
The Composition Theorem for Differential Privacy

Peter Kairouz

ECE Department, University of Illinois at Urbana-Champaign, Urbana, IL 61801 USA

KAIROUZ2@ILLINOIS.EDU

Sewoong Oh

IESE Department, University of Illinois at Urbana-Champaign, Urbana, IL 61801 USA

SWOH@ILLINOIS.EDU

Pramod Viswanath

ECE Department, University of Illinois at Urbana-Champaign, Urbana, IL 61801 USA

PRAMODV@ILLINOIS.EDU

Abstract

Sequential querying of differentially private mechanisms degrades the overall privacy level. In this paper, we answer the fundamental question of characterizing the level of overall privacy degradation as a function of the number of queries and the privacy levels maintained by each privatization mechanism. Our solution is complete: we prove an upper bound on the overall privacy level and construct a sequence of privatization mechanisms that achieves this bound. The key innovation is the introduction of an operational interpretation of differential privacy (involving hypothesis testing) and the use of new data processing inequalities. Our result improves over the state-of-the-art and has immediate applications to several problems studied in the literature.

1. Introduction

Differential privacy is a formal framework to quantify to what extent individual privacy in a statistical database is preserved while releasing useful aggregate information about the database. It provides strong privacy guarantees by requiring the indistinguishability of whether or not an individual is in the database based on the released information. Denoting the database when the individual is present as D and as D' when the individual is not, a differentially private mechanism provides indistinguishability guarantees with respect to the pair (D, D') . The databases D and D' are referred to as “neighboring” databases.

Definition 1.1 (Differential Privacy (Dwork et al.,

2006b;a)). A randomized mechanism M over a set of databases is (ϵ, δ) -differentially private if for all pairs of neighboring databases D and D' , and for all sets S in the output space of the mechanism \mathcal{X} ,

$$\mathbb{P}(M(D) \in S) \leq e^\epsilon \mathbb{P}(M(D') \in S) + \delta.$$

A basic problem in differential privacy is how does the overall privacy level degrade under the *composition* of interactive queries where each query meets a certain differential privacy guarantee. A routine argument shows that the composition of k queries, each of which is (ϵ, δ) -differentially private, is at least $(k\epsilon, k\delta)$ -differentially private (Dwork et al., 2006b;a; Dwork & Lei, 2009; Dwork et al., 2010). A tighter bound of $(\tilde{\epsilon}_\delta, k\delta + \tilde{\delta})$ -differential privacy under k -fold adaptive composition is provided, using more sophisticated arguments, in Dwork et al. (2010) for the case when each of the individual queries is (ϵ, δ) -differentially private. Here $\tilde{\epsilon}_\delta = O\left(k\epsilon^2 + \epsilon\sqrt{k\log(1/\tilde{\delta})}\right)$. On the other hand, it was not known if this bound could be improved until this work.

Our main result is the *exact* characterization of the privacy guarantee under k -fold composition. Any k -fold adaptive composition of (ϵ, δ) -differentially private mechanisms satisfies this privacy guarantee, stated as Theorem 3.3. Further, we construct a specific sequence of privacy mechanisms which under (in fact, nonadaptive) composition actually degrade privacy to the level guaranteed. Our result entails a strict improvement over the state-of-the-art: this can be seen immediately in the following approximation – using the same notation as above, the value of $\tilde{\epsilon}_\delta$ is now reduced to $\tilde{\epsilon}_\delta = O\left(k\epsilon^2 + \epsilon\sqrt{k\log(e + (\epsilon\sqrt{k}/\tilde{\delta}))}\right)$. Since a typical choice of $\tilde{\delta}$ is $\tilde{\delta} = \Theta(k\delta)$, in the regime where $\epsilon = \Theta(\sqrt{k}\delta)$, this improves the existing guarantee by a logarithmic factor. The gain is especially significant when both ϵ and δ are small.

We start with the view of differential privacy as providing certain guarantees for the two error types (false alarm and missed detection) in a binary hypothesis testing problem (involving two neighboring databases), as in previous work by Wasserman & Zhou (2010). We bring two benefits of this *operational* interpretation of the privacy definition to bear on the problem at hand.

- The first is conceptual: the operational setting directs the logic of the steps of the proof, makes the arguments straightforward and readily allows generalizations such as heterogeneous compositions. While the state-of-the-art work (Dwork et al., 2010) needed recourse to sophisticated mathematical results in Reingold et al. (2008); Tao & Ziegler (2008); Green & Tao (2004) to derive their results, our strengthening is arrived at using relatively elementary techniques.
- The second is technical: the operational interpretation of hypothesis testing brings both the natural data processing inequality, and the strong converse to the data processing inequality. These inequalities, while simple by themselves, lead to surprisingly strong technical results. As an aside, we mention that there is a strong tradition of such derivations in the information theory literature: the Fisher information inequality (Blachman, 1965; Zamir, 1998), the entropy power inequality (Stam, 1959; Blachman, 1965; Verdú & Guo, 2006), an extremal inequality involving mutual informations (Liu & Viswanath, 2007), matrix determinant inequalities (Cover & Thomas, 1988), the Brunn-Minkowski inequality and its functional analytic variants (Dembo et al., 1991) – Chapter 17 of Cover & Thomas (2012) enumerates a detailed list – were all derived using operational interpretations of mutual information and corresponding data processing inequalities.

One special case of our results, the strengthening of the state-of-the-art result in Dwork et al. (2010), could also have been arrived at directly by using stronger technical methods than used in Dwork et al. (2010). Specifically, we use a direct expression for the privacy region (instead of an upper bound) to arrive at our strengthened result.

The optimal composition theorem (Theorem 3.3) provides a fundamental limit on how much privacy degrades under composition. Such a characterization is a basic result in differential privacy and has been used widely in the literature (Dwork et al., 2010; Hardt et al., 2010; Blocki et al., 2012; Gupta et al., 2012; Muthukrishnan & Nikolov, 2012; Hardt & Roth, 2013). In each of these instances, the optimal composition theorem derived here (or the simpler characterization of Theorem 3.4) could be “cut-and-pasted”, allowing for corresponding strengthening of their conclusions. We

demonstrate this strengthening for two instances: variance of noise adding mechanisms in Section 4.1 and Gaussian projection in Section 4.2. We further show that a variety of existing noise adding mechanisms ensure the same level of privacy with similar variances. This implies that there is nothing special about the popular choice of adding a Gaussian noise when composing multiple queries, and the same utility as measured through the noise variance can be obtained using other known mechanisms. We start our discussions by operationally introducing differential privacy as certain guarantees on the error probabilities in a binary hypothesis testing problem.

