statistical analyses (including specific statistical analyses, query systems, and/or synthetic data) satisfying formal privacy requirements such as differential privacy (both upper and lower bounds). These would be aimed at analyses relevant to BOC work and would include, in particular, statistical inference tasks.
- New measures of statistical confidence incorporating the noise added for privacy.
- An understanding of the practical performance and usability of a variety of algorithms for analyzing and sharing privacy-sensitive data.

- A methodology for setting initial privacy parameters in formal privacy models.
- The transfer of information on current interdisciplinary research in privacy to BOC personnel and training of postdocs, and graduate students on issues related to BOC work.

The outcome of our work would be recorded in scientific and legal papers, memorandums, and reports. While we are optimistic that formal models of privacy will prove to be practical on some of the datasets collected by BOC, we note that most of the outcomes above are not dependent on this, and the lessons learned in trying to bring formal privacy models to practice will inform the next generation of privacy definitions and tools.

## 5 Timelines

The following timeline describes the anticipated activities at each stage of the project. Naturally, most aspects of the project will not be confined to a particular year, but the following reflects our expected emphases during different stages. In addition, we expect our work to generate new research questions, which may lead to some adjustments in our priorities over the course of the project.

<b>Year 1:</b>	<ul style="list-style-type: none"> <li>• Initiate engagement with BOC team.</li> <li>• Begin legal analysis of Title 13 and related policies and documents. We expect to have initial findings to support our work on bridging legal and mathematical definitions and developing law-based formal privacy models around mid-year.</li> <li>• Begin work on statistical inference for finite sample size, focusing on differential privacy. Map out limits of what can be achieved in terms of privacy-utility-efficiency tradeoffs, including design of new algorithms and theoretical bounds.</li> <li>• Begin research on confidence intervals for least-squares regression and other statistical inference tasks.</li> <li>• Survey and prioritize applicable laws and policies, past BOC reports, and accepted best practices relevant to informing the tuning of privacy parameters.</li> </ul>
<b>Year 2:</b>	<ul style="list-style-type: none"> <li>• Refine legal analysis of Title 13 as needed to support the other subgoals.</li> <li>• Begin formal legal-technical modeling of BOC’s privacy needs.</li> <li>• Continue work on the following, with the choice of privacy models depending on legal analysis and focusing on statistical utility most relevant for BOC data: <ul style="list-style-type: none"> <li>– Bridging mathematical and legal definitions of privacy.</li> <li>– Selecting initial privacy parameters.</li> <li>– Statistical inference procedures.</li> <li>– Measures of statistical confidence.</li> </ul> </li> <li>• Begin implementing and testing selected algorithms with BOC datasets.</li> <li>• Begin exploring the feasibility of interactive query models and synthetic data generation for BOC data releases.</li> </ul>
<b>Year 3:</b>	<ul style="list-style-type: none"> <li>• Iteratively refine the utility metrics and algorithmic tools according to the results of the experiments.</li> <li>• Iteratively refine statistical inference procedures to incorporate law-based formal models and tailoring to BOC’s needs.</li> <li>• Continue implementation and testing of algorithms.</li> <li>• Optimize methods for interactive queries and/or synthetic data generation.</li> </ul>

