# Privacy: From Database Reconstruction to Legal Theorems

Kobbi Nissim
kobbi.nissim@georgetown.edu
Georgetown University
Washington, D.C., USA

## ABSTRACT

There are significant gaps between legal and technical thinking around data privacy. Technical standards are described using mathematical language whereas legal standards are not rigorous from a mathematical point of view and often resort to concepts which they only partially define. As a result, arguments about the adequacy of technical privacy measures for satisfying legal privacy often lack rigor, and their conclusions are uncertain. The uncertainty is exacerbated by a litany of successful privacy attacks on privacy measures thought to meet legal expectations but then shown to fall short of doing so.

As computer systems manipulating individual privacy-sensitive data become integrated in almost every aspect of society, and as such systems increasingly make decisions of legal significance, the need to bridge the diverging, and sometimes conflicting legal and technical approaches becomes urgent.

We formulate and prove formal claims – "legal theorems" – addressing legal questions such as whether the use of technological measures satisfies the requirements of a legal privacy standard. In particular, we analyze the notion of singling out from the GDPR and whether technologies such as $k$-anonymity and differential privacy prevent singling out.

Our long-term goal is to develop concepts which are on one hand technical, so they can be integrated in the design of computer systems, and can be used in legal reasoning and for policymaking on the other hand.

## CCS CONCEPTS

• **Security and privacy** → **Privacy protections**.

## KEYWORDS

data privacy; GDPR; singling out; differential privacy; $k$-anonymity

## 1 INTRODUCTION

A 2010 article by Ohm drew a grim picture of how advances in the study of privacy by computer scientists have "undermined our faith in the privacy-protecting power of anonymization". Ohm referred to the belief that the practice of redacting personal identifiers and other context dependent identifying information from personally identifiable data provides a "best-of-both-worlds" compromise where "Analysts will still find the data useful, but unscrupulous marketers and malevolent identity thieves will find it impossible to identify the people tracked". The promise that "[a]nonymization ensures privacy" so "[s]ociety will be able to turn its collective attention to other problems because technology will have solved this one" was broken [36]. Central to Ohm's analysis were cases where supposedly anonymized data was very effectively re-identified, such as was demonstrated by Sweeney [41] and Narayanan and Shmatikov [32, 33].

It is illustrative to recall what led to the revelation that expectations of privacy were not met and promises were broken. Sweeney has re-identified very sensitive personal data, the medical records of state employees in Massachusetts' Group Insurance Commission (GIC). Publishers of the GIC data were aware of its sensitivity, and hence redacted fields containing directly identifying information such as patient names, addresses, and social security numbers before putting it in the public sphere. It turned out that this was not enough for keeping the published records anonymous.

At the heart of Sweeney's re-identification attack was the crucial observation that the seemingly innocuous combination of ZIP code, birth date, and sex – attributes which were not redacted from the GIC data – is unique for a vast majority of the US population. In her attack, Sweeney used such triples as *quasi-identifiers* which she matched with information available in public records – the Cambridge MA Voter Vegistration – hence linking patients' medical records in the GIC data with directly identifying information.

Narayanan and Shmatikov demonstrated that records containing user's ratings of movies which were released by Netflix for its "Netflix Prize" could be re-identified. Movie viewing information can reflect sensitive personal characteristics such as a person's sexual orientation. To protect its customers data, Netflix has applied disclosure limitation procedures to anonymize the dataset. Nevertheless, Narayanan and Shmatikov found that, analogously to the combination of ZIP code, birth date, and sex, the movies rated by a subscriber and the approximate times of their rating often makes the subscriber unique in the dataset. Even more so, little partial knowledge about a subscriber's viewings and ratings, when matched with publicly available movie ratings from the Internet Movie Database (IMDb), can lead to the exact re-identification of the subscriber (or to a small number of candidate identities, one of which is correct).

On the theoretical side, Dinur and Nissim presented in 2003 a new way for promises of privacy to be broken (at least in theory) – reconstruction attacks [16].[1] Consider a dataset over the binary data domain $X = \{0, 1\}$,

$$\mathbf{x} = (x_1, \ldots, x_n) \in \{0, 1\}^n,$$

containing information whether each of $n$ individuals has a sensitive trait, e.g., $x_i = 1$ if person $i$ is diabetic and $x_i = 0$ otherwise. Access to the dataset $\mathbf{x}$ is provided to an analyst via a mechanism answering statistical queries. If the analyst wishes to learn the number of diabetic individuals in a specific sub-population $q \subseteq [n]$, the analyst can issue the query $q \subseteq [n]$ and the mechanism would return as answer an estimate

$$a_q \approx \sum_{i \in q} x_i.$$

What if a privacy attacker has access to the mechanism? We say that a mechanism is *blatantly non-private* [16] if an attacker having query access to the mechanism can *reconstruct* a very accurate rendering of $\mathbf{x}$, say, $\tilde{\mathbf{x}} \in \{0, 1\}$ which agrees with $\mathbf{x}$ on all but at most 5% of the entries. It turns out that such reconstruction is possible unless either the mechanism introduces sufficiently large error in its answers or it limits the number of queries asked (or both). Assume query answers are guaranteed to be within error $\alpha$, i.e., for all queries $q \subseteq [n]$ we have that $|a_q - \sum_{i \in q} x_i| \leq \alpha$.

THEOREM 1.1 ([16] INFORMAL). *There exist constants $c, c' > 0$ such that reconstruction is possible in the following settings: (i) $\alpha = cn$ and an attacker can make all $2^n$ possible subset queries , or (ii) $\alpha = c'\sqrt{n}$ and an attacker can make polynomially many queries.*

These results have been strengthened and extended in a number of works, e.g., [18, 21, 31]. They demonstrate what was coined by Dwork and Roth as the *Fundamental Law of Information Recovery* [19]:

> overly accurate answers to too many questions will destroy privacy in a spectacular way.

Re-identification attacks and reconstruction attacks are just two examples out of a fast growing litany of failures to meet expectations of privacy and we mention a few other notable cases.[2] Backstrom, Dwork, and Kleinberg extended re-identification to the setting of social graphs [10]; Homer et al. introduced membership attacks on aggregate genomic data, allowing to infer whether a person's data was included in the aggregate [26]; Shokri, Stronati, Song, and Shmatikov developed membership attacks against machine learning models, allowing to infer whether a person's data was included in the training set [40]; and Carlini, Liu, Erlingsson, Kos, and Song identified that inadvertent memorization of training data can lead to the revealing of secret personal information, such as the exposure of a person's Social Security Number as an auto-complete for the sentence "my social-security number is . . ." [11].

