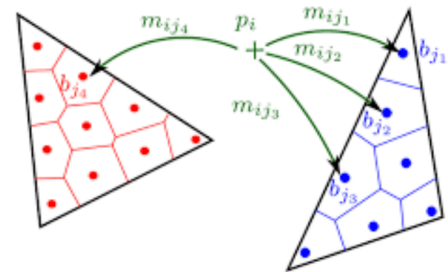




# REPAS

RELIABLE AND  
PRIVACY-AWARE  
SOFTWARE SYSTEMS



Deliverable D3.b

# 1 Introduction

The enormous growth in the use of internet-connected devices and the big-data revolution have created serious privacy concerns, and motivated an intensive area of research aimed at devising methods to protect the users' sensitive information. During the last decade, two main frameworks have emerged in this area: Differential Privacy (DP) and Quantitative Information Flow (QIF).

Differential privacy (DP) [10] was originally developed in the area of statistical databases, and it aims at protecting the individuals' data while allowing the release of aggregate information through queries. This is obtained by *obfuscating* the result of the query via the addition of controlled noise. Naturally, we need to assume that the *curator*, namely the entity collecting and storing the data and handling the queries, is honest and capable of protecting the data from security breaches. Since this assumption cannot always be guaranteed, a variant has been proposed: local differential privacy (LDP) [9], where the data are obfuscated individually before they are collected.

Both DP and LDP are subsumed by **d**-privacy [6], and in this paper we will use the latter as a unifying framework. The definition of **d**-privacy assumes an underlying metric structure on the data domain  $\mathcal{X}$ . An obfuscation mechanism  $K$  for  $\mathcal{X}$  is a probabilistic mapping from  $\mathcal{X}$  to some output domain  $\mathcal{Y}$ , namely a function from  $\mathcal{X}$  to probabilistic distributions over  $\mathcal{Y}$ . We will use the notation  $K_{x,y}$  to represent the probability that  $K$  on input  $x$  gives output  $y$ . The mechanism  $K$  is  $\varepsilon$ -**d**-private, where  $\varepsilon$  is a parameter representing the privacy level, if

$$K_{x_1,y} \leq e^{\varepsilon \mathbf{d}(x_1,x_2)} K_{x_2,y} \quad \text{for all } x_1, x_2 \in \mathcal{X}, y \in \mathcal{Y}. \quad (1)$$

Standard DP is obtained from this definition by assuming  $\mathcal{X}$  to be a set of all datasets and **d** the Hamming distance between two datasets, seen as vectors or records (i.e., the number of positions in which the two datasets differ)<sup>1</sup>. LDP is obtained by considering a trivial metric (distance 0 between identical elements, and 1 otherwise).

The other framework for the protection of sensitive information, quantitative information flow (QIF), focuses on the potentialities and the goals of the attacker, and the research on this area has developed rigorous foundations based on information theory [15, 21]. The idea is that a system processing some sensitive data from a random variable  $X$  and releasing some observable data as a random variable  $Y$  can be modelled as an information-theoretic channel with input  $X$  and output  $Y$ . The leakage is then measured in terms of correlation between  $X$  and  $Y$ . There are, however, many different ways to define such correlation, depending on the notion of adversary. In order to provide a unifying approach, [2] has proposed the theory of  $g$ -leakage, in which an adversary is characterized by a functional parameter  $g$  representing its *gain* for each possible outcomes of the attack.

---

<sup>1</sup>The more common definition of differential privacy assumes that  $x_1, x_2$  are adjacent, i.e. their Hamming distance is 1, and requires  $K_{x_1,y} \leq e^\varepsilon K_{x_2,y}$ . It is easy to prove that the two definitions are equivalent.

$P_1$	$P_2$	$P_3$	$P_4$
<pre> if H mod 8 = 0 then   L := H else   L := 1 </pre>	<pre> if H mod 4 = 0 then   L := H else   L := 1 </pre>	<pre> L := H &amp; 0<sup>24</sup>1<sup>8</sup> </pre>	<pre> L := H &amp; 0<sup>28</sup>1<sup>4</sup> </pre>

Table 1: Programs that take in input a secret  $H$  and leak information about  $H$  into the output  $L$ .

One issue that arises in both frameworks is how to compare systems from the point of view of their privacy guarantees. It is important to have a rigorous and effective way to establish whether a mechanism is better or worse than another one, in order to guide the design and the implementation of mechanisms for information protection. This is not always an obvious task. To illustrate the point, consider the following examples.

**Example 1.** Let  $P_1, P_2, P_3$  and  $P_4$  be the programs illustrated in Table 1, where  $H$  is a “high” (i.e., secret) input and  $L$  is a “low” (i.e., public) output. We assume that  $H$  is a uniformly distributed 32-bit integer with range  $0 \leq H < 2^{32}$ . All these programs leak information about  $H$  via  $L$ , in different ways:  $P_1$  reveals  $H$  whenever it is a multiple of 8 ( $H \bmod 8$  represents the integer division of  $H$  by 8), and reveals nothing otherwise.  $P_2$  does the same thing whenever  $H$  is a multiple of 4.  $P_3$  reveals the last 8 bits of  $H$  ( $H \& 0^{24}1^8$  represents the bitwise and between  $H$  and a string of 24 bits “0” followed by 8 bits “1”). Analogously,  $P_4$  reveals the last 4 bits of  $H$ . Now, it is clear that  $P_2$  leaks more than  $P_1$ , and that  $P_4$  leaks more than  $P_3$ , but how to compare, for instance,  $P_1$ , and  $P_3$ ? It is debatable which one is worse, because their behavior is very different:  $P_1$  reveals nothing in most cases, but when it does reveal something, it reveals everything.  $P_3$ , on the other hand, always reveals part of the secret. Clearly, we cannot decide which situation is worse, unless we have some more information about the goals and the capabilities of the attacker. For instance, if the adversary has only one attempt at his disposal (and no extra information), then the program  $P_3$  is better, because even after the output of  $L$  there are still 24 bits of  $H$  that are unknown. On the other hand, if the adversary can repeat the attacks on program similar to  $P_3$ , then eventually it will uncover the secret entirely all the times.

**Example 2.** Consider a geometric mechanism (cfr. Definition 9) with  $\epsilon^\epsilon = 27/25$ , and a randomized response one mechanism (cfr. Definition 13) with  $\epsilon^\epsilon = 22/11$ . The two mechanisms are illustrated in Figure ???. Clearly, it does not make sense to compare them on the basis of their respective privacy parameters  $\epsilon$ , because they represent different privacy properties, and it is not obvious how to compare them in general: The geometric mechanism tends to make the true value indistinguishable from his immediate neighbors, but it separates it from

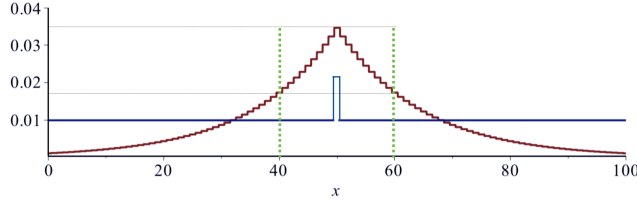


Figure 1: Comparison between the geometric (red) and the randomized response (blue) mechanisms.

the values far away. The randomized response introduces the same level of confusion between the true value and any other value of the domain. Thus, which mechanism is more private depends on the kind of attack we want to mitigate: if the attacker is trying to guess an approximation of the value, then the randomized response is better. If the attacker is only interested in identifying the true value among the immediate neighbors, then the geometric is better.

In this respect, the QIF approach has lead to an elegant theory of refinement (pre)order<sup>2</sup>  $\sqsubseteq_{\mathbb{G}}^{\text{avg}}$ , which provides strong guarantees:  $A \sqsubseteq_{\mathbb{G}}^{\text{avg}} B$  means that  $B$  is safer than  $A$  in all circumstances, in the sense that the *expected gain* of an attack on  $B$  is less than on  $A$ , for whatever kind of gain the attacker may be seeking. This means that we can always substitute the component  $A$  by  $B$  without compromising the security of the system. An appealing aspect of this particular refinement order is that it is characterized by a precise structural relation between the stochastic channels associated to  $A$  and  $B$  [2, 19], which makes it easy to reason about, and relatively efficient to verify. It is important to remark that this order is based on an *average* notion of adversarial gain (*vulnerability*), defined by mediating over all possible observations and their probabilities. We call this perspective *average-case*.

At the other end of the spectrum, DP, LDP and **d**-privacy are *max-case* measures. In fact, by applying the Bayes theorem to (1) we obtain:

$$\frac{p(x_1 | y)}{p(x_2 | y)} \leq e^{\varepsilon \mathbf{d}(x_1, x_2)} \frac{\pi(x_1)}{\pi(x_2)} \quad \text{for all } x_1, x_2 \in \mathcal{X}, y \in \mathcal{Y}. \quad (2)$$

where, for  $i \in \{1, 2\}$ ,  $\pi(x_i)$  is the *prior* probability of  $x_i$  and  $p(x_i | y)$  is the *posterior* probability of  $x_i$  given  $y$ . We can interpret  $\pi(x_1)/\pi(x_2)$  and  $p(x_1|y)/p(x_2|y)$  as knowledge about  $\mathcal{X}$ : they represent how much more likely  $x_1$  is with respect to  $x_2$ , *before* (prior) and *after* (posterior) observing  $y$ , respectively. Thus the property expresses a bound on how much the adversary can learn from each individual outcome of the mechanism<sup>3</sup>.

<sup>2</sup>In this paper we call  $\sqsubseteq_{\mathbb{G}}^{\text{avg}}$  and the other refinement relations “orders”, although, strictly speaking they are preorders.

<sup>3</sup>The property (2) is also the semantic interpretation of the guarantees of the Pufferfish framework (cfr. [14], Section 3.1). The ratio  $\frac{p(x_1|y)}{p(x_2|y)} / \frac{\pi(x_1)}{\pi(x_2)}$  is known as *odds ratio*.

In the literature of DP, LDP and  $\mathbf{d}$ -privacy, mechanisms are usually compared on the basis of their  $\varepsilon$ -value<sup>4</sup>, which controls a bound on the log-likelihood ratio of an observation  $y$  given two “secrets”  $x_1$  and  $x_2$ : smaller  $\varepsilon$  means more privacy. In DP and LDP the bound is  $\varepsilon$  itself, while in  $\mathbf{d}$ -privacy it is  $\varepsilon \times \mathbf{d}(s_1, s_2)$ . We remark that the relation induced by  $\varepsilon$  in  $\mathbf{d}$ -privacy is fragile, in the sense that the definition of  $\mathbf{d}$ -privacy assumes an underlying metric structure  $\mathbf{d}$  on the data, and whether a mechanism  $B$  is “better” than  $A$  depends in general on the metric considered.

Average-case and max-case are different principles, suitable for different scenarios: the former represent the point of view of an organization, for instance an insurance company providing coverage for risks related to credit cards, which for the cost-benefit analysis is interested in reasoning in terms of expectation (expected cost of an attack). The max-case represents the point of view of an individual, who is interested in limiting the cost of *any* attack. As such, the max-case seems particularly suitable for the domain of privacy.

In this paper, we combine the max-case perspective with the robustness of the QIF approach, and we introduce two refinement orders:

- $\sqsubseteq_{\mathbb{Q}}^{\max}$ , based on the max-case leakage introduced in [1]. This order takes into account all possible privacy breaches caused by any observable (like in the DP world), but it quantifies over all possible quasi-convex vulnerability functions (in the style of the QIF world).
- $\sqsubseteq_{\mathbb{M}}^{\text{priv}}$ , based on  $\mathbf{d}$ -privacy (like in the DP world), but quantified over all metrics  $\mathbf{d}$ .

To underline the importance of a robust order, let us consider the case of the oblivious mechanisms for differential privacy: These mechanisms are of the form  $K = H \circ f$ , where  $f : \mathcal{X} \rightarrow \mathcal{Y}$  is a query, namely a function from datasets in  $\mathcal{X}$  to some answer domain  $\mathcal{Y}$ , and  $H$  is a probabilistic mechanism implementing the noise. The idea is that the system first computes the result  $y \in \mathcal{Y}$  of the query (*true answer*), and then it applies  $H$  to  $y$  to obtain a *reported answer*  $z$ . In general, if we want  $K$  to be  $\varepsilon$ -DP, we need to tune the mechanism  $H$  so to take into account the *sensitivity* of  $f$ , which is the maximum distance between the results of  $f$  on two adjacent databases, and as such it depends on the metric on  $\mathcal{Y}$ . However, if we know that  $K = H \circ f$  is  $\varepsilon$ -DP, and that  $H \sqsubseteq_{\mathbb{M}}^{\text{priv}} H'$  for some other mechanism  $H'$ , then we can safely substitute  $H$  by  $H'$  as it is, because one of our results (cfr. Theorem 1 at Page 33) guarantees that  $K' = H' \circ f$  is also  $\varepsilon$ -DP. In other words,  $H \sqsubseteq_{\mathbb{M}}^{\text{priv}} H'$  implies that we can substitute  $H$  by  $H'$  in an oblivious mechanism for whatever query  $f$  and whatever metric on  $\mathcal{Y}$ , without the need to know the sensitivity of  $f$  and without the need to do any tuning of  $H'$ . Thanks to Theorem 3 (Page 13) and Theorem 5 (Page 17), we know that this is the case also for  $\sqsubseteq_{\mathbb{G}}^{\text{avg}}$  and  $\sqsubseteq_{\mathbb{Q}}^{\max}$ . We illustrate this with the following example.

---

<sup>4</sup>In DP and LDP  $\varepsilon$  is a parameter that usually appears explicitly in the definition of the mechanism. In  $\mathbf{d}$ -privacy it is an implicit scaling factor.

**Example 3.** Consider datasets  $x \in \mathcal{X}$  of records containing the age of people, expressed as natural numbers from 0 to 100, and assume that each dataset in  $\mathcal{X}$  contains at least 100 records. Consider two queries,  $f(x)$  and  $g(x)$ , which give the rounded average age and the minimum age of the people in  $x$ , respectively. Finally, consider the truncated geometric mechanism  $TG^\varepsilon$  (cfr. Definition 10, Page 21), and the randomized response mechanism  $R^\varepsilon$  (cfr. Definition 13, Page 22). It is easy to see that  $K_1 = TG^\varepsilon \circ f$  is  $\varepsilon$ -DP, and it is possible to prove that  $TG^\varepsilon \sqsubseteq_{\mathbb{M}}^{\text{prv}} R^\varepsilon$  (cfr. Theorem 15, Page 15). We can then conclude that  $K_2 = R^\varepsilon \circ f$  is  $\varepsilon$ -DP as well, and that in general it is safe to replace  $TG^\varepsilon$  by  $R^\varepsilon$  for whatever query. On the other hand,  $R^\varepsilon \not\sqsubseteq_{\mathbb{M}}^{\text{prv}} TG^\varepsilon$ , so we cannot expect that it is safe to replace  $R^\varepsilon$  by  $TG^\varepsilon$  in any context. In fact,  $K_3 = R^\varepsilon \circ g$  is  $\varepsilon$ -DP, but  $K_4 = TG^\varepsilon \circ g$  is not  $\varepsilon$ -DP, despite the fact that both mechanisms are constructed using the same privacy parameter  $\varepsilon$ . Hence we can conclude that a refinement relation based only on the comparison of the  $\varepsilon$  parameters would not be robust, at least not for a direct replacement in an arbitrary context. Note that  $K_4$  is  $100 \times \varepsilon$ -DP. In order to make it  $\varepsilon$ -DP we should divide the parameter  $\varepsilon$  by the sensitivity of  $g$  (with respect to the ordinary distance on natural numbers), which is 100, i.e. use  $TG^{\varepsilon/100}$ . For  $R^\varepsilon$  this is not necessary because it is defined using the discrete metric on  $\{0, \dots, 100\}$ , and the sensitivity of  $g$  with respect to this metric is 1.

The robust orders allow us to take into account different kinds of adversaries. The following example shows what is the idea.

**Example 4.** Consider the following three LDP mechanisms, represented by their stochastic matrices (where each element is the conditional probability of the outcome of the mechanism, given the secret value). The secrets are three possible economic situations of an individual,  $p$ ,  $a$  and  $r$ , standing for poor, average and rich, respectively. The observable outcomes are  $n$  and  $r$ , standing for rich and not rich.

$$\begin{array}{|c|c|c|} \hline A & n & r \\ \hline p & 3/4 & 1/4 \\ a & 1/2 & 1/2 \\ r & 1/4 & 3/4 \\ \hline \end{array}
 \quad
 \begin{array}{|c|c|c|} \hline B & n & r \\ \hline p & 2/3 & 1/3 \\ a & 2/3 & 1/3 \\ r & 1/3 & 2/3 \\ \hline \end{array}
 \quad
 \begin{array}{|c|c|c|} \hline C & n & r \\ \hline p & 2/3 & 1/3 \\ a & 1/2 & 1/2 \\ r & 1/3 & 2/3 \\ \hline \end{array}
 \tag{3}$$

Let us assume that the prior distribution  $\pi$  on the secrets is uniform. We note that  $A$  is  $(\log 3)$ -LDP while  $B$  is  $(\log 2)$ -LDP. Hence, if we only look at the value of  $\varepsilon$ , we would think that  $B$  is better than  $A$  from the privacy point of view. However, there are attackers that gain more from  $B$  than from  $A$  (which means that, with respect to those attackers, the privacy of  $B$  is worse). For instance, this is the case when the attacker is only interested in discovering whether the person is rich or not. In fact, if we consider a gain 1 when the attacker guesses the right class ( $r$  versus (either  $p$  or  $a$ )) and 0 otherwise, we have that the highest possible gain in  $A$  is  $(3/4)\pi(p) + (1/2)\pi(a) = 5/12$ , while in  $B$  is  $(2/3)\pi(p) + (2/3)\pi(a) = 4/9$ , which is higher than  $5/12$ . This is consistent with our orders: it is possible to show that none of the three orders hold between

$A$  and  $B$ , and that therefore we should not expect  $B$  to be better (for privacy) than  $A$  with respect to all possible adversaries.

On the other hand, the mechanism  $C$  is also  $(\log 2)$ -LDP, and in this case we have that the relation  $A \sqsubseteq_{\mathbb{Q}}^{\max} C$  holds, implying that we can safely replace  $A$  by  $C$ . We can also prove that the reverse does not hold, which means that  $C$  is strictly better than  $A$ .

A fundamental issue is how to prove that these robust orders hold: Since  $\sqsubseteq_{\mathbb{Q}}^{\max}$  and  $\sqsubseteq_{\mathbb{M}}^{\text{prv}}$  involve universal quantifications, it is important to devise finitary methods to verify them. To this purpose, we will study their characterizations as structural relations between stochastic matrices (representing the mechanisms to be compared), along the lines of what was done for  $\sqsubseteq_{\mathbb{G}}^{\text{avg}}$ .

We will also study the relation between the three orders (the two above and  $\sqsubseteq_{\mathbb{G}}^{\text{avg}}$ ), and their algebraic properties. Finally, we will analyze various mechanisms for DP, LDP, and  $\mathbf{d}$ -privacy to see in which cases the order induced by  $\varepsilon$  is consistent with the three orders above.

