
Privately Counting Unique Elements

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We study the problem of counting the number of unique elements in a dataset sub-
2 ject to the constraint of differential privacy. We consider the challenging setting
3 of person-level DP (a.k.a. user-level DP) where each person may contribute an
4 unbounded number of items and hence the sensitivity is unbounded.

5 Our approach is to compute a bounded-sensitivity version of this query, which
6 reduces to solving a max-flow problem. The sensitivity bound is optimized to
7 balance the noise we must add to privatize the answer against the error of the
8 approximation of the bounded-sensitivity query to the true number of unique ele-
9 ments.

10 1 Introduction

11 An elementary data analysis task is to count the number of unique elements occurring in a dataset.
12 The dataset may contain private data and even simple statistics can be combined to leak sensitive
13 information about people [Dinur and Nissim, 2003]. Our goal is to release (an approximation to)
14 this count in a way that ensures the privacy of the people who contributed their data. As a motivating
15 example, consider a collection of internet browsing histories, in which case the goal is to compute
16 the total number of websites that have been visited by at least one person.

17 Differential privacy (DP) [Dwork et al., 2006b] is a formal privacy standard. The simplest method
18 for ensuring DP is to add noise (from either a Laplace or Gaussian distribution) to the true answer,
19 where the scale of the noise corresponds to the sensitivity of the true answer – i.e., how much one
20 person’s data can change the true value.

21 If each person contributes a single element to the dataset, then the sensitivity of the number of
22 unique elements is one. However, a person may contribute multiple elements to the dataset and our
23 goal is to ensure privacy for all of these contributions simultaneously. That is, we seek to provide
24 person-level DP (a.k.a. user-level DP).

25 This is the problem we study: We have a dataset $D = (u_1, u_2, \dots, u_n)$ of person records. Each
26 person $i \in [n]$ contributes a finite dataset $u_i \in \Omega^*$, where Ω is some (possibly infinite) universe of
27 potential elements (e.g., all finite-length binary strings) and $\Omega^* := \bigcup_{\ell \in \mathbb{N}} \Omega^\ell$ denotes all subsets of
28 Ω of finite size. Informally, our goal is to compute the number of unique elements

$$\text{DC}(D) := \left| \bigcup_{i \in [n]} u_i \right| \tag{1}$$

29 in a way that preserves differential privacy. A priori, the sensitivity of this quantity is infinite, as a
30 single person can contribute an unbounded number of unique elements.

31 In particular, it is not possible to give a meaningful upper bound on the number of distinct elements
32 subject to differential privacy. However, it is possible to give a lower bound. Thus our formal goal

33 is to compute a high-confidence lower bound on the number of distinct elements that is as large as
 34 possible and which is computed in a differentially private manner.

35 1.1 Our Contributions

36 Given a dataset $D = (u_1, \dots, u_n) \in (\Omega^*)^n$ and an integer $\ell \geq 1$, we define

$$DC(D; \ell) := \max \left\{ \left| \bigcup_{i \in [n]} v_i \right| : \forall i \in [n] \ v_i \subset u_i \wedge |v_i| \leq \ell \right\}. \quad (2)$$

37 That is, $DC(D; \ell)$ is the number of distinct element if we restrict each person's contribution to ℓ
 38 elements. We take the maximum over all possible restrictions.

39 It is immediate that $DC(D; \ell) \leq DC(D)$ for all $\ell \geq 1$. Thus we obtain a lower bound on the true
 40 number of unique elements. The advantage of $DC(D; \ell)$ is that its sensitivity is bounded by ℓ and,
 41 hence, we can estimate it in a differentially private manner. Specifically,

$$\mathcal{M}_{\ell, \varepsilon}(D) := DC(D; \ell) + \text{Lap}(\ell/\varepsilon)$$

42 defines an ε -DP algorithm $\mathcal{M}_{\ell, \varepsilon} : (\Omega^*)^n \rightarrow \mathbb{R}$, where $\text{Lap}(b)$ denotes Laplace noise scaled to have
 43 mean 0 and variance $2b^2$. This forms the basis of our algorithm. Two challenges remain: Setting the
 44 sensitivity parameter ℓ and computing $DC(D; \ell)$ efficiently.

45 **Choosing the sensitivity parameter ℓ .** Any choice of $\ell \geq 1$ gives us a lower bound: $DC(D; \ell) \leq$
 46 $DC(D)$. Since $\forall D \lim_{\ell \rightarrow \infty} DC(D; \ell) = DC(D)$, this lower bound can be arbitrarily tight. How-
 47 ever, the larger ℓ is, the larger the sensitivity of $DC(D; \ell)$ is. That is, the noise we add scales linearly
 48 with ℓ .