2. Differential Privacy as Hypothesis Testing

Given a random output Y of a database access mechanism M , consider the following hypothesis testing experiment. We choose a null hypothesis as database D_0 and alternative hypothesis as D_1 :

$$\begin{aligned} H_0 &: Y \text{ came from a database } D_0, \\ H_1 &: Y \text{ came from a database } D_1. \end{aligned}$$

For a choice of a rejection region S , the probability of false alarm (type I error), when the null hypothesis is true but rejected, is defined as $P_{\text{FA}}(D_0, D_1, M, S) \equiv \mathbb{P}(M(D_0) \in S)$, and the probability of missed detection (type II error), when the null hypothesis is false but retained, is defined as $P_{\text{MD}}(D_0, D_1, M, S) \equiv \mathbb{P}(M(D_1) \in \bar{S})$ where \bar{S} is the complement of S . The differential privacy condition on a mechanism M is equivalent to the following set of constraints on the probability of false alarm and missed detection. Wasserman and Zhou proved that $(\epsilon, 0)$ -differential privacy implies the conditions (1) for a special case when $\delta = 0$ (Wasserman & Zhou, 2010). The same proof technique can be used to prove a similar result for general $\delta \in [0, 1]$, and to prove that the conditions (1) imply (ϵ, δ) -differential privacy as well. We refer to the supplementary material for a proof.

Theorem 2.1. *For any $\epsilon \geq 0$ and $\delta \in [0, 1]$, a database mechanism M is (ϵ, δ) -differentially private if and only if the following conditions are satisfied for all pairs of neighboring databases D_0 and D_1 , and all rejection regions $S \subseteq \mathcal{X}$:*

$$\begin{aligned} P_{\text{FA}}(D_0, D_1, M, S) + e^\epsilon P_{\text{MD}}(D_0, D_1, M, S) &\geq 1 - \delta, \\ e^\epsilon P_{\text{FA}}(D_0, D_1, M, S) + P_{\text{MD}}(D_0, D_1, M, S) &\geq 1 - \delta. \end{aligned} \quad (1)$$

This operational perspective of differential privacy relates the privacy parameters ϵ and δ to a set of conditions on probability of false alarm and missed detection. This shows that it is impossible to get both small P_{MD} and P_{FA} from data obtained via a differentially private mechanism, and

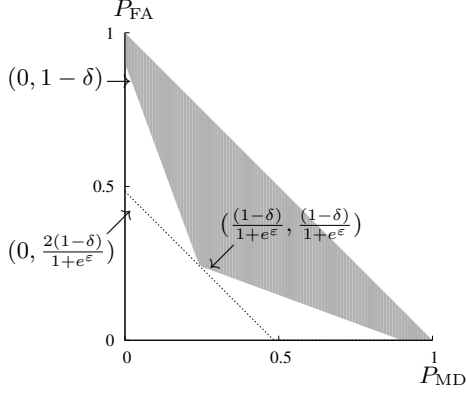


Figure 1. Privacy region for (ε, δ) -differential privacy. For simplicity, we only show the privacy region below the line $P_{\text{FA}} + P_{\text{MD}} \leq 1$, since the whole region is symmetric w.r.t. the line $P_{\text{FA}} + P_{\text{MD}} = 1$.

that the converse is also true. This operational interpretation of differential privacy suggests a graphical representation of differential privacy as illustrated in Figure 1.

$$\mathcal{R}(\varepsilon, \delta) \equiv \left\{ (P_{\text{MD}}, P_{\text{FA}}) \mid P_{\text{FA}} + e^\varepsilon P_{\text{MD}} \geq 1 - \delta, \text{ and } e^\varepsilon P_{\text{FA}} + P_{\text{MD}} \geq 1 - \delta \right\}. \quad (2)$$

Similarly, we define the *privacy region* of a database access mechanism M with respect to two neighboring databases D_0 and D_1 as

$$\begin{aligned} \mathcal{R}(M, D_0, D_1) \equiv & \\ \text{conv} \left(\left\{ (P_{\text{MD}}(D_0, D_1, M, S), P_{\text{FA}}(D_0, D_1, M, S)) \mid \right. \right. & \\ \left. \left. \text{for all } S \subseteq \mathcal{X} \right\} \right), & \quad (3) \end{aligned}$$

where $\text{conv}(\cdot)$ is the convex hull of a set. Operationally, by taking the convex hull, the region includes the pairs of false alarm and missed detection probabilities achieved by soft decisions that might use internal randomness in the hypothesis testing. Precisely, let $\gamma : \mathcal{X} \rightarrow \{H_0, H_1\}$ be any decision rule where we allow probabilistic decisions. For example, if the output is in a set S_1 we can accept the null hypothesis with a certain probability p_1 , and for another set S_2 accept with probability p_2 . In full generality, a decision rule γ can be fully described by a partition $\{S_i\}$ of the output space \mathcal{X} , and corresponding accept probabilities $\{p_i\}$. The probabilities of false alarm and missed detection for a decision rule γ is defined as $P_{\text{FA}}(D_0, D_1, M, \gamma) \equiv \mathbb{P}(\gamma(M(D_0)) = H_1)$ and $P_{\text{MD}}(D_0, D_1, M, \gamma) \equiv \mathbb{P}(\gamma(M(D_1)) = H_0)$.

Remark 2.2. For all neighboring databases D_0 and D_1 , and a database access mechanism M , the pair of a false alarm and a missed detection probabilities achieved by any decision rule γ is included in the privacy region:

$$(P_{\text{MD}}(D_0, D_1, M, \gamma), P_{\text{FA}}(D_0, D_1, M, \gamma)) \in \mathcal{R}(M, D_0, D_1)$$

for all decision rules γ .

Let $D_0 \sim D_1$ denote that the two databases are neighbors. The union over all neighboring databases define the *privacy region of the mechanism*.

$$\mathcal{R}(M) \equiv \bigcup_{D_0 \sim D_1} \mathcal{R}(M, D_0, D_1).$$

The following corollary, which follows immediately from Theorem 2.1, gives a necessary and sufficient condition on the privacy region for (ε, δ) -differential privacy.