## 1.1 Contribution

The main contributions of this paper are the following:

- We introduce two refinement orders for the max case,  $\sqsubseteq_{\mathbb{Q}}^{\max}$  and  $\sqsubseteq_{\mathbb{M}}^{\text{prv}}$ , that are robust with respect to a large class of adversaries.
- We give structural characterizations of both  $\sqsubseteq_{\mathbb{Q}}^{\max}$  and  $\sqsubseteq_{\mathbb{M}}^{\text{prv}}$  in terms of relations on the stochastic matrices of the mechanisms under comparison. These relations help the intuition and open the way to verification.
- We study efficient methods to verify the structural relations above. Furthermore, these methods are such that, when the verification fails, they produce counterexamples. In this way it is possible to pin down what is the problem and try to correct it.
- We show that  $\sqsubseteq_{\mathbb{G}}^{\text{avg}} \subset \sqsubseteq_{\mathbb{Q}}^{\max} \subset \sqsubseteq_{\mathbb{M}}^{\text{prv}}$ .
- We apply the three orders ( $\sqsubseteq_{\mathbb{G}}^{\text{avg}}$ ,  $\sqsubseteq_{\mathbb{Q}}^{\max}$ , and  $\sqsubseteq_{\mathbb{M}}^{\text{prv}}$ ) to the comparison of some well-known families of  $\mathbf{d}$ -private mechanisms: geometric, exponential and randomised response. We show that, in general,  $A \sqsubseteq_{\mathbb{G}}^{\text{avg}} B$  (and thus all the refinement orders between  $A$  and  $B$ ) holds within the same family whenever the  $\varepsilon$  of  $B$  is smaller than that of  $A$ .
- We show that if  $A$  and  $B$  are mechanisms from different families, then, even if the  $\varepsilon$  of  $B$  is smaller than that of  $A$ , the relations  $A \sqsubseteq_{\mathbb{G}}^{\text{avg}} B$  and  $A \sqsubseteq_{\mathbb{Q}}^{\max} B$  do not hold, and in most cases  $A \sqsubseteq_{\mathbb{M}}^{\text{prv}} B$  does not hold either. We conclude that a comparison based only on the value of the  $\varepsilon$ 's is not robust across different families, at least not for the purposes illustrated above.

- We study lattice-properties of these orders. In contrast to  $\sqsubseteq_{\mathbb{G}}^{\text{avg}}$ , which was shown to not be a lattice, we prove that suprema and infima exist for  $\sqsubseteq_{\mathbb{Q}}^{\text{max}}$  and  $\sqsubseteq_{\mathbb{M}}^{\text{prv}}$ , and that therefore these orders form lattices.

## 1.2 Related work

We are not aware of many studies on refinement relations for QIF. Yasuoka and Terauchi [23] and Malacaria [18] have explored strong orders on *deterministic* mechanisms, focusing on the fact that such mechanisms induce *partitions* on the space of secrets. They showed that the orders produced by min-entropy leakage [21] and Shannon leakage [8, 17] are the same and, moreover, they coincide with the *partition refinement* order in the *Lattice of Information* [16]. This order was extended to the probabilistic case in [2], resulting in the relation  $\sqsubseteq^{\text{avg}}$  mentioned in Section 2. The same paper [2] proposed the theory of *g*-leakage and introduced the corresponding order  $\sqsubseteq_{\mathbb{G}}^{\text{avg}}$ . Furthermore, [2] proved that  $\sqsubseteq^{\text{avg}} \subseteq \sqsubseteq_{\mathbb{G}}^{\text{avg}}$  and conjectured that also the reverse should hold. This conjecture was then proved valid in [19]. The max-case leakage, on which the relation  $\sqsubseteq_{\mathbb{Q}}^{\text{max}}$  is based, was introduced in [1], but  $\sqsubseteq_{\mathbb{Q}}^{\text{max}}$  and its properties were not investigated. Finally,  $\sqsubseteq^{\text{prv}}$  is a novel notion introduced in this paper.

In the field of differential privacy, on the other hand, there have been various works aimed at trying understand the operational meaning of the privacy parameter  $\varepsilon$  and at providing guidelines for the choice of its values. We mention for example [13] and [12], which consider the value of  $\varepsilon$  from an economical point of view, in terms of cost. We are not aware, however, of studies aimed at establishing orders between the level of privacy of different mechanisms, except the one based on the comparison of the  $\varepsilon$ 's.

The relation between QIF and DP, LDP, and **d**-privacy is based on the so-called *semantic interpretation* of the privacy notions, that regards these properties as expressing a bounds on the increase of knowledge (from prior to posterior) due to the answer reported by the mechanism. For **d**-privacy the semantic interpretation is expressed by (2). To the best of our knowledge, this interpretation was first pointed out (for the location privacy instance) in [3]. The seminal paper on **d**-privacy, [6], also proposed a semantic interpretation, with a rather different flavor, although formally equivalent. As for DP, as explained in the introduction, (2) instantiated to databases and Hamming distance corresponds to the odds ratio on which is based the semantics interpretation provided in [14]. Before that, another version of semantic interpretation was presented in [10] and proved equivalent to a form of DP called  $\varepsilon$ -indistinguishability. Essentially, in this version an adversary that queries the database, and knows all the database except one record, cannot infer too much about this record from the answer to the query reported by the mechanism. Later on, an analogous version of semantic interpretation was reformulated in [4] and proved equivalent to DP. A different interpretation of DP, called *semantic privacy*, was proposed by [20]. This interpretation is based on a comparison between two posteriors (rather between the posterior and the prior), and the authors show that, within certain limits, it is equivalent to DP.



A short version of this paper, with only part of the proofs, appeared in [7].

A short version of this paper appeared in the proceedings of the IEEE Computer Security Foundations Symposium [7].

### 1.3 Plan of the paper

In the next three sections, 2, 3 and 4, we define the order refinements  $\sqsubseteq_{\mathbb{G}}^{\text{avg}}$ ,  $\sqsubseteq_{\mathbb{Q}}^{\text{max}}$  and  $\sqsubseteq_{\mathbb{M}}^{\text{prv}}$  respectively, and we study their properties. In Section 5.1 we investigate methods to verify them. In Section 6 we consider various mechanisms for DP and its variants, and we investigate the relation between the parameter  $\varepsilon$  and the orders introduced in this paper. In Section 7 we show that  $\sqsubseteq_{\mathbb{Q}}^{\text{max}}$  and  $\sqsubseteq_{\mathbb{M}}^{\text{prv}}$  form a lattice. Finally, Section 8 concludes.

Note: In order to make the reading of the paper more fluid, we have moved all the proofs to the appendix at the end of the paper.

## 2 Average-case refinement

### 2.1 Vulnerability, channels and leakage

Quantitative Information Flow studies the problem of quantifying the *information leakage* of a system (eg. a program, or an anonymity protocol). A common model in this area is to consider that the user has a secret  $x$  from a finite set of possible secrets  $\mathcal{X}$ , about which the adversary has some probabilistic knowledge  $\pi: \mathbb{D}\mathcal{X}$  ( $\mathbb{D}\mathcal{X}$  denoting the set of probability distributions over  $\mathcal{X}$ ). A function  $V: \mathbb{D}\mathcal{X} \rightarrow \mathbb{R}_{\geq 0}$  is then employed to measure the *vulnerability* of our system:  $V(\pi)$  quantifies the adversary’s success in achieving some desired goal, when his knowledge about the secret is  $\pi$ .

Various such functions can be defined (eg. employing well-known notions of entropy), but it quickly becomes apparent that no single vulnerability function is meaningful for all systems. The family of *g-vulnerabilities* [2] tries to address this issue by parametrizing  $V$  in an operational scenario: first, the adversary is assumed to possess a set of *actions*  $\mathcal{W}$ ; second, a gain function  $g(w, x)$  models the adversary’s gain when choosing action  $w$  and the real secret is  $x$ . *g-vulnerability* can be then defined as the expected gain of an optimal guess:  $V_g(\pi) = \max_{w: \mathcal{W}} \sum_{x: \mathcal{X}} \pi_x g(w, x)$ . Different adversaries can be modelled by proper choices of  $\mathcal{W}$  and  $g$ . We denote by  $\mathbb{G}\mathcal{X}$  the set of all gain functions.

A *system* is then modelled as a *channel*: a probabilistic mapping from the (finite) set of secrets  $\mathcal{X}$  to a finite set of observations  $\mathcal{Y}$ , described by a stochastic matrix  $C$ , where  $C_{x,y}$  is the probability that secret  $x$  produces the observation  $y$ . When the adversary observes  $y$ , he can transform his initial knowledge  $\pi$  into a *posterior* knowledge  $\delta^y: \mathbb{D}\mathcal{X}$ . Since each observation  $y$  is produced with some probability  $a_y$ , it is sometimes conceptually useful to consider that the *result of running a channel*  $C$ , on the initial knowledge  $\pi$ , is a “hyper” distribution  $[\pi, C]$ : a probability distribution on posteriors  $\delta^y$ , each having probability  $a_y$ .

Defn	Vulnerabilities
Equation (4)	$V[\pi, C] := \sum_y a_y V(\delta^y)$
Section 3	$V^{\max}[\pi, C] := \max_y V(\delta^y)$
Definition 8	$V_{\mathbf{d}}(\pi) := \inf\{\varepsilon \geq 0 \mid \forall x, x' \in \mathcal{X}, \pi_x \leq e^{\varepsilon \cdot \mathbf{d}(x, x')} \pi_{x'}\}$
Defn	Leakage Measures
Definition 4	$\text{Priv}_{\mathbf{d}}(C) := \inf\{\varepsilon \geq 0 \mid C \text{ satisfies } \varepsilon \cdot \mathbf{d}\text{-privacy}\}$
Equation (25)	$\mathcal{L}_{\mathbf{d}}^{+, \max}(\pi, C) := V_{\mathbf{d}}^{\max}[\pi, C] - V_{\mathbf{d}}(\pi)$
Equation (26)	$\mathcal{ML}_{\mathbf{d}}^{+, \max}(C) := \max_{\pi} \mathcal{L}_{\mathbf{d}}^{+, \max}(\pi, C)$
Defn	Refinement Orders
Equation (7)	$A \sqsubseteq^{\text{avg}} B$ iff $AR = B$ for some channel $R$
Definition 2	$A \sqsubseteq^{\max} B$ iff $R\tilde{A} = \tilde{B}$ for some channel $R$
Definition 7	$A \sqsubseteq^{\text{prv}} B$ iff $B$ satisfies $\mathbf{d}_A$ -privacy
Defn	Leakage Orders
Section 2.2	$A \sqsubseteq_{\mathbb{G}}^{\text{avg}} B$ iff $\forall g: \mathbb{G}\mathcal{X}, \forall \pi: \mathbb{D}\mathcal{X}, V_g[\pi, A] \geq V_g[\pi, B]$
Definition 1	$A \sqsubseteq_{\mathbb{Q}}^{\max} B$ iff $\forall V: \mathbb{Q}\mathcal{X}, \forall \pi: \mathbb{D}\mathcal{X}, V^{\max}[\pi, A] \geq V^{\max}[\pi, B]$
Definition 6	$A \sqsubseteq_{\mathbb{M}}^{\text{prv}} B$ iff $\forall \mathbf{d} \in \mathbb{M}\mathcal{X}, A \sqsubseteq_{\mathbf{d}}^{\text{prv}} B$
Definition 5	$A \sqsubseteq_{\mathbf{d}}^{\text{prv}} B$ iff $\text{Priv}_{\mathbf{d}}(A) \geq \text{Priv}_{\mathbf{d}}(B)$

Table 2: Definitions and symbols used in this paper.

It is then natural to define the (average-case) *posterior vulnerability* of the system by applying  $V$  to each posterior  $\delta^y$ , then averaging by its probability  $a_y$  of being produced:

$$V[\pi, C] := \sum_y a_y V(\delta^y) . \quad (4)$$

When defining vulnerability in this way, it can be shown [1] that  $V$  has to be *convex on  $\pi$* , otherwise fundamental properties (such as the data processing inequality) are violated. Any continuous and convex function  $V$  can be written as  $V_g$  for a properly chosen  $g$ , so when studying average-case leakage we can safely restrict to using  $g$ -vulnerability.

*Leakage* can be finally defined by comparing the prior and posterior vulnerabilities, eg. as  $\mathcal{L}_g^+(\pi, C) = V_g[\pi, C] - V_g(\pi)$ .<sup>5</sup>

## 2.2 Refinement

A fundamental question arises in the study of leakage: can we guarantee that a system  $B$  is no less safe than a system  $A$ ? Having a family of vulnerability functions, we can naturally define a strong order  $\sqsubseteq_{\mathbb{G}}^{\text{avg}}$  on channels by explicitly requiring that  $B$  leaks<sup>6</sup> no more than  $A$ , for all priors  $\pi$  and all gain functions  $g: \mathbb{G}\mathcal{X}$ :<sup>7</sup>

$$A \sqsubseteq_{\mathbb{G}}^{\text{avg}} B \quad \text{iff} \quad V_g[\pi, A] \geq V_g[\pi, B] \quad (5)$$

$$\text{for all } g: \mathbb{G}\mathcal{X}, \pi: \mathbb{D}\mathcal{X} . \quad (6)$$

Although  $\sqsubseteq_{\mathbb{G}}^{\text{avg}}$  is intuitive and provides clear leakage guarantees, the explicit quantification over vulnerability functions makes it hard to reason about and verify. Thankfully, this order can be characterized in a “structural” way, that is as a direct property of the channel matrix. We first define the *refinement* order  $\sqsubseteq^{\text{avg}}$  on channels by requiring that  $B$  can be obtained by post-processing  $A$  by some other channel  $R$ , that is:

$$A \sqsubseteq^{\text{avg}} B \quad \text{iff} \quad AR = B \text{ for some channel } R . \quad (7)$$

A fundamental result [2, 19] states that  $\sqsubseteq^{\text{avg}}$  and  $\sqsubseteq_{\mathbb{G}}^{\text{avg}}$  coincide.

We read  $A \sqsubseteq^{\text{avg}} B$  as “ $A$  is refined by  $B$ ”, or “ $B$  is as safe as  $A$ ”. When  $A \sqsubseteq^{\text{avg}} B$  holds we have a strong privacy guarantee: we can safely replace  $A$  by  $B$  without decreasing the privacy of the system, independently from the adversary’s goals and his knowledge. But refinement can be also useful in case  $A \not\sqsubseteq^{\text{avg}} B$ ; namely, we can conclude that there must exist some adversary, modelled by a gain function  $g$ , and some initial knowledge  $\pi$ , such that the adversary actually prefers to interact with  $A$  rather than interacting with  $B$ .<sup>8</sup>

<sup>5</sup>Comparing vulnerabilities “multiplicatively” is also possible, but is orthogonal to the goals of this paper.

<sup>6</sup>Note that comparing the leakage of  $A, B$  is equivalent to comparing their posterior vulnerability, so we choose the latter for simplicity.

<sup>7</sup>Note also that quantifying over  $g: \mathbb{G}\mathcal{X}$  is equivalent to quantifying over all continuous and convex vulnerabilities.

<sup>8</sup>Whether this adversary is of practical interest or not is a different issue, but we know that there exists one.

Moreover, we can actually *construct* such a “counter-example” gain function; this is discussed in Section 5.

### 3 Max-case refinement

Although  $\sqsubseteq^{\text{avg}}, \sqsubseteq_{\mathbb{G}}^{\text{avg}}$  provide a strong and precise way of comparing systems, one could argue that average-case vulnerability might underestimate the threat of a system. More precisely, imagine that there is a certain observation  $y$  such that the corresponding posterior  $\delta^y$  is highly vulnerable (eg. the adversary can completely infer the real secret), but  $y$  happens with very small probability  $a_y$ . In this case the average-case posterior vulnerability  $V[\pi, C]$  can be relatively small, although  $V(\delta^y)$  is large for that particular  $y$ .

If such a scenario is considered problematic, we can naturally quantify leakage using a max-case<sup>9</sup> variant of posterior vulnerability, where all observations are treated equally regardless of their probability of being produced:

$$V^{\max}[\pi, C] := \max_y V(\delta^y). \quad (8)$$

Under this definition, it can be shown [1] that  $V$  has to be *quasi-convex* on  $\pi$  (instead of convex), in order to satisfy fundamental properties (such as the data processing inequality). Hence, in the max-case, we no longer restrict to  $g$ -vulnerabilities (which are always convex), but we can use any vulnerability  $V : \mathbb{Q}\mathcal{X}$ , where  $\mathbb{Q}\mathcal{X}$  denotes the set of all continuous quasi-convex functions  $\mathbb{D}\mathcal{X} \rightarrow \mathbb{R}_{\geq 0}$ .

Inspired by  $\sqsubseteq_{\mathbb{G}}^{\text{avg}}$ , we can now define a corresponding max-case leakage order.

**Definition 1.** *The max-case leakage order is defined as*

$$A \sqsubseteq_{\mathbb{Q}}^{\max} B \quad \text{iff} \quad V^{\max}[\pi, A] \geq V^{\max}[\pi, B] \quad \text{for all } V : \mathbb{Q}\mathcal{X}, \pi : \mathbb{D}\mathcal{X}. \quad (9)$$

Similarly to its average-case variant,  $\sqsubseteq_{\mathbb{Q}}^{\max}$  provides clear privacy guarantees by explicitly requiring that  $B$  leaks no more than  $A$  for all adversaries (modelled as a vulnerability  $V$ ). But this explicit quantification make the order hard to reason about and verify. We would thus like to characterize  $\sqsubseteq_{\mathbb{Q}}^{\max}$  by a refinement order that depends only on the structure of the two channels.

Given a channel  $C$  from  $\mathcal{X}$  to  $\mathcal{Y}$ , we denote by  $\tilde{C}$  the channel obtained by normalizing<sup>10</sup>  $C$ ’s columns and then transposing:

$$\tilde{C}_{y,x} := \frac{C_{x,y}}{\sum_x C_{x,y}}. \quad (10)$$

Note that the row  $y$  of  $\tilde{C}$  can be seen as the posterior distribution  $\delta^y$  obtained by  $C$  under the uniform prior. Note also that  $\tilde{C}$  is non-negative and its rows sum up to 1, so it is a valid channel from  $\mathcal{Y}$  to  $\mathcal{X}$ . The average-case refinement order required that  $B$  can be obtained by post-processing  $A$ . We define the *max-case refinement* order by requiring that  $B$  can be obtained by *pre-processing*  $\tilde{A}$ .

<sup>9</sup>Also called “worse”-case in some contexts, although the latter is more ambiguous, “worse” can refer to a variety of factors.

<sup>10</sup>If a column consists of only zeroes it is simply removed.

**Definition 2.** The max-case refinement order is defined as  $A \sqsubseteq^{\max} B$  iff  $R\tilde{A} = \tilde{B}$  for some channel  $R$ .

Our goal now is to show that  $\sqsubseteq_{\mathbb{Q}}^{\max}$  and  $\sqsubseteq^{\max}$  are different characterizations of the same order. To do so, we start by giving a “semantic” characterization of  $\sqsubseteq^{\max}$ , that is, expressing it, not in terms of the channel matrices  $A$  and  $B$ , but in terms of the *posterior distributions* that they produce. Thinking of  $[\pi, C]$  as a (“hyper”) distribution on the posteriors produced by  $\pi$  and  $C$ , its *support*  $\text{supp}[\pi, C]$  is the set of all posteriors produced with non-zero probability. We also denote by  $\text{ch} S$  the *convex hull* of  $S$ .

**Theorem 1.** Let  $\pi: \mathbb{D}\mathcal{X}$ . If  $A \sqsubseteq^{\max} B$  then the posteriors of  $B$  (under  $\pi$ ) are convex-combinations of those of  $A$ , that is

$$\text{supp}[\pi, B] \subseteq \text{ch supp}[\pi, A]. \quad (11)$$

Moreover, if (11) holds and  $\pi$  is full support then  $A \sqsubseteq^{\max} B$ .

Note that if (11) holds for any full-support prior, then it must hold for all priors.

Theorem 1 has a nice geometric intuition (cfr. Figure 2) that we are going to illustrate in the following example.

**Example 5.** Consider the following systems.

$A$	$y_1$	$y_2$	$y_3$	$y_4$	$B$	$y_1$	$y_2$	$y_3$	$y_4$	·	
$x_1$	$1/3$	$2/9$	$2/9$	$2/9$	$x_1$	$1/3$	$2/9$	$2/9$	$1/9$		
$x_2$	$1/9$	$1/3$	$2/9$	$1/3$	$x_2$	$2/9$	$1/3$	$2/9$	$1/9$		
$x_3$	$1/9$	$2/9$	$1/3$	$1/3$	$x_3$	$2/9$	$2/9$	$1/3$	$1/9$		(12)

Consider the prior  $\pi = (1/2, 1/4, 1/4)$ . The set of the posterior distributions generated by  $A$  under  $\pi$  are:

$$\text{supp}[\pi, A] = \{(3/4, 1/8, 1/8), (4/9, 1/3, 2/9), (4/9, 2/9, 1/3), (2/5, 3/10, 3/10)\} \quad (13)$$

while those generated by  $B$  are:

$$\text{supp}[\pi, B] = \{(3/5, 1/5, 1/5), (4/9, 1/3, 2/9), (4/9, 2/9, 1/3), (1/2, 1/4, 1/4)\} \quad (14)$$

These posteriors, and the convex hulls that they generate, are illustrated in Figure 2. The pink area is the  $\text{ch supp}[\pi, A]$  and the purple area is the  $\text{ch supp}[\pi, B]$ . We can see that  $\text{supp}[\pi, B] \subseteq \text{ch supp}[\pi, A]$ , or equivalently,  $\text{ch supp}[\pi, B] \subseteq \text{ch supp}[\pi, A]$ .