## Literature Cited

- [1] Confidential Information Protection and Statistical Efficiency Act, 44 U.S.C. § 3501 note.
- [2] E-Government Act of 2002, Pub. L. No. 107-347, 116 Stat. 2899.
- [3] Family Educational Rights and Privacy Act (FERPA), 20 U.S.C. § 1232g; 34 C.F.R. Part 99.
- [4] Health Insurance Portability and Accountability Act (HIPAA) Privacy Rule, 45 C.F.R. Part 160 and Subparts A and E of Part 164.
- [5] Privacy Act of 1974, 5 U.S.C. § 552a.
- [6] 13 U.S.C. et seq.
- [7] 13 U.S.C. § 8.
- [8] 13 U.S.C. § 9.
- [9] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang. Deep learning with differential privacy. arXiv:1607.00133 [stat.ML], July 2016.
- [10] J. T. Alexander, M. Davern, and B. Stevenson. Inaccurate age and sex data in the census pums files: Evidence and implications. Working Paper 15703, National Bureau of Economic Research, January 2010. URL <http://www.nber.org/papers/w15703>.
- [11] M. Altman, D. O’Brien, S. Vadhan, and A. Wood. Re: Big data study; request for information. Submitted to the White House Office of Science and Technology Policy (OSTP), March 2014. on behalf of the *Privacy Tools for Sharing Research Data Project*.
- [12] M. Altman, D. O’Brien, S. Vadhan, and A. Wood. Re: Proposed rule: Improve tracking of workplace injuries and illnesses; extension of comment period. Submitted to the Occupational Safety and Health Administration, March 2014. on behalf of the *Privacy Tools for Sharing Research Data Project*.
- [13] M. Altman, A. Wood, D. R. O’Brien, S. Vadhan, and U. Gasser. Towards a modern approach to a privacy-aware government data releases. *Berkeley Technology Law Journal*, 30(3), 2015.
- [14] L. Backstrom, C. Dwork, and J. M. Kleinberg. Wherefore art thou r3579x?: anonymized social networks, hidden patterns, and structural steganography. *Commun. ACM*, 54(12): 133–141, 2011. doi: 10.1145/2043174.2043199. URL <http://doi.acm.org/10.1145/2043174.2043199>.
- [15] B. Barak, K. Chaudhuri, C. Dwork, S. Kale, F. McSherry, and K. Talwar. Privacy, accuracy, and consistency too: a holistic solution to contingency table release. In *PODS*, pages 273–282, 2007.
- [16] G. Barthe, B. Köpf, F. Olmedo, and S. Z. Béguelin. Probabilistic relational reasoning for differential privacy. In J. Field and M. Hicks, editors, *Proceedings of the 39th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2012, Philadelphia, Pennsylvania, USA, January 22-28, 2012*, pages 97–110. ACM, 2012. ISBN 978-1-4503-1083-3. URL <http://dl.acm.org/citation.cfm?id=2103656>.
- [17] G. Barthe, M. Gaboardi, E. J. G. Arias, J. Hsu, C. Kunz, and P. Strub. Proving differential privacy in hoare logic. In *27<sup>th</sup> CSF’14 IEEE*, 2014. doi: 10.1109/CSF.2014.36.
- [18] R. Bassily, A. Smith, and A. Thakurta. Private empirical risk minimization, revisited. *CoRR*, abs/1405.7085, 2014. URL <http://arxiv.org/abs/1405.7085>.
- [19] R. Bassily, K. Nissim, A. D. Smith, T. Steinke, U. Stemmer, and J. Ullman. Algorithmic stability for adaptive data analysis. In D. Wichs and Y. Mansour, editors, *Proceedings of the 48th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2016, Cambridge, MA, USA, June 18-21, 2016*, pages 1046–1059. ACM, 2016. ISBN 978-1-4503-4132-5.

- doi: 10.1145/2897518.2897566. URL <http://doi.acm.org/10.1145/2897518.2897566>.
- [20] A. Beimel, S. P. Kasiviswanathan, and K. Nissim. Bounds on the sample complexity for private learning and private data release. In *TCC*, pages 437–454, 2010.
- [21] A. Beimel, K. Nissim, and U. Stemmer. Characterizing the sample complexity of private learners. In *ITCS*, pages 97–110, 2013.
- [22] A. Beimel, K. Nissim, and U. Stemmer. Private learning and sanitization: Pure vs. approximate differential privacy. In *APPROX-RANDOM’13*, pages 363–378, 2013.
- [23] A. Beimel, K. Nissim, and U. Stemmer. Characterizing the sample complexity of private learners. *CoRR*, abs/1402.2224, 2014. URL <http://arxiv.org/abs/1402.2224>.
- [24] A. Beimel, K. Nissim, and U. Stemmer. Learning privately with labeled and unlabeled examples. In *Proc. 26<sup>th</sup> SODA’15 ACM-SIAM*, pages 461–477, 2015. doi: 10.1137/1.9781611973730.32. URL <http://dx.doi.org/10.1137/1.9781611973730.32>.
- [25] A. Blum, C. Dwork, F. McSherry, and K. Nissim. Practical privacy: the SuLQ framework. In *PODS’05 ACM*, pages 128–138, 2005. doi: 10.1145/1065167.1065184. URL <http://doi.acm.org/10.1145/1065167.1065184>.
- [26] A. Blum, K. Ligett, and A. Roth. A learning theory approach to non-interactive database privacy. In *STOC*, pages 609–618, 2008.
- [27] M. Bun, J. Ullman, and S. Vadhan. Fingerprinting codes and the price of approximate differential privacy. In *Proceedings of the 46th Annual ACM Symposium on Theory of Computing (STOC ’14)*, pages 1–10, New York, NY, USA, 2014. ACM. doi: 10.1145/2591796.2591877. Full version posted as arXiv:1311.3158 [cs.CR]. Invited to *SIAM J. Computing* Special Issue on STOC ’14.
- [28] M. Bun, K. Nissim, U. Stemmer, and S. Vadhan. Differentially private release and learning of threshold functions. In *Proceedings of the 56th Annual IEEE Symposium on Foundations of Computer Science (FOCS ’15)*. IEEE, 18–20 October 2015. To appear. Full version posted as arXiv:1504.07553.
- [29] Center for Research on Computation and Society (CRCS). Privacy tools for sharing research data. <http://privacytools.seas.harvard.edu/>.
- [30] K. Chaudhuri and C. Monteleoni. Privacy-preserving logistic regression. In *NIPS*, 2008.
- [31] K. Chaudhuri, C. Monteleoni, and A. Sarwate. Differentially private empirical risk minimization. *J. Mach. Learn. Res.*, 12:1069–1109, July 2011. ISSN 1532-4435. URL <http://dl.acm.org/citation.cfm?id=1953048.2021036>.
- [32] K. Chaudhuri, A. D. Sarwate, and K. Sinha. Near-optimal differentially private principal components. In *Advances in Neural Information Processing Systems 25: 26th Annual Conference on Neural Information Processing Systems 2012. Proceedings of a meeting held December 3-6, 2012, Lake Tahoe, Nevada, United States.*, pages 998–1006, 2012. URL <http://papers.nips.cc/paper/4565-near-optimal-differentially-private-principal-components>.
- [33] Y. Chen, O. Sheffet, and S. P. Vadhan. Privacy games. In T. Liu, Q. Qi, and Y. Ye, editors, *Web and Internet Economics - 10th International Conference, WINE 2014, Beijing, China, December 14-17, 2014. Proceedings*, volume 8877 of *Lecture Notes in Computer Science*, pages 371–385. Springer, 2014. ISBN 978-3-319-13128-3. doi: 10.1007/978-3-319-13129-0\_30. URL [http://dx.doi.org/10.1007/978-3-319-13129-0\\_30](http://dx.doi.org/10.1007/978-3-319-13129-0_30).