Furthermore, reconstruction attacks, initially viewed as primarily a theoretical tool motivating and helping the development of

formal privacy models, have proved effective when applied to a commercial system [13] and, more importantly, when applied to the statistical tables published by the US Census Bureau following the 2010 Decennial Census [24]. The reconstruction of the 2010 Decennial data yielded exact sex, race, ethnicity, and location to the block-level and age up to one year difference for 71% of the US population. Furthermore, via matching with commercial databases that were available in 2010, records were accurately reconstructed and re-identified for 52 million people (17% of the US population) [7]. These numbers are very different from what the US Census expected in 2010. For comparison, the prior re-identification risk estimate for the 2010 Census release was merely 0.003% [8], i.e., lower by a factor of about 4500. Note that Title 13 of the US code, which outlines the role of the US Census Bureau, prohibits the Census Bureau from making "any publication whereby the data furnished by any particular establishment or individual under this title can be identified" [1].

## 1.1 New technical privacy concepts

The examples above demonstrate how some of the commonly used methods for ensuring privacy can fail. From a technical computer science point of view, a possible remedy for this failure is via the introduction of new paradigms for ensuring privacy, two of which we describe here briefly.

*k-anonymity*, introduced by Samarati and Sweeney [37, 42], is a restriction on the output of an anonymization process.[3] In this framework, a dataset $\mathbf{x}$ is anonymized via the application of suppression and generalization of potentially identifying attributes in the input dataset towards preventing the possibility that any of the records would be uniquely matched with a publicly available directly identifying dataset.[4] The $k$-anonymizer produces a dataset $\mathbf{x}'$ which is identical to $\mathbf{x}$ except for the suppressed or generalized attributes subject to the requirement that in $\mathbf{x}'$ every record is identical to at least $k$-1 other records.

As a toy example, see below a dataset $\mathbf{x}$ with four records (on the left) and a 2-anonymized version $\mathbf{x}'$ of the same dataset (on the right), The derived from $\mathbf{x}$ via suppression and hierarchical generalization operations.

| ZIP | Age | Sex | Disease | ZIP | Age | Sex | Disease |
|-------|-----|-----|---------|-------|-------|-----|---------|
| 23456 | 55 | F | COVID | 23456 | * | F | COVID |
| 23456 | 42 | F | COVID | 23456 | * | F | COVID |
| 12345 | 30 | M | CF | 1234* | 30-39 | * | PULM |
| 12346 | 33 | F | Asthma | 1234* | 30-39 | * | PULM |

Minimizing the number of suppressed attributes is an NP-hard problem [30] and a rich algorithmic literature exists, providing methods for $k$-anonymizing datasets while (approximately) maximizing some measure of information content in the anonymized dataset.

It is important to note that while the syntactic restriction imposed by $k$-anonymity prevents the type of attacks performed by Sweeney on the GIC data via finding unique matching with identifiable public data, this restriction is *syntactic* and does not imply

---

[1]We use the following notation throughout: A dataset $\mathbf{x}$ containing personal information of $n$ individuals is a vector containing $n$ records from a data domain $X$, i.e., $\mathbf{x} = (x_1, \ldots, x_n) \in X^n$. Each record $x_i$ in $\mathbf{x}$ corresponds to the personal information of a single individual $i \in [n]$.

[2]See also the survey of attacks on private data complied by Dwork, Smith, Steinke, and Ullman [20].

[3]The analysis of $k$-anonymity throughout also holds for variants of $k$-anonymity such as $\ell$-diversity [29] and $t$-closeness [28].

[4]Generalization is typically done in a hierarchical manner, e.g., by suppressing the last digit(s) of a ZIP code or replacing a geographic unit (e.g., a person's county of residence) with a coarser geographic unit (the person's state of residence).

that a $k$-anonymized dataset cannot be post-processed so as to infer personal data. In particular, in a $k$-anonymizer decisions on suppression and generalization are typically data dependent, and hence the pattern of suppression and generalization in the outcome of a $k$-anonymizer potentially leaks information which a privacy attacker can make use of. In this vein, a recent result of Cohen provides a reconstruction attack of generalization-based $k$-anonymized datasets. Cohen's attack applies even in cases where every attribute combination is treated as a potential quasi-identifier. The attack relies on knowledge of the underlying distribution but does not require the attacker to consult any other dataset beyond the $k$-anonymized dataset, and hence is purely post-processing [12].

Furthermore, $k$-anonymity is not closed under composition, i.e., it may well be that the combination of two or more $k$-anonymized datasets derived from the same (or similar) collection of personal information allows for uniquely identifying individuals in the data [12, 23].

*Differential privacy*, introduced by Dwork, McSherry, Nissim, and Smith in 2006 [17], is a definition of privacy capturing the desiderata that any information-related risk to a person should not change significantly as a result of that person's information being included, or not, in an analysis. Informally, an analysis (referred to in the literature as a mechanism) satisfies differential privacy if its outcome distribution is insensitive to any change in a single individual's record. More formally, let $X$ be an arbitrary data domain, and $Y$ be an arbitrary output domain and let $\varepsilon > 0$.

*Definition 1.2 ([17]).* A mechanism $M : X^n \rightarrow Y$ is $\varepsilon$-differentially private if for all datasets $\mathbf{x}, \mathbf{x}' \in X^n$ which differ on a single entry and for all events $T \subseteq Y$,

$$\Pr[M(\mathbf{x}) \in T] \leq e^{\varepsilon} \Pr[M(\mathbf{x}') \in T],$$

where the probability is over the randomness of the mechanism $M$.

A simple example of a differentially private mechanism is the Laplace mechanism for counting. This mechanism first counts the number of individuals in the databse with a sensitive trait (e.g., the number of individuals suffering diabetes) and then introduces statistical noise to the outcome, sampled from the Laplace distribution.

THEOREM 1.3 (LAPLACE MECHANISM [17]). *Let $M_{Lap}$ be the mechanism which on input $\mathbf{x} \in \{0, 1\}^n$ outputs*

$$\sum_{i=1}^{n} x_i + Y$$

*where $Y \sim Lap(1/\varepsilon)$.[5] Then, $M_{Lap}$ is $\varepsilon$-differentially private.*

Besides count queries, differentially private computations were developed for a large variety of tasks, including the computation of statistical estimates and machine learning.

Note that, in contrast with $k$-anonymity, differential privacy directly limits the information content in the outcome of the mechanism so as to limit the ability of any attacker to distinguish whether an individual's records holds one value or another.