We are now ready to give the main result of this section.

**Theorem 2.** The orders  $\sqsubseteq^{\max}$  and  $\sqsubseteq_{\mathbb{Q}}^{\max}$  coincide.

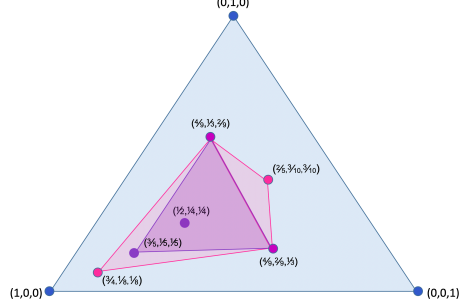


Figure 2: The simplex and the convex hulls of the posterior distributions in Example 5.

Similarly to the average case,  $A \sqsubseteq^{\max} B$  gives us a strong leakage guarantee: can safely replace  $A$  by  $B$ , knowing that for any adversary, the max-case leakage of  $B$  can be no-larger than that of  $A$ . Moreover, in case  $A \not\sqsubseteq^{\max} B$ , we can always find an adversary, modelled by a vulnerability function  $V$ , who prefers (wrt the max-case) interacting with  $A$  than with  $B$ . Such a function is discussed in Section 5.

Finally, we resolve the question of how  $\sqsubseteq^{\max}$  and  $\sqsubseteq^{\text{avg}}$  are related.

**Theorem 3.**  $\sqsubseteq^{\text{avg}}$  is strictly stronger than  $\sqsubseteq^{\max}$ .

This result might appear counter-intuitive at first; one might expect  $A \sqsubseteq^{\max} B$  to imply  $A \sqsubseteq^{\text{avg}} B$ . To understand why it does not, note that the former only requires that, for each output  $y_B$  of  $B$ , there exists some output of  $y_A$  that is at least as vulnerable, regardless of how likely  $y_A$  and  $y_B$  are to happen (this is max-case, after all). We illustrate this with the following example.

**Example 6.** Consider the following systems:

$$\begin{array}{c|cccc} A & y_1 & y_2 & y_3 & y_4 \\ \hline x_1 & 3/4 & 0 & 1/4 & 0 \\ x_2 & 3/4 & 1/4 & 0 & 0 \\ x_3 & 0 & 1/4 & 1/4 & 1/2 \end{array} \quad \cdot \quad \begin{array}{c|ccc} B & y_1 & y_2 & y_3 \\ \hline x_1 & 1/2 & 0 & 1/2 \\ x_2 & 1/2 & 1/2 & 0 \\ x_3 & 0 & 1/2 & 1/2 \end{array} \quad (15)$$

Under the uniform prior, the  $y_1, y_2, y_3$  posteriors for both channels are the same, namely  $(1/2, 1/2, 0)$ ,  $(0, 1/2, 1/2)$  and  $(1/2, 0, 1/2)$  respectively. So the knowledge that can be obtained by each output of  $B$  can be also obtained by some output of  $A$  (albeit with a different probability). Hence, from Theorem. 1 we get that  $A \sqsubseteq^{\max} B$ . However, we can check (see Section 5.1) that  $B$  cannot be obtained by post-processing  $A$ , that is  $A \not\sqsubseteq^{\text{avg}} B$ .

The other direction might also appear tricky: if  $B$  leaks no more than  $A$  in the average-case, it must also leak no more than  $B$  in the max-case. The quantification over all gain functions in the average-case is powerful enough to

“detect” differences in max-case leakage. The above result also means that  $\sqsubseteq^{\text{avg}}$  could be useful even if we are interested in the max-case, since it gives us  $\sqsubseteq^{\text{max}}$  for free.

## 4 Privacy-based refinement

So far we have compared systems based on their (average-case or max-case) leakage. In this section we turn our attention to the model of differential privacy, and discuss new ways of ordering mechanisms based on that model.

### 4.1 Differential privacy and d-privacy

Differential privacy relies on the observation that some pairs of secrets *need to be indistinguishable* from the point of view of the adversary in order to provide some meaningful notion of privacy; for instance, databases differing in a single individual should not be distinguishable, otherwise the privacy of that individual is violated. At the same, other pairs of secrets can be allowed to be distinguishable in order to provide some utility; for instance, distinguishing databases differing in many individuals allows us to answer a statistical query about those individuals.

This idea can be formalized by a *distinguishability metric*<sup>11</sup>  $\mathbf{d}$ . Intuitively,  $\mathbf{d}(x, x')$  models how *distinguishable* we allow these secrets to be. A value 0 means that we require  $x$  and  $x'$  to be completely indistinguishable to the adversary, while  $+\infty$  means that she can distinguish them completely.

In this context, a mechanism is simply a channel (the two terms will be used interchangeably), mapping secrets  $\mathcal{X}$  to some observations  $\mathcal{Y}$ . Denote by  $\mathbb{M}\mathcal{X}$  the set of all metrics on  $\mathcal{X}$ . Given  $\mathbf{d} \in \mathbb{M}\mathcal{X}$ , we define  $\mathbf{d}$ -privacy as follows.

**Definition 3.** A channel  $C$  satisfies  $\mathbf{d}$ -privacy iff

$$C_{x,y} \leq e^{\mathbf{d}(x,x')} C_{x',y} \quad \text{for all } x, x' \in \mathcal{X}, y \in \mathcal{Y}. \quad (16)$$

Intuitively, this definition requires that the closer  $x$  and  $x'$  are (as measured by  $\mathbf{d}$ ), the more similar (probabilistically) the output of the mechanism on these secrets should be.

**Remark 1.** Note that the definition of  $\mathbf{d}$ -privacy given in (1) is slightly different from the above one, because of the presence of  $\varepsilon$  in the exponent. Indeed, it is common to scale  $\mathbf{d}$  by a privacy parameter  $\varepsilon \geq 0$ , in which case  $\mathbf{d}$  can be thought of as the “kind” and  $\varepsilon$  as the “amount” of privacy. In other words, the structure determined by  $\mathbf{d}$  on the data specifies how we want to distinguish each pair of data, and  $\varepsilon$  specifies (uniformly) the degree of the distinction. Note that  $\varepsilon \cdot \mathbf{d}$  is itself a metric, so the two definitions are equivalent.

<sup>11</sup>To be precise, an *extended pseudo* metric, that is one in which distinct secrets can have distance 0, and distance  $+\infty$  is allowed.

Using a generic metric  $\mathbf{d}$  in this definition allows us to express different scenarios, depending on the domain  $\mathcal{X}$  on which the mechanism is applied and the choice of  $\mathbf{d}$ . For instance, in the standard model of differential privacy, the mechanism is applied to a database  $x$  (i.e.  $\mathcal{X}$  is the set of all databases), and produces some observation  $y$  (eg. a number). The *Hamming* metric  $\mathbf{d}_H$  – defined as the number of individuals in which  $x$  and  $x'$  differ – captures standard differential privacy.

## 4.2 Oblivious mechanisms

In the case of an *oblivious* mechanism, a query  $f: \mathcal{X} \rightarrow \mathcal{Y}$  is first applied to database  $x$ , and a noise mechanism  $H$  from  $\mathcal{Y}$  to  $\mathcal{Z}$  is applied to  $y = f(x)$ , producing an observation  $z$ . In this case, it is useful to study the privacy of  $H$  wrt some metric  $\mathbf{d}_Y$  on  $\mathcal{Y}$ . Then, to reason about the  $\mathbf{d}_X$ -privacy of the whole mechanism  $H \circ f$ , we can first compute the *sensitivity* of  $f$  wrt  $\mathbf{d}_X, \mathbf{d}_Y$ :

$$\Delta_{\mathbf{d}_X, \mathbf{d}_Y}^f = \max_{x, x'} \frac{\mathbf{d}_Y(f(x), f(x'))}{\mathbf{d}_X(x, x')} , \quad (17)$$

then apply the following property [6]:

$$\text{If } H \text{ satisfies } \mathbf{d}_Y\text{-privacy,} \quad (18)$$

$$\text{then } H \circ f \text{ satisfies } \Delta_{\mathbf{d}_X, \mathbf{d}_Y}^f \cdot \mathbf{d}_X\text{-privacy .} \quad (19)$$

For instance, the geometric mechanism  $G^\varepsilon$  satisfies  $\varepsilon \cdot \mathbf{d}_E$ -privacy (where  $\mathbf{d}_E$  denotes the Euclidean metric), hence it can be applied to any numeric query  $f$ : the resulting mechanism  $G^\varepsilon \circ f$  satisfies  $\Delta_{\mathbf{d}_H, \mathbf{d}_E}^f \cdot \varepsilon$ -differential privacy.<sup>12</sup>

## 4.3 Applying noise to the data of a single individual

There are also scenarios in which a mechanism  $C$  is applied directly to the data of a single individual (that is  $\mathcal{X}$  is the set of possible values). For instance, in the *local model* of differential privacy [9], the value of each individual is obfuscated before sending them to an untrusted curator. In this case,  $C$  should satisfy  $\mathbf{d}_D$ -privacy, where  $\mathbf{d}_D$  is the discrete metric, since *any change* in the individual's value should have negligible effects.

Moreover, in the context of *location-based services*, a user might want to obfuscate his location before sending it to the service provider. In this context, it is natural to require that locations that are geographically close are indistinguishable, while far away ones are allowed to be distinguished (in order to provide the service). In other words, we wish to provide  $\mathbf{d}_E$ -privacy, for the Euclidean metric on  $\mathbb{R}^2$ , called *geo-indistinguishability* in [3].

<sup>12</sup>The sensitivity wrt the Hamming and Euclidean metrics reduces to  $\Delta_{\mathbf{d}_H, \mathbf{d}_E}^f = \max_{x \sim x'} |f(x) - f(x')|$  where  $x \sim x'$  denotes  $\mathbf{d}_H(x, x') = 1$ .



#### 4.4 Comparing mechanisms by their “smallest $\varepsilon$ ” (for fixed $\mathbf{d}$ )

Scaling  $\mathbf{d}$  by a privacy parameter  $\varepsilon$  allows us to turn  $\mathbf{d}$ -privacy (for some fixed  $\mathbf{d}$ ) into a *quantitative “leakage” measure*, by associating each channel to the *smallest*  $\varepsilon$  by which we can scale  $\mathbf{d}$  without violating privacy.

**Definition 4.** *The privacy-based leakage (wrt  $\mathbf{d}$ ) of a channel  $C$  is defined as*

$$\text{Priv}_{\mathbf{d}}(C) := \inf\{\varepsilon \geq 0 \mid C \text{ satisfies } \varepsilon \cdot \mathbf{d}\text{-privacy}\} . \quad (20)$$

Note that  $\text{Priv}_{\mathbf{d}}(C) = +\infty$  iff there is no such  $\varepsilon$ ; also  $\text{Priv}_{\mathbf{d}}(C) \leq 1$  iff  $C$  satisfies  $\mathbf{d}$ -privacy.

It is then natural to compare two mechanisms  $A$  and  $B$  based on their “smallest  $\varepsilon$ ”.

**Definition 5.** *Define  $A \sqsubseteq_{\mathbf{d}}^{\text{prv}} B$  iff  $\text{Priv}_{\mathbf{d}}(A) \geq \text{Priv}_{\mathbf{d}}(B)$ .*

For instance,  $A \sqsubseteq_{\mathbf{d}_{\text{H}}}^{\text{prv}} B$  means that  $B$  satisfies standard differential privacy for  $\varepsilon$  at least as small as the one of  $A$ .

#### 4.5 Privacy-based leakage and refinement orders

When discussing the average- and max-case leakage orders  $\sqsubseteq_{\mathbb{G}}^{\text{avg}}, \sqsubseteq_{\mathbb{Q}}^{\text{max}}$ , we obtained strong leakage guarantees by quantifying over *all vulnerability functions*. It is thus natural to investigate a similar quantification in the context of  $\mathbf{d}$ -privacy. Namely, we define a stronger privacy-based “leakage” order, by comparing mechanisms not on a single metric  $\mathbf{d}$ , but on *all metrics* simultaneously.

**Definition 6.** *The privacy-based leakage order is defined as  $A \sqsubseteq_{\mathbb{M}}^{\text{prv}} B$  iff  $A \sqsubseteq_{\mathbf{d}}^{\text{prv}} B$  for all  $\mathbf{d} \in \mathbb{M}\mathcal{X}$ .*

Similarly to the other leakage orders, the drawback of  $\sqsubseteq_{\mathbb{M}}^{\text{prv}}$  is that it quantifies over an uncountable family of metrics. As a consequence, our first goal is to characterize it as a property of the channel matrix alone, which would make it much easier to reason about or verify.

To do so, we start by recalling an alternative way of thinking about  $\mathbf{d}$ -privacy. Consider the *multiplicative total variation* distance between probability distributions  $\mu, \mu' \in \mathbb{D}\mathcal{Y}$ , defined as:

$$\text{tv}_{\otimes}(\mu, \mu') := \max_{y: \mathcal{Y}} \left| \ln \frac{\mu_y}{\mu'_y} \right| . \quad (21)$$

If we think of  $C$  as a function  $\mathcal{X} \rightarrow \mathbb{D}\mathcal{Y}$  (mapping every  $x$  to the distribution  $C_{x,-}$ ),  $C$  satisfies  $\mathbf{d}$ -privacy iff  $\text{tv}_{\otimes}(C_{x,-}, C_{x',-}) \leq \mathbf{d}(x, x')$ , in other words iff  $C$  is non-expansive (1-Lipschitz) wrt  $\text{tv}_{\otimes}, \mathbf{d}$ .

Then, we introduce the concept of the *distinguishability metric*  $\mathbf{d}_C \in \mathbb{M}\mathcal{X}$  induced by the channel  $C$ , defined as

$$\mathbf{d}_C(x, x') := \text{tv}_{\otimes}(C_{x,-}, C_{x',-}) . \quad (22)$$

Intuitively,  $\mathbf{d}_C(x, x')$  expresses exactly how much the channel distinguishes (wrt  $\text{tv}_\otimes$ ) the secrets  $x, x'$ . It is easy to see that  $\mathbf{d}_C$  is the *smallest metric* for which  $C$  is private; in other words, for any  $\mathbf{d}$ :

$$C \text{ satisfies } \mathbf{d}\text{-privacy} \quad \text{iff} \quad \mathbf{d} \geq \mathbf{d}_C . \quad (23)$$

We can now give a refinement order on mechanisms, by comparing their corresponding induced metrics.

**Definition 7.** *The privacy-based refinement order is defined as  $A \sqsubseteq^{\text{prv}} B$  iff  $\mathbf{d}_A \geq \mathbf{d}_B$ , or equivalently iff  $B$  satisfies  $\mathbf{d}_A$ -privacy.*

This achieves our goal of goal of characterizing  $\sqsubseteq_{\mathbb{M}}^{\text{prv}}$ .

**Proposition 4.** *The orders  $\sqsubseteq_{\mathbb{M}}^{\text{prv}}$  and  $\sqsubseteq^{\text{prv}}$  coincide.*

We now turn our attention to the question of how these orders relate to each other.

**Theorem 5.**  *$\sqsubseteq^{\text{max}}$  is strictly stronger than  $\sqsubseteq^{\text{prv}}$ , which is strictly stronger than  $\sqsubseteq_{\mathbf{d}}^{\text{prv}}$ .*

The fact that  $\sqsubseteq^{\text{max}}$  is stronger than  $\sqsubseteq^{\text{prv}}$  is due to the fact that  $\text{Priv}_{\mathbf{d}}$  can be seen as a max-case information leakage, for a properly constructed vulnerability function  $V_{\mathbf{d}}$ . This is discussed in detail in Section 4.7. This implication means that  $\sqsubseteq^{\text{avg}}, \sqsubseteq^{\text{max}}$  can be useful even if we “only” care about  $\mathbf{d}$ -privacy.

## 4.6 Application to oblivious mechanisms

We conclude the discussion on privacy-based refinement by showing the usefulness of our strong  $\sqsubseteq^{\text{prv}}$  order in the case of oblivious mechanisms.

**Theorem 6.** *Let  $f : \mathcal{X} \rightarrow \mathcal{Y}$  be any query and  $A, B$  be two mechanisms on  $\mathcal{Y}$ . If  $A \sqsubseteq^{\text{prv}} B$  then  $A \circ f \sqsubseteq^{\text{prv}} B \circ f$ .*

This means that, replacing  $A$  by  $B$  in the context of an oblivious mechanism is always safe, regardless of the query (and its sensitivity) and regardless of the metric by which the privacy of the composed mechanism is evaluated.

Assume, for instance, that we care about standard differential privacy, and we have properly constructed  $A$  such that  $A \circ f$  satisfies  $\varepsilon$ -differential privacy for some  $\varepsilon$ . If we know that  $A \sqsubseteq^{\text{prv}} B$  (several such cases are discussed in Section 6) we can replace  $A$  by  $B$  without even knowing what  $f$  does. The mechanism  $B \circ f$  is guaranteed to also satisfy  $\varepsilon$ -differential privacy.

Note also that the above theorem *fails for the weaker order  $\sqsubseteq_{\mathbf{d}}^{\text{prv}}$* . Establishing  $A \sqsubseteq_{\mathbf{d}_{\mathcal{Y}}}^{\text{prv}} B$  for some metric  $\mathbf{d}_{\mathcal{Y}} : \mathbb{M}\mathcal{Y}$  gives no guarantees that  $A \circ f \sqsubseteq_{\mathbf{d}_{\mathcal{X}}}^{\text{prv}} B \circ f$  for some other metric of interest  $\mathbf{d}_{\mathcal{X}} : \mathbb{M}\mathcal{X}$ . It is possible that replacing  $A$  by  $B$  in that case is not safe (one would need to re-examine the behavior of  $B$ , and possibly reconfigure it to the sensitivity of  $f$ ).

Leakage orders		Refinement orders
$\sqsubseteq_{\mathbf{G}}^{\text{avg}}$	$\Leftrightarrow$	$\sqsubseteq^{\text{avg}}$
$\Downarrow$		$\Downarrow$
$\sqsubseteq_{\mathbf{Q}}^{\text{max}}$	$\Leftrightarrow$	$\sqsubseteq^{\text{max}}$
$\Downarrow$		$\Downarrow$
$\sqsubseteq_{\mathbf{M}}^{\text{prv}}$	$\Leftrightarrow$	$\sqsubseteq^{\text{prv}}$
$\Rightarrow$		$\Leftarrow$
		$\sqsubseteq_{\mathbf{d}}^{\text{prv}}$

Table 3: Comparison of leakage and refinement orders. All implications are strict.

#### 4.7 Privacy as max-case capacity

One way to check whether  $\sqsubseteq^{\text{max}}$  is stronger than  $\sqsubseteq^{\text{prv}}$  is to examine whether **d**-privacy can be expressed as a (max-case) information leakage. We start this by defining a suitable vulnerability function:

**Definition 8.** *The **d**-vulnerability function  $V_{\mathbf{d}}$  is defined as*

$$V_{\mathbf{d}}(\pi) := \inf\{\varepsilon \geq 0 \mid \forall x, x' \in \mathcal{X}, \pi_x \leq e^{\varepsilon \cdot \mathbf{d}(x, x')} \pi_{x'}\} . \quad (24)$$

Note the difference between  $V_{\mathbf{d}}(\pi)$  (a vulnerability function on *distributions*) and  $\text{Priv}_{\mathbf{d}}(C)$  (a “leakage” measure on *channels*).