49 Thus there is a bias-variance tradeoff in the choice of ℓ . To make this precise, suppose we want a
 50 lower bound on $DC(D)$ with confidence $1 - \beta \in [\frac{1}{2}, 1)$. We can obtain such a lower bound from
 51 $\mathcal{M}_{\ell}(D)$ using the cumulative distribution function (CDF) of the Laplace distribution:

$$\begin{aligned} \mathbb{P} \left[\underbrace{\mathcal{M}_{\ell, \varepsilon}(D) - \frac{\ell}{\varepsilon} \cdot \log \left(\frac{1}{2\beta} \right)}_{\text{lower bound}} \leq DC(D) \right] &= \mathbb{P} \left[\text{Lap}(\ell/\varepsilon) \leq \text{cdf}_{\text{Lap}(\ell/\varepsilon)}^{-1}(1 - \beta) + DC(D) - DC(D; \ell) \right] \\ &\geq \mathbb{P} \left[\text{Lap}(\ell/\varepsilon) \leq \text{cdf}_{\text{Lap}(\ell/\varepsilon)}^{-1}(1 - \beta) + 0 \right] = \underbrace{1 - \beta}_{\text{confidence}}. \end{aligned}$$

52 Thus, to obtain the tightest possible lower bound with confidence $1 - \beta$, we choose ℓ to maximize

$$q(D; \ell) := DC(D; \ell) - \frac{\ell}{\varepsilon} \cdot \log \left(\frac{1}{2\beta} \right).$$

53 We can use the exponential mechanism [McSherry and Talwar, 2007] to privately select ℓ that ap-
 54 proximately maximizes $q(D; \ell)$. However, directly applying the exponential mechanism is problem-
 55 atic because each score has a different sensitivity – the sensitivity of $q(\cdot; \ell)$ is ℓ . Instead, we apply
 56 the generalized exponential mechanism of Raskhodnikova and Smith [2015] (see Algorithm 3).

57 Our main algorithm attains the following guarantees.

58 **Theorem 1.1** (Theoretical Guarantees of Our Algorithm). *Let $\varepsilon > 0$ and $\beta \in (0, \frac{1}{2})$ and $\ell_{\max} \in \mathbb{N}$.
 59 Define $\mathcal{M} : (\Omega^*)^n \rightarrow \mathbb{N} \times \mathbb{R}$ to be $\mathcal{M}(D) = \text{DPDISTINCTCOUNT}(D; \ell_{\max}, \varepsilon, \beta)$ from Algorithm 1.
 60 Then \mathcal{M} satisfies all of the following properties.*

- 61 • **Privacy:** \mathcal{M} is ε -differentially private.
- 62 • **Lower bound:** For all $D \in (\Omega^*)^n$,

$$\mathbb{P}_{(\hat{\ell}, \hat{v}) \leftarrow \mathcal{M}(D)} [\hat{v} \leq DC(D)] \geq 1 - \beta. \quad (3)$$

63 • **Upper bound** For all $D \in (\Omega^*)^n$,

$$\mathbb{P}_{(\hat{\ell}, \hat{\nu}) \leftarrow \mathcal{M}(D)} \left[\hat{\nu} \geq \max_{\ell \in [\ell_{\max}]} \text{DC}(D; \ell) - \frac{10\ell + 18\ell_A^*}{\varepsilon} \log \left(\frac{\ell_{\max}}{\beta} \right) \right] \geq 1 - 2\beta, \quad (4)$$

64 where $\ell_A^* = \arg \max_{\ell \in [\ell_{\max}]} \text{DC}(D; \ell) - \frac{\ell}{\varepsilon} \log \left(\frac{1}{2\beta} \right)$.

65 • **Computational efficiency:** $\mathcal{M}(D)$ has running time $O(|D|^{1.5} \cdot \ell_{\max}^2)$, where $|D| :=$
66 $\sum_i |u_i|$.

67 In particular, if $D = (u_1, \dots, u_n) \in (\Omega^*)^n$ satisfies $\max_{i \in [n]} |u_i| \leq \ell_* \leq \ell_{\max}$, then combining
68 the upper and lower bounds of Theorem 1.1 gives

$$\mathbb{P}_{(\hat{\ell}, \hat{\nu}) \leftarrow \mathcal{M}(D)} \left[\text{DC}(D) \geq \hat{\nu} \geq \text{DC}(D) - \frac{28\ell_*}{\varepsilon} \log \left(\frac{\ell_{\max}}{\beta} \right) \right] \geq 1 - 3\beta. \quad (5)$$