The privacy loss parameter $\varepsilon$ (also referred to as the "privacy budget") quantifies and bounds the excessive risk to an individual due to her participation in a differentially private analysis. A lower value of $\varepsilon$ corresponds to a better privacy guarantee, but also restricts the utility that can be obtained from the data (e.g., in terms of accuracy of the analysis).

Two important properties of differential privacy are that post processing of the outcome of a differentially private analysis does not break the differential privacy guarantee – privacy loss cannot increase due to post processing. Furthermore, differential privacy is closed under composition. i.e., the result of applying two or more differentially private analysis on the same or related data preserves differential privacy (albeit with worse privacy loss parameter than that of each individual analysis).

Differential privacy provides a theoretical framework for reasoning about privacy in analysis of personal information, based on well defined concepts, well formed statements, and rigorous proofs, following a strong tradition established in Foundations of Cryptography. It is making significant steps towards being used as a standard for the processing of sensitive personal data in academia, industry, and government agencies.[6]

Potentially, $k$-anonymity and differential privacy improve over the state of affairs described in the introduction. However, technical concepts as well as mathematical rigor can only go a certain mileage towards ensuring that expectations of privacy are being met. It is important to keep in mind that privacy is not merely a technical concept. It is also a normative concept rooted in a variety of traditions, including philosophy, sociology, and economics. In particular, expectations of privacy are coded in numerous laws and regulations, raising the question of how to argue that the use of specific technologies meets these legal standards?

## 1.2 Legal privacy concepts

We provide a very brief review of legal privacy concepts which are most relevant to our discussion.[7] Legal standards of privacy (such as the Health Insurance Portability and Accountability Act (HIPAA) privacy rule [3], the Family Educational Rights and Privacy Act (FERPA) [2] and Title 13 [1], the EU General Data Protection Regulation (GDPR) [22]) and the recent California Consumer Privacy Act (CCPA) and California Privacy Rights Act (CPRA) typically focus on the protection of personally identifiable information (PII) or personal data. Definitions of PII vary between regulations, but PII generally includes information that could be linked to an individual.

Information which is not PII is generally excluded from protection. This includes information where PII has been removed by means of anonymization or de-identification [39]. As an example, the HIPAA privacy rule puts no restrictions on the use or disclosure of de-identified health information. The HIPAA de-identification standard provides two de-identification methods: (i) by expert determination, i.e., if a person with appropriate knowledge and experience determines that the identification risk is very small, and (ii) by using a safe-harbor method prescribed in the privacy rule where identifiers are redacted. In a little more detail, the HIPAA safe-harbor methods enumerates 18 identifiers to be redacted including

---

[5]The probability density function of the Laplace distribution $Lap(b)$ is $\frac{1}{2b} e^{-|x|/b}$.

[6]For a more in-depth non-technical introduction to differential privacy and its applications see [44, 45] . For a technical introduction see [19, 43].

[7]An attempt for a more thorough review will quickly hit the limits of the author's expertise.

name, geographic location at a resolution smaller than a state, telephone number, and medical record numbers. It is also required that the processor "has no actual knowledge that the remaining information could be used to identify the individual" [3].

In defining PII legal standards use a collection of related concepts, such as identification, linkability, and singling out. These concepts are not precisely defined from a technical point of view, and also seem to have a limited scope. As an example, while linkability may have a concrete meaning in a setting where a record can be linked to the identify of the individual whose information is in the records (such as was the case in re-identification attacks by Sweeney and Narayanan and Shmatikov described above), it is not clear how linkability should be interpreted in settings where the data is provided in other formats, such as when a statistical or a machine learning model is derived from the PII, or when PII is replaced with "synthetic data".

These meanings of legal concepts evolve in time. For example, the US Office of Management and Budget's guidance on protecting PII was updated over time to reflect evolving understanding of re-identification attacks in de-identified data. In particular, by their guidance, non-PII may eventually become PII [4].

### 1.3 Do technical concepts of privacy match legal expectations?

Attempts to argue rigorously that certain privacy measures satisfy legal requirements face a variety of issues and challenges to the point that it is not clear that technical concepts such as $k$-anonymity and differential privacy, which evolved as measures to protect data from re-identification and other privacy attacks, match specific legal requirements.[8] In particular, computer science and legal approaches to privacy have developed in parallel, leading to diverging concepts, to the point that legal expectations sometimes seem to contradict current scientific knowledge. Hence, computer science and legal approaches to privacy lack common grounds for a rigorous analysis [34, 35]. Furthermore, the two disciplines differ greatly in their definitions, their notions of what consists a formal argument, and in the values which are pursued.

The need to bridge these gaps is now more urgent than ever. Computer systems are deeply integrated in almost every aspect of society. They continually collect extremely fine-grained personal data. And, they make numerous decisions of legal significance and consequence, effectively interpreting, or even making up the law. With the fast growing number of such decisions made in computer systems, It is unlikely that a human judiciary system will be able to deliberate even a small fraction of these decisions.

Personal data which is collected in these socio-technical systems encompasses practically every aspect of our lives, small and large. This data is stored, linked with other data, analyzed, and shared in ways which are far from being transparent. There is an urgent if not desperate need for concepts which are on one hand technical, so they can be integrated in the design of computer systems manipulating personal information, and on the other hand shown to agree with legal (and ethical) privacy desiderata.

---

[8]And, similarly with ethical requirements and societal needs.

### 1.4 Overview

Our goal in the rest of this paper is to demonstrate how a principled legal-technical analysis can help answer the above question.

In Section 2 we give an example of such an analysis. We begin with a concept from EU privacy law – the notion of singling out from the EU General Data Protection Regulation (GDPR). We then present a definition of what we believe is a related concept – predicate singling out - which we examine mathematically. Finally we make formal claims – "legal theorems" – regarding whether the technologies we surveyed above, $k$-anonymity and differential privacy, meet the GDPR requirement for anonymizing personal data. Throughout, we motivate and justify our modeling choices while attempting to stay as close as possible to the spirit of the legal text.

We conclude in Section 3 with a short discussion of possible directions for this line of research to move forward.

## 2 AN EXAMPLE: FORMALIZING THE GDPR'S CONCEPT OF SINGLING OUT

As an example of how legal concepts can be modeled mathematically and then used in rigorous statements – "legal theorems" – whether the use of specific technologies satisfy the requirements of legal standards we focus on a recent analysis of the EU's General Data Protection Regulation (GDPR) concept of singling out by Altman, Cohen, Nissim, and Wood [9, 14]. Other related examples include the modeling of the Family Educational Rights and Privacy Act (FERPA) privacy standard towards examining whether the use of differentially private analyses satisfies the standard [34], the modeling of the forgone conclusion doctrine for compelled disclosure by the government and examining of when decryption is compellable [15, 38], and the modeling of the right to deletion of personal data (the *right to be forgotten*) [25].

