

---

# Extending the Foundations of Differential Privacy: Robustness and Flexibility

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

Differential Privacy (DP) is an area that has recently seen many direct and indirect applications to machine learning. In this work, we make foundational contributions to the area of DP.

Our first contribution is definitional. We define two complementary concepts that greatly enhance the applicability of DP, namely, robust privacy and flexible accuracy. Robust privacy requires that a mechanism provides the "best possible privacy" without further degrading accuracy guarantees, even if such privacy is not a priori anticipated based on input neighborhoods alone. Flexible accuracy allows small distortions in the input (e.g., dropping outliers) before measuring accuracy of the output. Along the way, we also extend the notion of DP to sampling (i.e. computation of randomized functions).

Our second contribution is in establishing versatile composition theorems that relate these notions.

Our third contribution is constructive: We present mechanisms that can help in achieving these notions, where previously no meaningful differentially private mechanisms were possible. In particular, we illustrate an application to differentially private histograms, which in turn yields mechanisms for revealing the support of a dataset or the extremal values in the data.

## 1 Introduction

### Blurb to connect DP and ML.

Below, we identify and address two limitations of the DP framework that seem to have evaded attention. At a high-level, these limitations follow from a seemingly natural choice: Accuracy guarantees of a mechanism are in terms of distances in the output space, and privacy demands are in terms of distances in the input space (neighboring inputs). Somewhat surprisingly, these choices turn out to be not always adequate. Our extensions can be seen as adding accuracy guarantees in terms of distances (or rather, distortions) in the input space, and privacy demands in terms of distances in the output space. Along the way, we extend the notion of DP to randomized functions over a metric space, for which distances are measured using a (generalization of) Wasserstein distance. These extensions greatly expand the scope of DP. Apart from the direct implications to privacy of training data, we anticipate that the implications of DP to generalization guarantees (as shown recently in [cite](#)) will also be strengthened by these extensions.

We start by discussing the limitations of DP.

**Lack of Flexibility.** Consider a simplistic learning task which tries to learn an upper bound on integer valued observations – say, ages of patients who recovered from a certain disease – presented to

35 it. For the sake of privacy, one may wish to apply a DP mechanism, rather than output the maximum  
 36 in the sample itself. Two possible datasets which differ in only one patient are considered neighbors  
 37 and a DP mechanism needs to make the outputs on these two samples indistinguishable from each  
 38 other. However, the function in question is *highly sensitive* – two neighboring datasets can have their  
 39 maxima differ by as much as the entire range of possible ages – and the standard DP mechanisms in  
 40 the literature will add so much noise that no useful information can be retained.<sup>1</sup>

41 As we shall see, the above limitation can be attributed to a rigidly defined notion of accuracy. This  
 42 same rigidity leads to another surprising limitation too. Consider the problem of reporting a *histogram*  
 43 (again, say, of patients’ ages). Here a standard DP mechanism, of adding a zero-mean Laplace noise  
 44 to each bar of the histogram is indeed reasonable, as the histogram function has low sensitivity  
 45 in each bar. Now, note that *maximum can be computed as a function of the histogram*. However,  
 46 even though the histogram mechanism was sufficiently accurate in the standard sense, the maximum  
 47 computed from its output is no longer accurate! This is because when a non-zero count is added to a  
 48 large-valued item which originally has a count of 0, the maximum can increase arbitrarily.

49 In this work we develop a more relaxed notion of accuracy, called *flexible accuracy*, that lets us address  
 50 both of the above issues. In particular, it not only enables new DP mechanisms for maximum, but  
 51 also allows one to derive the mechanism from a new DP mechanism for histograms. A composition  
 52 theorem enables us to transfer the accuracy guarantees on histogram to accuracy guarantees on the  
 53 maximum function.

54 **Lack of Robustness.** Differential Privacy focuses on making outputs from *neighboring* databases  
 55 indistinguishable, where neighborhood usually refers to databases obtained by adding or deleting  
 56 a small number of data items (or a single one). However, such a notion of neighborhood of the  
 57 databases may not capture all pairs of databases that *should be* indistinguishable from each other.

58 Consider training a machine learning model on either dataset  $D_1$  or dataset  $D_2$ , where the two  
 59 datasets are disjoint. Suppose both the datasets are representative and yield very similar models. In  
 60 this case, we may reasonably require that querying a model should not reveal whether it was trained  
 61 on  $D_1$  or  $D_2$ . Indeed, since the models are “similar,” one may expect them to yield results which  
 62 are indistinguishable from each other. Unfortunately, this is not generally true: Similarity of outputs  
 63 is measured in terms of a distance in the output space (or rather, the Wasserstein distance over that  
 64 space, since the output is probabilistic); but the extent of their indistinguishability is measured in  
 65 terms of total variation distance or the ratio of probabilities (as in DP), which are not influenced by  
 66 the metric space associated with the outputs. For instance, if the output from the model trained in  $D_1$   
 67 has an even value for the least significant digit, and the other has an odd value, the total variation  
 68 distance between the two output distributions is maximum, while the Wasserstein distance can be  
 69 very small.

70 In short, DP only guarantees indistinguishability between datasets which are close to each other  
 71 in the input space, whereas one may demand – without necessarily compromising on accuracy –  
 72 indistinguishability between datasets which result in outputs that are close to each other. Robustness  
 73 is a complementary notion defined for a mechanism that addresses this.

## 74 1.1 Our Contributions

75 **blurb**

76 **Flexibility.** Flexible accuracy is a notion that is designed to salvage the situation for functions like  
 77 maximum. The high-level idea is to allow for some *distortion of the input* when measuring accuracy.  
 78 **We shall require** distortion to be defined using a *quasi-metric* over the input space (a quasi-metric is  
 79 akin to a metric, but is not required to be symmetric). A typical form of distortion is to *drop a few*  
 80 *items* from the dataset; in this case, adding a data item is not considered low distortion. Referring  
 81 back to the example of reporting maximum, given a dataset with a single elderly patient and many  
 82 young patients, flexible accuracy with respect to this distortion allows a mechanism for maximum to  
 83 report the maximum age of the younger group.

<sup>1</sup>Indeed, *all datasets* with low maximum values have high sensitivity *locally*, by considering a neighboring dataset with a single additional data item with a large value. As such, mechanisms which add noise based on the local sensitivity rather than global sensitivity **cite** also do not fare any better.

84 Flexible accuracy also provides us with a means for *transferring accuracy guarantees* when composed  
 85 with other functions or mechanisms. Consider again the example of the histogram and maximum  
 86 functions from above. Recall that even a high (but less than perfect) accuracy of histograms under  
 87 a metric in the output space can result in maximum computed from the histogram to be wildly  
 88 inaccurate. But if the inaccuracy in the histogram can be entirely attributed to a distortion in the input,  
 89 computing maximum on this histogram does not amplify the inaccuracy at all.

90 **Robustness.** We define a mechanism whose output is in a metric space to be robust if, roughly, it  
 91 holds that whenever two input distributions result in output distributions that are close in Wasserstein  
 92 distance, then the output distributions are also indistinguishable in the sense of differential privacy.  
 93 Unlike in the definition of differential privacy, where an input neighborhood is specified, here the  
 94 neighborhood is implicitly defined by the mechanism itself.

95 **Composition Theorems.**

96 **Constructions.**

97 **Related Work.** Book [1]

## 98 2 Preliminaries

99 **Notations.** We denote by  $\mathbb{N}$  all the non-negative integers (including zero). For  $i, j \in \mathbb{N}$ , such that  
 100  $i \leq j$ , we write  $[i : j]$  to denote the set  $\{i, i + 1, \dots, j\}$ .

101 **Definition 1** (Total Variation Distance). *Let  $p$  and  $q$  be two discrete probability distributions on a*  
 102 *sample space  $\Omega$ . The total variation distance between  $p$  and  $q$ , denoted by  $\Delta(p, q)$ , is defined as*  
 103 *follows:*

$$\Delta(p, q) = \frac{1}{2} \sum_{\omega \in \Omega} |p(\omega) - q(\omega)|.$$

104 Let  $\theta \in [0, 1]$  be a constant, and let  $\Phi^\theta(p, q)$  denote the set of all joint distributions  $\phi^\theta$  with marginals  
 105  $\phi_1^\theta$  and  $\phi_2^\theta$  such that  $\Delta(\phi_1^\theta, p) + \Delta(\phi_2^\theta, q) \leq \theta$  hold. Note that at  $\theta = 0$ , all the joint distributions  
 106 in  $\Phi^0(p, q)$  have marginals exactly equal to  $p$  and  $q$ . In this case we will write  $\Phi(p, q)$  to denote  
 107  $\Phi^0(p, q)$ .

108 **Definition 2** (Wasserstein Distance). *Let  $p$  and  $q$  be two discrete probability distributions on  $\mathbb{R}$ . The*  
 109 *Wasserstein distance between  $p$  and  $q$  is defined as:*

$$W(p, q) = \inf_{\phi \in \Phi(p, q)} \mathbb{E}_{(x, y) \leftarrow \phi} [|x - y|]. \quad (1)$$

110 **Definition 3** ( $\theta$ -Wasserstein Distance). *Let  $p$  and  $q$  be two discrete probability distributions on  $\mathbb{R}$ .*  
 111 *Let  $\theta \in [0, 1]$ . The  $\theta$ -Wasserstein distance between  $p$  and  $q$  is defined as:*

$$W^\theta(p, q) = \inf_{\phi \in \Phi^\theta(p, q)} \mathbb{E}_{(x, y) \leftarrow \phi} [|x - y|]. \quad (2)$$

112 **Definition 4** ( $\infty$ -Wasserstein Distance). *Let  $p$  and  $q$  be two discrete probability distributions on  $\mathbb{R}$ .*  
 113 *The  $\infty$ -Wasserstein distance between  $p$  and  $q$  is defined as:*

$$W_\infty(p, q) = \inf_{\phi \in \Phi(p, q)} \max_{(x, y) \leftarrow \phi} |x - y|. \quad (3)$$

114 **Definition 5** ( $(\infty, \theta)$ -Wasserstein Distance). *Let  $p$  and  $q$  be two discrete probability distributions on*  
 115  *$\mathbb{R}$ . Let  $\theta \in [0, 1]$ . The  $(\infty, \theta)$ -Wasserstein distance between  $p$  and  $q$  is defined as:*

$$W_\infty^\theta(p, q) = \inf_{\phi \in \Phi^\theta(p, q)} \max_{(x, y) \leftarrow \phi} |x - y|. \quad (4)$$

116 **Claim 1.** *Let  $X$  and  $Y$  be a random variable with distributions  $\mu_1, \mu_2$ , respectively. Let  $Z$  be a an*  
 117 *independent noise random variable with a distribution  $\mu_3$ . Then we have,*

$$W^\gamma(X, Y) = W^\gamma(X + Z, Y + Z),$$

$$W_\infty^\gamma(X, Y) = W_\infty^\gamma(X + Z, Y + Z),$$

118  
 119 *i.e., convolution of distributions with the same independent noise does not change Wasserstein*  
 120 *distance.*

121 *Proof.* First, we show that applying convolution of distributions with another does not increase the  
 122 Wasserstein distance between them.

123 Let  $p, q, r$  be three distributions, with wasserstein distance between  $p, q$  as  $\beta$  and  $\pi$  be the optimal  
 124 transfer which achieves this. Let  $I_1, I_2$  be two distributions obtained by convolving  $p, q$  with  $r$   
 125 respectively. i.e.,

$$Pr(I_1 = x) = \int_Z p(x - z)r(z)$$

$$Pr(I_2 = y) = \int_Z q(y - z)r(z)$$

127 From the above equations, we can construct a transfer from  $I_1, I_2$  as follows. Let  $z$  be drawn from  
 128  $r$  with a probability  $p(z)$ . In this event, we transfer  $I_1$  to  $I_2$  using the policy,  $\pi + z$ . This policy  
 129 also transfers wasserstein distance of  $\beta$  for all  $z$ . Hence, it's expectation over  $z$  is also  $\beta$ . Since,  
 130 wasserstein distance is an infimum of all transfers, we have, the wasserstein distance between  $I_1, I_2$   
 131 to be at most  $\beta$ . Hence, proved.