Corollary 2.3. A mechanism M is (ε, δ) -differentially private if and only if $\mathcal{R}(M) \subseteq \mathcal{R}(\varepsilon, \delta)$.

Consider two database access mechanisms $M(\cdot)$ and $M'(\cdot)$. Let X and Y denote the random outputs of mechanisms M and M' respectively. We say M dominates M' if $M'(D)$ is conditionally independent of the database D conditioned on the outcome of $M(D)$. In other words, the database D , $X = M(D)$ and $Y = M'(D)$ form the following Markov chain: $D-X-Y$ (Cover & Thomas, 1988).

Theorem 2.4 (Data processing inequality for differential privacy). If a mechanism M dominates a mechanism M' , then for all pairs of neighboring databases D_1 and D_2 ,

$$\mathcal{R}(M', D_1, D_2) \subseteq \mathcal{R}(M, D_1, D_2).$$

We refer to the supplementary material for a proof. Wasserman & Zhou (2010) have proved that, for a special case when M is $(\varepsilon, 0)$ -differentially private, M' is also $(\varepsilon, 0)$ -differentially private, which is a corollary of the above theorem. Perhaps surprisingly, the converse is also true.

Theorem 2.5 (Corollary of Theorem 10 from (Blackwell, 1953)). Fix a pair of neighboring databases D_1 and D_2 and let X and Y denote the random outputs of mechanisms M and M' , respectively. If M and M' satisfy

$$\mathcal{R}(M', D_1, D_2) \subseteq \mathcal{R}(M, D_1, D_2),$$

then there exists a coupling of the random outputs X and Y such that they form a Markov chain $D-X-Y$.

When the privacy region of M' is included in M , then there exists a stochastic transformation T that operates on X and produce a random output that has the same marginal distribution as Y conditioned on the database D . We can consider this mechanism T as a privatization mechanism that takes a (privatized) output X and provides even further privatization. The above theorem was proved in Blackwell (1953) in the context of comparing two experiments, where a statistical *experiment* denotes a mechanism in the context of differential privacy.

3. Composition of Differentially Private Mechanisms

In this section, we address how differential privacy guarantees compose: when accessing databases multiple times via differentially private mechanisms, each of which having its own privacy guarantees, how much privacy is still guaranteed on the union of those outputs? To formally define composition, we consider the following scenario known as the ‘composition experiment’, proposed in [Dwork et al. \(2010\)](#).

A composition experiment takes as input a parameter $b \in \{0, 1\}$, and an adversary \mathcal{A} . From the hypothesis testing perspective proposed in the previous section, b can be interpreted as the hypothesis: null hypothesis for $b = 0$ and alternative hypothesis for $b = 1$. At each time i , a database $D^{i,b}$ is accessed depending on b . For example, $D^{1,0}$ could be medical records including a particular individual and $D^{1,1}$ does not include the person, and $D^{2,0}$ could be voter registration database with the same person present and $D^{2,1}$ with the person absent. An adversary \mathcal{A} is trying to break privacy (and figure out whether or not the particular individual is in the database) by testing the hypotheses on the output of k sequential access to those databases via differentially private mechanisms. In full generality, we allow the adversary to have full control over which database to access, which query to ask, and which mechanism to be used at each repeated access. Further, the adversary is free to make these choices adaptively based on the previous outcomes. The only restrictions are the differentially private mechanisms belong to a family \mathcal{M} (e.g., the family of all (ε, δ) -differentially private mechanisms), the internal randomness of the mechanisms are independent at each repeated access, and the hypothesis b is not known to the adversary. The outcome of this k -fold composition experiment is the *view of the adversary* \mathcal{A} : $V^b \equiv (R, Y_1^b, \dots, Y_k^b)$, where

$$Y_i^b = M_i(D^{i,b}, q_i),$$

q_i is the i^{th} query, $M_i \in \mathcal{M}$ is the i^{th} privatization mechanism, and R is the internal randomness of \mathcal{A} .

3.1. Optimal privacy region under composition

In terms of testing whether a particular individual is in the database ($b = 0$) or not ($b = 1$), we want to characterize how much privacy degrades after a k -fold composition experiment. It is known that the privacy degrades under composition by at most the ‘sum’ of the differential privacy parameters of each access.

Theorem 3.1 ([\(Dwork et al., 2006b;a; Dwork & Lei, 2009; Dwork et al., 2010\)](#)). *For any $\varepsilon > 0$ and $\delta \in [0, 1]$, the class of (ε, δ) -differentially private mechanisms satisfy*

$(k\varepsilon, k\delta)$ -differential privacy under k -fold adaptive composition.

In general, one can show that if M_i is $(\varepsilon_i, \delta_i)$ -differentially private, then the composition satisfies $(\sum_{i \in [k]} \varepsilon_i, \sum_{i \in [k]} \delta_i)$ -differential privacy. If we do not allow any slack in the δ , this bound cannot be tightened. Precisely, there are examples of mechanisms which under k -fold composition violate $(\varepsilon, \sum_{i \in [k]} \delta_i)$ -differential privacy for any $\varepsilon < \sum_{i \in [k]} \varepsilon_i$. We can prove this by providing a set S such that the privacy condition is met with equality: $\mathbb{P}(V^0 \in S) = e^{\sum_{i \in [k]} \varepsilon_i} \mathbb{P}(V^1 \in S) + \sum_{i \in [k]} \delta_i$. However, if we allow for a slightly larger value of δ , then [Dwork et al. \(2010\)](#) showed that one can gain a significantly higher privacy guarantee in terms of ε .

Theorem 3.2 (Theorem III.3 from [Dwork et al. \(2010\)](#)). *For any $\varepsilon > 0$, $\delta \in [0, 1]$, and $\tilde{\delta} \in (0, 1]$, the class of (ε, δ) -differentially private mechanisms satisfies $(\tilde{\varepsilon}_{\tilde{\delta}}, k\delta + \tilde{\delta})$ -differential privacy under k -fold adaptive composition, for*

$$\tilde{\varepsilon}_{\tilde{\delta}} = k\varepsilon(e^\varepsilon - 1) + \varepsilon\sqrt{2k \log(1/\tilde{\delta})}. \quad (4)$$

By allowing a slack of $\tilde{\delta} > 0$, one can get a higher privacy of $\tilde{\varepsilon}_{\tilde{\delta}} = O(k\varepsilon^2 + \sqrt{k\varepsilon^2})$, which is significantly smaller than $k\varepsilon$. This is the best known guarantee so far, and has been used whenever one requires a privacy guarantee under composition (e.g. [\(Dwork et al., 2010; Blocki et al., 2012; Hardt & Roth, 2013\)](#)). However, the important question of optimality has remained open. Namely, is there a composition of mechanisms where the above privacy guarantee is tight? In other words, is it possible to get a tighter bound on differential privacy under composition?