69 In addition to proving the above theoretical guarantees, we perform an experimental evaluation of
70 our algorithm.

Algorithm 1 Distinct Count Algorithm

```

1: procedure SENSITIVEDISTINCTCOUNT( $D = (u_1, \dots, u_n) \in (\Omega^*)^n; \ell \in \mathbb{N}$ ) ▷ DC( $D; \ell$ )
2:   Let  $U_\ell = \bigcup_{i \in [n]} (\{i\} \times [\min\{\ell, |u_i|\}]) \subset [n] \times [\ell]$ .
3:   Let  $V = \bigcup_{i \in [n]} u_i \subset \Omega$ .
4:   Define  $E_\ell \subseteq U \times V$  by  $((i, j), v) \in E \iff v \in u_i$ .
5:   Let  $G_\ell$  be a bipartite graph with vertices partitioned into  $U_\ell$  and  $V$  and edges  $E_\ell$ .
6:    $m_\ell \leftarrow \text{MAXIMUMMATCHINGSIZE}(G)$ . ▷ [Hopcroft and Karp, 1973, Karzanov, 1973]
7:   return  $m_\ell \in \mathbb{N}$ 
8: end procedure
9: procedure DPDISTINCTCOUNT( $D = (u_1, \dots, u_n) \in (\Omega^*)^n; \ell_{\max} \in \mathbb{N}, \varepsilon > 0, \beta \in (0, \frac{1}{2})$ )
10:  for  $\ell \in [\ell_{\max}]$  do
11:    Define  $q_\ell(D) := \text{SENSITIVEDISTINCTCOUNT}(D; \ell) - \frac{2\ell}{\varepsilon} \cdot \log \left( \frac{1}{2\beta} \right)$ .
12:  end for
13:   $\hat{\ell} \leftarrow \text{GEM}(D; \{q_\ell\}_{\ell \in [\ell_{\max}]}, \{\ell\}_{\ell \in [\ell_{\max}]}, \varepsilon/2, \beta)$ . ▷ Algorithm 3
14:   $\hat{\nu} \leftarrow q_{\hat{\ell}}(D) + \text{Lap} \left( 2\hat{\ell}/\varepsilon \right)$ .
15:  return  $(\hat{\ell}, \hat{\nu}) \in [\ell_{\max}] \times \mathbb{R}$ .
16: end procedure

```

71 **Efficient computation.** The main computational task for our algorithm is to compute $\text{DC}(D; \ell)$.
72 By definition (2), this is an optimization problem. For each person $i \in [n]$, we must select a subset
73 v_i of that person's data u_i of size at most ℓ so as to maximize the size of the union of the subsets
74 $\left| \bigcup_{i \in [n]} v_i \right|$.

75 We can view the dataset $D = (u_1, \dots, u_n) \in (\Omega^*)^n$ as a bipartite graph. On one side we have the n
76 people and on the other side we have the elements of the data universe Ω .¹ There is an edge between
77 $i \in [n]$ and $x \in \Omega$ if and only if $x \in u_i$.

78 We can reduce computing $\text{DC}(D; \ell)$ to a max-flow problem: Each edge in the bipartite graph has
79 capacity one. We add a source vertex s which is connected to each person $i \in [n]$ by an edge with
80 capacity ℓ . Finally we add a sink t that is connected to each $x \in \Omega$ by an edge with capacity 1. The
81 max flow through this graph is precisely $\text{DC}(D; \ell)$.

82 Alternatively, we can reduce computing $\text{DC}(D; \ell)$ to bipartite maximum matching. For $\ell = 1$,
83 $\text{DC}(D; 1)$ is exactly the maximum cardinality of a matching in the bipartite graph described above.

¹The data universe Ω may be infinite, but we can restrict the computation to the finite set $\bigcup_{i \in [n]} u_i$. Thus there are at most $n + \text{DC}(D) \leq n + |D|$ item vertices in the graph.

84 For $\ell \geq 2$, we simply create ℓ copies of each person vertex $i \in [n]$ and then $\text{DC}(D; \ell)$ is the
 85 maximum cardinality of a matching in this new bipartite graph.²

86 Using this reduction, standard algorithms for bipartite maximum matching [Hopcroft and Karp,
 87 1973, Karzanov, 1973] allow us to compute $\text{DC}(D; \ell)$ with $O(|D|^{1.5} \cdot \ell)$ operations. We must
 88 repeat this computation for each $\ell \in [\ell_{\max}]$.

Algorithm 2 Linear-Time Approximate Distinct Count Algorithm

```

1: procedure APPROXDPDISTINCTCOUNT( $D=(u_1, \dots, u_n) \in (\Omega^*)^n; \ell_{\max} \in \mathbb{N}, \varepsilon > 0, \beta \in (0, \frac{1}{2})$ )
2:    $S \leftarrow \emptyset$ .
3:   for  $\ell \in [\ell_{\max}]$  do
4:     for  $i \in [n]$  with  $u_i \setminus S \neq \emptyset$  do
5:       Choose lexicographically first  $v \in u_i \setminus S$ . ▷ Match  $(i, \ell)$  to  $v$ .
6:       Update  $S \leftarrow S \cup \{v\}$ .
7:     end for
8:     Define  $q_\ell(D) := |S| - \frac{2\ell}{\varepsilon} \cdot \log\left(\frac{1}{2\beta}\right)$ . ▷ This loop computes  $\{q_\ell(D)\}_{\ell \in [\ell_{\max}]}$ .
9:   end for
10:   $\hat{\ell} \leftarrow \text{GEM}(D; \{q_\ell\}_{\ell \in [\ell_{\max}]}, \{\ell\}_{\ell \in [\ell_{\max}]}, \varepsilon/2, \beta)$ . ▷ Algorithm 3
11:   $\hat{v} \leftarrow q_{\hat{\ell}}(D) + \text{Lap}\left(2\hat{\ell}/\varepsilon\right)$ .
12:  return  $(\hat{\ell}, \hat{v}) \in [\ell_{\max}] \times \mathbb{R}$ .
13: end procedure

```

89 **Linear-time algorithm.** Our algorithm above is polynomial-time. However, for many applica-
 90 tions the dataset size $|D|$ is enormous. Thus we also propose a linear-time variant of our algorithm.
 91 However, we must trade accuracy for efficiency.