A fundamental notion in QIF is that of *capacity* : the maximization of leakage over all priors. It turns out that for  $V_{\mathbf{d}}$ , the capacity-realizing prior is the uniform one. In the following,  $\mathcal{L}_{\mathbf{d}}^{+, \text{max}}$  denotes the additive max-case **d** leakage, namely:

$$\mathcal{L}_{\mathbf{d}}^{+, \text{max}}(\pi, C) = V_{\mathbf{d}}^{\text{max}}[\pi, C] - V_{\mathbf{d}}(\pi). \quad (25)$$

and  $\mathcal{ML}_{\mathbf{d}}^{+, \text{max}}$  denotes the additive max-case **d**-capacity, namely:

$$\mathcal{ML}_{\mathbf{d}}^{+, \text{max}}(C) = \max_{\pi} \mathcal{L}_{\mathbf{d}}^{+, \text{max}}(\pi, C). \quad (26)$$

**Theorem 7.**  $\mathcal{ML}_{\mathbf{d}}^{+, \text{max}}$  is always achieved on a uniform prior  $\pi^u$ . Namely

$$\max_{\pi} \mathcal{L}_{\mathbf{d}}^{+, \text{max}}(\pi, C) = \mathcal{L}_{\mathbf{d}}^{+, \text{max}}(\pi^u, C) = V_{\mathbf{d}}^{\text{max}}[\pi^u, C] . \quad (27)$$

This finally brings us to our goal of expressing  $\text{Priv}_{\mathbf{d}}$  in terms of information leakage (for a proper vulnerability function).

**Theorem 8.** *[DP as max-case capacity]  $C$  satisfies  $\varepsilon \cdot \mathbf{d}$ -privacy iff  $\mathcal{ML}_{\mathbf{d}}^{+, \text{max}}(C) \leq \varepsilon$ . In other words:  $\mathcal{ML}_{\mathbf{d}}^{+, \text{max}}(C) = \text{Priv}_{\mathbf{d}}(C)$ .*

## 5 Verifying the refinement orders

We now turn our attention to the problem of checking whether the various orders hold, given two *explicit representations* of channels  $A$  and  $B$  (in terms of their matrices). We show that, for all orders, this question can be answered in time polynomial in the size of the matrices. Moreover, when one of the order fails, we discuss how to obtain a *counter-example* (eg. a gain function  $g$  or a vulnerability function  $V$ ), demonstrating this fact. All the methods discussed in the section have been implemented in a publicly available library, and have been used in the experimental results of Section 6.

### 5.1 Average-case refinement

Verifying  $A \sqsubseteq^{\text{avg}} B$  can be done in polynomial time (in the size of  $A, B$ ) by solving the system of equations  $AR = B$ , with variables  $R$ , under the linear constraints that  $R$  is a channel matrix (non-negative and rows sum up to 1). However, if the system has no solution (i.e.  $A \not\sqsubseteq^{\text{avg}} B$ ), this method does not provide us with a counter-example gain function  $g$ .

We now show that there is an alternative efficient method: define  $C^\uparrow = \{CR \mid R \text{ is a channel}\}$ , the set of all channels obtainable by post-processing  $C$ . The idea is to compute the projection of  $B$  on  $A^\uparrow$ . Clearly, the projection is  $B$  itself iff  $A \sqsubseteq^{\text{avg}} B$ ; otherwise, the projection can be directly used to construct a counter-example  $g$ .

**Theorem 9.** *Let  $B^*$  be the projection of  $B$  on  $A^\uparrow$ .*

1. *If  $B = B^*$  then  $A \sqsubseteq^{\text{avg}} B$ .*
2. *Otherwise, let  $G = B - B^*$ . The gain function  $g(w, x) = G_{x,w}$  provides a counter-example to  $A \sqsubseteq^{\text{avg}} B$ , that is  $V_g(\pi^u, A) < V_g(\pi^u, B)$ , for uniform  $\pi^u$ .*

Since  $\|x - y\|_2^2 = x^T x - 2x^T y + y^T y$ , the projection of  $y$  to a convex set can be written as  $\min_x x^T x - 2x^T y$  for  $Ax \leq b$ . This is a quadratic program with  $Q$  being the identity matrix, which is positive definite, hence it can be solved in polynomial time.

Note that the proof that  $\sqsubseteq_{\mathbb{G}}^{\text{avg}}$  is stronger than  $\sqsubseteq^{\text{avg}}$  (the “coriaceous” theorem of [19]) uses the hyperplane-separation theorem to show the existence of a counter example  $g$  in case  $A \not\sqsubseteq^{\text{avg}} B$ . The above technique essentially computes such a separating hyperplane.

### 5.2 Max-case refinement

Similarly to  $\sqsubseteq^{\text{avg}}$ , we can verify  $A \sqsubseteq^{\text{max}} B$  directly using its definition, by solving the system  $R\tilde{A} = \tilde{B}$  under the constraint that  $R$  is a channel.

In contrast to  $\sqsubseteq^{\text{avg}}$ , when  $A \not\sqsubseteq^{\text{max}} B$ , the proof of Theorem. 2 directly gives us a counter-example:

$$V(\sigma) := \min_{\sigma' \in S} \|\sigma - \sigma'\|_2. \quad (28)$$

where  $S = \mathbf{ch\,supp}[\pi, A]$  and  $\pi$  is any full-support prior. For this vulnerability function it holds that  $V^{\max}(\pi, A) < V^{\max}(\pi, B)$ .

### 5.3 Privacy-based refinement

The  $\sqsubseteq^{\text{prv}}$  order can be verified directly from its definition, by checking that  $\mathbf{d}_A \geq \mathbf{d}_B$ . This can be done in time  $O(|\mathcal{X}|^2|\mathcal{Y}|)$ , by computing  $\text{tv}_{\otimes}(C_{x,-}, C_{x',-})$  for each pair of secrets. If  $A \not\sqsubseteq^{\text{prv}} B$ , then  $\mathbf{d} = \mathbf{d}_B$  provides an immediate counter-example metric, since  $B$  satisfies  $\mathbf{d}_B$ -privacy, but  $A$  does not.

## 6 Application: comparing DP mechanisms

In differential privacy it is common to compare the privacy guarantees provided by different mechanisms by ‘comparing the epsilons’. But it is interesting to ask to what extent  $\varepsilon$ -equivalent mechanisms are comparable wrt the other leakage measures defined here. Or we might want to know whether reducing  $\varepsilon$  in a mechanism also corresponds to a *refinement* of it. This could be useful if, for example, it is important to understand the privacy properties of a mechanism with respect to *any* max-case leakage measure, and not just the DP measure given by  $\varepsilon$ .

Since the  $\varepsilon$ -based order given by  $\sqsubseteq_{\mathbf{d}}^{\text{prv}}$  is (strictly) the weakest of the orders considered here, it cannot be the case that we *always* get a refinement (wrt other orders). But it may be true that for particular *families* of mechanisms some (or all) of the refinement orders hold.

We investigate 3 families of mechanisms commonly used in DP or LDP: *geometric*, *exponential* and *randomized response* mechanisms.

### 6.1 Preliminaries

We define each family of mechanisms in terms of their channel construction. We assume that mechanisms operate on a set of inputs (denoted by  $\mathcal{X}$ ) and produce a set of outputs (denoted by  $\mathcal{Y}$ ). In this sense our mechanisms can be seen as oblivious (as in standard DP) or as LDP mechanisms. (We use the term ‘mechanism’ in either sense). We denote by  $M^\varepsilon$  a mechanism parametrized by  $\varepsilon$ , where  $\varepsilon$  is defined to be the same as  $\text{Priv}_{\mathbf{d}}(M)$ .<sup>13</sup> In order to compare mechanisms, we restrict our input and output domains of interest to (possibly infinite) sequences of non-negative integers.<sup>14</sup> (We assume  $\mathcal{X}, \mathcal{Y}$  are finite unless specified.) Also, as we are operating in the framework of  $\mathbf{d}$ -privacy, it is necessary to provide an appropriate metric defined over  $\mathcal{X}$ ; here it makes sense to use the Euclidean distance metric  $\mathbf{d}_E$ .

<sup>13</sup>We note that the exponential mechanism *under-reports* its  $\varepsilon$ , thus for the purposes of comparison we make sure that we use the best possible  $\varepsilon$  for each mechanism.

<sup>14</sup>Our results hold for sequences of quantized integers  $q[0..]$  but we use integer sequences to simplify presentation.

**Definition 9.** A geometric mechanism is a channel  $(\mathcal{X}, \mathbb{Z}, G^\varepsilon)$ , parametrized by  $\varepsilon \geq 0$  constructed as follows:

$$G_{x,y}^\varepsilon = \frac{(1 - \alpha) \cdot \alpha^{\mathbf{d}_E(x,y)}}{1 + \alpha} \quad \text{for all } x \in \mathcal{X}, y \in \mathbb{Z} \quad (29)$$

$$(30)$$

where  $\alpha = e^{-\varepsilon}$  and  $\mathbf{d}_E(x, y) = \|x - y\|$ . Such a mechanism satisfies  $\varepsilon \cdot \mathbf{d}_E$ -privacy.

In practice, the truncated geometric mechanism is preferred to the infinite geometric. We define the truncated geometric mechanism as follows.

**Definition 10.** A truncated geometric mechanism is a channel  $(\mathcal{X}, \mathcal{Y}, TG^\varepsilon)$ , parametrized by  $\varepsilon \geq 0$  with  $\mathcal{X} \subseteq \mathcal{Y}$  constructed as follows:

$$TG_{x,y}^\varepsilon = \frac{(1 - \alpha) \cdot \alpha^{\mathbf{d}_E(x,y)}}{1 + \alpha} \quad \text{for all } y \neq \min \mathcal{Y}, \max \mathcal{Y} \quad (31)$$

$$TG_{x,y}^\varepsilon = \frac{\alpha^{\mathbf{d}_E(x,y)}}{1 + \alpha} \quad \text{for } y = \min \mathcal{Y}, \max \mathcal{Y} \quad (32)$$

where  $\alpha = e^{-\varepsilon}$  and  $\mathbf{d}_E(x, y) = \|x - y\|$ . Such a mechanism satisfies  $\varepsilon \cdot \mathbf{d}_E$ -privacy.

It is also possible to define the ‘over-truncated’ geometric mechanism whose input space is not entirely included in the output space.

**Definition 11.** An over-truncated geometric mechanism is a channel  $(\mathcal{X}, \mathcal{Y}, OTG^\varepsilon)$ , parametrized by  $\varepsilon \geq 0$  with  $\mathcal{X} \not\subseteq \mathcal{Y}$  constructed as follows:

1. Start with the truncated geometric mechanism  $(\mathcal{X}, \mathcal{X} \cup \mathcal{Y}, TG^\varepsilon)$ .
2. Sum up the columns at each end until the output domain is reached.

Such a mechanism satisfies  $\varepsilon \cdot \mathbf{d}_E$ -privacy.

For example, the set of inputs to an over-truncated geometric mechanism could be integers in the range  $[0 \dots 100]$  but the output space may have a range of  $[0 \dots 50]$  or perhaps  $[-50 \dots 50]$ . In either of these cases, the mechanism has to ‘over-truncate’ the inputs to accommodate the output space.

We remark that we do not consider the over-truncated mechanism a particularly useful mechanism in practice. However, we provide results on this mechanism for completeness since its construction is possible, if unusual.

**Definition 12.** An exponential mechanism is a channel  $(\mathcal{X}, \mathcal{Y}, E^\alpha)$ , parametrized by  $\varepsilon \geq 0$  constructed as follows:

$$E_{x,y}^\alpha = \lambda_x \cdot e^{-\frac{\varepsilon}{2} \mathbf{d}_E(x,y)} \quad \text{for all } x \in \mathcal{X}, y \in \mathcal{Y} \quad (33)$$

where  $\lambda_x$  are normalizing constants ensuring  $\sum_y E_{x,y}^\alpha = 1$ . Such a mechanism satisfies  $\alpha \cdot \mathbf{d}_E$ -privacy where  $\alpha \geq \frac{\varepsilon}{2}$  (which can be calculated exactly from the channel construction).<sup>15</sup>

The exponential mechanism was designed for arbitrary domains, thus its parameter  $\varepsilon$  does not correspond to the true (best-case)  $\varepsilon$ -DP guarantee that it provides. We will denote by  $E^\varepsilon$  the exponential mechanism with ‘true’ privacy parameter  $\varepsilon$  rather than the reported one, as our intention is to capture the privacy guarantee provided by the channel in order to make reasonable comparisons.

**Definition 13.** A randomized response mechanism is a channel  $(\mathcal{X}, \mathcal{Y}, R^\varepsilon)$ , parametrized by  $\varepsilon \geq 0$  constructed as follows:

$$R_{x,y}^\varepsilon = \frac{e^{\varepsilon(1-\mathbf{d}_D(x,y))}}{e^\varepsilon + n} \quad \text{for all } x, y \in \mathcal{Y} \quad (34)$$

$$R_{x,y}^\varepsilon = \frac{1}{n+1} \quad \text{for all } x \notin \mathcal{Y} \quad (35)$$

where  $n = \|\mathcal{Y}\| - 1$  and  $\mathbf{d}_D$  is the discrete metric (that is,  $\mathbf{d}_D(x, x) = 0$  and  $\mathbf{d}_D(x, y) = 1$  for  $x \in \mathcal{Y}, x \neq y$ ). Such a mechanism satisfies  $\varepsilon \cdot \mathbf{d}_D$ -privacy.

We note that the randomized response mechanism also satisfies  $\varepsilon \cdot \mathbf{d}_E$ -privacy.

Intuitively, the randomized response mechanism returns the true answer with high probability and all other responses with equal probability. In the case where the input  $x$  lies outside  $\mathcal{Y}$  (that is, in ‘over-truncated’ mechanisms), all of the outputs (corresponding to the outlying inputs) have equal probability.

**Example 7.** The following are examples of each of the mechanisms described above, represented as channel matrices. For this example, we set  $\varepsilon = \log(2)$  for the geometric and randomized response mechanisms, while for the exponential mechanism we use  $\varepsilon = \log(4)$ .

<i>TG</i>	$x_1$	$x_2$	$x_3$
$x_1$	$2/3$	$1/6$	$1/6$
$x_2$	$1/3$	$1/3$	$1/3$
$x_3$	$1/6$	$1/6$	$2/3$

<i>OTG</i>	$x_1$	$x_2$
$x_1$	$2/3$	$1/3$
$x_2$	$1/3$	$2/3$
$x_3$	$1/6$	$5/6$

(36)

<i>E</i>	$x_1$	$x_2$	$x_3$
$x_1$	$4/7$	$2/7$	$1/7$
$x_2$	$1/4$	$1/2$	$1/4$
$x_3$	$1/7$	$2/7$	$4/7$

<i>R</i>	$x_1$	$x_2$	$x_3$
$x_1$	$1/2$	$1/4$	$1/4$
$x_2$	$1/4$	$1/2$	$1/4$
$x_3$	$1/4$	$1/4$	$1/2$

Note that the exponential mechanism here actually satisfies  $\log(\frac{16}{7}) \cdot \mathbf{d}_E$ -privacy even though it is specified by  $\varepsilon = \log(4)$ .

<sup>15</sup>Note that the construction presented here uses the Euclidean distance metric since we only consider integer domains. The general construction of the exponential mechanism uses an arbitrary metric.

We now have 3 families of mechanisms which we can characterize by channels, and which satisfy  $\varepsilon \cdot \mathbf{d}_E$ -privacy. For the remainder of this section we will refer only to the  $\varepsilon$  parameter and take  $\mathbf{d}_E$  as given, as we wish to understand the effect of changing  $\varepsilon$  (for a fixed metric) on the various leakage measures.

## 6.2 Refinement order within families of mechanisms

We first ask which refinement orders hold within a family of mechanisms. That is, when does reducing  $\varepsilon$  for a particular mechanism produce a refinement? Since we have the convenient order  $\sqsubseteq^{\text{avg}} \subset \sqsubseteq^{\text{max}} \subset \sqsubseteq^{\text{prv}}$  it is useful to first check if  $\sqsubseteq^{\text{avg}}$  holds as we get the other refinements ‘for free’.

For the (infinite) geometric mechanism we have the following result.

**Theorem 10.** *Let  $G^\varepsilon, G^{\varepsilon'}$  be geometric mechanisms. Then  $G^\varepsilon \sqsubseteq^{\text{avg}} G^{\varepsilon'}$  iff  $\varepsilon \geq \varepsilon'$ . That is decreasing  $\varepsilon$  produces a refinement of the mechanism.*

This means that reducing  $\varepsilon$  in an infinite geometric mechanism is safe against *any* adversary that can be modelled using, for example, max-case or average-case vulnerabilities.

For the truncated geometric mechanism we get the same result.

**Theorem 11.** *Let  $TG^\varepsilon, TG^{\varepsilon'}$  be truncated geometric mechanisms. Then  $TG^\varepsilon \sqsubseteq^{\text{avg}} TG^{\varepsilon'}$  iff  $\varepsilon \geq \varepsilon'$ . That is, decreasing  $\varepsilon$  produces a refinement of the mechanism.*

However, the over-truncated geometric mechanism does not behave so well.

**Theorem 12.** *Let  $OTG^\varepsilon, OTG^{\varepsilon'}$  be over-truncated geometric mechanisms. Then  $OTG^\varepsilon \not\sqsubseteq^{\text{avg}} OTG^{\varepsilon'}$  for any  $\varepsilon \neq \varepsilon'$ . That is, decreasing  $\varepsilon$  does **not** produce a refinement.*

We can think of this last class of geometrics as ‘skinny’ mechanisms, that is, corresponding to a channel with a smaller output space than input space.

Intuitively, this theorem means that we can *always* find some (average-case) adversary who prefers the over-truncated geometric mechanism with the smaller  $\varepsilon$ .

We remark that the gain function we found can be easily calculated by treating the columns of channel  $A$  as vectors, and finding a vector orthogonal to both of these. This follows from the results in Section 5.1. Since the columns of  $A$  cannot span the space  $\mathbb{R}^3$  it is always possible to find such a vector, and when this vector is not orthogonal to the ‘column space’ of  $B$  it can be used to construct a gain function preferring  $B$  to  $A$ .

Even though the  $\sqsubseteq^{\text{avg}}$  refinement does not hold, we can check whether the other refinements are satisfied.

**Theorem 13.** *Let  $OTG^\varepsilon$  be an over-truncated geometric mechanism. Then reducing  $\varepsilon$  does **not** produce a  $\sqsubseteq^{\text{max}}$  refinement, however it **does** produce a  $\sqsubseteq^{\text{prv}}$  refinement.*



This means that although a smaller  $\varepsilon$  does not provide safety against all max-case adversaries, it *does* produce a safer mechanism wrt  $d$ -privacy for *any* choice of metric we like.

Intuitively, the  $\sqsubseteq^{\text{prv}}$  order relates mechanisms based on how they distinguish *inputs*. Specifically, if  $A \sqsubseteq^{\text{prv}} B$  then for any pair of inputs  $x, x'$ , the corresponding output distributions are ‘further apart’ in channel  $A$  than in channel  $B$ , and thus the inputs are more distinguishable using channel  $A$ . When  $\sqsubseteq^{\text{prv}}$  fails to hold, it means that there are some inputs in  $A$  which are more distinguishable than in  $B$ , and vice versa. This means an adversary who is interested in distinguishing some particular pair of inputs would prefer one mechanism to the other.

We now consider the exponential mechanism. In this case we do not have a theoretical result, but experimentally it appears that the exponential mechanism respects refinement, so we present the following conjecture.

**Conjecture 1.** *Let  $E^\varepsilon$  be an exponential mechanism. Then decreasing  $\varepsilon$  in  $E$  produces a refinement. That is,  $E^\varepsilon \sqsubseteq^{\text{avg}} E^{\varepsilon'}$  iff  $\varepsilon \geq \varepsilon'$ .*

Finally we consider the randomized response mechanism.

**Theorem 14.** *Let  $R^\varepsilon$  be a randomized response mechanism. Then decreasing  $\varepsilon$  in  $R$  produces a refinement. That is,  $R^\varepsilon \sqsubseteq^{\text{avg}} R^{\varepsilon'}$  iff  $\varepsilon \geq \varepsilon'$ .*

In conclusion we can say that, in general, the usual DP families of mechanisms are ‘well-behaved’ wrt all of the refinement orders. This means that it is safe (wrt *any* adversary we model here) to replace a mechanism from a particular family with another mechanism from the *same* family with a lower  $\varepsilon$ .