132 Now, for the random variables  $X, Y, Z$ . The distributions of  $X + Z, Y + Z$  are given by the following  
 133 convolutions.

$$Pr(A + C = x) = \int_Z \mu_1(x - z)\mu_3(z)$$

$$Pr(B + C = y) = \int_Z \mu_2(y - z)\mu_3(z)$$

135 From the above reasoning, we conclude that the wasserstein distance between  $X + Z, Y + Z$  is atmost  
 136 that of  $X, Y$ .

137 Note that the distributions of  $X$  and  $Y$  can be obtained by inverse convolution of  $Z$  with that of  $X + Z,$   
 138  $Y + Z$ . This is equivalent to convolution with an appropriate distribution  $Z'$ . **This is true for discrete**  
 139 **variables, convolution is matrix multiplication with a kernel, so inverse convolution is just multiplying**  
 140 **with inverse of the kernel. How does it go for continuous variables ?**. Again, as convolving with a  
 141 distribution doesn't increase the wasserstein distance, we conclude the wasserstein distance between  
 142  $X + Z, Y + Z$  is atmost that of  $X, Y$ .

143 Combining both these results, we have that, wasserstein distance between  $X, Y$  is same as wasserstein  
 144 distance between  $X + Z, Y + Z$ .  $\square$

145 **Claim 2.** We have the following triangle inequality for the  $\theta$ -Wasserstein distance, where  $p, q, r$  are  
 146 any three distributions defined over the same support, and  $\gamma_1, \gamma_2 \geq 0$ .

$$W_{\infty}^{\gamma_1 + \gamma_2}(p, r) \leq W_{\infty}^{\gamma_1}(p, q) + W_{\infty}^{\gamma_2}(q, r).$$

$$W^{\gamma_1 + \gamma_2}(p, r) \leq W^{\gamma_1}(p, q) + W^{\gamma_2}(q, r).$$

148 *Proof.* We will prove the result for  $W_{\infty}$ . The result for  $W$  is exactly the same.

149 We first prove the following result. Let  $p, q$  be two given distribution. Let  $p'$  be a distribution which  
 150 is  $\delta$  statistical distance apart from  $p$ . Then there exists another distribution  $q'$ , which is atmost  $\delta$   
 151 statistical distance apart from  $q$ , such that

$$W_{\infty}(p', q) = W_{\infty}(p, q')$$

152 Intuitively, this can be proved as follow. Let  $\phi$  be any joint distribution of  $p'$  and  $q$ . We will  
 153 convert this joint distribution into another joint distribution  $\phi'$  whose marginals are  $p$  and  $q'$  such  
 154 that  $\max_{(x,y) \leftarrow \phi} |x - y| = \max_{(x,y) \leftarrow \phi'} |x - y|$ .  $p'$  has extra mass at some places and less mass  
 155 at other places as compared to  $p$ . In  $\phi$ , this extra mass (which is equal to  $\delta$ ) would be joined with  
 156 some  $\delta$  mass of  $q$ . Now we restore  $p$  by moving this  $\delta$  extra part so as to convert  $p'$  into  $p$  but we  
 157 also move the corresponding  $\delta$  mass of  $q$  parallelly (keeping the distance same). The resulting joint  
 158 distribution is the required  $\phi'$ . Its easy to see that the total mass of  $q$  moved is  $\delta$  so the  $q'$  we get  
 159 has atmost  $\delta$  statistical difference with  $q$ . And since we moved the masses of  $p'$  and  $q$  parallelly, no  
 160 distance changed. This implies  $\max_{(x,y) \leftarrow \phi} |x - y| = \max_{(x,y) \leftarrow \phi'} |x - y|$ .

161 This allows us to say that  $W_{\infty}(p', q) \geq W_{\infty}(p, q')$ . We can use argument to prove that  $W_{\infty}(p', q) \leq$   
 162  $W_{\infty}(p, q')$ . Hence we proved that  $W_{\infty}(p', q) \geq W_{\infty}(p, q')$ .

163 Let us use  $W_{\infty}^{\gamma_a, \gamma_b}(p, q)$  to represent  $\inf_{\substack{p', q': \\ \Delta(p, p') \leq \gamma_a, \\ \Delta(q, q') \leq \gamma_b}} W_{\infty}(p', q')$ . The above result implies that,  $W_{\infty}^{\gamma}(p, q)$   
 164 is equal to  $W_{\infty}^{\gamma_a, \gamma_b}(p, q) \forall \gamma_a, \gamma_b \text{ such that } \gamma_a + \gamma_b = \gamma$ .

165 Now using our above proof, we can say that

$$\begin{aligned} W_{\infty}^{\gamma_1 + \gamma_2}(p, r) &= W_{\infty}^{\gamma_1, \gamma_2}(p, r) \\ &\leq W_{\infty}^{\gamma_1, 0}(p, q) + W_{\infty}^{0, \gamma_2}(q, r) \quad (\text{since wasserstein distance is a metric}) \\ &= W_{\infty}^{\gamma_1}(p, q) + W_{\infty}^{\gamma_2}(q, r) \end{aligned}$$

166 □

167 Let  $\mathcal{X}$  be a finite set, and let  $\mathbf{x} := (x_1, x_2, \dots, x_n) \in \mathcal{X}^n$  represent a database with  $n$  elements,  
 168 where  $x_i \in \mathcal{X}$  is the  $i$ -th element. For any two databases  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}^n$ , we say that  $\mathbf{x} \sim \mathbf{x}'$ , if  $\mathbf{x}'$  is  
 169 obtained from  $\mathbf{x}$  by adding/removing a single element. Let  $f : \mathcal{X}^n \rightarrow \mathcal{Y}$  be a randomized function,  
 170 where  $\mathcal{Y} = \{y_1, y_2, \dots, y_k\}$  is a discrete set for some finite  $k \in \mathbb{N}$ . For every  $\mathbf{x} \in \mathcal{X}^n$ , we have that  
 171  $\sum_{i=1}^k \Pr[f(\mathbf{x}) = y_i] = 1$ . **our results are also for continuous functions.**

172 **Definition 6** (Parameterized Sensitivity of a Randomized Query). For  $\theta \in [0, 1]$ , we define  $\theta$ -  
 173 sensitivity of a randomized query  $f$ , denoted by  $S^{\theta}(f)$ , as follows:

$$S^{\theta}(f) := \max_{\substack{\mathbf{x}, \mathbf{x}' \in \mathcal{X}^n: \\ \mathbf{x} \sim \mathbf{x}'}} W_{\infty}^{\theta}(f(\mathbf{x}), f(\mathbf{x}')). \quad (5)$$

174 **Definition 7** (Laplace Distribution). Let  $b$  be a positive real number. Laplace distribution with respect  
 175 to  $b$  ( $b$  is called the scaling parameter) and mean  $\mu$ , denoted by  $\text{Lap}(x|\mu, b)$ , is defined as:

$$\text{Lap}(x|\mu, b) := \frac{1}{2b} e^{-\frac{|x-\mu|}{b}}, \quad x \in \mathbb{R}.$$

176 We denote a random variable that is distributed as the Laplace distribution with the scaling parameter  
 177  $b$  and mean  $\mu$  by  $\text{Lap}(b, \mu)$ . If mean  $\mu$  is zero, then we will simply denote it by  $\text{Lap}(b)$ .

178 **Definition 8** ( $(\alpha, \beta, \gamma)$ -accuracy). Let  $d : \mathcal{X}^n \times \mathcal{X}^n \rightarrow [0, \infty)$  be a distortion function, not necessarily  
 179 a distance metric. A mechanism  $\mathcal{M}$  for computing a given randomized function  $f : \mathcal{X}^n \rightarrow \mathcal{Y}$  is said  
 180 to be  $(\alpha, \beta, \gamma)$ -accurate with respect to  $d$ , if for every  $\mathbf{x} \in \mathcal{X}^n$ , there exists a random variable  $X'$ ,  
 181 which is a distribution on all  $\mathbf{x}'$  such that  $d(\mathbf{x}, \mathbf{x}') \leq \alpha$  and that  $W^{\gamma}(f(X'), \mathcal{M}(\mathbf{x})) \leq \beta$ .

182 **Definition 9** (Robustness). Let  $\rho, \theta, \epsilon, \delta$  be non-negative real numbers with  $\delta, \theta \leq 1$ . Let  $\mathcal{RM} : A \rightarrow B$   
 183 be a randomized mechanism and  $p_x$  denote the output distribution of  $\mathcal{RM}$  when the input  
 184 is taken from a distribution  $x$  over  $A$ . The mechanism  $\mathcal{RM}$  is called  $(\theta, \rho, \epsilon, \delta)$ -robust if  $\forall x, x', S$   
 185 where  $S \subseteq \mathbb{R}$  and  $x$  and  $x'$  are two distributions over  $A$  such that  $W_{\infty}^{\theta}(p_x, p_{x'}) \leq \rho$ ,

$$\Pr_{y \leftarrow p_x} [y \in S] \leq e^{\epsilon} \Pr_{y \leftarrow p_{x'}} [y \in S] + \delta \quad (6)$$

186 **Definition 10** (distortion sensitivity). Let  $M : A \rightarrow B$  be a mechanism for computing a function  $f$   
 187 and  $d_1, d_2$  are distortion functions on  $A, B$ , respectively. We define distortion-sensitivity function  $\sigma_f$   
 188 for the randomized function  $f$  w.r.t.  $d_1, d_2$  as  $\sigma_f(\alpha)$  to be the minimum number such that for every  
 189  $x \in A$  and  $y \in B$ , s.t.  $d_2(f(x), y) \leq \alpha$ , there exists an  $x' \in A$  such that  $d_1(x, x') \leq \sigma_f(\alpha)$  and  
 190  $f(x') = y$ .

191 Note that this may not be well defined for every randomized function.

192 **Definition 11** ( $\theta$ -error sensitivity). Let  $M : A \rightarrow B$  be a mechanism for computing a function  
 193  $f$  and  $d_1, d_2$  are distance functions on  $A, B$ , respectively. Additionally, suppose that the input to  
 194 the mechanism  $M$  is drawn from a distribution. We define  $\theta$ -error-sensitivity function  $\tau$  for the  
 195 mechanism  $M$  w.r.t.  $d_1, d_2$  as

$$\tau_M^{\theta}(\beta) = \max_{\substack{\mathbf{x} \sim p, \mathbf{x}' \sim q: \\ W^{\theta}(p, q) \leq \beta}} W^{\theta}(M(\mathbf{x}), M(\mathbf{x}')) \quad (7)$$

### 3 Adding Robustness using Robust Mechanisms for Identity Functions

#### 3.1 Robust Mechanism for Identity Function over $\mathbb{R}$

We will now show a mechanism  $\mathcal{M}_{\text{rob}} : \mathbb{R} \rightarrow \mathbb{R}$  which is a robust mechanism for identity function, i.e., the output corresponding to any input is the same as input. Now Suppose the input distribution is  $x$ . Let  $\rho$  be any non-negative real number and  $\epsilon$  be any positive real number. Consider the following mechanism for the identity function:

**Mechanism 1.**  $\mathcal{M}_{\text{rob}}$ : On input  $y$ , sample  $z$  according to the probability distributions  $\text{Lap}(\frac{\rho}{\epsilon})$  and output  $y + z$ .

**Lemma 1.** For a fixed constant  $\theta \in [0, 1]$ ,  $\mathcal{M}_{\text{rob}}$  achieves  $(\theta, \rho, \epsilon, \delta)$ -robustness where  $\delta = 2e^\epsilon \theta$ .