*Note.* Sections 2.1-2.4 are based on [9, 14] adapted, simplified, and presented in varying levels of (in)formality. The interested reader is referred to [9, 14] for the detailed legal and technical analysis.

### 2.1 Singling out in the GDPR

We begin our analysis with the text of the GDPR [22]. Article 1 of the GDPR states that the regulation turns on the processing of personal data:

> This Regulation lays down rules relating to the protection of natural persons with regard to the processing of personal data and rules relating to the free movement of personal data.

The notion of personal data is defined in Article 4:

> '[P]ersonal data' means any information relating to an identified or identifiable natural person ('data subject'); an identifiable natural person is one who can be identified, directly or indirectly

On the other hand, Recital 26 clarifies that the regulation does not apply to the processing of non-personal data, including personal data that has been *rendered anonymous'*:

> The principles of data protection should therefore not apply to anonymous information, namely information which does not relate to an identified or identifiable

natural person or to personal data rendered anonymous in such a manner that the data subject is not or no longer identifiable.

It is hence essential to understand what processes, when applied to personal data, render it "anonymous" to understand when data can be excepted from the data protection requirements in the GDPR.

Finally, Recital 26 gives more details to further clarify the notion of identifiability:

To determine whether a natural person is identifiable, account should be taken of all the means reasonably likely to be used, such as singling out, either by the controller or by another person to identify the natural person directly or indirectly.

From Recital 26 we learn that singling out is one of the *means reasonably likely to be used'* to identify a person in data. Hence, for data to be excepted from the regulation by means of rendering the personal data anonymous, it is necessary to prevent singling out in that data.[9]

Further insight into what singling out means is provided in official opinion documents prepared by the Article 29 Data Protection Working Party, a working party established by the EU Data Protection Directive, which preceded the GDPR: the Opinion on the Concept of Personal Data [5] and the Opinion on Anonymisation Techniques [6]. Referring to the notion of indirect identification, the Article 29 Working Party writes:

As regards indirectly identified or identifiable persons, this category typically relates to the phenomenon of unique combinations, whether small or large in size. . . . A name may itself not be necessary in all cases to identify an individual. This may happen when other identifiers are used to single someone out.

They define singling out as

the possibility to isolate some or all records which identify an individual in the dataset.

We interpret the a collection of attributes to be a predicate, i.e., a function assigning truth values in $\{0, 1\}$ to records. Using this formulation we define isolation as follows:

*Definition 2.1 (isolation).* A predicate $p : X \rightarrow \{0, 1\}$ isolates in the database $\mathbf{x} = (x_1, \ldots, x_n) \in X^n$ if $p(x_i) = 1$ for exactly one record $i \in [n]$, i.e., $\sum_{i=1}^{n} p(x_i) = 1$.

In Section 2.4.3 we will return to the opinion documents of the Article 29 Working parties, when we will compare their conclusions with respect to $k$-anonymity and differential privacy with the results of our analysis.

## 2.2 Formalizing singling out mathematically

Following the text cited above from the Article 29 Working Party opinion documents, it seems natural to define singling out as isolation. We now present some notation to help formalize this idea. Let $X$ be a data domain and $Y$ an arbitrary range. An anonymization mechanism $M : X^n \rightarrow Y$ takes a dataset $\mathbf{x} = (x_1, \ldots, x_n) \in X^n$ as

input and produces an output $y$ in the range $Y$. A privacy attacker $A$ observes $y$ and outputs a predicate $p : X \rightarrow \{0, 1\}$. Isolation occurs if the predicate $p(x_i) = 1$ for exactly one record in $\mathbf{x}$. Note that the predicate produced by $A$ acts on the records of the original dataset $\mathbf{x}$ and not the output $y$. Furthermore, the above formulation rules out isolation by reference to a record's position in $\mathbf{x}$ (such as "the first record").

*Definition 2.2 (isolation, rough sketch).* A mechanism $M$ prevents isolation if there does not exist an attacker $A$ that isolates in the dataset $\mathbf{x}$, except for very low probability.

There are some design decisions to be made before turning Definition 2.2 into a formal mathematical definition. To be useful, a formalization of singling out should help decide whether the application of technologies like $k$-anonymity and differential privacy suffices for rendering data anonymous. Ideally, we would like to be able to develop a mathematical formulation that exactly captures the GDPR dichotomy so as to exactly distinguish between personal data and anonymous data. This task, however, seems too ambitious if only because Recital 26 requires considering "all the means reasonably likely to be used . . . to identify [a] person" but the Recital does not make clear what these means are – only singling out is explicitly mentioned. Furthermore, even exactly capturing singling out mathematically seems (to the least) problematic. While formal from a legal point of view, the definitions cited in Section 2.1 from the GDPR as well as the opinion documents of the Article 29 Working Party leave much room for interpretation. The legal standard does not exhibit the precision required for a mathematical formulation.

Instead of attempting to exactly capture singling out with a mathematical formulation, we choose to develop a formulation of security against a notion of singling out which is potentially weaker (but not stronger) than what is intended by the GDPR. Preventing a weaker than intended notion of security against singling out is *necessary* but potentially insufficient for rendering data anonymous. Hence, the most significant consequence of our analysis would be identifying those technologies which are insufficient for satisfying the GDPR standard of anonymization (as they do not provide even a weakened notion of preventing singling out). For technologies which do withstand the weakened requirement, further inquiry would be needed for determining whether they meet the GDPR anonymization standard.

We weaken the requirements for preventing singling out in two important ways. First, we only consider attackers who do not have access to any auxiliary information. The GDPR regulators were likely familiar with the use of auxiliary information (in form of an identified dataset which can be matched with quasi-identifiers) in re-identification attacks, and likely intended that singling out would be prevented even in the presence of some auxiliary information. Second, we only consider a simple data generation process where data is sampled i.i.d. from a fixed probability distribution (which may not known to the attacker). In a more realistic settings, an attacker can utilize knowledge of the underling distribution as well as dependencies between data items.