92 There are two key ideas that differentiate our linear-time algorithm (Algorithm 2) from our first
 93 algorithm (Algorithm 1) above: First, we compute a maximal bipartite matching instead of a
 94 maximum bipartite matching. This can be done using a linear-time greedy algorithm and gives a
 95 2-approximation to the maximal matching. (Experimentally we find that the approximation is better
 96 than a factor of 2.) Second, rather than repeating the computation from scratch for each $\ell \in [\ell_{\max}]$,
 97 we incrementally update our a maximal matching while increasing ℓ . The main challenge here is
 98 ensuring that the approximation to $\text{DC}(D; \ell)$ has low sensitivity – i.e., we must ensure that our
 99 approximation algorithm doesn't inflate the sensitivity. Note that $\text{DC}(D; \ell)$ having low sensitivity
 100 does not automatically ensure that the approximation has low sensitivity.

101 **Theorem 1.2** (Theoretical Guarantees of Our Linear-Time Algorithm). *Let $\varepsilon \geq 0$ and*
 102 *$\beta \in (0, \frac{1}{2})$ and $\ell_{\max} \in \mathbb{N}$. Define $\mathcal{M} : (\Omega^*)^n \rightarrow \mathbb{N} \times \mathbb{R}$ to be $\widehat{\mathcal{M}}(D) =$*
 103 *$\text{APPROXDPDISTINCTCOUNT}(D; \ell_{\max}, \varepsilon, \beta)$ from Algorithm 2. Then $\widehat{\mathcal{M}}$ satisfies all of the fol-*
 104 *lowing properties.*

105 • **Privacy:** $\widehat{\mathcal{M}}$ is ε -differentially private.

106 • **Lower bound:** For all $D \in (\Omega^*)^n$,

$$\mathbb{P}_{(\hat{\ell}, \hat{v}) \leftarrow \widehat{\mathcal{M}}(D)} [\hat{v} \leq \text{DC}(D)] \geq 1 - \beta. \quad (6)$$

107 • **Upper bound:** If $D = (u_1, \dots, u_n) \in (\Omega^*)^n$ satisfies $\max_{i \in [n]} |u_i| \leq \ell_* \leq \ell_{\max}$, then

$$\mathbb{P}_{(\hat{\ell}, \hat{v}) \leftarrow \widehat{\mathcal{M}}(D)} \left[\hat{v} \geq \frac{1}{2} \text{DC}(D) - O\left(\frac{\ell_*}{\varepsilon} \log\left(\frac{\ell_{\max}}{\beta}\right)\right) \right] \geq 1 - 2\beta. \quad (7)$$

108 • **Computational efficiency:** $\mathcal{M}(D)$ has running time $O(|D|)$, where $|D| := \sum_i |u_i|$.

²We need only create $\min\{\ell, |u_i|\}$ copies of the person $i \in [n]$. Thus the number of person vertices is at most $\min\{n\ell, |D|\}$.