*Proof.* Fix a constant  $\theta \geq 0$ . Suppose  $q$  and  $q'$  denote two input distributions, and let  $p$  and  $p'$  denote the corresponding output probability distributions of  $\mathcal{M}_{\text{rob}}$ . Suppose  $W_\infty^\theta(p, p') \leq \rho$ . Then by **Claim 1**,  $W_\infty^\theta(q, q') \leq \rho$ . Let  $b := \frac{\rho}{\epsilon}$  and  $Z \sim \text{Lap}(b)$ . We need to show that  $\forall S \subseteq \mathbb{R}$ ,  $\Pr_{y \leftarrow p}[y \in S] \leq e^\epsilon \Pr_{y \leftarrow p'}[y \in S] + \delta$ .

$$\Pr_{y \leftarrow p}[y \in S] = \int_{t \in S} p(t) \cdot dt = \int_{t \in S} \left[ \int_{-\infty}^{\infty} q(s) \cdot p_Z(t - s) \cdot ds \right] \cdot dt \quad (8)$$

Let  $\phi^\theta \in \Phi^\theta(q, q')$  be a joint distribution. Let  $M^\phi$  denote  $\max_{(x, y) \leftarrow \phi^\theta} |x - y|$  for the distribution  $\phi^\theta$ . For  $i \in \mathbb{R}$ , let  $\theta_1(i) := |\phi_1^\theta(i) - q(i)|$  and  $\theta_2(i) := |\phi_2^\theta(i) - q'(i)|$ . Note that, by definition of the distribution  $\phi^\theta$ , we have  $\Delta(\phi_1^\theta, q) + \Delta(\phi_2^\theta, q') \leq \theta$  which implies  $\int_{-\infty}^{\infty} \theta_1(i) \cdot di + \int_{-\infty}^{\infty} \theta_2(i) \cdot di \leq \theta$ . Using these in (8) we get the following:

$$\begin{aligned} \Pr_{y \leftarrow p}[y \in S] &\leq \int_{t \in S} \left[ \frac{1}{2b} \int_{-\infty}^{\infty} (\phi_1^\theta(i) + \theta_1(i)) \cdot e^{\frac{-|t-i|}{b}} \cdot di \right] \cdot dt \\ &= \int_{t \in S} \left[ \frac{1}{2b} \int_{-\infty}^{\infty} \phi_1^\theta(i) \cdot e^{\frac{-|t-i|}{b}} \cdot di \right] \cdot dt + \int_{t \in S} \left[ \frac{1}{2b} \int_{-\infty}^{\infty} \theta_1(i) \cdot e^{\frac{-|t-i|}{b}} \cdot di \right] \cdot dt \\ &= \int_{t \in S} \left[ \frac{1}{2b} \int_{-\infty}^{\infty} \phi_1^\theta(i) \cdot e^{\frac{-|t-i|}{b}} \cdot di \right] \cdot dt + \underbrace{\int_{-\infty}^{\infty} \theta_1(i) \cdot \int_{t \in S} \frac{1}{2b} e^{\frac{-|t-i|}{b}} \cdot dt \cdot di}_{\leq 1} \quad (9) \end{aligned}$$

By properties of joint distributions, we have  $\phi_1^\theta(i) = \int_{-\infty}^{\infty} \phi^\theta(i, j) \cdot dj$  and  $\phi_2^\theta(j) = \int_{-\infty}^{\infty} \phi^\theta(i, j) \cdot di$ ; and by triangle inequality we have  $|t - i| \geq |t - j| - |i - j|$ . Let  $\int_{-\infty}^{\infty} \theta_l(i) \cdot di$  be  $\omega_l$ ,  $l \in \{1, 2\}$ . Substituting all these in (9) we get the following:

$$\Pr_{y \leftarrow p}[y \in S] \leq \int_{t \in S} \left[ \frac{1}{2b} \int_{-\infty}^{\infty} \left[ \int_{-\infty}^{\infty} \phi^\theta(i, j) \cdot e^{\frac{-|t-j|+|i-j|}{b}} \cdot dj \right] \cdot di \right] \cdot dt + \omega_1. \quad (10)$$

Observe that  $\phi^\theta(i, j)$  is non-zero only if  $|i - j| \leq M^\phi$  (by definition of  $M^\phi$ ). Using this for every  $i, j \in \mathbb{R}$ , the integrand in (10) can be upper-bounded as follows:

$$\phi^\theta(i, j) \cdot e^{\frac{-|t-j|}{b}} e^{\frac{|i-j|}{b}} \leq \phi^\theta(i, j) \cdot e^{\frac{-|t-j|}{b}} e^{\frac{M^\phi}{b}} = \phi^\theta(i, j) \cdot e^{\frac{-|t-j|}{b}} e^{\frac{M^\phi}{b}}.$$

Substituting this in (10) and with some algebraic manipulations done in **Appendix B**, we get  $\Pr_{y \leftarrow p}[y \in S] \leq e^\epsilon \Pr_{y \leftarrow p'}[y \in S] + 2\theta e^\epsilon$ . This completes the proof of **Lemma 1**.  $\square$

**Lemma 2.** For every constant  $\gamma' \geq 0$ ,  $\mathcal{M}_{\text{rob}}$  is  $(0, \beta', \gamma')$ -accurate, where  $\beta' = \frac{\rho}{\epsilon(1-\gamma')}$   $\left(1 - \gamma' \left[1 + \ln\left(\frac{1}{\gamma'}\right)\right]\right)$ . Note that if  $\gamma' = 0$ ,  $\beta' = \frac{\rho}{\epsilon}$ .

*Proof.* Fix a constant  $\gamma \geq 0$  and any input  $x \in \mathbb{R}$ . Instead of treating  $x$  as a real, for this proof, we will treat  $x$  as a point distribution over the input space from which we are sampling the inputs. Clearly, this is equivalent to treating  $x$  as a deterministic input. Let  $q$  denote the output distribution of



226  $\mathcal{M}_{\text{rob}}$  when the input is drawn from  $\mathbf{x}$ . Let  $b$  denote  $\frac{\rho}{\epsilon}$ . We want to show that  $W^\gamma(x, q) \leq \beta$ , for the  
 227 above-mentioned  $\beta$ . By definition of  $W^\gamma$  from (2) we have  $W^\gamma(x, q) = \inf_{\phi \in \Phi^\gamma(x, q)} \mathbb{E}_{(y, t) \leftarrow \phi} [|y - t|]$ .

228 Consider the following  $\phi^*$ :

$$\phi^*(i, t) = \begin{cases} 0 & \text{if } t < -b \ln(\frac{1}{\gamma}) + i \text{ or } t > b \ln(\frac{1}{\gamma}) + i; \\ \frac{1}{1-\gamma} \text{Lap}(t|b, i)x(i) & \text{if } t \in [-b \ln(\frac{1}{\gamma}) + i, b \ln(\frac{1}{\gamma}) + i]. \end{cases} \quad (11)$$

229 It can be verified that  $\Delta(\phi_1^*, x) = 0$  and  $\Delta(\phi_2^*, q) \leq \gamma$ , which implies that  $\phi^* \in \Phi^\gamma(x, q)$ . This in  
 230 turn implies that  $W^\gamma(x, Q) \leq \mathbb{E}_{(y, t) \leftarrow \phi^*} [|y - t|]$ . We show in **Appendix B** that  $\mathbb{E}_{(y, t) \leftarrow \phi^*} [|y - t|] \leq$   
 231  $\frac{\rho}{\epsilon(1-\gamma)} \left(1 - \gamma[1 + \ln(\frac{1}{\gamma})]\right)$ . This completes the proof of **Lemma 2**.  $\square$

232 **Theorem 1.** For any  $\theta \in [0, 1]$  and  $\gamma \geq 0$ ,  $\mathcal{M}_{\text{rob}}$  is  $(\theta, \rho, \epsilon, \delta)$ -robust and  $(0, \beta, \gamma)$ -accurate, where  
 233  $\delta = 2e^\epsilon \theta$  and  $\beta = \frac{\rho}{\epsilon(1-\gamma)} \left(1 - \gamma[1 + \ln(\frac{1}{\gamma})]\right)$ .

234 *Proof.* Using **Lemma 1** and **Lemma 2**, the theorem trivially holds.  $\square$

### 235 3.2 Adding Robustness to a Mechanism

236 Consider a randomized mechanism  $\mathcal{M} : A \rightarrow B$  which has some privacy and accuracy bounds and a  
 237 randomized mechanism  $\mathcal{RM} : B \rightarrow B$  which is robust and has some accuracy bounds. We want to  
 238 construct a new mechanism which is robust, at least as private as  $\mathcal{M}$ , and doesn't compromise much  
 239 on accuracy as compared to  $\mathcal{M}$ . Let  $x \in A$  be the input. Consider the following mechanism:

241 **Mechanism 2.** Run  $\mathcal{M}$  on  $x$  to get an output  $y \in B$  then run  $\mathcal{M}_{\text{rob}}$  on  $y$  to get  $z \in B$ , output  $z$ .

242 **Theorem 2.** Let  $\rho, \theta, \epsilon, \delta$  be non-negative real numbers with  $\delta, \theta \leq 1$ . Then **Mechanism 2** achieves (i)  
 243  $(\theta, \rho, \epsilon, \delta)$ -robustness, where  $\delta = 2e^\epsilon \theta$ , (ii)  $(\epsilon_p, \delta_p)$ -differential privacy, if  $\mathcal{M}$  is  $(\epsilon_p, \delta_p)$ -differentially  
 244 private, and (iii)  $(\alpha_p, \beta_p, \gamma_p)$ -accuracy, if  $\mathcal{M}$  is  $(\alpha, \beta, \gamma)$ -accurate, where  $\alpha_p = \alpha$ ,  $\beta_p = \beta + \beta'$ ,  $\gamma_p =$   
 245  $\gamma + \gamma'$  and  $\beta', \gamma'$  are such that  $\gamma' \geq 0$  is arbitrary and  $\beta' = \frac{\rho}{\epsilon(1-\gamma')} \left(1 - \gamma'[1 + \ln(\frac{1}{\gamma'})]\right)$ .

## 246 4 Differential Privacy for Randomized Queries

247 Suppose the database is  $\mathbf{x} \in \mathcal{X}^n$  and  $\theta \in [0, 1], \epsilon > 0$  be fixed constants. Consider the following  
 248 randomized mechanism for a randomized query  $f : \mathcal{X}^n \rightarrow \mathbb{R}$ :

250 **Mechanism 3.** Sample  $y$  and  $z$  independently and according to the probability distributions  $f(\mathbf{x})$   
 251 and  $\text{Lap}(\frac{S^\theta(f)}{\epsilon})$ , respectively, and output  $y + z$ .

252 **Theorem 3.** For the fixed constants  $\theta' \in [0, 1], \epsilon' > 0, \gamma \geq 0$ , **Mechanism 3** achieves  
 253  $(\theta', S^\theta(f), \epsilon', \delta')$ -robustness,  $(\epsilon, \delta)$ -differential privacy and  $(0, \beta, \gamma)$ -accuracy, where  $\delta' = 2e^{\epsilon'} \theta'$ ,  
 254  $\delta = 2e^\epsilon \theta$  and  $\beta = \frac{S^\theta(f)}{\epsilon(1-\gamma)} \left(1 - \gamma[1 + \ln(\frac{1}{\gamma})]\right)$ .