Let $D \in \Delta(X)$ be a probability distribution over the data universe $X$. We will assume that each individual's record in the dataset is a

value sampled i.i.d. from $D$, i.e.,

$$\mathbf{x} = (x_1, \ldots, x_n) \sim D^n.$$

We can now rewrite Definition 2.2:

*Definition 2.3 (isolation).* A mechanism $M$ *prevents isolation* if for every attacker $A$,

$$\Pr[\mathbf{x} \sim D^n; y := M(\mathbf{x}); p := A(y) \text{ s.t. } \sum_{i=1}^{n} p(x_i) = 1]$$

is a negligible function of $n$, where the probability is taken over the i.i.d. drawing of $\mathbf{x}$ from $D$ and the randomness of $M$ and $A$.[10]

Unfortunately, Definition 2.3 is impossible to achieve. There exists trivial attackers, that do not even look at the outcome $y$ of the mechanism, and yet isolate with high probability! For example, let $\mathbf{x}$ be a dataset of birthdates, i.e., the records in $\mathbf{x}$ are from the data domain is $X = \{\text{Jan-1}, \ldots, \text{Dec-31}\}$ and assume $n = 365$. Assume further that $D$ is uniform over $X$. An attacker may choose an arbitrary date in $X$ to form a predicate such as

$$p(x) = \begin{cases} 1 & \text{if } x = \text{Apr-30} \\ 0 & \text{otherwise} \end{cases}$$

The probability that $p$ isolates is far from being negligible:

$$\binom{365}{1} \cdot \frac{1}{365} \cdot (1 - \frac{1}{365})^{364} \approx 37\%.$$

This is a general phenomena, if $D$ has moderate min-entropy (which would typically be the case), then one can construct a predicate $p$ such that $\Pr_{x \sim D}[p(x) = 1] = 1/n$ by applying the Leftover Hash Lemma [27].

Definition 2.3 hence needs to be modified so as to take into account the existence of trivial attackers. We would like to do that while keeping as close as possible to the spirit of the definition by the Article 29 Working party.

Define the weight of a predicate $p$ under distribution $D$ to be

$$w_D(p) = \Pr_{x \sim D}[p(x) = 1].$$

A predicate $p$ of weight $w$ which is chosen independently of the dataset succeeds in isolating with probability

$$nw(1 - w)^{n-1} \approx nwe^{-nw}.$$

This probability is negligible when either $w = negl(n)$ or $w = \omega(\log n / n)$ – in all other cases, a predicate with weight $w$ isolates with non-negligible probability. We can now modify Definition 2.3 to take this fact into account:[11]

*Definition 2.4 (Security against predicate singling out, simplified).* A mechanism $M$ is prevents predicate singling out if for every attacker $A$,

$$\Pr[\mathbf{x} \sim D^n; y := M(\mathbf{x}); p := A(y) \text{ s.t. } w_D(p) = negl(n) \wedge \sum_{i=1}^{n} p(x_i) = 1]$$

is a negligible function of $n$, where the probability is taken over the i.i.d. drawing of $\mathbf{x}$ from $D$ and the randomness of $M$ and $A$.

[10]The quantification over attackers can be modified to only consider polynomial-time attackers. A negligible function is a function approaches zero faster than any inverse polynomial, i.e., $f(n) = negl(n)$ if $f(n) = n^{-\omega(1)}$.

[11]We will onlt focus on predicates $p$ where $w_D(p)$ is negligible. The case of "heavy" predicates with $w_D(p) = \omega(\log n / n)$ can be treated analogously but seems less natural.

## 2.3 Making use of the formalization

We can now reap the benefits of faving a formal mathematical definition of security against predicate singling out. We begin by demonstrating a family of useful mechanisms which are secure against predicate singling out. We then investigate properties of the concept itself. Finally, we ask whether $k$-anonymity and differential privacy provide security against predicate singling out.

*2.3.1 Mechanisms secure against predicate singling out.* We first demonstrate that the notion of predicate singling out is not vacuous. In fact, it allows a very natural and useful family of statistical computations, namely count queries.

Let $q : X \rightarrow \{0, 1\}$ be a predicate. Let $M_{\#q}$ be the mechanism that on input $\mathbf{x} \in X^n$ returns the count of records in $\mathbf{x}$ satisfying $q$, i.e.,

$$M_{\#q}(\mathbf{x}) = \sum_{i=1}^{n} q(x_i).$$

THEOREM 2.5 ([14]). *$M_{\#q}$ prevents predicate singling out.*

*2.3.2 Properties of the definition.* Next, we examine whether security against predicate singling out is robust to post-processing and composition.

THEOREM 2.6 (ROBUSTNESS OF SECURITY AGAINST PREDICATE SINGLING OUT TO POST PROCESSING). *If a mechanism $M$ prevents predicate singling out, then for any function $f$ the mechanism which on input $\mathbf{x} \in X^n$ outputs $f(M(\mathbf{x}))$ also prevents predicate singling out.*

THEOREM 2.7 (INCOMPOSABILITY OF SECURITY AGAINST PREDICATE SINGLING OUT [14]). *There exists mechanisms $M_1 : X^n \rightarrow Y_1$ and $M_2 : X^n \rightarrow Y_2$, both preventing predicate singling out whereas the mechanism that on input $\mathbf{x} \in X^n$ outputs $(M_1(\mathbf{x}), M_2(\mathbf{x}))$ does not prevent predicate singling out.*

The proof of Theorem 2.7 provides an explicit construction of mechanisms $M_1, M_2$ satisfying the claim. However, these two mechanisms are not providing a "natural" "useful" functionality. A weaker incomposibility result can be proved using a larger number of mechanisms providing a "natural" functionality, namely the counting mechanisms of Theorem 2.5:

THEOREM 2.8 (INCOMPOSABILITY OF SECURITY AGAINST PREDICATE SINGLING OUT [14]). *There exists $\ell = \omega(\log n)$ count mechanisms $M_{\#q_1}, \ldots, M_{\#q_\ell}$ such that the mechanism that on input $\mathbf{x} \in X^n$ outputs $(M_{\#q_1}(\mathbf{x}), \ldots, M_{\#q_\ell}(\mathbf{x}))$ does not prevent predicate singling out.*

The proof of Theorem 2.8 demonstrates how count queries can be used to learn sufficiently many bits of a single record so as to isolate it with a predicate of negligible weight. This use of count queries in the proof implies that *any* formalization of security against singling out will fail to compose as long as count mechanisms would be deemed secure under that formalization.

*2.3.3 Differential privacy and security against predicate singling out.* The count mechanism $M_{\#q}$ does not satisfy differential privacy, and hence differential privacy is not necessary for security against predicate singling out. However, differential privacy does provide PSO security.

THEOREM 2.9 (INFORMAL [14]). *If M is ε-differentially private for some constant ε then M prevents predicate singling out.*

*2.3.4 k-anonymity and security against predicate singling out.* In contrast with differential privacy, it turns out that *k*-anonymity does not prevent predicate singling out attacks.