We give a complete answer to this fundamental question in the following theorems. We prove a tighter bound on the privacy under composition. Further, we also prove the achievability of the privacy guarantee: we provide a set of mechanisms such that the privacy region under k -fold composition is exactly the region defined by the conditions in (5). Hence, this bound on the privacy region is tight and cannot be improved upon.

Theorem 3.3. *For any $\varepsilon \geq 0$ and $\delta \in [0, 1]$, the class of (ε, δ) -differentially private mechanisms satisfies*

$$((k - 2i)\varepsilon, 1 - (1 - \delta)^k(1 - \delta_i))\text{-differential privacy} \quad (5)$$

under k -fold adaptive composition, for all $i = \{0, 1, \dots, \lfloor k/2 \rfloor\}$, where

$$\delta_i = \frac{\sum_{\ell=0}^{i-1} \binom{k}{\ell} (e^{(k-\ell)\varepsilon} - e^{(k-2i+\ell)\varepsilon})}{(1 + e^\varepsilon)^k}. \quad (6)$$

Hence, the privacy region of k -fold composition is an intersection of k regions, each of which is $((k -$

$2i)\varepsilon, \delta_i$)-differentially private: $\mathcal{R}(\{(k-2i)\varepsilon, \delta_i\}_{i \in [k/2]}) \equiv \bigcap_{i=0}^{\lfloor \frac{k}{2} \rfloor} \mathcal{R}((k-2i)\varepsilon, \delta_i)$. We refer to the supplementary material for a proof, where we give an explicit mechanism that achieves this region under composition. Hence, this bound on the privacy region is tight, and gives the exact description of how much privacy can degrade under k -fold adaptive composition. This settles the question left open in Dwork et al. (2006b;a); Dwork & Lei (2009); Dwork et al. (2010) by providing, for the first time, the fundamental limit of composition, and proving a matching mechanism with the worst-case privacy degradation.

To prove the optimality in Theorem 3.3, namely that it is impossible to have privacy worse than (5), we rely on the operational interpretation of the privacy as hypothesis testing. Precisely, for a given pair (ε, δ) , we define a *dominant mechanism* whose output can be used to simulate any (ε, δ) -differentially private mechanism. Therefore, any k -fold composition of (ε, δ) -differentially private mechanisms can be simulated from the output of k -fold composition of the dominant mechanisms. The analysis of this dominant mechanisms gives the exact composition theorem in Eq.(5), and the new analysis tools (Theorem 2.4 and Theorem 2.5) proves its optimality.

Figure 2 illustrates how much the privacy region of Theorem 3.3 degrades as we increase the number of composition k . Figure 3 provides a comparison of the three privacy guarantees in Theorems 3.1, 3.2 and 3.3 for 30-fold composition of $(0.1, 0.001)$ -differentially private mechanisms. Smaller region gives a tighter bound, since it guarantees the higher privacy.

3.2. Simplified privacy region under composition

In many applications of the composition theorems, a closed form expression of the composition privacy guarantee is required. The privacy guarantee in (5) is tight, but can be difficult to evaluate. The next theorem provides a simpler form expression which is an outer bound of the exact region described in (5). Comparing to (4), the privacy guarantee is significantly improved from $\tilde{\varepsilon}_{\tilde{\delta}} = O(k\varepsilon^2 + \sqrt{k\varepsilon^2 \log(1/\tilde{\delta})})$ to $\tilde{\varepsilon}_{\tilde{\delta}} = O(k\varepsilon^2 + \min\{\sqrt{k\varepsilon^2 \log(1/\tilde{\delta})}, \varepsilon \log(\varepsilon/\tilde{\delta})\})$, especially when composing a large number k of interactive queries. Further, the δ -approximate differential privacy degradation of $(1 - (1 - \delta)^k(1 - \tilde{\delta}))$ is also strictly smaller than the previous $(k\delta + \tilde{\delta})$. We discuss the significance of this improvement in the next section using examples from existing differential privacy literature.

Theorem 3.4. *For any $\varepsilon > 0$, $\delta \in [0, 1]$, and $\tilde{\delta} \in [0, 1]$, the class of (ε, δ) -differentially private mechanisms satisfies $(\tilde{\varepsilon}_{\tilde{\delta}}, 1 - (1 - \delta)^k(1 - \tilde{\delta}))$ -differential privacy under*

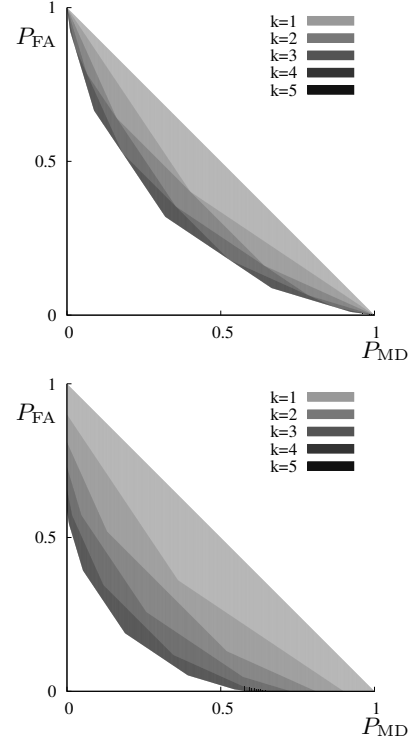


Figure 2. Privacy region $\mathcal{R}(\{(k-2i)\varepsilon, \delta_i\})$ for the class of $(\varepsilon, 0)$ -differentially private mechanisms (top) and (ε, δ) -differentially private mechanisms (bottom) under k -fold adaptive composition.

k -fold adaptive composition, for

$$\tilde{\varepsilon}_{\tilde{\delta}} = \min \left\{ k\varepsilon, \frac{(e^\varepsilon - 1)\varepsilon k}{e^\varepsilon + 1} + \varepsilon \sqrt{2k \log \left(e + \frac{\sqrt{k\varepsilon^2}}{\tilde{\delta}} \right)}, \frac{(e^\varepsilon - 1)\varepsilon k}{e^\varepsilon + 1} + \varepsilon \sqrt{2k \log \left(\frac{1}{\tilde{\delta}} \right)} \right\}. \quad (7)$$