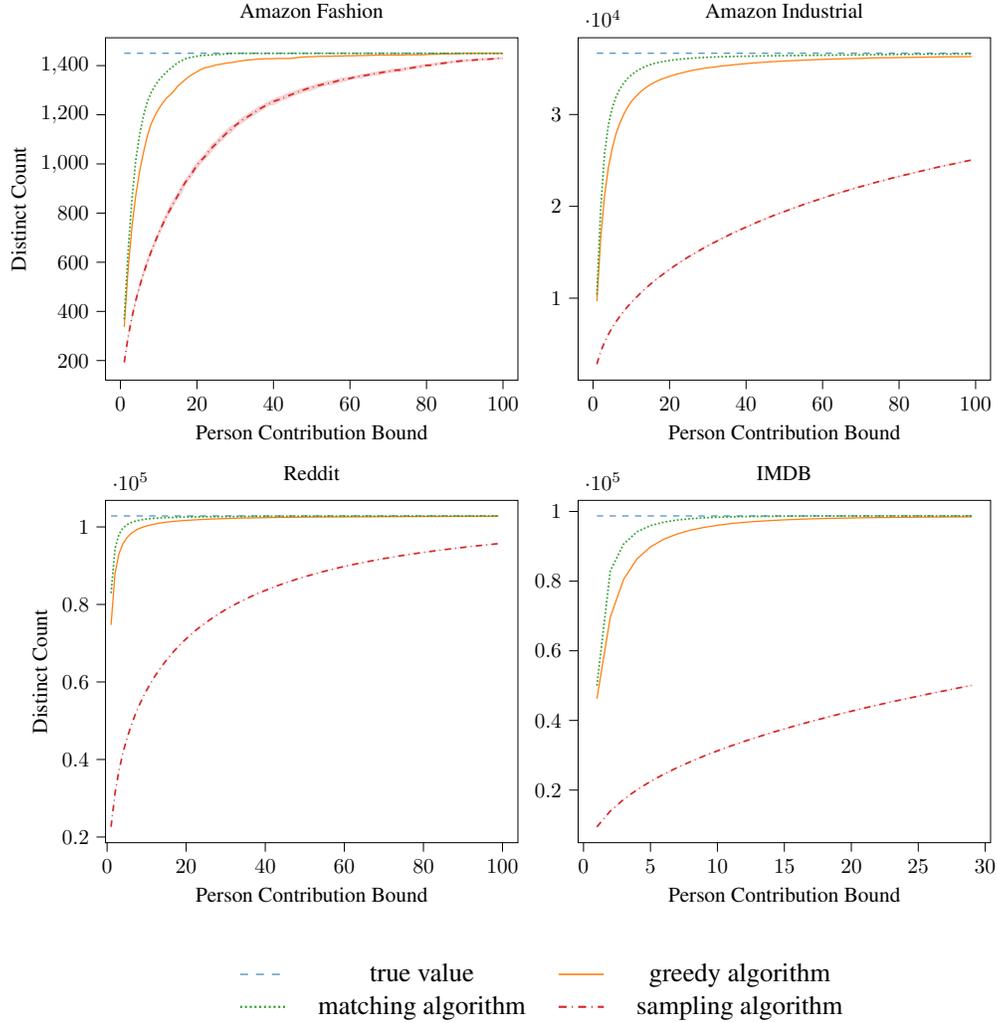


Figure 1: Performance of different algorithms estimating distinct count assuming that each person can contribute at most ℓ elements.

109 **2 Related Work**

110 Counting the number of distinct elements in a collection is one of the most fundamental operations.
 111 Hence, unsurprisingly, the problem of computing the number of unique elements in a differentially
 112 private way has been extensively investigated.

113 In the case where we assume each person contributes only one element (a.k.a. event-level privacy),
 114 the number of distinct elements has sensitivity 1 and, hence, we can simply use Laplace (or Gaus-
 115 sian) noise addition to release. However, it may not be possible to compute the number of distinct
 116 elements exactly (e.g. in the local model of DP [Kasiviswanathan et al., 2011]).

117 Most efforts have been focused on creating differentially private approximation schemes for count-
 118 ing distinct elements. Desfontaines et al. [2019] proved that a number of existing approximate
 119 algorithms allow an attacker to test whether a particular individual is in the collection; therefore,
 120 creation of a differentially private scheme requires care. Nonetheless, Smith et al. [2020] proved
 121 that Flajolet-Martin Sketch is private by itself and Dickens et al. [2022] proved that several other
 122 cardinality estimators can be tweaked to make them private. In case of local and shuffle models
 123 the only known results are communication complexity bounds [Chen et al., 2021]. Counting unique
 124 elements has been considered in the streaming setting [Dwork et al., 2010, Ghazi et al., 2023].

125 A closely related problem is that of identifying as many elements as possible (rather than just count-
 126 ing them); this is known as “partition selection,” “set union,” and “key selection” [Swanberg et al.,
 127 2023, Desfontaines et al., 2022, Korolova et al., 2009, Carvalho et al., 2022, Rivera Cardoso and
 128 Rogers, 2022, Gopi et al., 2020, Zhang et al., 2023]. Note that, by design, DP prevents us from
 129 identifying elements that only appear once in the dataset, or only a few times. Thus we can only
 130 output items that appear frequently.

131 For our problem of counting the number of unique elements under person-level/user-level privacy,
 132 the only known algorithm is the algorithm where each user independently samples a subset of their
 133 elements to reduce the sensitivity. We use this as a baseline in our experiments and show that our
 134 algorithm outperforms it.

135 3 Technical Background on Differential Privacy

136 For detailed background on differential privacy, see the survey by Vadhan [2017] or the book by
 137 Dwork and Roth [2014]. We briefly define pure DP and some basic mechanisms and results.

Algorithm 3 Generalized Exponential Mechanism [Raskhodnikova and Smith, 2015]

1: **procedure** GEM($D \in \mathcal{X}^*$; $q_i : \mathcal{X}^* \rightarrow \mathbb{R}$ for $i \in [m]$, $\Delta_i > 0$ for $i \in [m]$, $\varepsilon > 0$, $\beta > 0$)
 2: **Require:** q_i has sensitivity $\sup_{x, x' \in \mathcal{X}^*} |q(x) - q(x')| \leq \Delta_i$ for all $i \in [m]$.

3: Let $t = \frac{2}{\varepsilon} \log \left(\frac{m}{\beta} \right)$.

4: **for** $i \in [m]$ **do**

5: $s_i \leftarrow \min_{j \in [m]} \frac{(q_i(D) - t\Delta_i) - (q_j(D) - t\Delta_j)}{\Delta_i + \Delta_j}$.

6: **end for**

7: Sample $\hat{i} \in [m]$ from the Exponential Mechanism using the normalized scores s_i ; i.e.,

$$\forall i \in [m] \quad \mathbb{P}[\hat{i} = i] = \frac{\exp\left(\frac{1}{2}\varepsilon s_i\right)}{\sum_{k \in [m]} \exp\left(\frac{1}{2}\varepsilon s_k\right)}.$$

8: **return** $\hat{i} \in [m]$.