- [34] C. Dimitrakakis, B. Nelson, A. Mitrokotsa, and B. Rubinstein. Robust and private bayesian inference. In *ALT 2014*, 2014.
- [35] I. Dinur and K. Nissim. Revealing information while preserving privacy. In *Proc. 22<sup>nd</sup> ACM PODS*, pages 202–210. Winner ACM PODS Alberto O. Mendelzon Test-of-Time Award., 2003. doi: 10.1145/773153.773173. URL <http://doi.acm.org/10.1145/773153.773173>.
- [36] I. Dinur and K. Nissim. Revealing information while preserving privacy. In *Proc. 22<sup>nd</sup> PODS ACM*, pages 202–210. ACM, 2003.
- [37] C. Dwork and J. Lei. Differential privacy and robust statistics. In *Proc. Forty-first Annual ACM Symposium on Theory of Computing*, pages 371–380, New York, NY, 2009. ACM. ISBN 978-1-60558-506-2. doi: 10.1145/1536414.1536466. URL <http://doi.acm.org/10.1145/1536414.1536466>.
- [38] C. Dwork and K. Nissim. Privacy-preserving datamining on vertically partitioned databases. In *Advances in Cryptology - CRYPTO'04*, pages 528–544, 2004. doi: 10.1007/978-3-540-28628-8\_32. URL [http://dx.doi.org/10.1007/978-3-540-28628-8\\_32](http://dx.doi.org/10.1007/978-3-540-28628-8_32).
- [39] C. Dwork and A. Roth. The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science*, 9(3-4):211–407, 2014. doi: 10.1561/0400000042. URL <http://dx.doi.org/10.1561/0400000042>.
- [40] C. Dwork and S. Yekhanin. New efficient attacks on statistical disclosure control mechanisms. In D. Wagner, editor, *Advances in Cryptology - CRYPTO 2008, 28th Annual International Cryptology Conference, Santa Barbara, CA, USA, August 17-21, 2008. Proceedings*, volume 5157 of *Lecture Notes in Computer Science*, pages 469–480. Springer, 2008. ISBN 978-3-540-85173-8. doi: 10.1007/978-3-540-85174-5\_26. URL [http://dx.doi.org/10.1007/978-3-540-85174-5\\_26](http://dx.doi.org/10.1007/978-3-540-85174-5_26).
- [41] C. Dwork, A. D. Smith, T. Steinke, and J. Ullman. Hiding in plain sight: A survey of attacks on private data. Forthcoming (2017).
- [42] C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private data analysis. In *3<sup>rd</sup> Theory of Crypt. Conf.*, pages 265–284, 2006. doi: 10.1007/11681878\_14. URL [http://dx.doi.org/10.1007/11681878\\_14](http://dx.doi.org/10.1007/11681878_14).
- [43] C. Dwork, F. McSherry, and K. Talwar. The price of privacy and the limits of LP decoding. In *Proc. 39<sup>th</sup> STOC'07 ACM*, pages 85–94, 2007. doi: 10.1145/1250790.1250804. URL <http://doi.acm.org/10.1145/1250790.1250804>.
- [44] C. Dwork, M. Naor, O. Reingold, G. Rothblum, and S. Vadhan. On the complexity of differentially private data release: Efficient algorithms and hardness results. In *Proceedings of the 41st Annual ACM Symposium on Theory of Computing (STOC '09)*, pages 381–390, 31 May–2 June 2009.
- [45] C. Dwork, G. N. Rothblum, and S. P. Vadhan. Boosting and differential privacy. In *FOCS*, pages 51–60, 2010.
- [46] C. Dwork, V. Feldman, M. Hardt, T. Pitassi, O. Reingold, and A. L. Roth. Preserving statistical validity in adaptive data analysis. In R. A. Servedio and R. Rubinfeld, editors, *Proceedings of the Forty-Seventh Annual ACM on Symposium on Theory of Computing, STOC 2015, Portland, OR, USA, June 14-17, 2015*, pages 117–126. ACM, 2015. ISBN 978-1-4503-3536-2. doi: 10.1145/2746539.2746580. URL <http://doi.acm.org/10.1145/2746539.2746580>.
- [47] C. Dwork, A. Smith, T. Steinke, J. Ullman, and S. Vadhan. Robust traceability from trace