255 *Proof.* Observe that  $f$  can be treated as a mechanism for  $f$  with  $(0, 0, 0)$ -accuracy and  $(\infty, 0)$ -  
 256 privacy and that **Mechanism 3** is equivalent to  $\mathcal{M}_{\text{rand}} := \mathcal{M}_{\text{rob}} \circ f$ , where in  $\mathcal{M}_{\text{rob}}$ , the Laplace  
 257 noise has the parameter  $\rho = S^\theta(f)$ . The robustness and accuracy bounds are obtained from  
 258 **Theorem 2**. To show that  $\mathcal{M}_{\text{rand}}$  is  $(\epsilon, \delta)$ -differentially private, consider any two neighbouring  
 259 databases  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}^n$ . Since  $\mathbf{x} \sim \mathbf{x}'$ , it follows from **Definition 6** that  $W_\infty^\theta(f(\mathbf{x}), f(\mathbf{x}')) \leq$   
 260  $S^\theta(f)$ . This, together with the fact that **Mechanism 3** is  $(\theta', S^\theta(f), \epsilon', \delta')$ -robust, implies that  
 261  $\Pr[\mathcal{M}_{\text{rand}}(\mathbf{x}) \in S] \leq e^\epsilon \Pr[\mathcal{M}_{\text{rand}}(\mathbf{x}') \in S] + \delta$ , where  $\delta = (e^\epsilon + 1)2\theta$ . But the definition of  
 262 robustness requires that  $W_\infty^\theta(\mathcal{M}_{\text{rand}}(\mathbf{x}), \mathcal{M}_{\text{rand}}(\mathbf{x}')) \leq S^\theta(f)$ , not  $W_\infty^\theta(f(\mathbf{x}), f(\mathbf{x}')) \leq S^\theta(f)$ .  
 263 How do we fix this? Note that Claim 1 is for  $W^\theta$ , not  $W_\infty^\theta$ !  $\square$

## 5 Histogram Mechanism

In this section, we propose a new mechanism to get a histogram from a given (correct) histogram in a differentially private manner.

**Mechanism 4.** For each bar of the histogram, independently add a noise according to the following probability distribution:

$$\pi_q(x) = \begin{cases} \frac{1}{1-e^{-\frac{\sqrt{q}}{2}}} \text{Lap}(x \mid -\frac{q}{2}, \sqrt{q}) & -q < x < 0 \\ 0 & \text{elsewhere} \end{cases}$$

and then round each of the bar to the nearest integer. If any of the bar becomes negative, then set it to 0. Output the resulting histogram.

To prove that the above mechanism is  $(\epsilon, \delta)$ -differentially private, first we show below in **Lemma 5** that, if in **Mechanism 4**, we output the histogram before rounding and setting negative values to 0, then we get  $(\epsilon, \delta)$ -differential privacy.

**Lemma 3.** For any  $q \geq 1$ , if in **Mechanism 4** we output the histogram before rounding and setting negative values to 0, then we get  $(\epsilon, \delta)$ -differential privacy, where  $\epsilon = \frac{1}{\sqrt{q}}$  and  $\delta = \frac{4}{\sqrt{q}} e^{-\frac{\sqrt{q}}{2}}$ .

**Lemma 5** is proved in **Appendix C**. Observe that rounding and setting negative values to 0 is a post-processing step. Since **Lemma 5** is  $(\epsilon, \delta)$ -differentially private and post-processing preserves differential privacy [1, Proposition 2.1], it follows that **Mechanism 4** is also  $(\epsilon, \delta)$  private. The following theorem establishes  $(qt, 0, 0)$ -accuracy of **Mechanism 4**.

**Theorem 4.** **Mechanism 4** achieves  $(\alpha, 0, 0)$  accuracy with  $\alpha = qt$ .

## 6 Using the Histogram Mechanism for General Statistical Queries

We now show how histogram mechanism can be used for answering a general query. For a query which desires some statistic of the database, we construct an appropriate function,  $f'$  which takes the histogram given by the above mechanism and outputs the desired statistic within some error. The input database is first bucketed, i.e. elements are rounded off to a certain number of elements and modified according to the above mentioned histogram mechanism. The modified histogram given to the above function. The output thus obtained is guaranteed to be accurate and private. To find the accuracy and private parameters we prove the following theorems.

**Theorem 5.** Let  $M_1 : A \rightarrow B$  be a mechanism which is  $(\alpha_1, \beta_1, \gamma_1)$  accurate for a function  $f_1 : A \rightarrow B$  and let  $M_2 : B \rightarrow C$  be a mechanism which is  $(\alpha_2, \beta_2, \gamma_2)$  accurate for a function  $f_2 : B \rightarrow C$ . Let  $d_1, d_2$  be distortion functions on  $A, B$ , respectively. Then the composite mechanism  $M_2 \circ M_1 : A \rightarrow C$  is  $(\alpha, \beta, \gamma)$ -accurate for the function  $f_2 \circ f_1$ , where  $\alpha = \alpha_1 + \sigma_{f_1}(\alpha_2)$ ,  $\beta = \beta_2 + \tau_{M_2}^{\gamma_1}(\beta_1)$ ,  $\gamma = \gamma_1 + \gamma_2$ .

*Proof.* We prove the theorem using a hybrid argument. Consider the hybrid mechanism  $M_2 \circ f_1$ , along with  $f_2 \circ f_1$  and  $M_2 \circ M_1$ . Below, we compute the Wasserstein distances between  $M_2 \circ M_1$  &  $M_2 \circ f_1$  and between  $M_2 \circ f_1$  &  $f_2 \circ f_1$ . Then we use the triangle inequality from **Claim 2** to find the Wasserstein distance between  $M_2 \circ M_1$  and  $f_2 \circ f_1$ .

For a given database  $\mathbf{x}$ , since  $M_1$  is  $(\alpha_1, \beta_1, \gamma_1)$ -accurate mechanism for  $f_1$ , there exists a random variable  $X'$  which is a distribution over the  $\alpha_1$ -distorted databases from  $\mathbf{x}$ , such that

$$W^{\gamma_1}(f_1(X'), M_1(\mathbf{x})) \leq \beta_1. \quad (12)$$

Now, applying the mechanism  $M_2$  over the distributions  $M_1(\mathbf{x})$ ,  $f_1(X')$  may increase the error by at most  $\tau_{M_2}^{\gamma_1}(\beta_1)$  (see **Definition 11**), which gives

$$W^{\gamma_1}(M_2(f_1(X')), M_2(M_1(\mathbf{x}))) \leq \tau_{M_2}^{\gamma_1}(\beta_1). \quad (13)$$

Now, we bound the Wasserstein distance between  $M_2 \circ f_1$  and  $f_2 \circ f_1$ . For any database  $\mathbf{x}_1$  drawn from  $X'$ , since  $M_2$  is an  $(\alpha_2, \beta_2, \gamma_2)$ -accurate mechanism for  $f_2$ , there exists a random variable  $X''$ , which



is a distribution over  $\alpha_2$ -distorted databases from  $f_1(\mathbf{x}_1)$ , such that  $W^{\gamma_2}(f_2(X''), M_2(f_1(\mathbf{x}_1))) \leq \beta_2$ . Since this holds for every  $\mathbf{x}_1$  in the support of  $X'$  (note that  $X''$  may depend on  $\mathbf{x}_1$ ), we have

$$W^{\gamma_2}(f_2(X''), M_2(f_1(X'))) \leq \beta_2. \quad (14)$$

Combining (14) and (14) and using Claim 2, we get

$$W^{\gamma_1+\gamma_2}(f_2(X''), M_2(M_1(\mathbf{x}))) \leq \beta_2 + \tau_{M_2}^{\gamma_1}(\beta_1). \quad (15)$$

Note that  $X''$  is a distribution over  $B$  and we need to find a distribution over the  $\alpha = (\alpha_1 + \sigma_{f_1}(\alpha_2))$ -distorted databases from the database  $\mathbf{x}$ . For this  $\mathbf{x} \in A$ , consider the corresponding distribution of  $X'$  guaranteed by the mechanism  $M_1$ , which is defined over all the elements  $\mathbf{x}_1 \in A$  such that  $d_1(\mathbf{x}, \mathbf{x}_1) \leq \alpha_1$ . For every  $\mathbf{x}_1$  drawn from  $X'$ , consider the distribution of  $X''$  guaranteed by  $M_2$ , which is defined over all the elements in  $\mathbf{x}'_1 \in B$  such that  $d_2(f_1(\mathbf{x}_1), \mathbf{x}'_1) \leq \alpha_2$ .

Since  $d_2(f_1(\mathbf{x}_1), \mathbf{x}'_1) \leq \alpha_2$ , it follows from Definition 10 that there exists a database  $\mathbf{x}_2 \in A$ , such that  $f_1(\mathbf{x}_2) = \mathbf{x}'_1$  and that  $d_1(\mathbf{x}_1, \mathbf{x}_2) \leq \sigma_{f_1}(\alpha_2)$ . Let  $X_2$  denote the distribution over all such  $\mathbf{x}_2$ 's in  $A$ . It follows from the triangular inequality of  $d_1$  that  $d_1(\mathbf{x}, \mathbf{x}_2) \leq \alpha_1 + \sigma_{f_1}(\alpha_2)$ . Thus,  $X_2$  is a distribution over databases in  $A$ , which are at most  $(\alpha_1 + \sigma_{\alpha, f_1}(\alpha_2))$ -distorted from  $\mathbf{x}$  and satisfy

$$W^{\gamma_1+\gamma_2}(f_2(f_1(X_2)), M_2(M_1(\mathbf{x}))) \leq \beta_2 + \tau_{M_2}^{\gamma_1}(\beta_1). \quad (16)$$

This implies that the composite mechanism  $M_2 \circ M_1$  is  $(\alpha, \beta, \gamma)$ -accurate, where  $\alpha = \alpha_1 + \sigma_{f_1}(\alpha_2)$ ,  $\beta = \beta_1 + \tau_{M_2}^{\gamma_1}(\beta_1)$ ,  $\gamma = \gamma_1 + \gamma_2$ . This completes the proof of Theorem 5.  $\square$

Now, we present another composition theorem, which we prove in Appendix D.

**Theorem 6.** Let  $M_1 : A \rightarrow B$  be a neighbourhood preserving mechanism for a function  $f_1 : A \rightarrow B$ , and let  $M_2 : B \rightarrow C$  be a mechanism which is  $(\epsilon, \delta)$ -differential private for a function  $f_2 : B \rightarrow C$ . Then the composite mechanism,  $M_2 \circ M_1 : A \rightarrow C$  is  $(\epsilon, \delta)$ -differential private.

Now, we show how appropriate mechanisms can be composed to achieve good accuracy and privacy guarantees for computing a statistic on a database. We consider databases which consists of positive real number numbers. The bucketing is done with  $t$  number of buckets, i.e, rounding each element of database to the nearest multiple of  $B/t$ .

**Lemma 4.** The bucketing followed by Mechanism 4 is an  $(\alpha_p, \beta_p, 0)$ -accurate and  $(\epsilon_p, \delta_p)$ -differentially private mechanism for the identity function, with  $\alpha_p = qt$ ,  $\beta_p = \frac{B}{2t}$ ,  $\epsilon_p = \frac{1}{\sqrt{q}}$ ,  $\delta_p = \frac{4}{\sqrt{q}} e^{-\frac{\sqrt{q}}{2}}$ .

## 6.1 Case Study

We now show the application of the histogram for computing Max and support of a database in an accurate and private manner.

### 6.1.1 Maximum

Let  $d$  be a distance metric on  $\mathcal{X}$ . We define the function  $f : \mathcal{N}^{\mathcal{X}} \rightarrow \mathcal{X}$ , as one which takes a database as input and outputs the largest element. Let's denote the distance metric over input space by  $d_1$  and on output space by  $d_2$ . For any two databases  $x_1$  and  $x_2$ ,  $d_1(x_1, x_2)$  = the maximum  $d$ -distance each element of  $x_1$  is to be modified to get  $x_2$  or vice-versa (both are equivalent) and  $d_2(x_1, x_2)$  =  $d$ -distance between  $x_1$  and  $x_2$  in  $\mathcal{X}$  space. Note that  $f$  is a perfectly accurate mechanism for maximum. Also, if for any two databases  $x_1$  and  $x_2$ ,  $d_1(x_1, x_2) = l$  then  $d_2(f(x_1), f(x_2)) = l$  because we don't need to modify the max of  $x_1$  by more than  $l$  to get max of  $x_2$ . Using this we get that the error-sensitivity function for  $f$  is identity, i.e,  $\tau_f(\beta') = \beta'$ . Because if two distribution over databases have  $d_1$ -wasserstein distance  $\beta$ , then the output distribution will have  $d_2$ -wasserstein  $\beta$ .