THEOREM 2.10 (INFORMAL [14]). *Typical implementations of k-anonymity, which try to optimize on the information content of the k-anonymized dataset enable an attacker to predicate single out with probability approximately 37%, i.e., to isolate in the data using predicates of negligible weight.*

A short explanation of why *k*-anonymity typically does not prevent predicate singling out is in place. Let $\mathbf{x}'$ be a *k*-anonymized version of a dataset $\mathbf{x}$. The data in $\mathbf{x}'$ can viewed as a collection "equivalence classes" each of *k* or more records, and a predicate can be assigned to each of the equivalence classes in $\mathbf{x}'$ based on their data. For the toy example 2-anonymized dataset in Section 1.1 the top two records form one equivalence class, and, similarly, the bottom two records, and the predicate corresponding to the bottom records evaluates to 1 on a record *x* if $x[\text{ZIP}] \in \{12340, \dots, 12349\} \wedge x[\text{Age}] \in \{30, \dots, 39\} \wedge x[\text{Disease}] \in PULM$, where $PULM$ is the set of pulmonary diseases. A typical dataset would be include many more attributes than that of our toy example, and *k*-anonymizers attempt to retain as much as possible information in the *k*-anonymized data and hence suppress and generalize as little attributes as possible. Hence it typically the case that the predicates corresponding the equivalence classes in a *k*-anonymized datasets would have negligible weights.

Let *p* such a predicate obtained from $\mathbf{x}'$. Since *p* has negligible weight, *p* does not evaluate to 1 for any entry of $\mathbf{x}$ outside its equivalence class. We hence get

$$\sum_{i=1}^{n} p(x_i) = \sum_{i=1}^{n} p(x_i') = k' \geq k.$$

It remains to choose a predicate $p'$ of weight $1/k'$ over the equivalence class.[12] Noting that $p'$ isolates in the equivalence class, and that the weight of $p \wedge p'$ is bounded by the weight of *p* and hence negligible, we get that an attacker outputting the predicate $p \wedge p'$ succeeds in predicate singling out with probability $\approx 37\%$.

Recent attacks by Cohen strengthen the results of Theorem 2.10 in the case of generalization-based *k*-anonymization, as they allow isolation with a negligible weight predicate with probability approaching 100% [12]

*2.3.5 Is predicate singling out the only modeling possible?* Before ending this subsection, we note that other formulations of singling out may emerge from the very same text of the GDPR, which is a rather incomplete description of the concept. The emergence of such concepts can be of great benefit to studying how to bridge between legal and technical privacy concepts.

## 2.4 Legal theorems

Based on theorems 2.9 and 2.10, we now wish to make rigorous statements of legal implications whether or not the use of a specific

privacy technology results in satisfying a legal standard. In particular, we return to the question whether *k*-anonymity and differential privacy satisfy the GDPR standard of preventing singling out.

*2.4.1 Differential privacy and the GDPR anonymization standard.* As discussed in Section 2.2, preventing singling out is a necessary but not a sufficient condition for satisfying the GDPR standard of anonymization, let alone providing a porentially weaker notion than that contemplated by the GDPR. We hence get that differential privacy may provide the right level of anonymization required by the GDPR, but further analysis is needed for making such a determination.

*2.4.2 k-anonymity and the GDPR anonymization standard.* In contrast, Theorem 2.10 has legal implications. As *k*-anonymity fails to prevent predicate singling out, and, as by the design choices made in Section 2.2, failure to prevent predicate singling out implies failure to prevent the GDPR notion of singling out, we can conclude a "legal theorem" and a "legal corollary":

LEGAL THEOREM 2.1. *k-anonymity (similarly, ℓ-diversity and t-closeness) fails to prevent singling out as required by the GDPR.*

LEGAL COROLLARY 2.1. *k-anonymity (similarly, ℓ-diversity and t-closeness) does not meet the GDPR standard for anonymization.*

*2.4.3 Comparing our conclusions with those of the Article 29 Working Party.* In its Opinion on Anonymisation Techniques [6], the Article 29 Working Party explores a number of privacy technologies, including *k*-anonymity, ℓ-diversity, and differential privacy and reach different conclusions than ours. Asking "Is Singling out still a risk?" they answer "no" for *k*-anonymity and for ℓ-diversity while they answer "may not" for differential privacy. Such a conflict between our technical analysis and what may be the current most authoritative legal interpretation of the concept we study puts our approach under stress. We explore two very different attempts to resolve it.

One possibility is to claim that the opinions provided by the Working Party are ground truth, and any modeling of singling out must agree with their determination that *k*-anonymity eliminates the risk of singling out while differential privacy may not. With this view, *k*-anonymity captures the meaning of preventing singling out (partially or in full), and the claim that *k*-anonymity eliminates the risk of singling out is *unfalsifiable*. We believe that this approach should be rejected if only because defining privacy implicitly by stipulating that a specific technology achieves privacy is an approach that has proved to fail [36].

Another possibility, which we strongly advocate, is that statements that a given privacy technology satisfies a given legal standard for anonymity should be *mathematically falsifiable*. We believe that the work presented above demonstrates how a principled analysis supported by mathematical argument can and should play an important role in articulating and informing public policy at the interface between law and technology. Returning to the Working Party Opinion on Anonymisation Technologies, to our best understanding their assessments were not substantiated using such a rigorous approach. We suggest and hope that the European Data Protection Board (EDPB), which in May 2018 replaced the Article 29 Working Party, will reconsider the Working Party's recommendations regarding *k*-anonymity and its variants.

---

[12]Technically, to see that such a predicate exists, one need to prove that sufficient min-entropy remains given that $p(x) = 1$, hence the leftover hash lemma can be used for constructing $p'$.

# 3 DISCUSSION

At the time of writing, it is not clear whether this work will have any real-world influence. In particular, it is not clear whether the European Data Protection Board (EDPB) (which replaced in 2018 the Article 29 Working Party) would find our reasoning compelling, adopt its conclusions, and update their Opinion on Anonymisation Techniques accordingly. Still, we believe that the value of a rigorous analysis has significant value, from both legal and technical points of view. The problem of bridging between legal and technical concepts is magnanimous, and rigor – our only tool for dispelling wrong intuitions while promoting reliable well-founded understanding and solutions – is probably the only direction towards tackling this extremely complex task.