### 6.3 Refinement order between families of mechanisms

Now we explore whether it is possible to compare mechanisms from different families. We first ask, can we compare mechanisms which have the same  $\varepsilon$ ? We assume that the input and output domains are the same, and the intention is to decide whether to replace one mechanism with another.

**Theorem 15.** *Let  $R$  be a randomized response mechanism,  $E$  an exponential mechanism and  $TG$  a truncated geometric mechanism. Then  $TG^\varepsilon \sqsubseteq^{\text{prv}} R^\varepsilon$  and  $TG^\varepsilon \sqsubseteq^{\text{prv}} E^\varepsilon$ . However  $\sqsubseteq^{\text{prv}}$  does not hold between  $E^\varepsilon$  and  $R^\varepsilon$ .*

*Proof.* We present a counter-example to show  $E^\varepsilon \not\sqsubseteq^{\text{prv}} R^\varepsilon$  and  $R^\varepsilon \not\sqsubseteq^{\text{prv}} E^\varepsilon$ . The remainder of the proof is in the appendix.

Consider the following channels:

$A$	$x_1$	$x_2$	$x_3$	$x_4$	$B$	$x_1$	$x_2$	$x_3$	$x_4$
$x_1$	$8/15$	$4/15$	$2/15$	$1/15$	$x_1$	$4/9$	$5/27$	$5/27$	$5/27$
$x_2$	$2/9$	$4/9$	$2/9$	$1/9$	$x_2$	$5/27$	$4/9$	$5/27$	$5/27$
$x_3$	$1/9$	$2/9$	$4/9$	$2/9$	$x_3$	$5/27$	$5/27$	$4/9$	$5/27$
$x_4$	$1/15$	$2/15$	$4/15$	$8/15$	$x_4$	$5/27$	$5/27$	$5/27$	$4/9$

(37)

Channel  $A$  represents an exponential mechanism and channel  $B$  a randomized response mechanism. Both have (true)  $\varepsilon$  of  $\log(12/5)$ .<sup>16</sup> However  $\mathbf{d}_A(x_1, x_3) > \mathbf{d}_B(x_1, x_3)$  and  $\mathbf{d}_A(x_2, x_3) < \mathbf{d}_B(x_2, x_3)$ . Thus  $A$  does not satisfy  $\mathbf{d}_B$ -privacy, nor does  $B$  satisfy  $\mathbf{d}_A$ -privacy.  $\square$

Intuitively, the randomized response mechanism maintains the same ( $\varepsilon$ ) distinguishability level between inputs, whereas the exponential mechanism causes some inputs to be *less* distinguishable than others. This means that, for the same (true)  $\varepsilon$ , an adversary who is interested in certain inputs could learn more from the randomized response than the exponential. In the above counter-example, points  $x_2, x_3$  in the exponential mechanism of channel  $A$  are *less* distinguishable than the corresponding points in the randomized response mechanism  $B$ .

As an example, let's say the mechanisms are to be used in geo-location privacy and the inputs represent adjacent locations (such as addresses along a street). Then an adversary (your boss) may be interested in how far you are from work, and therefore wants to be able to distinguish between points distant from  $x_1$  (your office) and points within the vicinity of your office, without requiring your precise location. Your boss chooses channel  $A$  as the most informative. However, another adversary (your suspicious partner) is more concerned about where exactly you are, and is particularly interested in distinguishing between your expected position ( $x_2$ , the boulangerie) versus your suspected position ( $x_3$ , the brothel). Your partner chooses channel  $B$  as the most informative.

Regarding the other refinements, we find (experimentally) that *none* of them hold (in general) between families of mechanisms.<sup>17</sup> We present these results in Table 5 (Section ??).

We next check what happens when we compare mechanisms with *different* epsilons. We note the following.

**Theorem 16.** *For any (truncated geometric, randomized response, exponential) mechanisms  $M_1^{\varepsilon_1}, M_2^{\varepsilon_2}$ , if  $M_1^{\varepsilon_1} \sqsubseteq M_2^{\varepsilon_2}$  for any of our refinements ( $\sqsubseteq^{\text{avg}}, \sqsubseteq^{\text{max}}, \sqsubseteq^{\text{prv}}$ ) then  $M_1^{\varepsilon_1} \sqsubseteq M_2^{\varepsilon'_2}$  for  $\varepsilon'_2 < \varepsilon_2$ .*

*Proof.* This follows directly from transitivity of the refinement relations, and our results on refinement with families of mechanisms.<sup>18</sup>  $\square$

This tells us that once we have a refinement between mechanisms, it continues to hold for reduced  $\varepsilon$  in the refining mechanism.

**Corollary 17.** *Let  $G, TG, R, E$  be the geometric, truncated geometric, randomized response and exponential mechanisms respectively. Then for all  $\varepsilon' \leq \varepsilon$  we have that  $TG^\varepsilon \sqsubseteq^{\text{prv}} R^{\varepsilon'}$ ,  $TG^\varepsilon \sqsubseteq^{\text{prv}} E^{\varepsilon'}$ ,  $G^\varepsilon \sqsubseteq^{\text{prv}} R^{\varepsilon'}$  and  $G^\varepsilon \sqsubseteq^{\text{prv}} E^{\varepsilon'}$ .*

<sup>16</sup>Channel A was generated using  $\varepsilon = \log(4)$ . However, as noted earlier, this corresponds to a lower *true*  $\varepsilon$ .

<sup>17</sup>Recall that we only need to produce a single counter-example to show that a refinement doesn't hold, and this can be done using the methods presented in Section 5.

<sup>18</sup>We recall however that our result for the exponential mechanism is only a conjecture.

So it is safe to ‘compare epsilons’ wrt  $\sqsubseteq^{\text{prv}}$  if we want to replace a geometric mechanism with either a randomized response or exponential mechanism.<sup>19</sup> What this means is that if, for example, we have a geometric mechanism  $TG$  that operates on databases with distance measured using the Hamming metric  $\mathbf{d}_H$  and satisfying  $\varepsilon \cdot \mathbf{d}_H$ -privacy, then any randomized response mechanism  $R$  parametrized by  $\varepsilon' \leq \varepsilon$  will also satisfy  $\varepsilon \cdot \mathbf{d}_H$ -privacy. Moreover, if we decide we’d rather use the Manhattan metric  $\mathbf{d}_M$  to measure distance between the databases, then we only need to check that  $TG$  also satisfies  $\varepsilon \cdot \mathbf{d}_M$ -privacy, as this implies that  $R$  will too.

The following tables 4, 5, and 6 summarize the refinement relations with respect to the various families of mechanisms.

Mechanism	Are these valid for decreasing $\varepsilon$ ?		
	$\sqsubseteq^{\text{avg}}$	$\sqsubseteq^{\text{max}}$	$\sqsubseteq^{\text{prv}}$
Geometric	Y	Y	Y
Truncated Geometric	Y	Y	Y
Over-Truncated Geometric	N	N	Y
Exponential	Y	Y	Y
Randomized Response	Y	Y	Y

Table 4: The refinements respected by families of mechanisms for decreasing  $\varepsilon$ .

## 6.4 Asymptotic behavior

We now consider the behavior of the relations when  $\varepsilon$  approximates 0, which represents the absence of leakage. We start with the following result:

**Theorem 18.** *Every (truncated geometric, randomized response, exponential) mechanism is ‘the safest possible mechanism’ when parametrized by  $\varepsilon = 0$ . That is  $L^\varepsilon \sqsubseteq^{\text{avg}} M^0$  for all mechanisms  $L, M$  (possibly from different families) and  $\varepsilon > 0$ .*

While this result may be unsurprising, it means that we know that refinement must *eventually* occur when we reduce  $\varepsilon$ . It is interesting then to ask just *when* this refinement occurs. We examine this question experimentally by considering different mechanisms and investigating for which values of  $\varepsilon$  average-case refinement holds. For simplicity of presentation, we show results for  $5 \times 5$  matrices, noting that we observed similar results for experiments across different matrix dimensions.<sup>20</sup> The results are plotted in Figure 3.

The plots show the relationship between  $\varepsilon_1$  (x-axis) and  $\varepsilon_2$  (y-axis) where  $\varepsilon_1$  parametrizes the mechanism being refined and  $\varepsilon_2$  parametrizes the refining

<sup>19</sup>As with the previous theorem, note that the results for the exponential mechanism are stated as conjecture only, and this conjecture is assumed in the statement of this corollary.

<sup>20</sup>By similar results, we mean wrt the coarse-grained comparison of plots that we do here.

Refinements across families with same $\varepsilon$		
$TG \not\sqsubseteq^{\text{avg}} R$	$TG \not\sqsubseteq^{\text{max}} R$	$TG \sqsubseteq^{\text{prv}} R$
$R \not\sqsubseteq^{\text{avg}} TG$	$R \not\sqsubseteq^{\text{max}} TG$	$R \not\sqsubseteq^{\text{prv}} TG$
$TG \not\sqsubseteq^{\text{avg}} E$	$TG \not\sqsubseteq^{\text{max}} E$	$TG \sqsubseteq^{\text{prv}} E$
$E \not\sqsubseteq^{\text{avg}} TG$	$E \not\sqsubseteq^{\text{max}} TG$	$E \not\sqsubseteq^{\text{prv}} TG$
$G \not\sqsubseteq^{\text{avg}} R$	$G \not\sqsubseteq^{\text{max}} R$	$G \sqsubseteq^{\text{prv}} R$
$R \not\sqsubseteq^{\text{avg}} G$	$R \not\sqsubseteq^{\text{max}} G$	$R \not\sqsubseteq^{\text{prv}} G$
$G \not\sqsubseteq^{\text{avg}} E$	$G \not\sqsubseteq^{\text{max}} E$	$G \sqsubseteq^{\text{prv}} E$
$E \not\sqsubseteq^{\text{avg}} G$	$E \not\sqsubseteq^{\text{max}} G$	$E \not\sqsubseteq^{\text{prv}} G$
$R \not\sqsubseteq^{\text{avg}} E$	$R \not\sqsubseteq^{\text{max}} E$	$R \not\sqsubseteq^{\text{prv}} E$
$E \not\sqsubseteq^{\text{avg}} R$	$E \not\sqsubseteq^{\text{max}} R$	$E \not\sqsubseteq^{\text{prv}} R$

Table 5: Comparing different families of mechanisms with respect to the different refinements under the same  $\varepsilon$ .

Comparison of refinements with $\varepsilon_1 > \varepsilon_2$ .		
$TG^{\varepsilon_1} \not\sqsubseteq^{\text{avg}} R^{\varepsilon_2}$	$TG^{\varepsilon_1} \not\sqsubseteq^{\text{max}} R^{\varepsilon_2}$	$TG^{\varepsilon_1} \sqsubseteq^{\text{prv}} R^{\varepsilon_2}$
$R^{\varepsilon_1} \not\sqsubseteq^{\text{avg}} TG^{\varepsilon_2}$	$R^{\varepsilon_1} \not\sqsubseteq^{\text{max}} TG^{\varepsilon_2}$	$R^{\varepsilon_1} \not\sqsubseteq^{\text{prv}} TG^{\varepsilon_2}$
$TG^{\varepsilon_1} \not\sqsubseteq^{\text{avg}} E$	$TG^{\varepsilon_1} \not\sqsubseteq^{\text{max}} E^{\varepsilon_2}$	$TG^{\varepsilon_1} \sqsubseteq^{\text{prv}} E^{\varepsilon_2}$
$E^{\varepsilon_1} \not\sqsubseteq^{\text{avg}} TG$	$E^{\varepsilon_1} \not\sqsubseteq^{\text{max}} TG^{\varepsilon_2}$	$E^{\varepsilon_1} \not\sqsubseteq^{\text{prv}} TG^{\varepsilon_2}$
$G^{\varepsilon_1} \not\sqsubseteq^{\text{avg}} R^{\varepsilon_2}$	$G^{\varepsilon_1} \not\sqsubseteq^{\text{max}} R^{\varepsilon_2}$	$G^{\varepsilon_1} \sqsubseteq^{\text{prv}} R^{\varepsilon_2}$
$R^{\varepsilon_1} \not\sqsubseteq^{\text{avg}} G^{\varepsilon_2}$	$R^{\varepsilon_1} \not\sqsubseteq^{\text{max}} G^{\varepsilon_2}$	$R^{\varepsilon_1} \not\sqsubseteq^{\text{prv}} G^{\varepsilon_2}$
$G^{\varepsilon_1} \not\sqsubseteq^{\text{avg}} E^{\varepsilon_2}$	$G^{\varepsilon_1} \not\sqsubseteq^{\text{max}} E^{\varepsilon_2}$	$G^{\varepsilon_1} \sqsubseteq^{\text{prv}} E^{\varepsilon_2}$
$E^{\varepsilon_1} \not\sqsubseteq^{\text{avg}} G^{\varepsilon_2}$	$E^{\varepsilon_1} \not\sqsubseteq^{\text{max}} G^{\varepsilon_2}$	$E^{\varepsilon_1} \not\sqsubseteq^{\text{prv}} G^{\varepsilon_2}$
$R^{\varepsilon_1} \not\sqsubseteq^{\text{avg}} E^{\varepsilon_2}$	$R^{\varepsilon_1} \not\sqsubseteq^{\text{max}} E^{\varepsilon_2}$	$R^{\varepsilon_1} \not\sqsubseteq^{\text{prv}} E^{\varepsilon_2}$
$E^{\varepsilon_1} \not\sqsubseteq^{\text{avg}} R^{\varepsilon_2}$	$E^{\varepsilon_1} \not\sqsubseteq^{\text{max}} R^{\varepsilon_2}$	$E^{\varepsilon_1} \not\sqsubseteq^{\text{prv}} R^{\varepsilon_2}$

Table 6: Comparing different families of mechanisms with differing  $\varepsilon$ .

mechanism. For example, the blue line on the top graph represents  $TG^{\varepsilon_1} \sqsubseteq^{\text{avg}} E^{\varepsilon_2}$ . We fix  $\varepsilon_1$  and ask for what value of  $\varepsilon_2$  do we get a  $\sqsubseteq^{\text{avg}}$  refinement. Notice that the line  $\varepsilon_1 = \varepsilon_2$  corresponds to the same mechanism in both axes (since every mechanism refines itself).

We can see that refining the randomized response mechanism requires much smaller values of epsilon in the other mechanisms. For example, from the middle graph we can see that  $R^4 \sqsubseteq^{\text{avg}} TG^1$  (approximately) whereas from the top graph we have  $TG^4 \sqsubseteq^{\text{avg}} R^4$ . This means that the randomized response mechanism is very ‘safe’ against average-case adversaries compared with the other mechanisms, as it is much more ‘difficult’ to refine than the other mechanisms.

We also notice that for ‘large’ values of  $\varepsilon_1$ , the exponential and geometric mechanisms refine each other for approximately the same  $\varepsilon_2$  values. This suggests that for these values, the epsilons are comparable (that is, the mechanisms are equally ‘safe’ for similar values of  $\varepsilon$ ). However, smaller values of  $\varepsilon_1$  require a (relatively) large reduction in  $\varepsilon_2$  to obtain a refinement.

## 6.5 Discussion

At the beginning of this section we asked whether it is safe to compare differential privacy mechanisms by ‘comparing the epsilons’. We found that it *is* safe to compare epsilons *within* families of mechanisms (except in the unusual case of the over-truncated geometric mechanism). However, when comparing different mechanisms it is *not* safe to just compare the epsilons, since none of the refinements hold in general. Once a ‘safe’ pair of epsilons has been calculated, then reducing epsilon in the refining mechanism is always safe. However, computing safe epsilons relies on the ability to construct a channel representation, which may not always be feasible.

## 7 Lattice properties

The orders  $\sqsubseteq^{\text{avg}}$ ,  $\sqsubseteq^{\text{max}}$  and  $\sqsubseteq^{\text{prv}}$  are all reflexive and transitive (i.e. preorders), but not anti-symmetric (i.e. not partial orders). This is due to the fact that there exist channels that have “syntactic” differences but the same semantics; eg. two channels having their columns swapped. However, if we are only interested in a specific type of leakage, then all channels such that  $A \sqsubseteq B \sqsubseteq A$  (where  $\sqsubseteq$  is one of  $\sqsubseteq^{\text{avg}}$ ,  $\sqsubseteq^{\text{max}}$ ,  $\sqsubseteq^{\text{prv}}$ ) have identical leakage, so we can view them as the “same channel” (either by working on the equivalence classes of  $\sqsubseteq \cup \supseteq$  or by writing all channels in some canonical form).

Seeing now  $\sqsubseteq$  as a partial order, the natural question is whether it forms a lattice, that is whether suprema and infima exist. If it exists, the supremum  $A \vee B$  has an interesting property: it is the “least safe” channel that is safer than both  $A$  and  $B$  (any channel  $C$  such that  $A \sqsubseteq C$  and  $B \sqsubseteq C$  would necessarily satisfy  $A \vee B \sqsubseteq C$ ). If we wanted a channel that is safer than both  $A$  and  $B$ ,  $A \vee B$  would be a natural choice.

In this section we briefly discuss this problem and show that – in contrast to  $\sqsubseteq^{\text{avg}}$  – both  $\sqsubseteq^{\text{max}}$  and  $\sqsubseteq^{\text{prv}}$  do have suprema and infima (i.e. they form a lattice).

## 7.1 Average-case refinement

In the case of  $\sqsubseteq^{\text{avg}}$ , “equivalent” channels are those producing the exact same hypers. But even if we identify such channels, it is known [19] that two channels  $A, B$  do not necessarily have a *least upper bound* wrt  $\sqsubseteq^{\text{avg}}$ , hence  $\sqsubseteq^{\text{avg}}$  does not form a lattice.

## 7.2 Max-case refinement

In the case of  $\sqsubseteq^{\text{max}}$ , “equivalent” channels are those producing the same posteriors (or more generally the same convex hull of posteriors). But, in contrast to  $\sqsubseteq^{\text{avg}}$ , if we identify such channels, that is if we represent a channel only by the convex hull of its posteriors, then  $\sqsubseteq^{\text{max}}$  becomes a lattice.

First, note that given a finite set of posteriors  $P = \{\delta^y|y\}$ , such that  $\pi \in \mathbf{ch}\{\delta^y\}_y$ , i.e. such that  $\pi = \sum_y a_y \delta^y$ , it is easy to construct a channel  $C$  producing each posterior  $\delta^y$  with output probability  $a_y$ . It suffices to take  $C_{x,y} := \delta_x^y a_y / \pi_x$ .

So  $A \vee^{\text{max}} B$  can be simply constructed by taking the intersection of the convex hulls of the posteriors of  $A, B$ . This intersection is a convex polytope itself, so it has (finitely many) extreme points, so we can construct  $A \vee^{\text{max}} B$  as the channel having exactly those as posteriors.  $A \wedge^{\text{max}} B$ , on the other hand, can be constructed as the channel having as posteriors the union of those of  $A$  and  $B$ .

Note that computing the intersection of polytopes is NP-hard in general [22], so  $A \vee^{\text{max}} B$  might be hard to construct. However, efficient special cases do exist [11]; we leave the study of the hardness of  $\vee^{\text{max}}$  as future work.

## 7.3 Privacy-based refinement

In the case  $\sqsubseteq^{\text{prv}}$ , “equivalent” channels are those producing the same induced metric, i.e.  $\mathbf{d}_A = \mathbf{d}_B$ . Representing channels only by their induced metric, we can use the fact that  $\mathbb{M}\mathcal{X}$  does form a lattice under  $\geq$ . We first show that any metric can be turned into a corresponding channel.

**Theorem 19.** *For any metric  $\mathbf{d} : \mathbb{M}\mathcal{X}$ , we can construct a channel  $C^{\mathbf{d}}$  such that  $\mathbf{d}_{C^{\mathbf{d}}} = \mathbf{d}$ .*

Then  $A \vee^{\text{prv}} B$  will be simply the channel whose metric is  $\mathbf{d}_A \vee \mathbf{d}_B$ , where  $\vee$  is the supremum in the lattice of metrics, and similarly for  $\wedge^{\text{prv}}$ .