In the high privacy regime, where $\varepsilon \leq 0.9$, this bound can be further simplified as

$$\tilde{\varepsilon}_{\tilde{\delta}} \leq \min \left\{ k\varepsilon, k\varepsilon^2 + \varepsilon \sqrt{2k \log \left(e + (\sqrt{k\varepsilon^2}/\tilde{\delta}) \right)}, k\varepsilon^2 + \varepsilon \sqrt{2k \log(1/\tilde{\delta})} \right\}.$$

We refer to the supplementary material for a proof. This privacy guarantee improves over the existing result of Theorem 3.2 when $\tilde{\delta} = \Theta(\sqrt{k\varepsilon^2})$. Typical regime of interest is the high-privacy regime for composition privacy guarantee, i.e. when $\sqrt{k\varepsilon^2} \ll 1$. The above theorem suggests that we only need the extra slack of approximate privacy $\tilde{\delta}$ of order $\sqrt{k\varepsilon^2}$, instead of $\Omega(1)$ as suggested by the existing results.

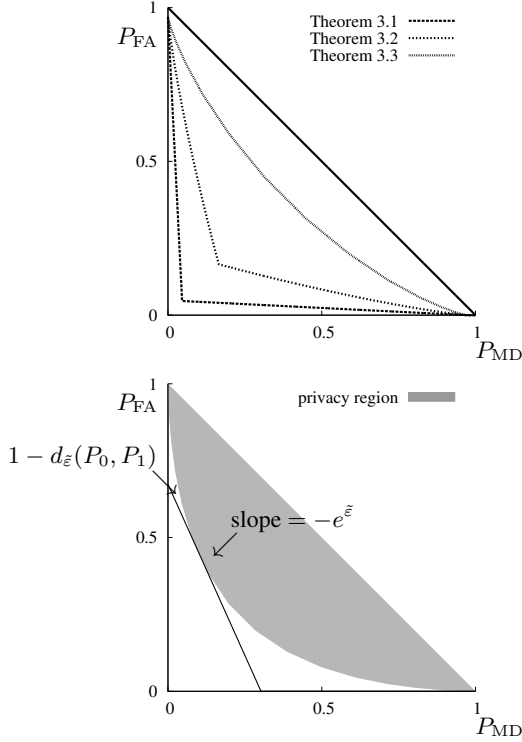


Figure 3. Theorem 3.3 provides the tightest bound (top). Given a mechanism M , the privacy region can be completely described by its boundary, which is represented by a set of tangent lines of the form $P_{\text{FA}} = -e^\epsilon P_{\text{MD}} + 1 - d_\epsilon(P_0, P_1)$ (bottom).

3.3. Composition Theorem for Heterogeneous Mechanisms

We considered homogeneous mechanisms, where all mechanisms are (ϵ, δ) -differentially private. Our analysis readily extends to heterogeneous mechanisms, where the ℓ -th query satisfies $(\epsilon_\ell, \delta_\ell)$ -differential privacy (we refer to such mechanisms as $(\epsilon_\ell, \delta_\ell)$ -differentially private mechanisms).

Theorem 3.5. *For any $\epsilon_\ell > 0$, $\delta_\ell \in [0, 1]$ for $\ell \in \{1, \dots, k\}$, and $\tilde{\delta} \in [0, 1]$, the class of $(\epsilon_\ell, \delta_\ell)$ -differentially private mechanisms satisfy $(\tilde{\epsilon}_{\tilde{\delta}}, 1 - (1 - \tilde{\delta}) \prod_{\ell=1}^k (1 - \delta_\ell))$ -differential privacy under k -fold adaptive composition, for $\tilde{\epsilon}_{\tilde{\delta}} =$*

$$\min \left\{ \sum_{\ell=1}^k \epsilon_\ell, \sum_{\ell=1}^k \frac{(e^{\epsilon_\ell} - 1)\epsilon_\ell}{e^{\epsilon_\ell} + 1} + \sqrt{\sum_{\ell=1}^k 2\epsilon_\ell^2 \log \left(e + \frac{\sqrt{\sum_{\ell=1}^k \tilde{\epsilon}_\ell^2}}{\tilde{\delta}} \right)}, \sum_{\ell=1}^k \frac{(e^{\epsilon_\ell} - 1)\epsilon_\ell}{e^{\epsilon_\ell} + 1} + \sqrt{\sum_{\ell=1}^k 2\epsilon_\ell^2 \log \left(\frac{1}{\tilde{\delta}} \right)} \right\}. \quad (8)$$

This tells us that the ϵ_ℓ 's *sum up under composition*: whenever we have $k\epsilon$ or $k\epsilon^2$ in (7) we can replace it by the summation to get the general result for heterogeneous case.

4. Applications

When analyzing a complex mechanism with multiple sub-mechanisms each with (ϵ_0, δ_0) -differential privacy guarantee, we can apply the composition theorem (Theorem 3.3 and Theorem 3.4). To ensure overall (ϵ, δ) -differential privacy for the whole complex mechanism, one chooses $\epsilon_0 = \epsilon / (2\sqrt{k \log(e + \epsilon/\delta)})$ and $\delta_0 = \delta/2k$, when there are k sub-mechanisms. The existing composition theorem guarantees the desired overall privacy. Then, the *utility* of the complex mechanism is calculated for the choice of ϵ_0 and δ_0 .

Following this recipe, we first provide a sufficient condition on the variance of noise adding mechanisms. This analysis shows that one requires smaller variance than what is previously believed, in the regime where $\epsilon = \Theta(\delta)$. Further, we show that a variety of known mechanisms achieve the desired privacy under composition with the *same* level of variance. Applying this analysis to known mechanisms for cut queries of a graph, we show that again in the regime where $\epsilon = \Theta(\delta)$, one can achieve the desired privacy under composition with improved utility.

For count queries with sensitivity one, the geometric noise adding mechanism is known to be universally optimal in a general cost minimization framework (Bayesian setting in Ghosh et al. (2012) and worst-case setting in Geng & Viswanath (2012)). Here we provide a new interpretation of the geometric noise adding mechanism as an optimal mechanism under *composition* for counting queries. In the course of proving Theorem 3.3, we show that a family of mechanisms are optimal under composition, in the sense that they achieve the largest false alarm and missed detection region. Larger region under composition implies that one can achieve smaller error rates, while ensuring the same level of privacy at each step of the composition. In this section, we show that the geometric mechanism is one of such mechanisms, thus providing the new interpretation to the optimality of the geometric mechanisms.