Now, using Theorem 5 and Theorem 6, we have that, the histogram mechanism for finding max is  $(\alpha_p, \beta_p, 0)$  accurate and  $(\epsilon_p, \delta_p)$  private where  $\alpha_p, \beta_p, \epsilon_p$  and  $\delta_p$  are defined in Lemma 4.

### 6.1.2 Support

Here we define the function  $f : \mathcal{N}^{\mathcal{X}} \rightarrow \mathcal{N}^{\mathcal{X}}$ , which takes a database as input and outputs another database. The function just removes duplicate entries from the input database and outputs the resulting database. Let  $d$  be a metric over  $X$ . The distance metrics  $d_1, d_2$  are defined as follows for any two databases  $x_1$  and  $x_2$ :  $d_1(x_1, x_2) = d_2(x_1, x_2) = l$  (according to  $d$ ) such that we can move each element of  $x_1$  and  $x_2$  by at most  $l$  and ensure that after moving there is no element which is in only  $x_1$  or only  $x_2$ . In other words, after moving, the set of elements in  $x_1$  must be equal to the set of elements in  $x_2$ . Again,  $f$  is a perfectly accurate mechanism for support. Also the error-sensitivity function for  $f$  is identity, i.e., Also, if for any two databases  $x_1$  and  $x_2$ ,  $d_1(x_1, x_2) = l$  then  $d_2(f(x_1), f(x_2)) = l$  because if we could move each element of input databases by a distance  $\beta'$  and satisfy the condition then the output is just the same entries just with duplicates removed so we can use the same distance to move them to satisfy the constraints. Now,  $\tau_f(\beta') = \beta'$  using the same argument as in the case of maximum. Now, using [Theorem 5](#) and [Theorem 6](#), we have that, the histogram mechanism for finding max is  $(\alpha_p, \beta_p, 0)$  accurate and  $(\epsilon_p, \delta_p)$  private where  $\alpha_p, \beta_p, \epsilon_p$  and  $\delta_p$  are defined in [Lemma 4](#).

356 **References**

- 357 [1] Cynthia Dwork and Aaron Roth. The algorithmic foundations of differential privacy. *Found.*  
358 *Trends Theor. Comput. Sci.*, 9(3&#8211;4):211–407, August 2014. 3, 8, 16, 18, 20

## 359 A Omitted Details from Section 2

**Claim 3.**

$$\inf_{\phi \in \Phi^\theta(p, q)} \max_{(x, y) \leftarrow \phi} |x - y| = \inf_{\substack{p', q': \\ \Delta(p, p') \leq \theta, \\ \Delta(q, q') \leq \theta}} W_\infty(p', q'),$$

360 where  $p'$  and  $q'$  are defined over the same alphabets where  $p$  and  $q$  are defined, respectively.

361 *Proof.* fill in. □

362 **We don't need the above claim, right?**

363 **Claim** (Restating Claim 1). *let  $X$  and  $Y$  be a random variable with distributions  $\mu_1, \mu_2$ .  $Z$  be a an*  
 364 *independent noise random variable with a distribution  $\mu_3$ . Then we have,*

$$W^\gamma(X, Y) = W^\gamma(X + Z, Y + Z)$$

365 *i.e., convolution of distributions with the same independent noise doesn't change wasserstein distance.*

366 *Proof.* First, we show that applying convolution of distributions with another doesn't increase the  
 367 wasserstein distance between them.

368 Let  $p, q, r$  be three distributions, with wasserstein distance between  $p, q$  as  $\beta$  and  $\pi$  be the optimal  
 369 transfer which achieves this . Let  $I_1, I_2$  be two distributions obtained by convolving  $p, q$  with  $r$   
 370 respectively. i.e.,

$$Pr(I_1 = x) = \int_Z p(x - z)r(z)$$

371

$$Pr(B + C = y) = \int_Z q(y - z)r(z)$$

372 From the above equations, we can construct a transfer from  $I_1, I_2$  as follows. Let  $z$  be drawn from  
 373  $r$  with a probability  $p(z)$ . In this event, we transfer  $I_1$  to  $I_2$  using the policy,  $\pi + z$ . This policy  
 374 also transfers wasserstein distance of  $\beta$  for all  $z$ . Hence, it's expectation over  $z$  is also  $\beta$ . Since,  
 375 wasserstein distance is an infimum of all transfers, we have, the wasserstein distance between  $I_1, I_2$   
 376 to be at most  $\beta$ . Hence, proved.

377 Now, for the random variables  $X, Y, Z$ . The distributions of  $X + Z, Y + Z$  are given by the following  
 378 convolutions.

$$Pr(A + C = x) = \int_Z \mu_1(x - z)\mu_3(z)$$

379

$$Pr(B + C = y) = \int_Z \mu_2(y - z)\mu_3(z)$$

380 From the above reasoning, we conclude that the wasserstein distance between  $X + Z, Y + Z$  is atmost  
 381 that of  $X, Y$ .

382 Note that the distributions of  $X$  and  $Y$  can be obtained by inverse convolution of  $Z$  with that of  $X + Z,$   
 383  $Y + Z$ . This is equivalent to convolution with an appropriate distribution  $Z'$ . **This is true for discrete**  
 384 **variables, convolution is matrix multiplication with a kernel, so inverse convolution is just multiplying**  
 385 **with inverse of the kernel. How does it go for continuous variables ?** . Again, as convolving with a  
 386 distribution doesn't increase the wasserstein distance, we conclude the wasserstein distance between  
 387  $X + Z, Y + Z$  is atmost that of  $X, Y$ .

388 Combining both these results, we have that, wasserstein distance between  $X, Y$  is same as wasserstein  
 389 distance between  $X + Z, Y + Z$ . □

390 **Claim** (Restating Claim 2).  *$\theta$ -Wasserstein distance follows the below triangular inequality. For*  
 391 *distributions,  $p, q$  and  $r$ , we have,*

$$W^{\gamma_1 + \gamma_2}(p, r) \leq W^{\gamma_1}(p, q) + W^{\gamma_2}(q, r)$$

392 *Proof.* fill in. □

Let  $\mathcal{X}$  be a finite set, and let  $\mathbf{x} := (x_1, x_2, \dots, x_n) \in \mathcal{X}^n$  represent a database with  $n$  elements, where  $x_i \in \mathcal{X}$  is the  $i$ -th element. Let  $f : \mathcal{X}^n \rightarrow \mathcal{Y}$  be a randomized function, where  $\mathcal{Y} = \{y_1, y_2, \dots, y_k\}$  is a discrete set for some finite  $k \in \mathbb{N}$ . For every  $\mathbf{x} \in \mathcal{X}^n$ , we have that  $\sum_{i=1}^k \Pr[f(\mathbf{x}) = y_i] = 1$ .

## B Omitted Details from Section 3

### B.1 Robust Mechanism for Identity Function over $\mathbb{R}$

We will now show a mechanism  $\mathcal{M}_{\text{rob}} : \mathbb{R} \rightarrow \mathbb{R}$  which is a robust mechanism for identity function, i.e., the output corresponding to any input is the same as input. Now Suppose the input distribution is  $x$ . Let  $\rho$  be any non-negative real number and  $\epsilon$  be any positive real number. Consider the following mechanism for the identity function:

**Mechanism** (Restating Mechanism 1).  $\mathcal{M}_{\text{rob}}$ : On input  $y$ , sample  $z$  according to the probability distributions  $\text{Lap}(\frac{\rho}{\epsilon})$  and output  $y + z$ .

**Lemma** (Restating Lemma 1). For a fixed constant  $\theta \in [0, 1]$ ,  $\mathcal{M}_{\text{rob}}$  achieves  $(\theta, \rho, \epsilon, \delta)$ -robustness where  $\delta = 2e^\epsilon \theta$ .

*Proof.* Fix a constant  $\theta \geq 0$ . Suppose  $q$  and  $q'$  denote two input distributions, and let  $p$  and  $p'$  denote the corresponding output probability distributions of  $\mathcal{M}_{\text{rob}}$ . Suppose  $W_\infty^\theta(p, p') \leq \rho$ . Then by Claim 1,  $W_\infty^\theta(q, q') \leq \rho$ . Let  $b := \frac{\rho}{\epsilon}$  and  $Z \sim \text{Lap}(b)$ . We need to show that  $\forall S \subseteq \mathbb{R}$ ,  $\Pr_{y \leftarrow p}[y \in S] \leq e^\epsilon \Pr_{y \leftarrow p'}[y \in S] + \delta$ .

$$\Pr_{y \leftarrow p}[y \in S] = \int_{t \in S} p(t) \cdot dt \quad (17)$$

$$= \int_{t \in S} \left[ \int_{-\infty}^{\infty} q(s) \cdot p_Z(t - s) \cdot ds \right] \cdot dt \quad (18)$$

Let  $\phi^\theta \in \Phi^\theta(q, q')$  be a joint distribution. Let  $M^\phi$  denote  $\max_{(x, y) \leftarrow \phi^\theta} |x - y|$  for the distribution  $\phi^\theta$ . For  $i \in \mathbb{R}$ , let  $\theta_1(i) := |\phi_1^\theta(i) - q(i)|$  and  $\theta_2(i) := |\phi_2^\theta(i) - q'(i)|$ . Note that, by definition of the distribution  $\phi^\theta$ , we have  $\Delta(\phi_1^\theta, q) + \Delta(\phi_2^\theta, q') \leq \theta$  which implies  $\int_{-\infty}^{\infty} \theta_1(i) \cdot di + \int_{-\infty}^{\infty} \theta_2(i) \cdot di \leq \theta$ . Using these in (18) we get the following:

$$\begin{aligned} \Pr_{y \leftarrow p}[y \in S] &\leq \int_{t \in S} \left[ \frac{1}{2b} \int_{-\infty}^{\infty} (\phi_1^\theta(i) + \theta_1(i)) \cdot e^{\frac{-|t-i|}{b}} \cdot di \right] \cdot dt \\ &= \int_{t \in S} \left[ \frac{1}{2b} \int_{-\infty}^{\infty} \phi_1^\theta(i) \cdot e^{\frac{-|t-i|}{b}} \cdot di \right] \cdot dt + \int_{t \in S} \left[ \frac{1}{2b} \int_{-\infty}^{\infty} \theta_1(i) \cdot e^{\frac{-|t-i|}{b}} \cdot di \right] \cdot dt \\ &= \int_{t \in S} \left[ \frac{1}{2b} \int_{-\infty}^{\infty} \phi_1^\theta(i) \cdot e^{\frac{-|t-i|}{b}} \cdot di \right] \cdot dt + \underbrace{\int_{-\infty}^{\infty} \theta_1(i) \cdot \int_{t \in S} \frac{1}{2b} e^{\frac{-|t-i|}{b}} \cdot dt \cdot di}_{\leq 1} \quad (19) \end{aligned}$$

By properties of joint distributions, we have  $\phi_1^\theta(i) = \int_{-\infty}^{\infty} \phi^\theta(i, j) \cdot dj$  and  $\phi_2^\theta(j) = \int_{-\infty}^{\infty} \phi^\theta(i, j) \cdot di$ ; and by triangle inequality we have  $|t - i| \geq |t - j| - |i - j|$ . Let  $\int_{-\infty}^{\infty} \theta_l(i) \cdot di$  be  $\omega_l$ ,  $l \in \{1, 2\}$ . Substituting all these in (19) we get the following:

$$\Pr_{y \leftarrow p}[y \in S] \leq \int_{t \in S} \left[ \frac{1}{2b} \int_{-\infty}^{\infty} \left[ \int_{-\infty}^{\infty} \phi^\theta(i, j) \cdot e^{\frac{-|t-j|+|i-j|}{b}} \cdot dj \right] \cdot di \right] \cdot dt + \omega_1. \quad (20)$$

Observe that  $\phi^\theta(i, j)$  is non-zero only if  $|i - j| \leq M^\phi$  (by definition of  $M^\phi$ ). Using this for every  $i, j \in \mathbb{R}$ , the integrand in (20) can be upper-bounded as follows:

$$\phi^\theta(i, j) \cdot e^{\frac{-|t-j|}{b}} e^{\frac{|i-j|}{b}} \leq \phi^\theta(i, j) \cdot e^{\frac{-|t-j|}{b}} e^{\frac{M^\phi}{b}} = \phi^\theta(i, j) \cdot e^{\frac{-|t-j|}{b}} e^{\frac{M^\phi}{b}}.$$