9: **end procedure**

138 **Definition 3.1** (Differential Privacy (DP) [Dwork et al., 2006b]). A randomized algorithm $M : \mathcal{X}^* \rightarrow \mathcal{Y}$ satisfies ε -DP if, for all inputs $D, D' \in \mathcal{X}^*$ differing only by the addition or removal of an
 139 element and for all measurable $S \subset \mathcal{Y}$, we have $\mathbb{P}[M(D) \in S] \leq e^\varepsilon \cdot \mathbb{P}[M(D') \in S]$.
 140

141 We refer to pairs of inputs that differ only by the addition or removal of one person’s data as *neigh-*
 142 *boring*. Note that it is common to also consider replacement of one person’s data; for simplicity,
 143 we do not do this. We remark that there are also variants of DP such as approximate DP [Dwork
 144 et al., 2006a] and concentrated DP [Dwork and Rothblum, 2016, Bun and Steinke, 2016], which
 145 quantitatively relax the definition, but these are not relevant in our application. A key property of
 146 DP is that it composes and is invariant under postprocessing.

147 **Lemma 3.2** (Composition & Postprocessing). Let $M_1 : \mathcal{X}^* \rightarrow \mathcal{Y}$ be ε_1 -DP. Let $M_2 : \mathcal{X}^* \times \mathcal{Y} \rightarrow \mathcal{Z}$
 148 be such that, for all $y \in \mathcal{Y}$, the restriction $M(\cdot, y) : \mathcal{X}^* \rightarrow \mathcal{Z}$ is ε_2 -DP. Define $M_{12} : \mathcal{X}^* \rightarrow \mathcal{Z}$ by
 149 $M_{12}(D) = M_2(D, M_1(D))$. Then M_{12} is $(\varepsilon_1 + \varepsilon_2)$ -DP.

150 A basic DP tool is the Laplace mechanism [Dwork et al., 2006b]. Note that we could also use the
 151 *discrete* Laplace mechanism [Ghosh et al., 2009, Canonne et al., 2020].

152 **Lemma 3.3** (Laplace Mechanism). Let $q : \mathcal{X}^* \rightarrow \mathbb{R}$. We say q has sensitivity Δ if $|q(D) - q(D')| \leq \Delta$
 153 for all neighboring $D, D' \in \mathcal{X}^*$. Define $M : \mathcal{X}^* \rightarrow \mathbb{R}$ by $M(D) = q(D) + \text{Lap}(\Delta/\varepsilon)$,
 154 where $\text{Lap}(b)$ denotes laplace noise with mean 0 and variance $2b^2$ – i.e., $\mathbb{P}_{\xi \leftarrow \text{Lap}(b)}[\xi > t] =$

155 $\mathbb{P}_{\xi \leftarrow \text{Lap}(b)}[\xi < -t] = \frac{1}{2} \exp\left(-\frac{t}{b}\right)$ for all $t > 0$. Then M is ε -DP.

156 Another fundamental tool for DP is the exponential mechanism [McSherry and Talwar, 2007]. It
 157 selects the approximately best option from among a set of options, where each option i has a quality
 158 function q_i with sensitivity Δ . The following result generalizes the exponential mechanism by
 159 allowing each of the quality functions to have a different sensitivity.

Data Set	Size		Words per Person			Vocabulary Size
	People	Records	Min	Median	Max	
Amazon Fashion	404	8533	1	14.0	139	1450
Amazon Industrial and Scientific	11041	1446031	0	86	2059	36665
Reddit	223388	7117494	0	18.0	1724	102835
IMDB	50000	6688844	5	110.0	925	98726

Table 1: Data sets details.

160 **Theorem 3.4** (Generalized Exponential Mechanism [Raskhodnikova and Smith, 2015, Theorem
161 1.4]). For each $i \in [m]$, let $q_i : \mathcal{X}^* \rightarrow \mathbb{R}$ be a query with sensitivity Δ_i . Let $\varepsilon, \beta > 0$. The
162 generalized exponential mechanism (GEM($;$ $\{q_i\}_{i \in [m]}, \{\Delta_i\}_{i \in [m]}, \varepsilon, \beta$) in Algorithm 3) is ε -DP
163 and has the following utility guarantee. For all $D \in \mathcal{X}^*$, we have

$$\mathbb{P}_{\hat{i} \leftarrow \text{GEM}(D; \{q_i\}_{i \in [m]}, \{\Delta_i\}_{i \in [m]}, \varepsilon, \beta)} \left[q_{\hat{i}}(D) \geq \max_{j \in [m]} q_j(D) - \Delta_j \cdot \frac{4}{\varepsilon} \log \left(\frac{m}{\beta} \right) \right] \geq 1 - \beta.$$

164 4 Experimental Results

165 We empirically validate the performance of our algorithms using data sets of various sizes from
166 different text domains. We focus on the problem of computing vocabulary size with person-level
167 DP. Section 4.1 describes the data sets and Section 4.2 discusses the algorithms we compare.