Note that the infimum of two metrics  $\mathbf{d}_1, \mathbf{d}_2$  is simply the max of the two (which is always a metric). The supremum, however, is more tricky, since the min of two metrics is not always a metric: the triangle inequality might be violated. So we first need to take the min of  $\mathbf{d}_1, \mathbf{d}_2$ , then compute its “triangle closure”, by finding the shortest path between all pairs of elements, for instance using the well-known Floyd-Warshall algorithm.

## 8 Conclusion

We have investigated various refinement orders for mechanisms for information protection, combining the max-case perspective typical of DP and its variants with the robustness of the QIF approach. We have provided structural characterizations of these preorders and methods to verify them efficiently. Then we have considered various DP mechanisms, and investigated the relation between the  $\varepsilon$ -based measurement of privacy and our orders. We have shown that, while within the same family of mechanisms a smaller  $\varepsilon$  implies the refinement order, this is almost never the case for mechanisms belonging to different families.

## Acknowledgement

This work has been supported by the European Research Council (ERC) under the European Unions Horizon 2020 research and innovation programme, grant agreement N. 835294, by the project ANR-16-CE25-0011 REPAS, and by the Equipe Associée LOGIS.

## A Proofs Omitted from Section 3

**Theorem 1.** *Let  $\pi: \mathbb{D}\mathcal{X}$ . If  $A \sqsubseteq^{\max} B$  then the posteriors of  $B$  (under  $\pi$ ) are convex-combinations of those of  $A$ , that is*

$$\mathbf{supp}[\pi, B] \subseteq \mathbf{ch} \mathbf{supp}[\pi, A]. \quad (11)$$

Moreover, if (11) holds and  $\pi$  is full support then  $A \sqsubseteq^{\max} B$ .

*Proof.* Note that seeing  $\pi$  as a row vector,  $\pi A$  and  $\pi B$  are the output distributions of  $A$  and  $B$  respectively. Denote by  $\alpha^y$  and  $\beta^z$  the posteriors of  $[\pi, A]$  and  $[\pi, B]$  respectively; we have as many posteriors as the elements in the support of the output distributions, that is for each  $y: \mathbf{supp} \pi A, z: \mathbf{supp} \pi B$ . (11) can be written as

$$\forall z: \mathbf{supp} \pi B. \left( \beta^z = \sum_y c_y^z \alpha^y \quad \text{where} \quad \sum_y c_y^z = 1 \right). \quad (38)$$

The proof consists of two parts: first, we show that (38) for uniform  $\pi$  is *equivalent* to  $A \sqsubseteq^{\max} B$ . Second, we show that (38) for full-support  $\pi$  *implies* (38) for any other prior.

For the first part, letting  $\pi$  be uniform, we show that (38) is equivalent to  $R\tilde{A} = \tilde{B}$ . This is easy to see since since the  $y$ -th row of  $\tilde{A}$  is  $\alpha^y$  and the  $z$ -th row of  $\tilde{B}$  is  $\beta^z$ . Hence we can construct  $R$  from the convex coefficients, and vice versa, as  $R_{z,y} = c_y^z$ .

For the second part, let  $\pi: \mathbb{D}\mathcal{X}$  be full-support and  $\hat{\pi}: \mathbb{D}\mathcal{X}$  be arbitrary. Since  $\mathbf{supp} \hat{\pi} \subseteq \mathbf{supp} \pi$ , we necessarily have  $\mathbf{supp} \hat{\pi} C \subseteq \mathbf{supp} \hat{\pi} C$  for any channel  $C$ . Assume that (11) holds for  $\pi$  and let  $c_y^z$  be the corresponding convex coefficients. Fixing an arbitrary  $z: \mathbf{supp} \hat{\pi} B$ , define

$$\hat{c}_y^z := c_y^z \frac{(\hat{\pi} A)_y (\pi B)_z}{(\pi A)_y (\hat{\pi} B)_z}. \quad (39)$$

We first show that

$$\begin{aligned}
& \sum_y c_y^z \frac{(\hat{\pi}A)_y}{(\pi A)_y} \\
&= \sum_x \frac{\hat{\pi}_x}{\pi_x} \sum_y c_y^z \frac{\pi_x A_{x,y}}{(\pi A)_y} && \text{Expand } \hat{\pi}A, \text{ rearrangement} \\
&= \sum_x \frac{\hat{\pi}_x}{\pi_x} \sum_y c_y^z \alpha_x^y && \text{Def. of } \alpha^y \\
&= \sum_x \frac{\hat{\pi}_x}{\pi_x} \beta_x^z && (11) \\
&= \sum_x \frac{\hat{\pi}_x}{\pi_x} \frac{\pi_x B_{x,z}}{(\pi B)_z} && \text{Def. of } \beta^z \\
&= \frac{(\hat{\pi}B)_z}{(\pi B)_z} . && \text{rearrangement}
\end{aligned}$$

From this it follows that  $\sum_y \hat{c}_y^z = 1$ .

Finally, denote by  $\hat{\alpha}^y$  and  $\hat{\beta}^z$  the posteriors of  $[\hat{\pi}, A]$  and  $[\hat{\pi}, B]$  respectively; we show that (11) holds for  $\hat{\pi}$ . Fixing  $x: \mathcal{X}$ , we have that

$$\begin{aligned}
& \sum_y \hat{c}_y^z \hat{\alpha}_x^y \\
&= \sum_y c_y^z \frac{(\hat{\pi}A)_y}{(\pi A)_y} \frac{(\pi B)_z}{(\hat{\pi}B)_z} \frac{\hat{\pi}_x A_{x,y}}{(\hat{\pi}A)_y} && \text{Def. of } c_y^z \text{ and } \hat{\alpha}^y \\
&= \frac{(\pi B)_z}{(\hat{\pi}B)_z} \frac{\hat{\pi}_x}{\pi_x} \sum_y c_y^z \frac{\pi_x A_{x,y}}{(\pi A)_y} && \text{Def. of } d_y^z, \text{ rearrangement} \\
&= \frac{(\pi B)_z}{(\hat{\pi}B)_z} \frac{\hat{\pi}_x}{\pi_x} \sum_y c_y^z \alpha_x^y && \text{Def. of } \alpha^y \\
&= \frac{(\pi B)_z}{(\hat{\pi}B)_z} \frac{\hat{\pi}_x}{\pi_x} \beta_x^z && (11) \\
&= \frac{(\pi B)_z}{(\hat{\pi}B)_z} \frac{\hat{\pi}_x}{\pi_x} \frac{\pi_x B_{x,z}}{(\pi B)_z} && \text{Def. of } \beta^y \\
&= \frac{\hat{\pi}_x B_{x,z}}{(\hat{\pi}B)_z} && \text{Rearrangement} \\
&= \hat{\beta}_x^z . && \text{Def. of } \hat{\beta}^z
\end{aligned}$$

□

**Theorem 2.** *The orders  $\sqsubseteq^{\max}$  and  $\sqsubseteq_{\mathbb{Q}}^{\max}$  coincide.*

*Proof.* Fix some arbitrary  $\pi$  and denote by  $\alpha^y$  and  $\beta^z$  the posteriors of  $[\pi, A]$  and  $[\pi, B]$  respectively. Assuming  $A \sqsubseteq^{\max} B$ , from Theorem. 1 we get that each  $\beta^z$  can be written as a convex combination  $\sum_y c_y^z \alpha^y$ . Hence

$$\begin{aligned}
& V^{\max}[\pi, B] \\
&= \max_z V(\beta^z) && \text{Def. of } V^{\max} \\
&= \max_z V(\sum_y c_y^z \alpha^y) && \text{Theorem. 1} \\
&\leq \max_z \max_y V(\alpha^y) && \text{quasi-convexity of } V \\
&= \max_y V(\alpha^y) \\
&= V^{\max}[\pi, A] ,
\end{aligned}$$



from which  $A \sqsubseteq_{\mathbb{Q}}^{\max} B$  follows.

Now assume that  $A \not\sqsubseteq_{\mathbb{Q}}^{\max} B$ , let  $S = \mathbf{ch\,supp}[\pi, A] \subseteq \mathbb{D}\mathcal{X}$  and define a vulnerability function  $V: \mathbb{Q}\mathcal{X}$  that maps every prior  $\sigma: \mathbb{D}\mathcal{X}$  to its Euclidean distance from  $S$ , that is

$$V(\sigma) := \min_{\sigma' \in S} \|\sigma - \sigma'\|_2 . \quad (40)$$

Since  $S$  is a convex set, it is well known that  $V(\sigma)$  is convex on  $\sigma$  (hence also quasi-convex). Note that  $V(\sigma) = 0$  for all  $\sigma \in S$  and strictly positive anywhere else.

By definition of  $S$  we have that  $\alpha^y \in S$  and hence  $V(\alpha^y) = 0$  for all posteriors of  $A$ , as a consequence  $V^{\max}[\pi, A] = 0$ . On the other hand, since  $A \not\sqsubseteq_{\mathbb{Q}}^{\max} B$ , from Theorem. 1 we get that there exists some posterior of  $B$  such that  $\delta^z \notin S$ . As a consequence  $V^{\max}[\pi, B] \geq V(\delta^z) > 0 = V^{\max}[\pi, A]$  which implies that  $A \not\sqsubseteq_{\mathbb{Q}}^{\max} B$ .  $\square$

**Theorem 3.**  $\sqsubseteq^{\text{avg}}$  is strictly stronger than  $\sqsubseteq^{\max}$ .

*Proof.* The “stronger” part is essentially the data-processing inequality for max-case vulnerability [1, Prop. 14]. To show it directly, assume that  $A \sqsubseteq^{\text{avg}} B$ , that is  $AR = B$  for some channel  $R$ , and define a channel  $S$  from  $\mathcal{Z}$  to  $\mathcal{Y}$  as

$$S_{z,y} := R_{y,z} \frac{\sum_x A_{x,y}}{\sum_x B_{x,z}} . \quad (41)$$

It is easy to check that  $S$  is a valid channel, i.e. that  $\sum_y S_{z,y} = 1$  for all  $z$ . Moreover, we have that

$$\begin{aligned} & (S\tilde{A})_{z,x} \\ &= \sum_y R_{y,z} \frac{\sum_x A_{x,y}}{\sum_x B_{x,z}} \frac{A_{x,y}}{\sum_x A_{x,y}} \text{Def. of } S, \tilde{A} \\ &= \frac{\sum_y A_{x,y} R_{y,z}}{\sum_x B_{x,z}} \text{Algebra} \\ &= \frac{B_{x,z}}{\sum_x B_{x,z}} \quad AR = B \\ &= \tilde{B}_{z,x} , \quad \text{Def. of } \tilde{B} \end{aligned}$$

hence  $A \sqsubseteq^{\max} B$ .

The “strictly” part has already been shown in the body of the paper: The two matrices  $A$  and  $B$  in 15 provide an example in which  $A \sqsubseteq^{\max} B$ , while  $B \not\sqsubseteq^{\max} A$ .  $\square$

## B Proofs Omitted from Section 4

**Proposition 4.** *The orders  $\sqsubseteq_{\mathbb{M}}^{\text{prv}}$  and  $\sqsubseteq^{\text{prv}}$  coincide.*

*Proof.* Assuming  $A \sqsubseteq_{\mathbb{M}}^{\text{prv}} B$ , recall that a channel  $C$  satisfies  $\mathbf{d}$ -privacy iff  $\text{Priv}_{\mathbf{d}}(C) \leq 1$ . Note also that  $\text{Priv}_{\mathbf{d}_C}(C) = 1$ . Setting  $\mathbf{d} = \mathbf{d}_A$  we get  $1 = \text{Priv}_{\mathbf{d}_A}(A) \geq \text{Priv}_{\mathbf{d}_A}(B)$ , which implies that  $B$  satisfies  $\mathbf{d}_A$ -privacy, hence  $A \sqsubseteq^{\text{max}} B$ .

Assuming  $A \sqsubseteq^{\text{prv}} B$ , to show that  $A \sqsubseteq_{\mathbb{M}}^{\text{prv}} B$  it is equivalent to show that  $A$  satisfies  $\mathbf{d}$ -privacy only if  $B$  also satisfies it. Let  $\mathbf{d} \in \mathbb{M}\mathcal{X}$ , if  $A$  satisfies  $\mathbf{d}$ -privacy then  $\mathbf{d} \geq \mathbf{d}_A \geq \mathbf{d}_B$ , hence  $B$  also satisfies  $\mathbf{d}$ -privacy, concluding the proof.  $\square$

**Theorem 5.**  $\sqsubseteq^{\text{max}}$  is strictly stronger than  $\sqsubseteq^{\text{prv}}$ , which is strictly stronger than  $\sqsubseteq_{\mathbf{d}}^{\text{prv}}$ .

*Proof.* The “stronger” part is a direct consequence of the fact that  $\text{Priv}_{\mathbf{d}}(C)$  can be expressed as max-case capacity for a suitable vulnerability measure  $V_{\mathbf{d}}: \mathbb{Q}\mathcal{X}$  (more concretely, a consequence of Theorems 2, 7 and 8). This is discussed in detail in Section 4.7.

For the “strictly” part consider the following counter-example:

$$\begin{array}{|c|c|c|} \hline A & y_1 & y_2 \\ \hline x_1 & 0.8 & 0.2 \\ x_2 & 0.4 & 0.6 \\ \hline \end{array} \quad \begin{array}{|c|c|c|} \hline B & y_1 & y_2 \\ \hline x_1 & 0.4 & 0.6 \\ x_2 & 0.8 & 0.2 \\ \hline \end{array} . \quad (42)$$

The only difference between  $A$  and  $B$  is that the two rows have been swapped. Hence  $\mathbf{d}_A = \mathbf{d}_B$  which implies  $A \sqsubseteq^{\text{prv}} B \sqsubseteq^{\text{prv}} A$ . However, the posteriors of  $A, B$  (for uniform prior) are (written in columns):

$$\begin{array}{|c|c|c|} \hline A & y_1 & y_2 \\ \hline x_1 & 2/3 & 1/4 \\ x_2 & 1/3 & 3/4 \\ \hline \end{array} \quad \begin{array}{|c|c|c|} \hline B & y_1 & y_2 \\ \hline x_1 & 1/3 & 3/4 \\ x_2 & 2/3 & 1/4 \\ \hline \end{array} . \quad (43)$$

Since  $(3/4, 1/4)$  cannot be written as a convex combination of  $(2/3, 1/3)$  and  $(1/4, 3/4)$ , and similarly  $(1/4, 3/4)$  cannot be written as a convex combination of  $(1/3, 2/3)$  and  $(3/4, 1/4)$ , from Theorem 1 we conclude that  $A \not\sqsubseteq^{\text{max}} B \not\sqsubseteq^{\text{max}} A$ .  $\square$

**Theorem 1.** *Let  $f: \mathcal{X} \rightarrow \mathcal{Y}$  be any query and  $A, B$  be two mechanisms on  $\mathcal{Y}$ . If  $A \sqsubseteq^{\text{prv}} B$  then  $A \circ f \sqsubseteq^{\text{prv}} B \circ f$ .*

*Proof.* Define  $\mathbf{d}_{A \circ f}(x_1, x_2) = \mathbf{d}_A(f(x_1), f(x_2))$ , and similarly for  $\mathbf{d}_{B \circ f}$ . Then we have:

$$\mathbf{d}_{A \circ f}(x_1, x_2) = \mathbf{d}_A(f(x_1), f(x_2)) \geq \mathbf{d}_B(f(x_1), f(x_2)) = \mathbf{d}_{B \circ f} . \quad (44)$$

$\square$

**Theorem 7.**  $\mathcal{ML}_{\mathbf{d}}^{+, \text{max}}$  is always achieved on a uniform prior  $\pi^u$ . Namely

$$\max_{\pi} \mathcal{L}_{\mathbf{d}}^{+, \text{max}}(\pi, C) = \mathcal{L}_{\mathbf{d}}^{+, \text{max}}(\pi^u, C) = V_{\mathbf{d}}^{\text{max}}[\pi^u, C] . \quad (27)$$

*Proof.* Fix  $\pi, C$ , and let  $(a, \delta^y)$  and  $(b, \rho^y)$  be the outer and inners of  $[\pi, V_{\mathbf{d}}]$  and  $[\pi^u, V_{\mathbf{d}}]$  respectively. Since  $\delta_x^y = C_{x,y} \pi_x a_y^{-1}$  and  $\rho_x^y = C_{x,y} |\mathcal{X}|^{-1} b_y^{-1}$  we have that:

$$V_{\mathbf{d}}(\delta^y) = \max_{x,x'} \mathbf{d}^{-1}(x, x') \left| \ln \frac{C_{x,y} \pi_x}{C_{x',y} \pi_{x'}} \right|, \quad \text{and} \quad (45)$$

$$V_{\mathbf{d}}(\rho^y) = \max_{x,x'} \mathbf{d}^{-1}(x, x') \left| \ln \frac{C_{x,y}}{C_{x',y}} \right|. \quad (46)$$

Moreover, it holds that:

$$\begin{aligned} & V_{\mathbf{d}}(\delta^y) \\ &= \max_{x,x'} \mathbf{d}^{-1}(x, x') \left| \ln \frac{C_{x,y} \pi_x}{C_{x',y} \pi_{x'}} \right| \\ &\leq \max_{x,x'} \mathbf{d}^{-1}(x, x') \left( \left| \ln \frac{C_{x,y}}{C_{x',y}} \right| + \left| \ln \frac{\pi_x}{\pi_{x'}} \right| \right) \quad \text{triangle inequality} \\ &\leq \max_{x,x'} (\mathbf{d}^{-1}(x, x') \left| \ln \frac{C_{x,y}}{C_{x',y}} \right|) + \max_{x,x'} (\mathbf{d}^{-1}(x, x') \left| \ln \frac{\pi_x}{\pi_{x'}} \right|) \quad \text{independent max} \\ &= V_{\mathbf{d}}(\rho^y) + V_{\mathbf{d}}(\pi) \end{aligned}$$

Finally, we have that:

$$\begin{aligned} & \mathcal{L}_{\mathbf{d}}^{+, \max}(\pi, C) \\ &= \max_y V_{\mathbf{d}}(\delta^y) - V_{\mathbf{d}}(\pi) \\ &\leq \max_y (V_{\mathbf{d}}(\rho^y) + V_{\mathbf{d}}(\pi)) - V_{\mathbf{d}}(\pi) \\ &= \max_y V_{\mathbf{d}}(\rho^y) \\ &= V_{\mathbf{d}}^{\max}[\pi^u, C] \\ &= \mathcal{L}_{\mathbf{d}}^{+, \max}(\pi^u, C) \quad \text{since } V_{\mathbf{d}}(\pi^u) = 0 \end{aligned}$$

which concludes the proof.  $\square$

**Theorem 8.** [DP as max-case capacity]  $C$  satisfies  $\varepsilon \cdot \mathbf{d}$ -privacy iff  $\mathcal{ML}_{\mathbf{d}}^{+, \max}(C) \leq \varepsilon$ . In other words:  $\mathcal{ML}_{\mathbf{d}}^{+, \max}(C) = \text{Priv}_{\mathbf{d}}(C)$ .