4.1. Variance of noise adding mechanisms under composition

In this section, we consider real-valued queries $q : \mathcal{D} \rightarrow \mathbb{R}$. The *sensitivity* of a real-valued query is defined as the maximum absolute difference of the output between two neighboring databases:

$$\Delta \equiv \max_{D \sim D'} |q(D) - q(D')|,$$

where \sim indicates that the pair of databases are neighbors. A common approach to privatize such a query output is to add noise to it, and the variance of the noise grows with sensitivity of the query and the desired level of privacy. A popular choice of the noise is Gaussian. It is previously known that it is sufficient to add Gaussian noise with variance $O(k\Delta^2 \log(1/\delta)/\varepsilon^2)$ to each query output in order to ensure (ε, δ) -differential privacy under k -fold composition. We improve the analysis of Gaussians under composition, and show that for a certain regime where $\varepsilon = \Theta(\delta)$, the sufficient condition can be improved by a log factor.

When composing real-valued queries, the *Gaussian mechanism* is a popular choice (Blocki et al., 2012; Hardt & Roth, 2013). However, we show that there is nothing special about Gaussian mechanisms for composition. We prove that the *Laplacian mechanism* or the *staircase mechanism* introduced in Geng & Viswanath (2012) can achieve the same level of privacy under composition with the same variance.

We can use Theorem 3.4 to find how much noise we need to add to each query output, in order to ensure (ε, δ) -differential privacy under k -fold composition. We know that if each query output is $(\varepsilon_0, \delta_0)$ -differentially private, then the composed outputs satisfy $(k\varepsilon_0^2 + \sqrt{2k\varepsilon_0^2 \log(e + \sqrt{k\varepsilon_0^2/\delta})}, k\delta_0 + \tilde{\delta})$ -differential privacy assuming $\varepsilon_0 \leq 0.9$. With the choice of $\delta_0 = \delta/2k$, $\tilde{\delta} = \delta/\sqrt{2}$, and $\varepsilon_0^2 = \varepsilon^2/4k \log(e + (\varepsilon/\delta))$, this ensures that the target privacy of (ε, δ) is satisfied under k -fold composition as described in the following corollary.

Corollary 4.1. *For any $\varepsilon \in (0, 0.9]$ and $\delta \in (0, 1]$, if the database access mechanism satisfies $(\sqrt{\varepsilon^2/4k \log(e + (\varepsilon/\delta))}, \delta/2k)$ -differential privacy on each query output, then it satisfies (ε, δ) -differential privacy under k -fold composition.*

One of the most popular noise adding mechanisms is the *Laplacian mechanism*, which adds Laplacian noise to real-valued query outputs. When the sensitivity is Δ , one can achieve $(\varepsilon_0, 0)$ -differential privacy with the choice of the distribution $\text{Lap}(\varepsilon_0/\Delta) = (\varepsilon_0/2\Delta)e^{-\varepsilon_0|x|/\Delta}$. The resulting variance of the noise is $2\Delta^2/\varepsilon_0^2$. The above corollary implies a certain sufficient condition on the variance of the Laplacian mechanism to ensure privacy under composition.

Corollary 4.2. *For real-valued queries with sensitivity $\Delta > 0$, the mechanism that adds Laplacian noise with variance $(8k\Delta^2 \log(e + (\varepsilon/\delta)))/\varepsilon^2$ satisfies (ε, δ) -differential privacy under k -fold adaptive composition for any $\varepsilon \in (0, 0.9]$ and $\delta \in (0, 1]$. The mean squared error achieved by the Laplacian mechanism is $O(k^2\Delta^2 \log(e + (\varepsilon/\delta)))/\varepsilon^2$.*

In terms of variance-privacy trade-off for real-valued queries, the optimal noise-adding mechanism known as the *staircase mechanism* was introduced in Geng &

Viswanath (2012). The probability density function of this noise is piecewise constant, and the probability density on the pieces decay geometrically. It is shown in Geng & Viswanath (2013) that that with variance of $O(\min\{1/\varepsilon^2, 1/\delta^2\})$, the staircase mechanism achieved (ε, δ) -differential privacy. Corollary 4.1 implies that with variance $O(k\Delta^2 \log(e + \varepsilon/\delta)/\varepsilon^2)$, the staircase mechanism satisfies (ε, δ) -differential privacy under k -fold composition.

Another popular mechanism known as the *Gaussian mechanism* privatizes each query output by adding a Gaussian noise with variance σ^2 . It is not difficult to show that when the sensitivity of the query is Δ , with a choice of $\sigma^2 \geq 2\Delta^2 \log(2/\delta_0)/\varepsilon_0^2$, the Gaussian mechanism satisfies $(\varepsilon_0, \delta_0)$ -differential privacy (e.g. (Dwork et al., 2006a)). The above corollary implies that the Gaussian mechanism with variance $O(k\Delta^2 \log(1/\delta) \log(e + (\varepsilon/\delta)))/\varepsilon^2$ ensures (ε, δ) -differential privacy under k -fold composition. However, we can get a tighter sufficient condition by directly analyzing how Gaussian mechanisms compose, and we refer to the supplementary material for a proof.

Theorem 4.3. *For real-valued queries with sensitivity $\Delta > 0$, the mechanism that adds Gaussian noise with variance $(8k\Delta^2 \log(e + (\varepsilon/\delta)))/\varepsilon^2$ satisfies (ε, δ) -differential privacy under k -fold adaptive composition for any $\varepsilon > 0$ and $\delta \in (0, 1]$. The mean squared error achieved by the Gaussian mechanism is $O(k^2\Delta^2 \log(e + (\varepsilon/\delta)))/\varepsilon^2$.*

It is previously known that it is sufficient to add i.i.d. Gaussian noise with variance $O(k\Delta^2 \log(1/\delta)/\varepsilon^2)$ to ensure (ε, δ) -differential privacy under k -fold composition (e.g. Theorem 2.7 from Hardt & Talwar (2010)). The above theorem shows that when $\delta = \Theta(\varepsilon)$, one can achieve the same privacy with smaller variance by a factor of $\log(1/\delta)$.