- amounts. In *Proceedings of the 56th Annual IEEE Symposium on Foundations of Computer Science (FOCS '15)*. IEEE, 18–20 October 2015. To appear.
- [48] C. Dwork, A. D. Smith, T. Steinke, and J. Ullman. Hiding in plain sight: A survey of attacks on private data. *Annual Review of Statistics and Its Application*, 4(1), 2016.
- [49] Ú. Erlingsson, V. Pihur, and A. Korolova. RAPPOR: randomized aggregatable privacy-preserving ordinal response. In G. Ahn, M. Yung, and N. Li, editors, *Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security, Scottsdale, AZ, USA, November 3-7, 2014*, pages 1054–1067. ACM, 2014. ISBN 978-1-4503-2957-6. doi: 10.1145/2660267.2660348. URL <http://doi.acm.org/10.1145/2660267.2660348>.
- [50] Federal Committee on Statistical Methodology, Office of Management and Budget. Statistical Policy Working Paper 22. Second Version, 2005. URL [http://www.fcsfm.gov/working-papers/SPWP22\\_rev.pdf](http://www.fcsfm.gov/working-papers/SPWP22_rev.pdf).
- [51] M. Gaboardi, A. Haeberlen, J. Hsu, A. Narayan, and B. C. Pierce. Linear dependent types for differential privacy. In *40<sup>th</sup> POPL'13 ACM*, pages 357–370, 2013. doi: 10.1145/2429069.2429113.
- [52] M. Gaboardi, E. J. G. Arias, J. Hsu, A. Roth, and Z. S. Wu. Dual query: Practical private query release for high dimensional data. In *31<sup>th</sup> ICML 2014*, pages 1170–1178, 2014. URL <http://jmlr.org/proceedings/papers/v32/gaboardi14.html>.
- [53] M. Gaboardi, J. Honaker, G. King, K. Nissim, J. Ullman, and S. Vadhan.  $\psi$  - a private data-sharing interface. Poster at 2nd Theory and Practice of Differential Privacy Workshop (TPDP 2016), 2016. URL <http://privacytools.seas.harvard.edu/publications/psipaper>.
- [54] M. Gaboardi, H. Lim, R. M. Rogers, and S. P. Vadhan. Differentially private chi-squared hypothesis testing: Goodness of fit and independence testing. In M. Balcan and K. Q. Weinberger, editors, *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016*, volume 48 of *JMLR Workshop and Conference Proceedings*, pages 2111–2120. JMLR.org, 2016. URL <http://jmlr.org/proceedings/papers/v48/rogers16.html>.
- [55] A. Ghosh and A. Roth. Selling privacy at auction. In *ACM EC'11*, pages 199–208, 2011. doi: 10.1145/1993574.1993605. URL <http://doi.acm.org/10.1145/1993574.1993605>.
- [56] A. Ghosh, T. Roughgarden, and M. Sundarajan. Universally utility-maximizing privacy mechanisms. Manuscript, November 2008.
- [57] A. Haeberlen, B. C. Pierce, and A. Narayan. Differential privacy under fire. In *Proceedings of the 20th USENIX Security Symposium*, Aug. 2011.
- [58] S. Haney, A. Machanavajjhala, J. M. Abowd, M. Graham, M. Kutzbach, and L. Vilhuber. The cost of provable privacy: A case study on linked employer-employee data. *Presented in Theory and Practice of Differential Privacy (TPDP)*, 2016.
- [59] M. Hardt and G. N. Rothblum. A multiplicative weights mechanism for privacy-preserving data analysis. In *FOCS*, pages 61–70, 2010.
- [60] M. Hardt, K. Ligett, and F. McSherry. A simple and practical algorithm for differentially private data release. In *NIPS*, pages 2348–2356, 2012.
- [61] M. Hay, V. Rastogi, G. Miklau, and D. Suci. Boosting the Accuracy of Differentially Private Histograms Through Consistency. *PVLDB*, 3(1):1021–1032, 2010.
- [62] N. Homer, S. Szlinger, M. Redman, D. Duggan, W. Tembe, J. Muehling, J. V. Pearson,