Substituting this in (20) gives

$$\Pr_{y \leftarrow p}[y \in S] \leq \int_{t \in S} \left[ \frac{e^{\frac{M^\phi}{b}}}{2b} \int_{-\infty}^{\infty} \left[ \int_{-\infty}^{\infty} \phi^\theta(i, j) \cdot e^{\frac{-|t-j|}{b}} \cdot dj \right] \cdot di \right] \cdot dt + \omega_1$$

$$\begin{aligned}
&= \int_{t \in S} \left[ \frac{e^{\frac{M^\phi}{b}}}{2b} \int_{-\infty}^{\infty} \left[ \int_{-\infty}^{\infty} \phi^\theta(i, j) \cdot e^{\frac{-|t-j|}{b}} \cdot dj \right] \cdot dt + \omega_1 \right. \\
&= \int_{t \in S} \left[ \frac{e^{\frac{M^\phi}{b}}}{2b} \int_{-\infty}^{\infty} \phi_2^\theta(j) \cdot e^{\frac{-|t-j|}{b}} \cdot dj \right] \cdot dt + \omega_1 \\
&\leq \int_{t \in S} \left[ \frac{e^{\frac{M^\phi}{b}}}{2b} \int_{-\infty}^{\infty} (q'_j + \theta_2(j)) \cdot e^{\frac{-|t-j|}{b}} \cdot dj \right] \cdot dt + \omega_1 \quad (\text{since } \theta_2(j) = |\phi_2^\theta(j) - q'(j)|) \\
&\leq \int_{t \in S} \left[ \frac{e^{\frac{M^\phi}{b}}}{2b} \int_{-\infty}^{\infty} q'_j \cdot e^{\frac{-|t-j|}{b}} \cdot dj \right] \cdot dt + \int_{t \in S} \left[ \frac{e^{\frac{M^\phi}{b}}}{2b} \int_{-\infty}^{\infty} \theta_2(j) \cdot e^{\frac{-|t-j|}{b}} \cdot dj \right] \cdot dt + \omega_1 \\
&\leq \int_{t \in S} \left[ \frac{e^{\frac{M^\phi}{b}}}{2b} \int_{-\infty}^{\infty} q'_j \cdot e^{\frac{-|t-j|}{b}} \cdot dj \right] \cdot dt + e^{\frac{M^\phi}{b}} \int_{-\infty}^{\infty} \theta_2(j) \underbrace{\left[ \int_{t \in S} \frac{1}{2b} \cdot e^{\frac{-|t-j|}{b}} \cdot dt \right]}_{\leq 1} \cdot dj + \omega_1 \\
&\leq \int_{t \in S} \left[ \frac{e^{\frac{M^\phi}{b}}}{2b} \int_{-\infty}^{\infty} q'_j \cdot e^{\frac{-|t-j|}{b}} \cdot dj \right] \cdot dt + e^{\frac{M^\phi}{b}} \cdot \omega_2 + \omega_1 \\
&= \int_{t \in S} \left[ e^{\frac{M^\phi}{b}} \int_{-\infty}^{\infty} p'(t) \cdot dj \right] \cdot dt + e^{\frac{M^\phi}{b}} \cdot \omega_2 + \omega_1 \\
&= e^{\frac{M^\phi}{b}} \Pr_{y \leftarrow p'}[y \in S] + e^{\frac{M^\phi}{b}} \cdot \omega_2 + \omega_1
\end{aligned}$$

421 Now, by the limiting case of  $\phi^\theta$  to be optimal,  $M^\phi$  tends to  $\rho$ , thus we have,

$$\begin{aligned}
\Pr_{y \leftarrow p}[y \in S] &\leq e^{\frac{\rho}{b}} \Pr_{y \leftarrow p'}[y \in S] + e^{\frac{\rho}{b}} \cdot \omega_2 + \omega_1 \\
&\leq e^{\frac{\rho}{b}} \Pr_{y \leftarrow p'}[y \in S] + e^{\frac{\rho}{b}} \cdot 2\theta \quad (\text{since } \omega_1 + \omega_2 \leq 2\theta, \text{ RHS is maximized at } \omega_2 = 2\theta) \\
&= e^\epsilon \Pr_{y \leftarrow p'}[y \in S] + \delta, \quad \text{where } \delta = 2\theta e^\epsilon
\end{aligned}$$

422

□

423 **Lemma (Restating Lemma 2).** For every constant  $\gamma' \geq 0$ ,  $\mathcal{M}_{\text{rob}}$  is  $(0, \beta', \gamma')$ -accurate, where

$$424 \quad \beta' = \frac{\rho}{\epsilon(1-\gamma')} \left( 1 - \gamma' [1 + \ln(\frac{1}{\gamma'})] \right). \text{ Note that if } \gamma' = 0, \beta' = \frac{\rho}{\epsilon}.$$

425 *Proof.* Fix a constant  $\gamma \geq 0$  and take an arbitrary distribution (over  $\mathbb{R}$ ) and let's denote it by  $x$ . Let  $q$   
426 denote the output distribution of  $\mathcal{M}_{\text{rob}}$  when the input is drawn from  $x$ . Let  $b$  denote  $\frac{\rho}{\epsilon}$ . We want  
427 to show that  $W^\gamma(x, q) \leq \beta$ , for the above-mentioned  $\beta$ . By definition of  $W^\gamma$  from (2) we have  
428  $W^\gamma(x, q) = \inf_{\phi \in \Phi^\gamma(x, q)} \mathbb{E}_{(y, t) \leftarrow \phi} [|y - t|].$

429 Consider the following  $\phi^*$ :

$$\phi^*(i, t) = \begin{cases} 0 & \text{if } t < -b \ln(\frac{1}{\gamma}) + i \text{ or } t > b \ln(\frac{1}{\gamma}) + i; \\ \frac{1}{1-\gamma} \text{Lap}(t|b, i)x(i) & \text{if } t \in [-b \ln(\frac{1}{\gamma}) + i, b \ln(\frac{1}{\gamma}) + i]. \end{cases}$$

430 It can be verified that  $\Delta(\phi_1^*, x) = 0$  and  $\Delta(\phi_2^*, q) \leq \gamma$ , which implies that  $\phi^* \in \Phi^\gamma(x, q)$ . This  
431 in turn implies that  $W^\gamma(x, q) \leq \mathbb{E}_{(y, t) \leftarrow \phi^*} [|y - t|]$ . We show below that  $\mathbb{E}_{(y, t) \leftarrow \phi^*} [|y - t|] \leq$

432  $\frac{b}{(1-\gamma)} \left( 1 - \gamma [1 + \ln(\frac{1}{\gamma})] \right)$ . This will prove Lemma 2.

$$\begin{aligned}
\mathbb{E}_{(y, t) \leftarrow \phi^*} [|y - t|] &= \int_{w=-\infty}^{w=\infty} |w| \cdot \Pr_{(y, t) \leftarrow \phi^*} [|y - t| = w] \cdot dw \\
&= \int_{w=0}^{w=\infty} w \cdot \Pr_{(y, t) \leftarrow \phi^*} [|y - t| = w] \cdot dw \\
&= \int_{w=0}^{w=\infty} w \cdot \left[ \Pr_{(y, t) \leftarrow \phi^*} [y - t = w] + \Pr_{(y, t) \leftarrow \phi^*} [y - t = -w] \right] \cdot dw
\end{aligned}$$



$$\begin{aligned}
&= \int_{w=0}^{w=\infty} w \cdot \left[ \int_{-\infty}^{\infty} \phi^*(i, i-w) \cdot di + \int_{-\infty}^{\infty} \phi^*(i, i+w) \cdot di \right] \cdot dw \\
&= \int_{w=0}^{w=\infty} w \cdot \left[ \int_{-\infty}^{\infty} 2\phi^*(i, i+w) \cdot di \right] \cdot dw \\
&\quad \text{(Since } \phi^*(i, i-w) = \phi^*(i, i+w) \text{)} \\
&= \int_{-\infty}^{\infty} \left[ \int_0^{\infty} w \cdot 2\phi^*(i, i+w) \cdot dw \right] \cdot di
\end{aligned}$$

433 Let  $d = b \ln(\frac{1}{\gamma})$ . By definition of  $\phi^*$ ,  $\phi^*(i, i+w)$  is non-zero only if  $i+w \in [i-d, i+d]$ , which  
434 implies that  $w \in [-d, d]$ . Using this in above gives:

$$\begin{aligned}
\mathbb{E}_{(y,t) \leftarrow \phi^*}[|y-t|] &= \int_{-\infty}^{\infty} \left[ \int_0^d w \cdot 2\phi^*(i, i+w) \cdot dw \right] \cdot di \\
&= \int_{-\infty}^{\infty} x_i \left[ \int_0^d w \cdot \frac{1}{(1-\gamma)} \left( \frac{1}{b} e^{\frac{-|i+w-i|}{b}} \right) \cdot dw \right] \cdot di \\
&= \frac{1}{b(1-\gamma)} \int_{-\infty}^{\infty} x(i) \left[ \int_0^d w \cdot \left( e^{\frac{-|w|}{b}} \right) \cdot dw \right] \cdot di \\
&= \frac{1}{b(1-\gamma)} \int_{-\infty}^{\infty} x(i) \left[ \int_0^d w \cdot \left( e^{\frac{-w}{b}} \right) \cdot dw \right] \cdot di \\
&= \frac{1}{b(1-\gamma)} \int_{-\infty}^{\infty} x(i) \left[ \int_0^{d/b} (bw) \cdot e^{-w} \cdot (b \cdot dw) \right] \cdot di \\
&= \frac{b}{(1-\gamma)} \int_{-\infty}^{\infty} x(i) \left[ \int_0^{\ln(\frac{1}{\gamma})} w \cdot e^{-w} \cdot dw \right] \cdot di \quad \text{(since } d = b \ln(\frac{1}{\gamma}) \text{)} \\
&= \frac{b}{(1-\gamma)} \int_{-\infty}^{\infty} x(i) \left[ (-we^{-w} - e^{-w}) \Big|_{w=0}^{w=\ln(\frac{1}{\gamma})} \right] \cdot di \\
&= \frac{b}{(1-\gamma)} \int_{-\infty}^{\infty} x(i) \left( 1 - \gamma \left[ 1 + \ln\left(\frac{1}{\gamma}\right) \right] \right) \cdot di \\
&= \frac{b}{(1-\gamma)} \left( 1 - \gamma \left[ 1 + \ln\left(\frac{1}{\gamma}\right) \right] \right) \\
&= \frac{\rho}{\epsilon(1-\gamma)} \left( 1 - \gamma \left[ 1 + \ln\left(\frac{1}{\gamma}\right) \right] \right) \quad \text{(since } b = \frac{\rho}{\epsilon} \text{)}
\end{aligned}$$

435 Note that at  $\gamma = 0$ , we have  $\mathbb{E}_{(y,t) \leftarrow \phi^*}[|y-t|] = \frac{\rho}{\epsilon}$ . □

436 **Theorem 7.** For any  $\theta \in [0, 1]$  and  $\gamma \geq 0$ ,  $\mathcal{M}_{\text{rob}}$  is  $(\theta, \rho, \epsilon, \delta)$ -robust and  $(0, \beta, \gamma)$ -accurate, where  
437  $\delta = 2e^\epsilon \theta$  and  $\beta = \frac{\rho}{\epsilon(1-\gamma)} \left( 1 - \gamma \left[ 1 + \ln\left(\frac{1}{\gamma}\right) \right] \right)$ .