168 4.1 Datasets

169 We used four publicly available datasets to assess the accuracy of our algorithms compared to base-
170 lines. Two small datasets were used: Amazon Fashion 5-core [Ni et al., 2019] (reviews of fashion
171 products on Amazon) and Amazon Industrial and Scientific 5-core [Ni et al., 2019] (reviews of in-
172 dustrial and scientific products on Amazon). Two large data sets were also used: Reddit [Shen,
173 2020] (a data set of posts collected from r/AskReddit) and IMDb [N, 2020, Maas et al., 2011] (a set
174 of movie reviews scraped from IMDb). See details of the datasets in Table 1.

175 4.2 Comparisons

176 Computing the number of distinct elements using a differentially private mechanism involves two
177 steps: selecting a contribution bound (ℓ in our algorithms) and counting the number of distinct
178 elements in a way that restricts each person to only contribute the given number of elements.

179 **Selection:** We examine three algorithms for determining the contribution limit:

- 180 1. Choosing the true maximum person contribution (due to computational restrictions this was
181 only computed for Amazon Fashion data set).
- 182 2. Choosing the 90th percentile of person contributions.
- 183 3. Choosing the person contribution that maximizes the utility function $q_\ell(D) = \text{DC}(D; \ell) -$
184 $\frac{\ell}{\varepsilon} \log(\frac{1}{2\beta})$, where $\varepsilon = 1$, and $\beta = 0.001$.
- 185 4. Choosing the person contribution that maximizes the utility function using generalized
186 exponential mechanism with $\varepsilon = 1$.

187 Note that only the last option is differentially private, but we consider the other comparison points
188 nonetheless.

189 **Counting:** We also consider three algorithms for estimating the number of distinct elements for a
190 given sensitivity bound ℓ :

- 191 1. For each person, we independently sample ℓ elements and count the number of distinct
192 elements in the union of the samples.

Selection	Counting	Person Contribution Bound			Distinct Count		
		10th PC	Median	90th PC	10th PC	Median	90th PC
Max Contrib	DP Sampling	–	139	–	1196.8	1407.5	1649.1
Max Contrib	DP Greedy	–	139	–	1174.2	1439.2	1646.5
Max Contrib	DP Matching	–	139	–	1222.2	1460.9	1631.0
90th PC Contrib	DP Sampling	–	48	–	1225.4	1296.2	1377.9
90th PC Contrib	DP Greedy	–	48	–	1367.0	1432.6	1516.3
90th PC Contrib	DP Matching	–	48	–	1365.3	1444.7	1524.8
Max Utility	Sampling	–	41	–	1247.0	1259.0	1270.0
Max Utility	Greedy	–	20	–	–	1376	–
Max Utility	Matching	–	17	–	–	1428	–
DP Max Utility	Sampling	8.9	16.0	28.0	661.6	892.5	1124.5
DP Max Utility	Greedy	8.0	11.0	17.0	1148.0	1241.0	1348.0
DP Max Utility	Matching	7.0	9.0	14.0	1252.0	1317.0	1400.0
DP Max Utility	DP Sampling	9.0	16.0	27.1	702.4	899.1	1145.1
DP Max Utility	DP Greedy	8.0	10.0	19.0	1128.5	1224.4	1370.8
DP Max Utility	DP Matching	6.9	9.0	13.1	1220.6	1319.1	1394.2

Table 2: Amazon Fashion: the comparison is for $\ell_{\max} = 100$.

Selection	Counting	Person Contribution Bound			Distinct Count		
		10th PC	Median	90th PC	10th PC	Median	90th PC
90th PC Contrib	DP Sampling	–	297	–	32458.1	32943.8	33452.6
90th PC Contrib	DP Greedy	–	297	–	36270.3	36669.5	37019.0
90th PC Contrib	DP Matching	–	297	–	36236.2	36651.7	37102.7
Max Utility	Sampling	–	99	–	24967.0	25039.0	25121.2
Max Utility	Greedy	–	79	–	–	36246	–
Max Utility	Matching	–	42	–	–	36364	–
DP Max Utility	Sampling	85.9	96.0	99.0	23852.8	24739.0	25049.8
DP Max Utility	Greedy	34.0	49.0	66.1	35393.0	35839.0	36116.9
DP Max Utility	Matching	22.9	30.5	43.2	36026.8	36243.5	36371.2
DP Max Utility	DP Sampling	87.0	95.0	99.0	23997.6	24701.1	25067.7
DP Max Utility	DP Greedy	32.9	47.5	68.0	35336.6	35776.2	36136.6
DP Max Utility	DP Matching	22.0	28.0	38.0	35970.5	36198.9	36326.7

Table 3: Amazon Industrial and Scientific: the comparison is for $\ell_{\max} = 100$.

193 2. The linear-time greedy algorithm (Algorithm 2) with $\varepsilon = 1$ and $\beta = 0.001$.

194 3. The matching-based algorithm (Algorithm 1) with $\varepsilon = 1$ and $\beta = 0.001$.

195 All of these can be converted into DP algorithms by adding Laplace noise to the result.

196 In all our datasets “true maximum person contribution” and “90th percentile of person contributions”
 197 output bounds that are much larger than necessary to obtain true distinct count; hence, we only
 198 consider DP versions of the estimation algorithm for these selection algorithms.