- D. A. Stephan, S. F. Nelson, and D. W. Craig. Resolving individuals contributing trace amounts of dna to highly complex mixtures using high-density snp genotyping microarrays. *PLoS Genetics*, 4(8), 2008.
- [63] J. Hsu, M. Gaboardi, A. Haeberlen, S. Khanna, A. Narayan, B. C. Pierce, and A. Roth. Differential privacy: An economic method for choosing epsilon. In *IEEE 27th Computer Security Foundations Symposium, CSF 2014, Vienna, Austria, 19-22 July, 2014*, pages 398–410. IEEE Computer Society, 2014. ISBN 978-1-4799-4290-9. doi: 10.1109/CSF.2014.35. URL <http://dx.doi.org/10.1109/CSF.2014.35>.
- [64] P. Kairouz, S. Oh, and P. Viswanath. The composition theorem for differential privacy. In F. R. Bach and D. M. Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, volume 37 of *JMLR Proceedings*, pages 1376–1385. JMLR.org, 2015. URL <http://jmlr.org/proceedings/papers/v37/kairouz15.html>.
- [65] A. Kantor and K. Nissim. Attacks on statistical databases: The highly noisy case. *Inf. Process. Lett.*, 113(12):409–413, 2013. doi: 10.1016/j.ipl.2013.03.005. URL <http://dx.doi.org/10.1016/j.ipl.2013.03.005>.
- [66] V. Karwa and S. Vadhan. Nearly optimal differentially private confidence intervals for a normal mean. In preparation, August 2016.
- [67] V. Karwa, D. Kifer, and A. B. Slavkovic. Private posterior distributions from variational approximations. *CoRR*, abs/1511.07896, 2015. URL <http://arxiv.org/abs/1511.07896>.
- [68] S. P. Kasiviswanathan and A. Smith. On the ‘semantics’ of differential privacy: A bayesian formulation. *Journal of Privacy and Confidentiality*, 6(1), 2014.
- [69] S. P. Kasiviswanathan, H. K. Lee, K. Nissim, S. Raskhodnikova, and A. Smith. What can we learn privately? *SIAM J. Comput.*, 40(3):793–826, 2011.
- [70] S. P. Kasiviswanathan, K. Nissim, S. Raskhodnikova, and A. D. Smith. Analyzing graphs with node differential privacy. In *TCC*, pages 457–476, 2013. doi: 10.1007/978-3-642-36594-2\_26. URL [http://dx.doi.org/10.1007/978-3-642-36594-2\\_26](http://dx.doi.org/10.1007/978-3-642-36594-2_26).
- [71] D. Kifer, A. D. Smith, and A. Thakurta. Private convex optimization for empirical risk minimization with applications to high-dimensional regression. In *25<sup>th</sup> COLT*, 2012. URL <http://www.jmlr.org/proceedings/papers/v23/kifer12/kifer12.pdf>.
- [72] J. Lee and C. Clifton. Differential identifiability. In *The 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '12, Beijing, China, August 12-16, 2012*, pages 1041–1049, 2012. doi: 10.1145/2339530.2339695. URL <http://doi.acm.org/10.1145/2339530.2339695>.
- [73] J. Lee, Y. Wang, and D. Kifer. Maximum likelihood postprocessing for differential privacy under consistency constraints. In L. Cao, C. Zhang, T. Joachims, G. I. Webb, D. D. Margineantu, and G. Williams, editors, *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Sydney, NSW, Australia, August 10-13, 2015*, pages 635–644. ACM, 2015. ISBN 978-1-4503-3664-2. doi: 10.1145/2783258.2783366. URL <http://doi.acm.org/10.1145/2783258.2783366>.
- [74] J. Lei, A.-S. Charest, A. Slavkovic, A. Smith, and S. Fienberg. Differentially private model selection with penalized and constrained likelihood. arXiv:1607.04204 [stat.ME], July 2016.