438 *Proof.* Using [Lemma 1](#) and [Lemma 2](#), the theorem trivially holds. □

## 439 B.2 Adding Robustness to a Mechanism

440 Consider a randomized mechanism  $\mathcal{M} : A \rightarrow B$  which has some privacy and accuracy bounds and a  
441 randomized mechanism  $\mathcal{M}_{\text{rob}} : B \rightarrow B$  which is robust and has some accuracy bounds. We want to  
442 construct a new mechanism which is robust, atleast as private as  $\mathcal{M}$  and doesn't compromise much  
443 on accuracy as compared to  $\mathcal{M}$ . Let  $x \in A$  be the input. Consider the following mechanism:

445 **Mechanism** (Restating [Mechanism 2](#)). Run  $\mathcal{M}$  on  $x$  to get an output  $y \in B$  then run  $\mathcal{M}_{\text{rob}}$  on  $y$  to  
446 get  $z \in B$ , output  $z$ .

**Theorem (Restating Theorem 2).** Let  $\rho, \theta, \epsilon, \delta$  be non-negative real numbers with  $\delta, \theta \leq 1$ . Then **Mechanism 2** achieves (i)  $(\theta, \rho, \epsilon, \delta)$ -robustness, where  $\delta = 2e^\epsilon \theta$ , (ii)  $(\epsilon_p, \delta_p)$ -differential privacy, if  $\mathcal{M}$  is  $(\epsilon_p, \delta_p)$ -differentially private, and (iii)  $(\alpha_p, \beta_p, \gamma_p)$ -accuracy, if  $\mathcal{M}$  is  $(\alpha, \beta, \gamma)$ -accurate, where  $\alpha_p = \alpha, \beta_p = \beta + \beta', \gamma_p = \gamma + \gamma'$  and  $\beta', \gamma'$  are such that  $\gamma' \geq 0$  is arbitrary and  $\beta' = \frac{\rho}{\epsilon(1-\gamma')} \left(1 - \gamma'[1 + \ln(\frac{1}{\gamma'})]\right)$ .

*Proof.* Observe that **Mechanism 2** is equivalent to  $\mathcal{M}_{\text{rob}} \circ \mathcal{M}$ . First we prove the accuracy guarantee. From **Lemma 2**, we have that  $\mathcal{M}_{\text{rob}}$  is  $(0, \beta', \gamma')$  accurate, where  $\gamma' \geq 0$  is arbitrary and  $\beta' = \frac{\rho}{\epsilon(1-\gamma')} \left(1 - \gamma'[1 + \ln(\frac{1}{\gamma'})]\right)$ . Now, let us compute the error-sensitivity function of  $\mathcal{M}_{\text{rob}}$ .

$$\begin{aligned} \tau_{\mathcal{M}_{\text{rob}}}^\theta(\beta) &= \max_{\substack{x \sim p, x' \sim q: \\ W^\theta(p, q) \leq \beta}} W^\theta(\mathcal{M}_{\text{rob}}(x), \mathcal{M}_{\text{rob}}(x')) \\ &\stackrel{(a)}{=} \max_{\substack{x \sim p, x' \sim q: \\ W^\theta(p, q) \leq \beta}} W^\theta(p, q) \\ &= \beta. \end{aligned}$$

Here (a) follows from **Claim 1** and that  $\mathcal{M}_{\text{rob}}$  adds independent laplacian noise. Thus, we have that the required error sensitivity function is identity. Now, if  $\mathcal{M}$  is  $(\alpha, \beta, \gamma)$ -accurate, then by **Theorem 5**, we have that  $\mathcal{M}_{\text{rob}} \circ \mathcal{M}$  is  $(\alpha_p, \beta_p, \gamma_p)$ -accurate with  $\alpha_p = \alpha, \beta_p = \beta' + \beta, \gamma_p = \gamma' + \gamma$ .

For the privacy guarantee, since  $\mathcal{M}$  is  $(\epsilon_p, \delta_p)$ -differentially private and post-processing preserves differential privacy [1, Proposition 2.1],  $\mathcal{M}_{\text{rob}} \circ \mathcal{M}$  is also  $(\epsilon_p, \delta_p)$ -differentially private.

The proof of robustness is straightforward. By **Definition 9**, if any two input distributions gives close (in terms of the Wasserstein distance) output distributions, then they are also close in terms of differential privacy sense. Since the requirement is only on the output distributions, it follows that if a mechanism  $A$  is robust, then  $A \circ B$  is also robust for every mechanism  $B$ , with exactly the same parameters. This implies that, since  $\mathcal{M}_{\text{rob}}$  is  $(\theta, \rho, \epsilon, 2e^\epsilon \theta)$ -robust (see **Lemma 1**),  $\mathcal{M}_{\text{rob}} \circ \mathcal{M}$  is also  $(\theta, \rho, \epsilon, 2e^\epsilon \theta)$ -robust.

This completes the proof of **Theorem 2**. □

## C Omitted Details from Section 5

We will now define a new mechanism to get a new histogram from a given histogram in a differential private manner.

**Mechanism (Restating Mechanism 4).** For each bar of the histogram, independently add a noise according to the following probability distribution:

$$\pi_q(x) = \begin{cases} \frac{1}{1-e^{-\frac{\sqrt{q}}{2}}} \text{Lap}(x \mid -\frac{q}{2}, \sqrt{q}) & -q < x < 0 \\ 0 & \text{elsewhere} \end{cases}$$

and then round each of the bar to the nearest integer. If any of the bar becomes negative then set it to 0. Output the resulting histogram.

**Lemma 5.** For any  $q \geq 1$ , if in **Mechanism 4** we output the histogram before rounding and setting negative values to 0 then we get  $(\epsilon, \delta)$ -differential privacy where  $\epsilon = \frac{1}{\sqrt{q}}$  and  $\delta = \frac{4}{\sqrt{q}} e^{-\frac{\sqrt{q}}{2}}$ .

*Proof.* Let  $t$  be the size of domain of the values in the database. Let  $\mathbf{x}$  and  $\mathbf{x}'$  be two neighbouring databases and let  $n_k$  and  $n'_k$  represent the number of values of  $k^{\text{th}}$  type in  $\mathbf{x}$  and  $\mathbf{x}'$  respectively. Let  $i^*$  be the index of the type in which  $\mathbf{x}$  and  $\mathbf{x}'$  differ. That means that  $n_k = n'_k$  if  $k \neq i^*$  and  $n_{i^*} = n'_{i^*} + 1$ .  $p_D(s)$  denotes the probability of outputting  $s$  from a distribution  $D$ .

$$\begin{aligned}\Pr[\mathcal{M}(\mathbf{x}) \in S] &= \int_{s \in S} p_{\mathcal{M}(\mathbf{x})}(s) \cdot ds \\ &= \int_{s \in S} \left[ \prod_{i=1}^t \pi_q(s_i - n_i) \right] \cdot ds\end{aligned}\quad (21)$$

$$\Pr[\mathcal{M}(\mathbf{x}') \in S] = \int_{s \in S} \left[ \prod_{i=1}^t \pi_q(s_i - n'_i) \right] \cdot ds \quad (22)$$

477 Using the fact that  $\forall k \neq i^*, n_k = n'_k$  and  $n_{i^*} = n'_{i^*} + 1$ , we can divide  $S$  into 6 sets as follows:

478  $S_0$  contains those outputs in which  $-\frac{q}{2} \leq s_{i^*} - n_{i^*}, s_{i^*} - n'_{i^*} < 0$ ,

479  $S_1$  contains those outputs in which  $-q < s_{i^*} - n_{i^*}, s_{i^*} - n'_{i^*} \leq -\frac{q}{2}$ ,

480  $S_2$  contains those outputs in which  $-q < s_{i^*} - n_{i^*} < -\frac{q}{2} < s_{i^*} - n'_{i^*} < 0$ ,

481  $S_3$  contains those outputs in which  $s_{i^*} - n_{i^*} \leq -q < s_{i^*} - n'_{i^*} < 0$ ,

482  $S_4$  contains those outputs in which  $-q < s_{i^*} - n_{i^*} < 0 \leq s_{i^*} - n'_{i^*}$  and

483  $S_5$  contains those outputs in which either  $0 \leq s_{i^*} - n_{i^*}, s_{i^*} - n'_{i^*}$  or  $s_{i^*} - n_{i^*}, s_{i^*} - n'_{i^*} \leq -q$ .

484

485 For each of the above set except  $S_4$ , we analyze the difference between  $a(= \pi_q(s_{i^*} - n'_{i^*}))$  and

486  $b(= \pi_q(s_{i^*} - n_{i^*}) = \pi_q(s_{i^*} - n'_{i^*} - 1))$ .

487 For  $S_0$ ,  $b = ae^{\frac{1}{\sqrt{q}}} \leq ae^\epsilon$

488 For  $S_1$ ,  $b = ae^{-\frac{1}{\sqrt{q}}} \leq ae^\epsilon$

489 For  $S_2$ ,  $b \leq ae^{\frac{1}{\sqrt{q}}} \leq ae^\epsilon$

490 For  $S_3$ ,  $b = 0 \leq ae^\epsilon$

491 For  $S_5$ ,  $b = 0 \leq ae^\epsilon$

492

493 Now let us analyze (21) for each  $S_i$ :

494

495 **Case1:**  $S_i, i = 0, 1, 2, 3, 5$

$$\begin{aligned}&= \int_{s \in S_i} \left[ \prod_{i=1}^t \pi_q(s_i - n_i) \right] \cdot ds \\ &= \int_{s \in S_i} \left[ \prod_{\substack{i=1 \\ i \neq i^*}}^t \pi_q(s_i - n_i) \right] \pi_q(s_{i^*} - n_{i^*}) \cdot ds \\ &= \int_{s \in S_i} \left[ \prod_{\substack{i=1 \\ i \neq i^*}}^t \pi_q(s_i - n'_i) \right] \pi_q(s_{i^*} - n'_{i^*} - 1) \cdot ds \\ &\leq \int_{s \in S_i} \left[ \prod_{\substack{i=1 \\ i \neq i^*}}^t \pi_q(s_i - n'_i) \right] \pi_q(s_{i^*} - n'_{i^*}) e^\epsilon \cdot ds \\ &= e^\epsilon \int_{s \in S_i} \left[ \prod_{i=1}^t \pi_q(s_i - n'_i) \right] \cdot ds\end{aligned}$$

496 **Case2:**  $S_4$ , the condition forces  $s_{i^*} - n_{i^*}$  to be in  $[-1, 0)$

$$\begin{aligned}&= \int_{s \in S_4} \left[ \prod_{i=1}^t \pi_q(s_i - n_i) \right] \cdot ds \\ &= \int_{s_1} \dots \int_{s_{i^*}} \dots \int_{s_t} \left[ \prod_{i=1}^t \pi_q(s_i - n_i) \right] \cdot ds_t \dots ds_{i^*} \dots ds_1\end{aligned}$$

$$\begin{aligned}
&= \int_{s_{i^*}} \pi_q(s_{i^*} - n_{i^*}) \left( \int_{s_1} \dots \int_{s_t} \left[ \prod_{\substack{i=1 \\ i \neq i^*}}^t \pi_q(s_i - n_i) \right] \cdot ds_t \dots ds_1 \right) ds_{i^*} \\
&\quad \text{(as } s_{i^*} \text{ varies from -1 to 0 independent of other } s'_k s) \\
&\leq \int_{s_{i^*}} \pi_q(s_{i^*} - n_{i^*}) (1) ds_{i^*} \\
&= \int_{-1}^0 \pi_q(s_{i^*} - n_{i^*}) ds_{i^*} \\
&= \frac{e^{1/\sqrt{q}} - 1}{2(1 - e^{-\sqrt{q}/2})} e^{-\sqrt{q}/2} \\
&\leq 2(e^{1/\sqrt{q}} - 1)e^{-\sqrt{q}/2} \quad (1 - e^{-\sqrt{q}/2} \geq 1 - e^{-1/2} \geq \frac{1}{4}) \\
&\leq \frac{4}{\sqrt{q}} e^{-\sqrt{q}/2} \quad (\text{For } x \leq 1, e^x \leq 2x + 1) \\
&= \delta
\end{aligned}$$