4.2. Cut queries of a graph and variance queries of a matrix

Blocki et al. (2012) showed that classical Johnson-Lindenstrauss transform can generate a differentially private version of a database. Further, they show that this achieves the best tradeoff between privacy and utility for two applications: cut queries of a graph and variance queries of a matrix. In this section, we show how the best known trade off can be further improved by applying Theorem 3.4.

First, Blocki et. al. provide a differentially private mechanism for cut queries $q(G, S)$: the number of edges crossing a (S, \bar{S}) -cut in a weighted undirected graph G . This mechanism produces a sanitized graph satisfying (ε, δ) -differential privacy, where two graphs are neighbors if they only differ on a single edge. The *utility* of the mechanism is measured via the additive error τ incurred by the

privatization. Precisely, a mechanism M is said to give a (η, τ, ν) -approximation for a *single* cut query $q(\cdot, \cdot)$, if for every graph G and every nonempty S it holds that

$$\mathbb{P}\left((1 - \eta)q(G, S) - \tau \leq M(G, S) \leq (1 + \eta)q(G, S) + \tau\right) \geq 1 - \nu. \quad (9)$$

For the proposed Johnson-Lindenstrauss mechanism satisfying (ε, δ) -differential privacy, it is shown that the additive error τ_0 incurred by querying the database k times is bounded by Theorem 3.2 from Blocki et al. (2012)¹

$$\tau_0 = O\left(|S| \frac{\sqrt{\log(1/\delta) \log(k/\nu)}}{\varepsilon} \log\left(\frac{\log(k/\nu)}{\eta^2 \delta}\right)\right). \quad (10)$$

Compared to other state-of-the-art privacy mechanisms such as the Laplace noise adding mechanism (Dwork, 2006), exponential mechanism (McSherry & Talwar, 2007), multiplicative weights (Hardt & Rothblum, 2010), and Iterative Database Construction (Gupta et al., 2012), it is shown in (Blocki et al., 2012) that the Johnson-Lindenstrauss mechanism achieves the best tradeoff between the additive error τ_0 and the privacy ε . This tradeoff in (10) is proved using the existing Theorem 3.2. We can improve this analysis using the optimal composition theorem of Theorem 3.4, which gives

$$\tau = O\left(|S| \frac{\sqrt{\log(e + \varepsilon/\delta) \log(k/\nu)}}{\varepsilon} \log\left(\frac{\log(k/\nu)}{\eta^2 \delta}\right)\right). \quad (11)$$

This is smaller than (10) by (a square root of) a logarithmic factor when $\varepsilon = \Theta(\delta)$. We refer to the supplementary material for a proof of the analysis in (11).

A similar technique has been used in Blocki et al. (2012) to provide a differentially private mechanism for variance queries $v(A, x) = x^T A^T A x$: the variance of a given matrix in a direction x . The proposed mechanism produces a sanitized covariance matrix that satisfy (ε, δ) -differential privacy, where two matrices are neighbors if they differ only in a single row and the difference is by Euclidean distance at most one. With the previous composition theorem in Theorem 3.2, Blocki et al. (2012) get an error bound $\tau_1 = O\left(\frac{\log(1/\delta) \log(k/\nu)}{\varepsilon^2 \eta} \log^2\left(\frac{\log(k/\nu)}{\eta^2 \delta}\right)\right)$. Using our tight composition theorem, this can be improved as $\tau = O\left(\frac{\log(e + \varepsilon/\delta) \log(k/\nu)}{\varepsilon^2 \eta} \log^2\left(\frac{\log(k/\nu)}{\eta^2 \delta}\right)\right)$. Again, for $\varepsilon = \Theta(\delta)$, this is an improvement of a logarithmic factor.

4.3. Optimality of geometric noise adding

In this section, we consider integer valued queries $q : \mathcal{D} \rightarrow \mathbb{Z}$ with sensitivity one, also called *counting queries*. Such

¹The original theorem is stated for a single query with $k = 1$. Here we state it more generally with arbitrary k . This requires scaling ν by $1/k$ to take into account the union bound over k query outputs in the *utility* guarantee in (9).

queries are common in practice, e.g. ‘‘How many individuals have income less than \$100,000?’’. Presence of absence of an individual record changes the output at most by one. Counting query is a well-studied topic in differential privacy (Dinur & Nissim, 2003; Dwork & Nissim, 2004; Blum et al., 2005; 2013) and they provide a primitive for constructing more complex queries (Blum et al., 2005).

The *geometric* noise adding mechanism is a discrete variant of the popular Laplacian mechanism. For integer-valued queries with sensitivity one, the mechanism adds a noise distributed according to a double-sided geometric distribution whose probability density function is $p(k) = ((e^\varepsilon - 1)/(e^\varepsilon + 1))e^{-\varepsilon|k|}$. This mechanism is known to be universally optimal in a general cost minimization framework (Bayesian setting in (Ghosh et al., 2012) and worst-case setting in (Geng & Viswanath, 2012)). In this section, we show that the geometric noise adding mechanism is also optimal under composition.

Consider the composition experiment for counting queries. For a pair of neighboring databases D_0 and D_1 , some of the query outputs differ by one, since sensitivity is one, and for other queries the output might be the same. Let k denote the number of queries whose output differs with respect to D_0 and D_1 . Then, we show that the privacy region achieved by geometric mechanism, that adds geometric noise for each integer-valued query output, is exactly described by the optimal composition theorem of (5). Further, since this is the largest privacy region under composition for the pair of database D_0 and D_1 that differ in k queries, no other mechanism can achieve a larger privacy region. Since the geometric mechanism does not depend on the particular choice of pairs of databases D_0 and D_1 , nor does it depend on the specific query being asked, the mechanism is *optimal* universally for every pair of neighboring databases simultaneously.

Here, optimality is with respect to the composed privacy region itself. Among the mechanisms guaranteeing the same level of privacy, one with larger privacy region under composition is considered better, in terms of allowing for smaller false alarm and missed detection rate in hypothesis testing whether the database contains a particular entry or not. In this sense, larger privacy degradation under composition has more utility. The geometric mechanism has the largest possible privacy degradation under composition, stated formally below; we refer to the supplementary material for a proof.