Looking forward, as socio-technical systems occupy a central place in society, there is an urgent need to develop new methodologies and frameworks for the design of such systems which are based on legally sound and mathematically sound principles. New *hybrid legal-technical concepts* that harmonize legal and technical aspects of privacy are needed for specifying requirements of socio-technical systems, reasoning about their properties, and for making sure their design meets legal-normative standards of protection to individuals, groups, and society at large. The development of such hybrid concepts would require a deep interdisciplinary collaboration, a legal-technical co-design taking into consideration the normative-legal desiderata as well as the result of a technical-scientific study of privacy. In turn, these concepts will serve as building blocks for both policymaking and for the design and implementation of accountable socio-technical systems. The analysis presented in section 2, as well as other recent work [15, 25, 34, 38] provides an encouraging signal significant progress can be made in this direction.

## ACKNOWLEDGMENTS

## REFERENCES

[1] 13 U.S.C. § 9.
[2] Family Educational Rights and Privacy Act (FERPA), 20 U.S.C. § 1232g.
[3] Health Insurance Portability and Accountability Act (HIPAA) Privacy Rule, 45 C.F.R. Part 160 and Subparts A and E of Part 164.
[4] Office of Management and Budget, Memorandum M-17-12. *Preparing for and Responding to a Breach of Personally Identifiable Information.* January 3, 2017.
[5] 2007. *Article 29 Data Protection Working Party Opinion 04/2007 on the Concept of Personal Data.* https://iapp.org/media/pdf/resource_center/wp216_Anonymisation-Techniques_04-2014.pdf
[6] 2014. *Article 29 Data Protection Working Party Opinion 05/2014 on Anonymisation Techniques.* https://iapp.org/media/pdf/resource_center/wp136_concept-of-personal-data_06-2007.pdf
[7] John Abowd. 2019. Stepping-up: The Census Bureau Tries to Be a Good Data Steward in the 21st Century, Presentation at the Simons Institute for the Theory of Computing. https://simons.berkeley.edu/talks/tba-30.
[8] John Abowd. 2019. Tweetorial: Reconstruction-abetted re-identification attacks and other traditional vulnerabilities (blogpost). http://blogs.cornell.edu/abowd/special-materials/245-2/.
[9] Micah Altman, Aloni Cohen, Kobbi Nissim, and Alexandra Wood. 2021 forthcoming. What a Hybrid Legal-Technical Analysis Teaches Us About Privacy Regulation: The Case of Singling Out. *Boston University Journal of Science & Technology Law* 27, 1 (2021 forthcoming).
[10] Lars Backstrom, Cynthia Dwork, and Jon M. Kleinberg. 2011. Wherefore art thou R3579X?: anonymized social networks, hidden patterns, and structural steganography. *Commun. ACM* 54, 12 (2011), 133–141. https://doi.org/10.1145/2043174.2043199
[11] Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Song. 2019. The Secret Sharer: Evaluating and Testing Unintended Memorization in Neural Networks. In *28th USENIX Security Symposium, USENIX Security 2019, Santa Clara, CA, USA, August 14-16, 2019*, Nadia Heninger and Patrick Traynor (Eds.). USENIX Association, 267–284. https://www.usenix.org/conference/usenixsecurity19/presentation/carlini
[12] Aloni Cohen. 2021. The Quasi-identifiers are the Problem: Attacking and Reidentifying $k$-Anonymous Datasets. Manuscript in preparation.
[13] Aloni Cohen and Kobbi Nissim. 2020. Linear Program Reconstruction in Practice. *Journal of Privacy and Confidentiality* 10, 1 (Jan. 2020). https://doi.org/10.29012/jpc.711
[14] Aloni Cohen and Kobbi Nissim. 2020. Towards formalizing the GDPR's notion of singling out. *Proceedings of the National Academy of Sciences* 117, 15 (2020), 8344–8352.
[15] Aloni Cohen and Sunoo Park. 2018. Compelled Decryption and the Fifth Amendment: Exploring the Technical Boundaries. *Harvard Journal of Law & Technology* 32 (2018), 169–234. https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3117984
[16] Irit Dinur and Kobbi Nissim. 2003. Revealing information while preserving privacy. In *Proc. 22nd ACM PODS.* 202–210. https://doi.org/10.1145/773153.773173
[17] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. 2006. Calibrating Noise to Sensitivity in Private Data Analysis. In *3rd Theory of Crypt. Conf.* 265–284. https://doi.org/10.1007/11681878_14
[18] Cynthia Dwork, Frank McSherry, and Kunal Talwar. 2007. The price of privacy and the limits of LP decoding. In *Proc. 39th STOC'07 ACM.* 85–94. https://doi.org/10.1145/1250790.1250804
[19] Cynthia Dwork and Aaron Roth. 2014. The Algorithmic Foundations of Differential Privacy. *Foundations and Trends in Theoretical Computer Science* 9, 3-4 (2014), 211–407. https://doi.org/10.1561/0400000042
[20] Cynthia Dwork, Adam Smith, Thomas Steinke, and Jonathan Ullman. 2017. Exposed! A Survey of Attacks on Private Data. *Annual Review of Statistics and Its Application* 4, 1 (2017), 61–84. https://doi.org/10.1146/annurev-statistics-060116-054123 arXiv:https://doi.org/10.1146/annurev-statistics-060116-054123
[21] Cynthia Dwork and Sergey Yekhanin. 2008. New Efficient Attacks on Statistical Disclosure Control Mechanisms. In *Advances in Cryptology - CRYPTO 2008, 28th Annual International Cryptology Conference, Santa Barbara, CA, USA, August 17-21, 2008. Proceedings.* 469–480. https://doi.org/10.1007/978-3-540-85174-5_26
[22] European Parliament and the Council of the European Union. 2016. Regulation (EU) 2016/679 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC. https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX%3A32016R0679.
[23] Srivatsava Ranjit Ganta, Shiva Prasad Kasiviswanathan, and Adam D. Smith. 2008. Composition attacks and auxiliary information in data privacy. In *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Las Vegas, Nevada, USA, August 24-27, 2008*, Ying Li, Bing Liu, and Sunita Sarawagi (Eds.). ACM, 265–273. https://doi.org/10.1145/1401890.1401926
[24] Simson Garfinkel, John M. Abowd, and Christian Martindale. 2018. Understanding Database Reconstruction Attacks on Public Data: These Attacks on Statistical Databases Are No Longer a Theoretical Danger. *ACM Queue* 16, 5 (Oct. 2018), 28–53. https://doi.org/10.1145/3291276.3295691
[25] Sanjam Garg, Shafi Goldwasser, and Prashant Nalini Vasudevan. 2020. Formalizing Data Deletion in the Context of the Right to Be Forgotten. In *Advances in Cryptology - EUROCRYPT 2020 - 39th Annual International Conference on the Theory and Applications of Cryptographic Techniques, Zagreb, Croatia, May 10-14, 2020, Proceedings, Part II (Lecture Notes in Computer Science, Vol. 12106)*, Anne Canteaut and Yuval Ishai (Eds.). Springer, 373–402. https://doi.org/10.1007/978-3-030-45724-2_13
[26] Nils Homer, Szabolcs Szelinger, Margot Redman, David Duggan, Waibhav Tembe, Jill Muehling, John V. Pearson, Dietrich A. Stephan, Stanley F. Nelson, and David W. Craig. 2008. Resolving Individuals Contributing Trace Amounts of DNA to Highly Complex Mixtures Using High-Density SNP Genotyping Microarrays. *PLoS Genetics* 4(8) (2008).
[27] Russell Impagliazzo, Leonid A. Levin, and Michael Luby. 1989. Pseudo-random Generation from one-way functions (Extended Abstracts). In *Proceedings of the*