497 Now we combine the above results with Equations (21) and (22) as follows:

$$\begin{aligned}
\Pr[\mathcal{M}(\mathbf{x}) \in S] &= \int_{s \in S} p_{\mathcal{M}(\mathbf{x})}(s) \cdot ds \\
&= \sum_{l=0}^5 \int_{s \in S_l} p_{\mathcal{M}(\mathbf{x})}(s) \cdot ds \\
&= \sum_{\substack{l=0 \\ l \neq 4}}^5 \int_{s \in S_l} p_{\mathcal{M}(\mathbf{x})}(s) \cdot ds + \int_{s \in S_4} p_{\mathcal{M}(\mathbf{x})}(s) \cdot ds \\
&\leq \sum_{\substack{l=0 \\ l \neq 4}}^5 e^\epsilon \int_{s \in S_l} p_{\mathcal{M}(\mathbf{x}')} (s) \cdot ds + \delta \\
&= e^\epsilon \left( \sum_{\substack{l=0 \\ l \neq 4}}^5 \int_{s \in S_l} p_{\mathcal{M}(\mathbf{x}')} (s) \cdot ds \right) + \delta \\
&= e^\epsilon \left( \sum_{\substack{l=0 \\ l \neq 4}}^5 \int_{s \in S_l} p_{\mathcal{M}(\mathbf{x}')} (s) \cdot ds + 0 \right) + \delta \\
&= e^\epsilon \left( \sum_{\substack{l=0 \\ l \neq 4}}^5 \int_{s \in S_l} p_{\mathcal{M}(\mathbf{x}')} (s) \cdot ds + \int_{s \in S_4} p_{\mathcal{M}(\mathbf{x}')} (s) \cdot ds \right) + \delta \\
&= e^\epsilon \left( \sum_{l=0}^5 \int_{s \in S_l} p_{\mathcal{M}(\mathbf{x}')} (s) \cdot ds \right) + \delta \\
&= e^\epsilon \Pr[\mathcal{M}(\mathbf{x}') \in S] + \delta
\end{aligned}$$

498 Similarly we can also prove the following:  $\Pr[\mathcal{M}(\mathbf{x}') \in S] \leq e^\epsilon \Pr[\mathcal{M}(\mathbf{x}) \in S] + \delta$ .  $\square$

499 **Theorem 8.** For any  $q \geq 1$ , **Mechanism 4** achieves  $(\epsilon, \delta)$ -differential privacy where  $\epsilon = \frac{1}{\sqrt{q}}$  and

$$500 \quad \delta = \frac{4}{\sqrt{q}} e^{-\frac{\sqrt{q}}{2}}.$$

501 *Proof.* Rounding and setting negative values to 0 is a post-processing step after adding noise. By  
502 **Lemma 5**, the addition of noise is  $(\epsilon, \delta)$ -private. Using the fact from [1] that post-processing doesn't  
503 decrease the privacy, we can say that **Mechanism 4** is also  $(\epsilon, \delta)$  private.  $\square$

504 **Theorem** (Restating **Theorem 4**). *Mechanism 4* achieves  $(\alpha, 0, 0)$  accuracy with  $\alpha = qt$ .

505 *Proof.* This can be easily seen. In essence the histogram mechanism deletes some elements from  
 506 the buckets. Now, the mechanism deletes almost  $q$  elements from a bucket, as the probability of  
 507 adding more negative noise is zero. Hence, the histogram returned by the mechanism can be obtained  
 508 by deleting at most  $qt$ , where  $t$  is the number of buckets from the original database. Thus, by the  
 509 considering the distorted random variable as the corresponding database for every instance of the  
 510 histogram mechanism noise, we get its distribution as that of output distributions of the mechanism.  
 511 Hence, the mechanism is  $(qt, 0, 0)$  accurate.  $\square$

## 512 D Omitted Details from Section 6

513 Now, we present our privacy composition theorem.

514 **Theorem** (Restating **Theorem 6**). *Let  $M_1 : A \rightarrow B$  be a neighbourhood preserving mechanism for*  
 515 *a function  $f_1 : A \rightarrow B$ , and let  $M_2 : B \rightarrow C$  be a mechanism which is  $(\epsilon, \delta)$ -differential private for*  
 516 *a function  $f_2 : B \rightarrow C$ . Then the composite mechanism,  $M_2 \circ M_1 : A \rightarrow C$  is  $(\epsilon, \delta)$ -differential*  
 517 *private.*

518 *Proof.* Note that we defined a randomized mechanism to be neighbourhood preserving, if it is a  
 519 convex combination of neighbourhood preserving deterministic functions.

520 First we consider the case when the mechanism  $M_1$  is deterministic. In this case, for neighbouring  
 521 databases  $\mathbf{x}, \mathbf{x}'$ ,  $M_1(\mathbf{x}), M_1(\mathbf{x}')$  are deterministic databases in  $B$ , which are also neighbours. Hence,  
 522 by the  $(\epsilon, \delta)$ -differential privacy of the mechanism  $M_2$ , for all subsets  $S \subseteq C$ , we have

$$P(M_2(M_1(\mathbf{x})) \in S) \leq e^\epsilon P(M_2(M_1(\mathbf{x}')) \in S) + \delta.$$

523 This proves the privacy guarantee of the mechanism in case of deterministic  $M_1$ .

524 For randomized mechanisms  $M_1$ , we use the fact that we can represent a randomized mechanism  
 525 as a convex combination of deterministic mechanisms, where the linear weights corresponds to the  
 526 randomness in the mechanism. Thus, we have  $M_1 = \sum_{i=1}^n p_i M_{1,i}$ , where  $\sum_{i=1}^n p_i = 1$ . Here,  $M_{1,i}$   
 527 are the deterministic mechanism, and  $M_1 = M_{1,i}$  with probability  $p_i$ . This implies that for every  
 528  $S \subseteq C$ , we have

$$P(M_1(\mathbf{x}) \in S) = \sum_{i=1}^n p_i P(M_{1,i}(\mathbf{x}) \in S).$$

529 Now, for two neighbouring databases  $\mathbf{x}, \mathbf{x}' \in A$ , and any subset  $S \subseteq C$ , we have

$$\begin{aligned} P(M_2(M_1(\mathbf{x})) \in S) &= \sum_{i=1}^n p_i P(M_2(M_{1,i}(\mathbf{x})) \in S) \\ &\stackrel{(a)}{\leq} \sum_{i=1}^n p_i (e^\epsilon P(M_2(M_{1,i}(\mathbf{x}')) \in S) + \delta) \\ &= e^\epsilon \sum_{i=1}^n p_i P(M_2(M_{1,i}(\mathbf{x}')) \in S) + \delta \\ &= e^\epsilon P(M_2(M_1(\mathbf{x}')) \in S) + \delta \end{aligned}$$

530 Here (a) follows because  $M_2$  is  $(\epsilon, \delta)$  private and  $M_{1,i}(\mathbf{x})$  and  $M_{1,i}(\mathbf{x}')$  are neighbours in  $B$ . This  
 531 proves that the  $(\epsilon, \delta)$  privacy guarantee of the composite mechanism  $M_2 \circ M_1$ .  $\square$

532 Now, we show how appropriate mechanisms can be composed to achieve good accuracy and privacy  
 533 guarantees for computing a statistic on a database. We consider databases which consists of positive  
 534 real number numbers. The bucketing is done with  $t$  number of buckets, i.e, rounding each element of  
 535 database to the nearest multiple of  $\frac{B}{t}$ .

536 **Lemma** (Restating **Lemma 4**). *The bucketing followed by Mechanism 4 is an  $(\alpha_p, \beta_p, 0)$ -accurate*  
 537 *and  $(\epsilon_p, \delta_p)$ -differentially private mechanism for the identity function, with  $\alpha_p = qt, \beta_p = \frac{B}{2t}, \epsilon_p =$   
 538  $\frac{1}{\sqrt{q}}, \delta_p = \frac{4}{\sqrt{q}} e^{-\frac{\sqrt{q}}{2}}.$*

539 *Proof.* Let us denote the Bucketing mechanism by  $\mathcal{M}_{buck}$  and Mechanism 4 by  $\mathcal{M}_{hist}$ . We compute  
 540 distance between histograms by considering them as distributions and taking the Wasserstein distance  
 541 between them. Now, using this distance metric, we have the bucketing mechanism to be  $(0, \frac{B}{2t}, 0)$  for  
 542 identity function, as each element is distorted by at most  $B/2t$ , which means that the corresponding  
 543 histogram will shift by at most the same distance.

544 Also note that bucketing mechanism is a deterministic mechanism that is neighbourhood preserving.  
 545 This is because, removing a single element changes the output of bucketing by at most one element.  
 546 Hence, neighbours remain neighbours after bucketing. This implies that the distortion sensitivity  
 547 function of bucketing mechanism is identity, i.e,  $\sigma_{\mathcal{M}_{buck}}(\alpha') = \alpha'$ .

548 We also have that the histogram mechanism is  $(qt, 0, 0)$ -accurate (see Theorem 4) and  $(\frac{1}{\sqrt{q}}, \frac{4}{\sqrt{q}}e^{-\frac{\sqrt{q}}{2}})$ -  
 549 differentially private (see Theorem 8). Note that the error sensitivity function of the histogram  
 550 mechanism is also the identity function, i.e,  $\tau_{\mathcal{M}_{hist}}(\beta') = \beta'$ . This follows from Claim 1, which  
 551 states that adding the same noise doesn't change the wasserstein distance between distributions.

552 Hence, by application of Theorem 5 and the fact that post-processing preserves differential privacy [1,  
 553 Proposition 2.1], we have that the bucketing followed by histogram mechanism is  $(qt, \frac{B}{2t}, 0)$ -accurate  
 554 and  $(\frac{1}{\sqrt{q}}, \frac{4}{\sqrt{q}}e^{-\frac{\sqrt{q}}{2}})$ -differentially private.

555 □

## 556 D.1 Case Study

557 We now show the application of the histogram for computing Max and support of a database in an  
 558 accurate and private manner.

### 559 D.1.1 Maximum

560 Let  $d$  be a distance metric on  $\mathcal{X}$  which has a well-defined order. We define the function  $f : \mathcal{N}^{\mathcal{X}} \rightarrow \mathcal{X}$ ,  
 561 as one which takes a database as input and outputs the largest element and take  $d_1$  as the  
 562 maximum  $d$ -distance each element of a histogram is to be modified to get the another,  $d_2$  as  $d$ -distance  
 563 in  $\mathcal{X}$  space. Note that  $f$  is a perfectly accurate mechanism for maximum. Also the error-sensitivity  
 564 function for  $f$  is identity, i.e,  $\tau_f(\beta') = \beta'$ . Because if the input distributions of histograms have a  
 565 wasserstein  $d_1$ -wasserstein  $\beta$ , implies that the maxima of these distributions are also distributions  
 566 which are atmost  $d_2$ -wasserstein  $\beta$  apart.

567 Now, using 5 and 6, we have that, the histogram mechanism for finding max is  $(\alpha_p, \beta_p, 0)$  accurate  
 568 and  $(\epsilon_p, \delta_p)$  private.

### 569 D.1.2 Support

570 Here we define the function  $f : \mathcal{N}^{\mathcal{X}} \rightarrow \mathcal{N}^{\mathcal{X}}$ , which takes a database as input and outputs the  
 571 database after removing duplicates and take  $d_1, d_2$  as wasserstein distance between inputs interpreting  
 572 the histograms as distributions. Again, this is a perfectly accurate and private mechanism for  
 573 maximum. Also the error-sensitivity function for  $f$  is identity, i.e,  $\tau_f(\beta') = \beta'$ . This is because, by  
 574 the same above analysis and if each element is modified by atmost  $\beta'$ , then the support elements also  
 575 changes by atmost  $\beta'$ .

576 Now, using 5 and 6, we have that, the histogram mechanism for finding max is  $(\alpha_p, \beta_p, 0)$  accurate  
 577 and  $(\epsilon_p, \delta_p)$  private.