Theorem 4.4. *Under the k -fold composition experiment of counting queries, the geometric mechanism achieves the largest privacy region among all $(\varepsilon, 0)$ -differentially private mechanisms, universally for every pair of neighboring databases simultaneously.*

References

- Blachman, N. The convolution inequality for entropy powers. *Information Theory, IEEE Transactions on*, 11(2): 267–271, 1965.
- Blackwell, David. Equivalent comparisons of experiments. *The Annals of Mathematical Statistics*, 24(2):265–272, 1953.
- Blocki, Jeremiah, Blum, Avrim, Datta, Anupam, and Sheffet, Or. The johnson-lindenstrauss transform itself preserves differential privacy. In *Foundations of Computer Science (FOCS), 2012 IEEE 53rd Annual Symposium on*, pp. 410–419. IEEE, 2012.
- Blum, Avrim, Dwork, Cynthia, McSherry, Frank, and Nissim, Kobbi. Practical privacy: the SuLQ framework. In *Proceedings of the twenty-fourth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pp. 128–138. ACM, 2005.
- Blum, Avrim, Ligett, Katrina, and Roth, Aaron. A learning theory approach to noninteractive database privacy. *Journal of the ACM (JACM)*, 60(2):12, 2013.
- Cover, Thomas M and Thomas, A. Determinant inequalities via information theory. *SIAM journal on Matrix Analysis and Applications*, 9(3):384–392, 1988.
- Cover, Thomas M and Thomas, Joy A. *Elements of information theory*. John Wiley & Sons, 2012.
- Dembo, Amir, Cover, Thomas M, and Thomas, Joy A. Information theoretic inequalities. *Information Theory, IEEE Transactions on*, 37(6):1501–1518, 1991.
- Dinur, Irit and Nissim, Kobbi. Revealing information while preserving privacy. In *Proceedings of the twenty-second ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pp. 202–210. ACM, 2003.
- Dwork, Cynthia. Differential privacy. In *Automata, languages and programming*, pp. 1–12. Springer, 2006.
- Dwork, Cynthia and Lei, Jing. Differential privacy and robust statistics. In *Proceedings of the 41st annual ACM symposium on Theory of computing*, pp. 371–380. ACM, 2009.
- Dwork, Cynthia and Nissim, Kobbi. Privacy-preserving datamining on vertically partitioned databases. In *Advances in Cryptology-CRYPTO 2004*, pp. 528–544. Springer, 2004.
- Dwork, Cynthia, Kenthapadi, Krishnam, McSherry, Frank, Mironov, Ilya, and Naor, Moni. Our data, ourselves: Privacy via distributed noise generation. In *Advances in Cryptology-EUROCRYPT 2006*, pp. 486–503. Springer, 2006a.
- Dwork, Cynthia, McSherry, Frank, Nissim, Kobbi, and Smith, Adam. Calibrating noise to sensitivity in private data analysis. In *Theory of Cryptography*, pp. 265–284. Springer, 2006b.
- Dwork, Cynthia, Rothblum, Guy N, and Vadhan, Salil. Boosting and differential privacy. In *Foundations of Computer Science (FOCS), 2010 51st Annual IEEE Symposium on*, pp. 51–60. IEEE, 2010.
- Geng, Quan and Viswanath, Pramod. Optimal noise-adding mechanism in differential privacy. *arXiv preprint arXiv:1212.1186*, 2012.
- Geng, Quan and Viswanath, Pramod. The optimal mechanism in (ϵ, δ) -differential privacy. *arXiv preprint arXiv:1305.1330*, 2013.
- Ghosh, Arpita, Roughgarden, Tim, and Sundararajan, Mukund. Universally utility-maximizing privacy mechanisms. *SIAM Journal on Computing*, 41(6):1673–1693, 2012.
- Green, Ben and Tao, Terence. The primes contain arbitrarily long arithmetic progressions. *arXiv preprint math/0404188*, 2004.
- Gupta, Anupam, Roth, Aaron, and Ullman, Jonathan. Iterative constructions and private data release. In *Theory of Cryptography*, pp. 339–356. Springer, 2012.
- Hardt, Moritz and Roth, Aaron. Beyond worst-case analysis in private singular vector computation. In *Proceedings of the 45th annual ACM symposium on Symposium on theory of computing*, pp. 331–340. ACM, 2013.
- Hardt, Moritz and Rothblum, Guy N. A multiplicative weights mechanism for privacy-preserving data analysis. In *Foundations of Computer Science (FOCS), 2010 51st Annual IEEE Symposium on*, pp. 61–70. IEEE, 2010.
- Hardt, Moritz and Talwar, Kunal. On the geometry of differential privacy. In *Proceedings of the 42nd ACM symposium on Theory of computing*, pp. 705–714. ACM, 2010.
- Hardt, Moritz, Ligett, Katrina, and McSherry, Frank. A simple and practical algorithm for differentially private data release. *arXiv preprint arXiv:1012.4763*, 2010.
- Liu, Tie and Viswanath, Pramod. An extremal inequality motivated by multiterminal information-theoretic problems. *Information Theory, IEEE Transactions on*, 53(5): 1839–1851, 2007.
- McSherry, Frank and Talwar, Kunal. Mechanism design via differential privacy. In *Foundations of Computer Science, 2007. FOCS'07. 48th Annual IEEE Symposium on*, pp. 94–103. IEEE, 2007.

- Muthukrishnan, S and Nikolov, Aleksandar. Optimal private halfspace counting via discrepancy. In *Proceedings of the 44th symposium on Theory of Computing*, pp. 1285–1292. ACM, 2012.
- Reingold, Omer, Trevisan, Luca, Tulsiani, Madhur, and Vadhan, Salil. Dense subsets of pseudorandom sets. In *Foundations of Computer Science, 2008. FOCS'08. IEEE 49th Annual IEEE Symposium on*, pp. 76–85. IEEE, 2008.
- Stam, AJ. Some inequalities satisfied by the quantities of information of fisher and shannon. *Information and Control*, 2(2):101–112, 1959.
- Tao, Terence and Ziegler, Tamar. The primes contain arbitrarily long polynomial progressions. *Acta Mathematica*, 201(2):213–305, 2008.
- Verdú, Sergio and Guo, Dongning. A simple proof of the entropy-power inequality. *IEEE Transactions on Information Theory*, 52(5):2165–2166, 2006.
- Wasserman, Larry and Zhou, Shuheng. A statistical framework for differential privacy. *Journal of the American Statistical Association*, 105(489):375–389, 2010.
- Zamir, Ram. A proof of the fisher information inequality via a data processing argument. *Information Theory, IEEE Transactions on*, 44(3):1246–1250, 1998.