*Proof.* Let  $\rho^y$  denote the inners of  $[\pi^u, V_{\mathbf{d}}]$ . From Theorem. 7 we have that

$$\mathcal{ML}_{\mathbf{d}}^{+, \max}(C) \leq \varepsilon \quad \text{iff} \quad V_{\mathbf{d}}(\rho^y) \leq \varepsilon \quad \text{for all } y$$

which, from the definition of  $V_{\mathbf{d}}$ , holds iff  $C_{x,y} \leq e^{\varepsilon \cdot \mathbf{d}(x, x')} C_{x',y}$  for all  $x, x', y$ .  $\square$

## C Proofs Omitted from Section 5

**Proposition 2** (Projection theorem, [5, Prop 1.1.9]). *Let  $C \subset \mathbb{R}^n$  be closed and convex and let  $z \in \mathbb{R}^n$ . There exists a unique  $z^* \in C$  that minimizes  $\|z - x\|_2$  over  $x \in C$ , called the projection of  $z$  on  $C$ . Moreover, a vector  $z^*$  is the projection of  $z$  on  $C$  iff*

$$(z - z^*) \cdot (x - z^*) \leq 0 \quad \text{for all } x \in C. \quad (47)$$

**Theorem 9.** *Let  $B^*$  be the projection of  $B$  on  $A^\uparrow$ .*

1. *If  $B = B^*$  then  $A \sqsubseteq^{\text{avg}} B$ .*
2. *Otherwise, let  $G = B - B^*$ . The gain function  $g(w, x) = G_{x,w}$  provides a counter-example to  $A \sqsubseteq^{\text{avg}} B$ , that is  $V_g(\pi^u, A) < V_g(\pi^u, B)$ , for uniform  $\pi^u$ .*

*Proof.* (1) is immediate from the definition of  $\sqsubseteq^{\text{avg}}$ . For (2), we first show that

$$B \cdot G > B^* \cdot G \geq X \cdot G \quad \text{for all } X \in A^\uparrow, \quad (48)$$

(in other words, that  $X \cdot G = B^* \cdot G$  is a hyperplane with normal  $G$ , separating  $B$  from  $A^\uparrow$ ). For the left-hand inequality we have  $B \cdot G - B^* \cdot G = G \cdot G = \|G\|_2^2 > 0$ . Moreover, since  $A^\uparrow$  is closed and convex, from the projection theorem (Proposition 2) we get that  $(B - B^*) \cdot (X - B^*) \leq 0$  for all  $X \in A^\uparrow$ , from which  $B^* \cdot G \geq X \cdot G$  directly follows.

The proof continues similarly to the one of [19, Thm. 9]. We write posterior vulnerability (for uniform prior) as a maximization over all remapping strategies  $S_A, S_B$  for  $A, B$  respectively, namely

$$V_g(\pi, A) = \frac{1}{|\mathcal{X}|} \max_{AS_A \in A^\uparrow} AS_A \cdot G, \quad (49)$$

$$V_g(\pi, B) = \frac{1}{|\mathcal{X}|} \max_{BS_B \in B^\uparrow} BS_B \cdot G. \quad (50)$$

Then  $V_g(\pi^u, A) < V_g(\pi^u, B)$  follows from (48) and the fact that  $B \in B^\uparrow$  (the identity is a remapping strategy).  $\square$

## D Proofs Omitted from Section 6

We call a truncated geometric mechanism ‘square’ if it has the same input and output space (that is, the channel representation is a square matrix).

We first show that geometric mechanisms and truncated geometric mechanisms are equivalent to square mechanisms under  $\sqsubseteq^{\text{avg}}$ .

**Lemma 3.** *Let  $G^\varepsilon$  be a geometric mechanism. Then the reduced (abstract) channel form of  $G^\varepsilon$  is the square channel  $(\mathcal{X}, \mathcal{X}, TG^\varepsilon)$ .*

*Proof.* First note that the square channel is obtained from the (infinite) geometric by summing up all the ‘extra’ columns (ie. those columns in  $\mathcal{Y} \setminus \mathcal{X}$ ). Now, note that these ‘extra’ columns are scalar multiples of each other, since each column has the form

$$\begin{bmatrix} \frac{(1-\alpha) \cdot \alpha^k}{1+\alpha} \\ \frac{(1-\alpha) \cdot \alpha^{k-1}}{1+\alpha} \\ \frac{(1-\alpha) \cdot \alpha^{k-2}}{1+\alpha} \\ \vdots \end{bmatrix} \quad (51)$$

for increasing values of  $k$ . Thus the ‘summing up’ operation is a valid reduction operation, and so the infinite geometric is reducible to the corresponding square channel.  $\square$

**Lemma 4.** *Let  $TG^\varepsilon$  be a truncated geometric mechanism. Then the reduced (abstract) channel form of  $TG^\varepsilon$  is the square channel  $(\mathcal{X}, \mathcal{X}, TG^\varepsilon)$ .*

*Proof.* First note that the truncated geometric is obtained from the infinite geometric by summing up columns at the ends of the matrix. This is exactly the ‘reduction’ step noted above. We can continue, as above, to sum up ‘extra’ columns until we get a square matrix.  $\square$

**Corollary 5.** *Any  $\sqsubseteq^{\text{avg}}$  refinement that holds for a square geometric mechanism  $(\mathcal{X}, \mathcal{X}, G^\varepsilon)$  also holds for any truncated geometric mechanism or (the) geometric mechanism  $G^\varepsilon$  having domain  $\mathcal{X}$ .*

Note that we only define truncation as far as the square geometric matrix, since at this point the columns of the matrix are linearly independent and can no longer be truncated via matrix reduction operations. We now show that refinement holds for the square geometric mechanisms.

**Lemma 6.** *Let  $TG^\varepsilon$  be a square geometric mechanism. Then decreasing  $\varepsilon$  produces a refinement of it. That is,  $TG^\varepsilon \sqsubseteq^{\text{avg}} TG^{\varepsilon'}$  iff  $\varepsilon \geq \varepsilon'$ .*

*Proof.* The square geometric mechanism  $TG^\varepsilon$  has the following form:

$TG^\varepsilon$	$x_1$	$x_2$	$\dots$	$x_n$
$x_1$	$\frac{1}{1+\alpha}$	$\frac{\alpha \cdot (1-\alpha)}{1+\alpha}$	$\dots$	$\frac{\alpha^{n-1}}{1+\alpha}$
$x_2$	$\frac{\alpha}{1+\alpha}$	$\frac{1-\alpha}{1+\alpha}$	$\dots$	$\frac{\alpha^{n-2}}{1+\alpha}$
$\dots$	$\dots$	$\dots$	$\dots$	$\dots$
$x_n$	$\frac{\alpha^{n-1}}{1+\alpha}$	$\frac{\alpha^{n-2} \cdot (1-\alpha)}{1+\alpha}$	$\dots$	$\frac{1}{1+\alpha}$

(52)

where  $\alpha = e^{-\varepsilon}$ , and similarly for  $TG^{\varepsilon'}$  with  $\alpha = e^{-\varepsilon'}$ .

Now, this matrix is invertible and the inverse has the following form:

$$(TG^\varepsilon)^{-1} = \begin{array}{c|ccccc} & x_1 & x_2 & x_3 & x_4 & \dots \\ \hline x_1 & \frac{1}{1-\alpha} & \frac{-\alpha}{1-\alpha^2} & 0 & 0 & \dots \\ x_2 & \frac{-\alpha}{(1-\alpha)^2} & \frac{1+\alpha^2}{(1-\alpha)^2} & \frac{-\alpha}{(1-\alpha)^2} & 0 & \dots \\ x_3 & 0 & \frac{-\alpha}{(1-\alpha)^2} & \frac{1+\alpha^2}{(1-\alpha)^2} & \frac{-\alpha}{(1-\alpha)^2} & \dots \\ x_4 & 0 & 0 & \frac{-\alpha}{(1-\alpha)^2} & \frac{1+\alpha^2}{(1-\alpha)^2} & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \end{array} \quad (53)$$

Recalling that

$$TG^\varepsilon \sqsubseteq^{\text{avg}} TG^{\varepsilon'} \text{ iff } TG^{\varepsilon'} = TG^\varepsilon R \quad (54)$$

for some channel  $R$ , we can construct a suitable  $R$  using  $(TG^\varepsilon)^{-1}$ , namely  $R = (TG^\varepsilon)^{-1} \cdot TG^{\varepsilon'}$ . It suffices to show that  $R$  is a valid channel.

It is clear that the rows of  $R$  sum to 1, since it is the product of matrices with rows summing to 1.

Multiplying out the matrix  $R$  yields:

$$R = \begin{array}{c|cccc} & x_1 & x_2 & x_3 & \dots \\ \hline x_1 & \frac{1-\alpha\beta}{(1-\alpha)(1+\beta)} & \frac{(\beta-\alpha)(1-\beta)}{(1-\alpha)(1+\beta)} & \frac{\beta(\beta-\alpha)(1-\beta)}{(1-\alpha)(1+\beta)} & \dots \\ x_2 & \frac{(1-\alpha\beta)(\beta-\alpha)}{(1-\alpha)^2(1+\beta)} & \frac{(1-2\alpha\beta+\alpha^2)(1-\beta)}{(1-\alpha)^2(1+\beta)} & \frac{(1-\alpha\beta)(1-\beta)(\beta-\alpha)}{(1-\alpha)^2(1+\beta)} & \dots \\ x_3 & \frac{\beta(1-\alpha\beta)(\beta-\alpha)}{(1-\alpha)^2(1+\beta)} & \frac{(1-\alpha\beta)(1-\beta)(\beta-\alpha)}{(1-\alpha)^2(1+\beta)} & \frac{(1-2\alpha\beta+\alpha^2)(1-\beta)}{(1-\alpha)^2(1+\beta)} & \dots \\ \dots & \dots & \dots & \dots & \dots \end{array} \quad (55)$$

where  $\alpha = e^{-\varepsilon}$  and  $\beta = e^{-\varepsilon'}$ . The only way that any of these matrix entries can be less than 0 is if  $\alpha > \beta$ , or  $\varepsilon < \varepsilon'$ . Thus  $R$  is a valid channel precisely when  $\varepsilon \geq \varepsilon'$  and so  $TG^\varepsilon \sqsubseteq^{\text{avg}} TG^{\varepsilon'}$  as required.  $\square$

The following theorems now follow from the previous lemmas.

**Theorem 10.** *Let  $G^\varepsilon, G^{\varepsilon'}$  be geometric mechanisms. Then  $G^\varepsilon \sqsubseteq^{\text{avg}} G^{\varepsilon'}$  iff  $\varepsilon \geq \varepsilon'$ . That is decreasing  $\varepsilon$  produces a refinement of the mechanism.*

*Proof.* Using Lemma 3 we can express  $G^\varepsilon$  as a square channel and from Lemma 6 it follows that the refinement holds.  $\square$

**Theorem 11.** *Let  $TG^\varepsilon, TG^{\varepsilon'}$  be truncated geometric mechanisms. Then  $TG^\varepsilon \sqsubseteq^{\text{avg}} TG^{\varepsilon'}$  iff  $\varepsilon \geq \varepsilon'$ . That is, decreasing  $\varepsilon$  produces a refinement of the mechanism.*

*Proof.* As above, using Lemma 4 and Lemma 6.  $\square$

**Theorem 7.** *Let  $OTG^\varepsilon, OTG^{\varepsilon'}$  be over-truncated geometric mechanisms. Then  $OTG^\varepsilon \not\sqsubseteq^{\text{avg}} OTG^{\varepsilon'}$  for any  $\varepsilon \neq \varepsilon'$ . That is, decreasing  $\varepsilon$  does **not** produce a refinement.*

*Proof.* Consider the following counter-example:

$A^\varepsilon$	$x_1$	$x_2$
$x_1$	$4/5$	$1/5$
$x_2$	$1/5$	$4/5$
$x_3$	$1/20$	$19/20$

$B^{\varepsilon'}$	$x_1$	$x_2$
$x_1$	$2/3$	$1/3$
$x_2$	$1/3$	$2/3$
$x_3$	$1/6$	$5/6$

(56)

Channels  $A^\varepsilon$  and  $B^{\varepsilon'}$  are over-truncated geometric mechanisms parametrized by  $\varepsilon = 2 \log 2$ ,  $\varepsilon' = \log 2$  respectively. We expect  $B^{\varepsilon'}$  to be safer than  $A^\varepsilon$ , that is,  $V_G[\pi^u, B^{\varepsilon'}] < V_G[\pi^u, A^\varepsilon]$ . However, under the uniform prior  $\pi^u$ , the gain function

$G$	$w_1$	$w_2$
$x_1$	$1/5$	$0$
$x_2$	$0$	$1$
$x_3$	$4/5$	$0$

(57)

yields  $V_G[\pi^u, A^\varepsilon] = 0.33$  and  $V_G[\pi^u, B^{\varepsilon'}] = 0.36$ , thus  $B^{\varepsilon'}$  leaks more than  $A^\varepsilon$  for this adversary. (In fact, for this gain function we have  $V_G[\pi^u, A^\varepsilon] = V_G(\pi^u)$  and so the adversary learns nothing from observing the output of  $A^\varepsilon$ ).  $\square$

**Theorem 13.** *Let  $OTG^\varepsilon$  be an over-truncated geometric mechanism. Then reducing  $\varepsilon$  does **not** produce a  $\sqsubseteq^{\max}$  refinement, however it **does** produce a  $\sqsubseteq^{\text{prv}}$  refinement.*

*Proof.* We show the first part using a counter-example. Consider the following channels:

$A^\varepsilon$	$x_1$	$x_2$
$x_1$	$4/5$	$1/5$
$x_2$	$1/5$	$4/5$
$x_3$	$1/20$	$19/20$

$B^{\varepsilon'}$	$x_1$	$x_2$
$x_1$	$2/3$	$1/3$
$x_2$	$1/3$	$2/3$
$x_3$	$1/6$	$5/6$

(58)

Channels  $A^\varepsilon$  and  $B^{\varepsilon'}$  are over-truncated geometric mechanisms using  $\varepsilon = 2 \log 2$ ,  $\varepsilon' = \log 2$  respectively. We can define the (prior) vulnerability  $V$  as the usual (convex)  $g$ -vulnerability. Then under a uniform prior  $\pi^u$ , the gain function given by:

$$g(w, x_1) = \frac{1}{5} \tag{59}$$

$$g(w, x_2) = -1 \tag{60}$$

$$g(w, x_3) = \frac{4}{5} \tag{61}$$

$$\tag{62}$$

yields  $V^{\max}[\pi^u, A] = 0$  and  $V^{\max}[\pi^u, B] = \frac{2}{55}$ . Thus  $A \not\sqsubseteq^{\max} B$ .

For the second part, we first note that for any square geometric channel  $A^\varepsilon$  we have  $\mathbf{d}_A(x, x') = \varepsilon$  exactly when  $x, x'$  are adjacent rows in the matrix (this can be seen from the construction of the square channel). Now, the over-truncated geometric is obtained by summing columns of the square geometric.

By construction, the square geometric  $A$  has adjacent elements  $A_{x,y}, A_{x',y}$  satisfying  $A_{x,y}/A_{x',y} = e^\varepsilon$  when  $x > x'$  and  $x$  is *above* (or on) the diagonal of the channel matrix, otherwise  $A_{x,y}/A_{x',y} = e^{-\varepsilon}$ . This means that each (over)-truncation step maintains the  $A_{x,y}/A_{x',y}$  ratio except when  $x, y$  and  $x', y'$  occur on diagonal elements, in which case their sum is between  $e^{-\varepsilon}$  and  $e^\varepsilon$ . Since this affects only 2 elements in each row, we still have that  $\mathbf{d}_A(x, x') = \varepsilon$  (until the final truncation step to produce a single 1 vector). Therefore since  $\sqsubseteq^{\text{prv}}$  holds for the square matrix, and it holds under truncation, we must have that it holds for over-truncated geometric mechanisms. Thus reducing  $\varepsilon$  corresponds to refinement under  $\sqsubseteq^{\text{prv}}$  as required.  $\square$

We show that the randomized response mechanism behaves well with respect to  $\sqsubseteq^{\text{avg}}$  by considering 3 cases.

Firstly, we consider the case where  $\mathcal{X} = \mathcal{Y}$ . We use the following lemmas:

**Lemma 8.** *Let  $R^\alpha, R^\beta$  be ‘square’ randomized response mechanisms. Then  $B = R^\alpha R^\beta$  is a randomized response mechanism with parameter  $\varepsilon = \log \frac{e^{\alpha+\beta} + k}{e^\alpha + e^\beta + k - 1}$  where  $k + 1$  is the dimension of  $R^\alpha, R^\beta$  and  $B$ .*

*Proof.* Observe that  $R^\alpha$  can be factorised as  $\frac{1}{e^\alpha + k} R$  where  $R$  has the form:

$$\begin{array}{c|cccc} R & x_1 & x_2 & \dots & x_{k+1} \\ \hline x_1 & e^\alpha & 1 & \dots & 1 \\ x_2 & 1 & e^\alpha & \dots & 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{k+1} & 1 & 1 & \dots & e^\alpha \end{array} \quad (63)$$

and similarly for  $R^\beta$ . Multiplying out gives the matrix:

$$\begin{array}{c|cccc} B & x_1 & x_2 & \dots & x_{k+1} \\ \hline x_1 & e^{\alpha+\beta} + k & e^\alpha + e^\beta + (k-1) & \dots & e^\alpha + e^\beta + (k-1) \\ x_2 & e^\alpha + e^\beta + (k-1) & e^{\alpha+\beta} + k & \dots & e^\alpha + e^\beta + (k-1) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{k+1} & e^\alpha + e^\beta + (k-1) & e^\alpha + e^\beta + (k-1) & \dots & e^{\alpha+\beta} + k \end{array} \quad (64)$$

(Note that the constant co-efficient factorised out the front does not affect the  $\varepsilon$  calculation for the channel). This is exactly the randomized response mechanism required.  $\square$

**Lemma 9.** *For any  $a \geq 1, b \geq 0$  the function*

$$f(x) = \frac{ae^x + b}{e^x + a + b - 1} \quad (65)$$

*defined for  $x \geq 0$  is increasing and has range  $[1, a)$ .*

*Proof.* We can see that  $f(x)$  is continuous for the given domain and the derivative  $f'(x) = \frac{e^x(a-1)(a+b)}{(e^x + a + b - 1)^2}$  is  $\geq 0$  for all  $a \geq 1, b \geq 0$ .



Additionally, at  $x = 0$  the function is defined and equal to 1. And

$$\lim_{x \rightarrow \infty} \frac{ae^x + b}{e^x + a + b - 1} = \lim_{x \rightarrow \infty} \frac{ae^x}{e^x} \quad (66)$$

$$= a \quad (67)$$

□

**Lemma 10.** *Let  $R^\varepsilon, R^{\varepsilon'}$  be randomized response mechanisms represented by square matrices (that is,  $\mathcal{X} = \mathcal{Y}$ ). Then  $R^\varepsilon \sqsubseteq^{\text{avg}} R^{\varepsilon'}$  iff  $\varepsilon \geq \varepsilon'$ .*

*Proof.* Note first that  $R^\varepsilon, R^{\varepsilon'}$  are in reduced (abstract channel) form and so the partial order of  $\sqsubseteq^{\text{avg}}$  holds.

From Lemma 8, we know that the composition of 2 randomized response mechanisms is another randomized response mechanism. Therefore, for the reverse direction, if  $\varepsilon > \varepsilon'$  then Lemma 9 tells us that we can find a randomized response mechanism  $R'$  such that  $R^\varepsilon R' = R^{\varepsilon'}$ . In the case of equality, we can choose the identity mechanism.

For the forward direct we show the contrapositive. If  $\varepsilon < \varepsilon'$  then we know there exists an  $R$  such that  $R^{\varepsilon'} R = R^\varepsilon$ . But this means that  $R^{\varepsilon'} \sqsubseteq^{\text{avg}} R^\varepsilon$ . Since the matrices are reduced (as channels), then  $\sqsubseteq^{\text{avg}}$  is a partial order and so this implies  $R^\varepsilon \not\sqsubseteq^{\text{avg}} R^{\varepsilon'}$ . Thus we must have  $R^\varepsilon \sqsubseteq^{\text{avg}} R^{\varepsilon'} \implies \varepsilon \geq \varepsilon'$ . □

Interestingly, we can use this result to show the second case where we consider ‘over-truncated’ mechanisms.

**Lemma 11.** *Let  $R^\varepsilon, R^{\varepsilon'}$  be randomized response mechanisms with  $\mathcal{Y} \subset \mathcal{X}$ . Then  $R^\varepsilon \sqsubseteq^{\text{avg}} R^{\varepsilon'}$  iff  $\varepsilon \geq \varepsilon'$ .*

*Proof.* Notice that this ‘over-truncated’ randomized response mechanism is just a square randomized response mechanism ‘glued’ onto a matrix containing all values  $\frac{1}{n+1}$ .

Denote the corresponding square mechanisms by  $S^\varepsilon, S^{\varepsilon'}$  (note the parameters are the same) and denote by  $N$  the matrix containing only  $\frac{1}{n+1}$  (note that this is the same matrix for both  $R^\varepsilon$  and  $R^{\varepsilon'}$ ).