- [75] B. Lin and D. Kifer. Information preservation in statistical privacy and bayesian estimation of unattributed histograms. In K. A. Ross, D. Srivastava, and D. Papadias, editors, *Proceedings of the ACM SIGMOD International Conference on Management of Data, SIGMOD 2013, New York, NY, USA, June 22-27, 2013*, pages 677–688. ACM, 2013. ISBN 978-1-4503-2037-5. doi: 10.1145/2463676.2463721. URL <http://doi.acm.org/10.1145/2463676.2463721>.
- [76] A. Machanavajjhala, D. Kifer, J. Abowd, J. Gehrke, and L. Vilhuber. Privacy Theory meets practice on the map. Proceedings of the IEEE 24th International Conference on Data Engineering, 2008. Available at [http://lehd.ces.census.gov/doc/help/ICDE08\\_conference\\_0768.pdf](http://lehd.ces.census.gov/doc/help/ICDE08_conference_0768.pdf).
- [77] F. McSherry. Privacy integrated queries: an extensible platform for privacy-preserving data analysis. In U. Çetintemel, S. B. Zdonik, D. Kossmann, and N. Tatbul, editors, *Proceedings of the ACM SIGMOD International Conference on Management of Data, SIGMOD 2009, Providence, Rhode Island, USA, June 29 - July 2, 2009*, pages 19–30. ACM, 2009. ISBN 978-1-60558-551-2. doi: 10.1145/1559845.1559850. URL <http://doi.acm.org/10.1145/1559845.1559850>.
- [78] F. McSherry and K. Talwar. Mechanism design via differential privacy. In *FOCS'07 IEEE*, pages 94–103, 2007. ISBN 0-7695-3010-9. doi: 10.1109/FOCS.2007.41. URL <http://dx.doi.org/10.1109/FOCS.2007.41>.
- [79] D. J. Mir, S. Isaacman, R. Cáceres, M. Martonosi, and R. N. Wright. DP-WHERE: differentially private modeling of human mobility. In X. Hu, T. Y. Lin, V. Raghavan, B. W. Wah, R. A. Baeza-Yates, G. Fox, C. Shahabi, M. Smith, Q. Yang, R. Ghani, W. Fan, R. Lempel, and R. Nambiar, editors, *Proceedings of the 2013 IEEE International Conference on Big Data, 6-9 October 2013, Santa Clara, CA, USA*, pages 580–588. IEEE, 2013. ISBN 978-1-4799-1292-6. doi: 10.1109/BigData.2013.6691626. URL <http://dx.doi.org/10.1109/BigData.2013.6691626>.
- [80] P. Mohan, A. Thakurta, E. Shi, D. Song, and D. E. Culler. GUPT: privacy preserving data analysis made easy. In K. S. Candan, Y. Chen, R. T. Snodgrass, L. Gravano, and A. Fuxman, editors, *Proceedings of the ACM SIGMOD International Conference on Management of Data, SIGMOD 2012, Scottsdale, AZ, USA, May 20-24, 2012*, pages 349–360. ACM, 2012. ISBN 978-1-4503-1247-9. doi: 10.1145/2213836.2213876. URL <http://doi.acm.org/10.1145/2213836.2213876>.
- [81] J. Murtagh and S. Vadhan. The complexity of computing the optimal composition of differential privacy. In E. Kushilevitz and T. Malkin, editors, *Proceedings of the 13th IACR Theory of Cryptography Conference (TCC '16-A)*, volume 9562 of *Lecture Notes in Computer Science*, pages 157–175. Springer-Verlag, 10–13 January 2016. ISBN 978-3-662-49095-2. doi: 10.1007/978-3-662-49096-9. URL <http://dx.doi.org/10.1007/978-3-662-49096-9>. Full version posted on *CoRR*, abs/1507.03113, July 2015.
- [82] A. Narayanan and V. Shmatikov. How to break anonymity of the netflix prize dataset. *CoRR*, abs/cs/0610105, 2006. URL <http://arxiv.org/abs/cs/0610105>.
- [83] A. Narayanan and V. Shmatikov. Robust de-anonymization of large sparse datasets. In *Security and Privacy*, pages 111–125. IEEE, 2008.
- [84] A. Narayanan and V. Shmatikov. De-anonymizing social networks. In *30th IEEE Symposium on Security and Privacy (S&P 2009), 17-20 May 2009, Oakland, California, USA*, pages 173–187, 2009. doi: 10.1109/SP.2009.22. URL <http://dx.doi.org/10.1109/SP.2009.22>.