*21st Annual ACM Symposium on Theory of Computing, May 14-17, 1989, Seattle, Washigton, USA*, David S. Johnson (Ed.). ACM, 12–24. https://doi.org/10.1145/73007.73009

[28] Ninghui Li, Tiancheng Li, and Suresh Venkatasubramanian. 2007. *t*-Closeness: Privacy Beyond *k*-Anonymity and *ℓ*-Diversity. In *Proceedings of the 23rd International Conference on Data Engineering, ICDE 2007, The Marmara Hotel, Istanbul, Turkey, April 15-20, 2007*, Rada Chirkova, Asuman Dogac, M. Tamer Özsu, and Timos K. Sellis (Eds.). IEEE Computer Society, 106–115. https://doi.org/10.1109/ICDE.2007.367856

[29] Ashwin Machanavajjhala, Daniel Kifer, Johannes Gehrke, and Muthuramakrishnan Venkitasubramaniam. 2007. *L*-diversity: Privacy beyond *k*-anonymity. *ACM Trans. Knowl. Discov. Data* 1, 1 (2007), 3. https://doi.org/10.1145/1217299.1217302

[30] Adam Meyerson and Ryan Williams. 2004. On the Complexity of Optimal *k*-Anonymity. In *Proceedings of the Twenty-third ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, June 14-16, 2004, Paris, France*, Catriel Beeri and Alin Deutsch (Eds.). ACM, 223–228. https://doi.org/10.1145/1055558.1055591

[31] S. Muthukrishnan and Aleksandar Nikolov. 2012. Optimal private halfspace counting via discrepancy. In *Proceedings of the 44th Symposium on Theory of Computing Conference, STOC 2012, New York, NY, USA, May 19 - 22, 2012*, Howard J. Karloff and Toniann Pitassi (Eds.). ACM, 1285–1292. https://doi.org/10.1145/2213977.2214090

[32] Arvind Narayanan and Vitaly Shmatikov. 2006. How To Break Anonymity of the Netflix Prize Dataset. *CoRR* abs/cs/0610105 (2006). arXiv:cs/0610105 http://arxiv.org/abs/cs/0610105

[33] Arvind Narayanan and Vitaly Shmatikov. 2008. Robust de-anonymization of large sparse datasets. In *Security and Privacy*. IEEE, 111–125.

[34] Kobbi Nissim, Aaron Bembenek, Alexandra Wood, Mark Bun, Marco Gaboardi, Urs Gasser, David R. O'Brien, Thomas Steinke, and Salil Vadhan. 2018. Bridging the gap between computer science and legal approaches to privacy. *Harvard Journal of Law & Technology* 31, 2 (2018), 687–780.

[35] Kobbi Nissim and Alexandra Wood. 2018. Is privacy *privacy?* *Philosophical Transactions of the Royal Society A* 376, 2128 (2018).

[36] Paul Ohm. 2010. Broken promises of privacy: Responding to the surprising failure of anonymization. *UCLA Law Review* 57 (2010), 1701.

[37] Pierangela Samarati and Latanya Sweeney. 1998. Generalizing Data to Provide Anonymity when Disclosing Information (Abstract). In *Proceedings of the Seventeenth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, June 1-3, 1998, Seattle, Washington, USA*, Alberto O. Mendelzon and Jan Paredaens (Eds.). ACM Press, 188. https://doi.org/10.1145/275487.275508

[38] Sarah Scheffler and Mayank Varia. 2020. Protecting Cryptography Against Compelled Self-Incrimination. *IACR Cryptol. ePrint Arch.* 2020 (2020), 862. https://eprint.iacr.org/2020/862

[39] P Schwartz and D Solove. 2011. The PII Problem: Privacy and a New Concept of Personally Identifiable Information. *New York University Law Review* 86, 4 (2011), 1814–1895.

[40] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. 2017. Membership Inference Attacks Against Machine Learning Models. In *2017 IEEE Symposium on Security and Privacy, SP 2017, San Jose, CA, USA, May 22-26, 2017*. IEEE Computer Society, 3–18. https://doi.org/10.1109/SP.2017.41

[41] Latanya Sweeney. 2000. *Uniqueness of simple demographics in the US population*. Technical Report. Technical report, Carnegie Mellon University.

[42] Latanya Sweeney. 2002. *k*-Anonymity: A Model for Protecting Privacy. *Int. J. Uncertain. Fuzziness Knowl. Based Syst.* 10, 5 (2002), 557–570. https://doi.org/10.1142/S0218488502001648

[43] Salil P. Vadhan. 2017. The Complexity of Differential Privacy. In *Tutorials on the Foundations of Cryptography*, Yehuda Lindell (Ed.). Springer International Publishing, 347–450. https://doi.org/10.1007/978-3-319-57048-8_7

[44] Alexandra Wood, Micah Altman, Aaron Bembenek, Mark Bun, Marco Gaboardi, James Honaker, Kobbi Nissim, David R O'Brien, Thomas Steinke, and Salil Vadhan. 2018. Differential privacy: A primer for a non-technical audience. *Vand. J. Ent. & Tech. L.* 21 (2018), 209.

[45] Alexandra Wood, Micah Altman, Kobbi Nissim, and Salil Vadhan. 2020. Designing Access with Differential Privacy. In *Handbook on Using Administrative Data for Research and Evidence-based Policy*, Shawn Cole, Iqbal Dhaliwal, Anja Sautmann, and Lars Vilhuber (Eds.). Abdul Latif Jameel Poverty Action Lab, Cambridge, MA.