From Lemma 10 we can find a square randomized response mechanism  $R$  satisfying  $S^\varepsilon R = S^{\varepsilon'}$  iff  $\varepsilon \geq \varepsilon'$ . Notice that  $R$  must be doubly symmetric, since its  $i$ th row is the same as its  $i$ th column. Thus the dot product of any column of  $R$  with any row vector containing only  $\frac{1}{n+1}$  must yield  $\frac{1}{n+1}$ . And so we must have  $N * R = N$ . Now we have that  $R^\varepsilon R$  is just  $S^\varepsilon R$  glued onto  $N$  which is the same as  $S^{\varepsilon'}$  glued onto  $N$ . And thus  $R$  also satisfies  $R^\varepsilon R = R^{\varepsilon'}$ . And so, following the same arguments as for Lem 10, we have  $R^\varepsilon \sqsubseteq^{\text{avg}} R^{\varepsilon'}$  iff  $\varepsilon \geq \varepsilon'$ . □

Finally, we consider the case where  $\mathcal{X} \subset \mathcal{Y}$ .

**Lemma 12.** *Let  $R^\varepsilon, R^{\varepsilon'}$  be randomized response mechanisms with  $\mathcal{X} \subset \mathcal{Y}$ . Then  $R^\varepsilon \sqsubseteq^{\text{avg}} R^{\varepsilon'}$  iff  $\varepsilon \geq \varepsilon'$ .*

*Proof.* The channel matrix  $R^\varepsilon$  is equivalent to the square randomized response mechanism  $S^\varepsilon$  of dimension  $\|\mathcal{Y}\| \times \|\mathcal{Y}\|$  with the bottom  $\|\mathcal{Y}\| - \|\mathcal{X}\|$  rows removed (and similarly for  $R^{\varepsilon'}$ ). This means any solution  $R$  for  $S^\varepsilon$  is a solution for  $R^\varepsilon$ . So, for the reverse direction, if  $\varepsilon \geq \varepsilon'$ , we can always find an  $R$  such that  $R^\varepsilon \sqsubseteq^{\text{avg}} R^{\varepsilon'}$ .

For the forward direction we prove the contrapositive. Note that in this case we cannot assume a partial order relation for  $\sqsubseteq^{\text{avg}}$ , since there may be columns of  $R^\varepsilon$  which are identical. If  $\varepsilon < \varepsilon'$  we want to show that  $R^\varepsilon \not\sqsubseteq^{\text{avg}} R^{\varepsilon'}$ . To do this we need to find a gain function and prior  $\pi$  such that  $V_g(\pi, R^\varepsilon) < V_g(\pi, R^{\varepsilon'})$ . The min-entropy leakage will do: this is simply the sum of the column maxima of the channel matrix. For the randomized response channels, this is given by

$$V(R^\varepsilon) = \frac{ae^\varepsilon + b}{e^\varepsilon + a + b - 1} \quad (68)$$

for a channel with dimensions  $a \times (a + b)$ . Since this is an increasing function of  $\varepsilon$  (for  $a \geq 1$ ),<sup>21</sup> we must have  $\varepsilon < \varepsilon' \implies V(R^\varepsilon) < V(R^{\varepsilon'})$ . Thus  $R^\varepsilon \not\sqsubseteq^{\text{avg}} R^{\varepsilon'}$ .  $\square$

We can now conclude the following theorem from the main body of the paper:

**Theorem 14.** *Let  $R^\varepsilon$  be a randomized response mechanism. Then decreasing  $\varepsilon$  in  $R$  produces a refinement. That is,  $R^\varepsilon \sqsubseteq^{\text{avg}} R^{\varepsilon'}$  iff  $\varepsilon \geq \varepsilon'$ .*

*Proof.* Follows from Lemma 10, Lemma 11 and Lemma 12.  $\square$

**Theorem 15.** *Let  $R$  be a randomized response mechanism,  $E$  an exponential mechanism and  $TG$  a truncated geometric mechanism. Then  $TG^\varepsilon \sqsubseteq^{\text{prv}} R^\varepsilon$  and  $TG^\varepsilon \sqsubseteq^{\text{prv}} E^\varepsilon$ . However  $\sqsubseteq^{\text{prv}}$  does not hold between  $E^\varepsilon$  and  $R^\varepsilon$ .*

*Proof.* We first show that  $TG^\varepsilon \sqsubseteq^{\text{prv}} R^\varepsilon$ , which is equivalent to showing that  $\mathbf{d}_R \leq \mathbf{d}_{TG}$ . For the geometric mechanism, we have  $\mathbf{d}_{TG}(x, x') = \varepsilon \mathbf{d}(x, x')$  for all  $x, x' \in \mathcal{X}$ . For the randomized response mechanism, we have  $\mathbf{d}_R(x, x') = \varepsilon$  or  $\mathbf{d}_R(x, x') = 0$  (when  $x, x' \notin \mathcal{Y}$ ). Thus  $\mathbf{d}_R \leq \mathbf{d}_{TG}$  and so  $TG^\varepsilon \sqsubseteq^{\text{prv}} R^\varepsilon$ .

We now show  $TG^\varepsilon \sqsubseteq^{\text{prv}} E^\varepsilon$ . Recall that we parametrize the exponential mechanism by the smallest possible  $\varepsilon$  such that it satisfies  $\varepsilon \mathbf{d}$ -privacy. In this case, we find that for any pair  $x, x'$  we have  $\mathbf{d}_E(x, x') \leq \varepsilon \mathbf{d}(x, x')$  whereas for the geometric mechanism we have  $\mathbf{d}_{TG}(x, x') = \varepsilon \mathbf{d}(x, x')$ . Therefore  $\mathbf{d}_E \leq \mathbf{d}_{TG}$  and so  $TG^\varepsilon \sqsubseteq^{\text{prv}} E^\varepsilon$ .

The proof of  $\sqsubseteq^{\text{prv}}$  not holding between  $E^\varepsilon$  and  $R^\varepsilon$  was provided in the main body of the paper.  $\square$

<sup>21</sup>Since the derivative is always positive for  $a > 1$

**Theorem 18.** *Every (truncated geometric, randomized response, exponential) mechanism is ‘the safest possible mechanism’ when parametrized by  $\varepsilon = 0$ . That is  $L^\varepsilon \sqsubseteq^{\text{avg}} M^0$  for all mechanisms  $L, M$  (possibly from different families) and  $\varepsilon > 0$ .*

*Proof.* The intuition is that all channels parametrized by  $\varepsilon = 0$  are equivalent to the **1** channel (that is, the  $m \times 1$  channel consisting only of 1s). Indeed, the exponential and randomized response mechanisms parametrized by  $\varepsilon = 0$  have every element equal to  $\frac{1}{n}$  where  $n = \|\mathcal{Y}\|$ . These clearly reduce to the **1** channel. The truncated geometric mechanism contains all 0s except for the first and last column which contain  $\frac{1}{2}$ . Again, this reduces to the **1** channel. Since the **1** channel refines everything (that is,  $L \sqsubseteq^{\text{avg}} \mathbf{1}$  for any channel  $L$ ), the result follows.  $\square$

## E Proofs Omitted from Section 7

**Theorem 19.** *For any metric  $\mathbf{d} : \mathbb{M}\mathcal{X}$ , we can construct a channel  $C^{\mathbf{d}}$  such that  $\mathbf{d}_{C^{\mathbf{d}}} = \mathbf{d}$ .*

*Proof.* Let  $\mathbf{d} : \mathbb{M}\mathcal{X}$ , we first show that for any  $x_0 : \mathcal{X}$ , we can construct a channel  $C^{x_0}$  whose induced metric is below  $\mathbf{d}$ , but coincides with it on all distances to  $x_0$ , that is:

$$\mathbf{d}_{C^{x_0}} \leq \mathbf{d} \quad \text{and} \quad \mathbf{d}_{C^{x_0}}(x_0, x) = \mathbf{d}(x_0, x) \quad \text{for all } x : \mathcal{X} . \quad (69)$$

To construct  $C^{x_0}$ , we use just two outputs (i.e.  $\mathcal{Y} = \{y_1, y_2\}$ ) and we use the fact that  $\text{tv}_\otimes$  on  $\mathbb{D}\mathcal{Y}$  admits a *geodesic*, that is a curve  $\gamma : [0, +\infty] \rightarrow \mathbb{D}\mathcal{Y}$  such that

$$\text{tv}_\otimes(\gamma(t), \gamma(t')) = |t - t'| \quad \text{for all } t, t' : [0, +\infty] . \quad (70)$$

For instance, we can check that  $\gamma(t) = (e^{-t-1}, 1 - e^{-t-1})$  is such a geodesic.

We can now use the geodesic, to assign probability distributions on each secret such that the properties (69) are satisfied. Concretely, define each row  $x$  of  $C^{x_0}$  as:

$$C_{x,-}^{x_0} := \gamma(\mathbf{d}(x_0, x)) . \quad (71)$$

We now check that the properties (69) are satisfied:

$$\begin{aligned} & \mathbf{d}_{C^{x_0}}(x_1, x_2) \\ &= \text{tv}_\otimes(\gamma(\mathbf{d}(x_0, x_1)), \gamma(\mathbf{d}(x_0, x_2))) && \text{Def. of } \mathbf{d}_C, C^{x_0} \\ &= |\mathbf{d}(x_0, x_1) - \mathbf{d}(x_0, x_2)| && \gamma \text{ is a geodesic} \\ &\leq \mathbf{d}(x_1, x_2) && \text{triangle ineq. for } \mathbf{d} \end{aligned}$$

and also

$$\begin{aligned}
& \mathbf{d}_{C^{x_0}}(x_0, x) \\
&= \text{tv}_{\otimes}(\gamma(\mathbf{d}(x_0, x_0)), \gamma(\mathbf{d}(x_0, x))) && \text{Def. of } \mathbf{d}_C, C^{x_0} \\
&= |\mathbf{d}(x_0, x_0) - \mathbf{d}(x_0, x)| && \gamma \text{ is a geodesic} \\
&= \mathbf{d}(x_0, x) . && \mathbf{d}(x_0, x_0) = 0
\end{aligned}$$

Finally,  $C^{\mathbf{d}}$  is constructed as the visible choice of all  $\{C^x\}_x$ . As a consequence,  $\mathbf{d}_{C^{\mathbf{d}}}$  will be the max of the corresponding induced metrics  $\{\mathbf{d}_{C^x}\}_x$ , from which and (69) we can easily conclude that  $\mathbf{d}_{C^{\mathbf{d}}} = \mathbf{d}$ .

Finally, note that the visible choice adds the columns of all mechanisms, so the constructed channel has  $2|\mathcal{X}|$  columns. However, the equality of distances in (69) is given by the *first column* of  $C^{\tilde{x}}$  (this is because of the way  $\gamma$  is constructed), hence we can merge all second columns together, giving finally a simple construction for  $C^{\mathbf{d}}$  with  $\mathcal{Y} = \mathcal{X} \cup \{\perp\}$  (i.e. having  $|\mathcal{X}| + 1$  columns)

$$C_{x,y}^{\mathbf{d}} = |\mathcal{X}|^{-1} e^{-\mathbf{d}(x,y)-1}, \quad x, y \in \mathcal{X}, \quad (72)$$

$$C_{x,\perp}^{\mathbf{d}} = 1 - |\mathcal{X}|^{-1} \sum_{y \in \mathcal{X}} e^{-\mathbf{d}(x,y)-1}, \quad x \in \mathcal{X}. \quad (73)$$

□

## References

- [1] Mário S. Alvim, Konstantinos Chatzikokolakis, Annabelle McIver, Carroll Morgan, Catuscia Palamidessi, and Geoffrey Smith. Axioms for information leakage. In *Proceedings of the 29th IEEE Computer Security Foundations Symposium (CSF)*, pages 77–92, 2016.
- [2] Mário S. Alvim, Konstantinos Chatzikokolakis, Catuscia Palamidessi, and Geoffrey Smith. Measuring information leakage using generalized gain functions. In *Proceedings of the 25th IEEE Computer Security Foundations Symposium (CSF)*, pages 265–279, 2012.
- [3] Miguel E. Andrés, Nicolás E. Bordenabe, Konstantinos Chatzikokolakis, and Catuscia Palamidessi. Geo-indistinguishability: differential privacy for location-based systems. In *Proceedings of the 20th ACM Conference on Computer and Communications Security (CCS 2013)*, pages 901–914. ACM, 2013.
- [4] Gilles Barthe and Boris Köpf. Information-theoretic bounds for differentially private mechanisms. In *Proceedings of the 24th IEEE Computer Security Foundations Symposium (CSF)*, pages 191–204. IEEE Computer Society, 2011.
- [5] Dimitri P Bertsekas. *Convex optimization theory*. Athena Scientific Belmont, 2009.

- [6] Konstantinos Chatzikokolakis, Miguel E. Andrés, Nicolás E. Bordenabe, and Catuscia Palamidessi. Broadening the scope of Differential Privacy using metrics. In Emiliano De Cristofaro and Matthew Wright, editors, *Proceedings of the 13th International Symposium on Privacy Enhancing Technologies (PETs 2013)*, volume 7981 of *Lecture Notes in Computer Science*, pages 82–102. Springer, 2013.
- [7] Konstantinos Chatzikokolakis, Natasha Fernandes, and Catuscia Palamidessi. Comparing systems: max-case refinement orders and application to differential privacy. In *Proceedings of the 32nd IEEE Computer Security Foundations Symposium*, pages 442–457, Hoboken, United States, 2019.
- [8] David Clark, Sebastian Hunt, and Pasquale Malacaria. Quantitative information flow, relations and polymorphic types. *J. of Logic and Computation*, 18(2):181–199, 2005.
- [9] John C. Duchi, Michael I. Jordan, and Martin J. Wainwright. Local privacy and statistical minimax rates. In *Proceedings of the 54th Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 429–438. IEEE Computer Society, 2013.
- [10] Cynthia Dwork, Frank Mcsherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In Shai Halevi and Tal Rabin, editors, *In Proceedings of the Third Theory of Cryptography Conference (TCC)*, volume 3876 of *Lecture Notes in Computer Science*, pages 265–284. Springer, 2006.
- [11] Komei Fukuda, Thomas M. Liebling, and Christine Lütolf. Extended convex hull. *Comput. Geom.*, 20:13–23, 2000.
- [12] Arpita Ghosh and Aaron Roth. Selling privacy at auction. In *Proceedings of the 12th ACM Conference on Electronic Commerce, EC ’11*, pages 199–208, New York, NY, USA, 2011. ACM.
- [13] Justin Hsu, Marco Gaboardi, Andreas Haeberlen, Sanjeev Khanna, Arjun Narayan, Benjamin C. Pierce, and Aaron Roth. Differential privacy: An economic method for choosing epsilon. In *IEEE 27th Computer Security Foundations Symposium, CSF 2014, Vienna, Austria, 19-22 July, 2014*, pages 398–410. IEEE Computer Society, 2014.
- [14] Daniel Kifer and Ashwin Machanavajjhala. Pufferfish: A framework for mathematical privacy definitions. *ACM Transactions on Database Systems*, 39(1):3:1–3:36, 2014.
- [15] Boris Köpf and David A. Basin. An information-theoretic model for adaptive side-channel attacks. In Peng Ning, Sabrina De Capitani di Vimercati, and Paul F. Syverson, editors, *Proceedings of the 2007 ACM Conference on Computer and Communications Security (CCS 2007)*, pages 286–296. ACM, 2007.

- [16] Jaisook Landauer and Timothy Redmond. A lattice of information. In *Proc. Computer Security Foundations Workshop VI*, pages 65–70, June 1993.
- [17] Pasquale Malacaria. Assessing security threats of looping constructs. In Martin Hofmann and Matthias Felleisen, editors, *Proceedings of the 34th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL 2007)*, pages 225–235. ACM, 2007.
- [18] Pasquale Malacaria. Algebraic foundations for quantitative information flow. *Mathematical Structures in Computer Science*, 25(2):404–428, 2015.
- [19] Annabelle McIver, Carroll Morgan, Geoffrey Smith, Barbara Espinoza, and Larissa Meinicke. Abstract channels and their robust information-leakage ordering. In Martin Abadi and Steve Kremer, editors, *Proceedings of the Third International Conference on Principles of Security and Trust (POST)*, volume 8414 of *Lecture Notes in Computer Science*, pages 83–102. Springer, 2014.
- [20] Shiva Prasad Kasiviswanathan and Adam Smith. On the ‘semantics’ of differential privacy: A bayesian formulation. *Journal of Privacy and Confidentiality*, 6, 03 2008.
- [21] Geoffrey Smith. On the foundations of quantitative information flow. In Luca de Alfaro, editor, *Proceedings of the 12th International Conference on Foundations of Software Science and Computation Structures (FOSSACS 2009)*, volume 5504 of *LNCS*, pages 288–302, York, UK, 2009. Springer.
- [22] Hans Raj Tiwary. On the hardness of computing intersection, union and minkowski sum of polytopes. *Discrete & Computational Geometry*, 40(3):469–479, Oct 2008.
- [23] Hirotoshi Yasuoka and Tachio Terauchi. Quantitative information flow - verification hardness and possibilities. In *Proceedings of the 23rd IEEE Computer Security Foundations Symposium, CSF 2010, Edinburgh, United Kingdom, July 17-19, 2010*, pages 15–27. IEEE Computer Society, 2010.

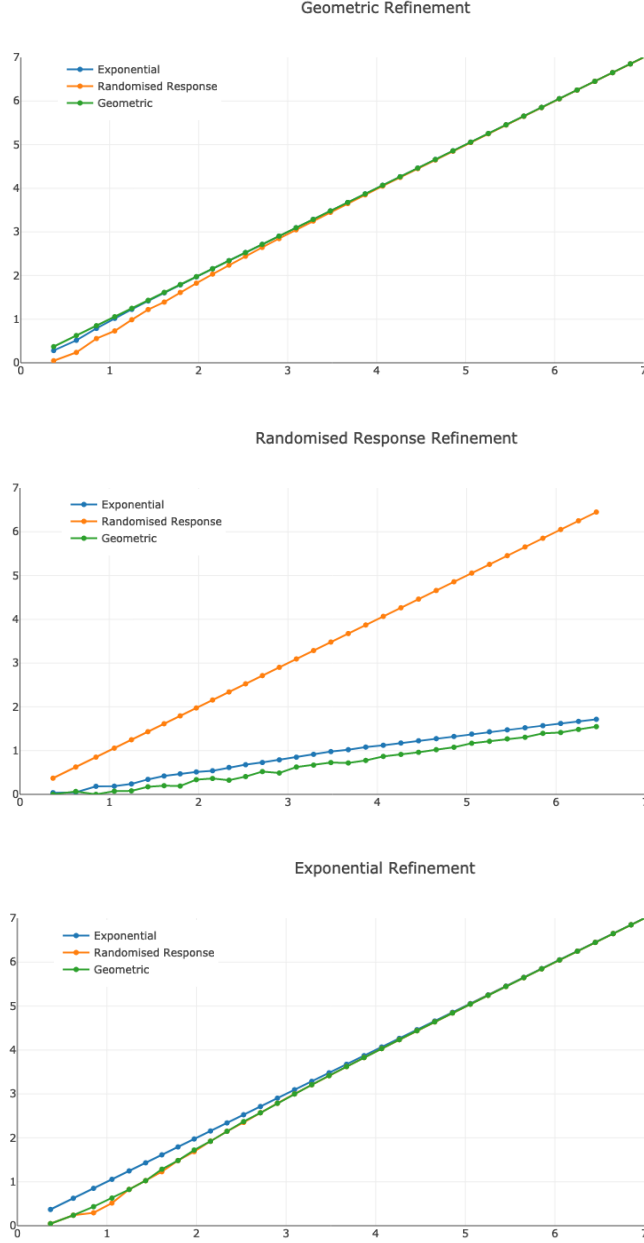


Figure 3: Refinement of mechanisms under  $\sqsubseteq^{avg}$  for  $5 \times 5$  channels. The x-axis represents the  $\varepsilon$  on the LHS of the relation, and the y-axis represents the one on the RHS. The top graph represents refinement of the truncated geometric mechanism (that is,  $TG \sqsubseteq^{avg}$ ), the middle graph is refinement of randomized response ( $R \sqsubseteq^{avg}$ ), and the bottom graph is refinement of the exponential mechanism ( $E \sqsubseteq^{avg}$ ).