- [85] A. Nikolov, K. Talwar, and L. Zhang. The geometry of differential privacy: the sparse and approximate cases. In *Proc. 45<sup>th</sup> STOC'13 ACM*, pages 351–360, 2013.
- [86] K. Nissim, A. Bembenek, A. Wood, M. Bun, M. Gaboardi, U. Gasser, D. R. O'Brien, T. Steinke, and S. Vadhan. Bridging the gap between computer science and legal approaches to privacy. URL <http://privacytools.seas.harvard.edu/publications/bridging-gap-between-computer-science-and-legal-approaches-privacy>. Working Paper (2016).
- [87] K. Nissim, S. Raskhodnikova, and A. Smith. Smooth sensitivity and sampling in private data analysis. In *Proc. 39<sup>th</sup> STOC'07 ACM*, 2007. doi: 10.1145/1250790.1250803. URL <http://doi.acm.org/10.1145/1250790.1250803>.
- [88] K. Nissim, C. Orlandi, and R. Smorodinsky. Privacy-aware mechanism design. In *ACM EC'12*, pages 774–789, 2012. doi: 10.1145/2229012.2229073. URL <http://doi.acm.org/10.1145/2229012.2229073>.
- [89] K. Nissim, S. P. Vadhan, and D. Xiao. Redrawing the boundaries on purchasing data from privacy-sensitive individuals. In *ITCS'14*, pages 411–422, 2014. doi: 10.1145/2554797.2554835. URL <http://doi.acm.org/10.1145/2554797.2554835>.
- [90] K. Nissim, U. Stemmer, and S. P. Vadhan. Locating a small cluster privately. In T. Milo and W. Tan, editors, *Proceedings of the 35th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, PODS 2016, San Francisco, CA, USA, June 26 - July 01, 2016*, pages 413–427. ACM, 2016. ISBN 978-1-4503-4191-2. doi: 10.1145/2902251.2902296. URL <http://doi.acm.org/10.1145/2902251.2902296>.
- [91] D. R. O'Brien, J. Ullman, M. Altman, U. Gasser, M. Bar-Sinai, K. Nissim, S. Vadhan, M. J. Wojcik, and A. Wood. Integrating approaches to privacy across the research lifecycle: When is information purely public? *Berkman Center Research Publication*, (2015-7), 2015.
- [92] P. Ohm. Broken promises of privacy: Responding to the surprising failure of anonymization. *UCLA Law Review*, 57:1701, 2010.
- [93] R. Rogers, A. Roth, A. Smith, and O. Thakkar. Max-information, differential privacy, and post-selection hypothesis testing. In *Foundations of Computer Science (FOCS)*, 2016. arXiv:1604.03924 [cs.LG].
- [94] I. Roy, S. T. V. Setty, A. Kilzer, V. Shmatikov, and E. Witchel. Airavat: Security and privacy for mapreduce. In *Proceedings of the 7th USENIX Symposium on Networked Systems Design and Implementation, NSDI 2010, April 28-30, 2010, San Jose, CA, USA*, pages 297–312. USENIX Association, 2010. ISBN 978-1-931971-73-7. URL [http://www.usenix.org/events/nsdi10/tech/full\\_papers/roy.pdf](http://www.usenix.org/events/nsdi10/tech/full_papers/roy.pdf).
- [95] B. I. P. Rubinstein, P. L. Bartlett, L. Huang, and N. Taft. Learning in a large function space: Privacy-preserving mechanisms for svm learning. *CoRR*, abs/0911.5708, 2009.
- [96] P. Schwartz and D. Solove. The PII problem: Privacy and a new concept of personally identifiable information. *New York University Law Review*, 86(4):1814–1895, 2011.
- [97] A. Smith. Privacy-preserving statistical estimation with optimal convergence rates. In *STOC*, 2011.
- [98] A. Smith and A. Thakurta. Differentially private feature selection via stability arguments, and the robustness of the lasso. In *COLT 2013 - The 26th Annual Conference on Learning Theory, June 12-14, 2013, Princeton University, NJ, USA*, pages 819–850, 2013. URL <http://jmlr.org/proceedings/papers/v30/Guha13.html>.
- [99] A. Smith and A. G. Thakurta. (nearly) optimal algorithms for private online learning in full-

- information and bandit settings. In *Advances in Neural Information Processing Systems*, pages 2733–2741, 2013.
- [100] L. Sweeney. Weaving technology and policy together to maintain confidentiality. *The Journal of Law, Medicine & Ethics*, 25(2-3):98–110, 1997.
- [101] L. Sweeney. Uniqueness of simple demographics in the us population. Technical report, Technical report, Carnegie Mellon University, 2000.
- [102] K. Talwar, A. Thakurta, and L. Zhang. Nearly optimal private LASSO. In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015*, pages 3025–3033, 2015.
- [103] C. Uhler, A. Slavkovic, and S. E. Fienberg. Privacy-preserving data sharing for genome-wide association studies. *Journal of Privacy and Confidentiality*, 5(1), 2013.
- [104] J. Ullman and S. P. Vadhan. PCPs and the hardness of generating private synthetic data. In *TCC*, pages 400–416, 2011.
- [105] S. Vadhan. The complexity of differential privacy. Unpublished survey, 2016.
- [106] E. Vayena, U. Gasser, A. Wood, D. R. O’Brien, and M. Altman. Elements of a new ethical framework for big data research. *Washington and Lee Law Review Online*, 72(420), 2016.
- [107] Y.-X. Wang, S. Fienberg, , and A. Smola. Privacy for free: Posterior sampling and stochastic gradient monte carlo. In *ICML 2015*, pages 2493—2502, 2015.
- [108] L. Wasserman and S. Zhou. A statistical framework for differential privacy. *Journal of the American Statistical Association*, 105(489):375–389, 2010. ISSN 01621459. URL <http://www.jstor.org/stable/29747034>.
- [109] O. Williams and F. McSherry. Probabilistic inference and differential privacy. In *NIPS*, 2010.
- [110] A. Wood, D. O’Brien, M. Altman, A. Karr, U. Gasser, M. Bar-Sinai, K. Nissim, J. Ullman, S. Vadhan, and M. J. Wojcik. Integrating approaches to privacy across the research lifecycle: Long-term longitudinal studies. *Berkman Center Research Publication*, (2014-12), 2014.
- [111] A. Wood, E. Airoidi, M. Altman, Y.-A. de Montjoye, U. Gasser, D. O’Brien, and S. Vadhan. Re: Federal policy for the protection of human subjects proposed rules (hhsophs20150008). Submitted to the Office for Human Research Protections, Department of Health and Human Services, January 2016. on behalf of the *Privacy Tools for Sharing Research Data Project*.
- [112] F. T. Wu. Defining privacy and utility in data sets. *University of Colorado Law Review*, 84: 1117–1177, 2013.
- [113] D. Xiao. Is privacy compatible with truthfulness? *IACR Cryptology ePrint Archive*, page 5, 2011. URL <http://eprint.iacr.org/2011/005>.
- [114] D. Xiao. Is privacy compatible with truthfulness? In *ITCS’13*, pages 67–86, 2013. doi: 10.1145/2422436.2422448. URL <http://doi.acm.org/10.1145/2422436.2422448>.
- [115] J. Zhang, G. Cormode, C. M. Procopiuc, D. Srivastava, and X. Xiao. Privbayses: Private data release via bayesian networks. In *ACM SIGMOD ’14*, pages 1423–1434, New York, NY, USA, 2014. ACM. ISBN 978-1-4503-2376-5. doi: 10.1145/2588555.2588573. URL <http://doi.acm.org/10.1145/2588555.2588573>.
- [116] J. Zheng. The differential privacy of bayesian inference. Undergraduate thesis, April 2015. <https://dash.harvard.edu/handle/1/14398533>.

# Data Sharing Plan

The PIs have a proven track record of releasing in a timely manner as much material (scientific and legal publications, memoranda, code, and data) in ways that are visible, accessible, and understandable by relevant scientific or legal communities. They will continue this policy of actively sharing information, within the constraints of feasibility and the law.

**Typical Data Types.** All data used in this research are already public-use data available from the Bureau of Census and other public data sites.

**Other Data.** We do not expect other data to be produced. In case such data will be generated, it will be archived. We will consider in particular using the Dataverse Project, an open source data repository developed by the Institute for Quantitative Social Science, and with an instance at Harvard. The Dataverse Project has already provided a Data Sharing solution for many scholars by providing easy access to data, standardized metadata, and archival and preservation support.