199 4.3 Results

200 Figure 1 shows the dependency of the result on the contribution bound for each of the algorithms for
 201 computing the number of distinct elements with fixed person contribution. It is clear that matching
 202 and greedy algorithms vastly outperform the sampling approach that is currently used in practice.

203 Tables 2 to 5 show the performance of algorithms for selecting optimal person contribution bounds
 204 on different data sets. For all bound selection algorithms and all data sets, the sampling approach to

Selection	Counting	Person Contribution Bound			Distinct Count		
		10th PC	Median	90th PC	10th PC	Median	90th PC
90th PC Contrib	DP Sampling	–	75	–	92480.7	92654.8	92812.1
90th PC Contrib	DP Greedy	–	75	–	102544.8	102665.7	102817.7
90th PC Contrib	DP Matching	–	75	–	102651.1	102784.1	102907.8
Max Utility	Sampling	–	99	–	95606.9	95692.0	95750.3
Max Utility	Greedy	–	52	–	–	102543	–
Max Utility	Matching	–	32	–	–	102685	–
DP Max Utility	Sampling	89.0	96.0	99.0	94549.9	95394.5	95656.5
DP Max Utility	Greedy	26.0	33.0	50.0	102015.0	102253.0	102527.0
DP Max Utility	Matching	14.0	18.5	30.0	102357.0	102501.5	102671.0
DP Max Utility	DP Sampling	88.8	96.0	99.0	94665.2	95375.5	95693.5
DP Max Utility	DP Greedy	27.0	34.0	53.0	102053.2	102289.6	102531.2
DP Max Utility	DP Matching	14.9	18.5	28.0	102379.7	102512.6	102643.9

Table 4: Reddit: the comparison is for $\ell_{\max} = 100$.

Selection	Counting	Person Contribution Bound			Distinct Count		
		10th PC	Median	90th PC	10th PC	Median	90th PC
90th PC Contrib	DP Sampling	–	238	–	95264.5	95593.5	95966.1
90th PC Contrib	DP Greedy	–	238	–	98411.0	98734.0	99120.0
90th PC Contrib	DP Matching	–	238	–	98354.2	98729.4	99164.2
Max Utility	Sampling	–	29	–	49907.8	50036.5	50195.3
Max Utility	Greedy	–	29	–	–	98459	–
Max Utility	Matching	–	19	–	–	98712	–
DP Max Utility	Sampling	29.0	29.0	29.0	49899.6	50070.5	50220.9
DP Max Utility	Greedy	22.0	25.0	29.0	98244.0	98364.0	98459.0
DP Max Utility	Matching	13.0	16.0	21.0	98586.0	98674.0	98721.0
DP Max Utility	DP Sampling	29.0	29.0	29.0	49924.2	50053.7	50211.9
DP Max Utility	DP Greedy	20.0	26.0	29.0	98126.7	98369.6	98451.8
DP Max Utility	DP Matching	12.0	16.0	21.0	98555.6	98670.4	98726.8

Table 5: IMDB: the comparison is for $\ell_{\max} = 30$.

205 estimating the distinct count performs much worse than the greedy and matching-based approaches.
206 The greedy approach performs worse than the matching-based approach, but the difference is about
207 10% for Amazon Fashion and is almost negligible for other data sets since they are much larger. As
208 for the matching-based algorithm, it performs as follows on all the data sets:

- 209 1. The algorithm that uses the bound equal to the maximal person contribution overestimates
210 the actual necessary bound. Therefore, we only consider the DP algorithms for counts
211 estimation. It is easy to see that while the median of the estimation is close to the actual
212 distinct count, the amount of noise is somewhat large.
- 213 2. The algorithm that uses the bound equal to the 99th percentile of person contributions
214 also overestimates the necessary bound and behaves similarly to the one we just described
215 (though the spread of the noise is a bit smaller).
- 216 3. The algorithms that optimize the utility function are considered: one non-private and one
217 private. The non-private algorithm with non-private estimation gives the answer that is
218 very close to the true number of distinct elements. The private algorithm with non-private
219 estimation gives the answer that is worse, but not too much. Finally, the private algorithm
220 with the private estimation gives answers very similar to the results of the non-private
221 estimation.

222 References

- 223 Mark Bun and Thomas Steinke. Concentrated differential privacy: Simplifications, extensions, and
224 lower bounds. In *Theory of Cryptography: 14th International Conference, TCC 2016-B, Beijing,*
225 *China, October 31-November 3, 2016, Proceedings, Part I*, pages 635–658. Springer, 2016. URL
226 <https://arxiv.org/abs/1605.02065>.
- 227 Clément L Canonne, Gautam Kamath, and Thomas Steinke. The discrete gaussian for differential
228 privacy. *Advances in Neural Information Processing Systems*, 33:15676–15688, 2020. URL
229 <https://arxiv.org/abs/2004.00010>.
- 230 Ricardo Silva Carvalho, Ke Wang, and Lovedeep Singh Gondara. Incorporating item frequency
231 for differentially private set union. In *Thirty-Sixth AAAI Conference on Artificial Intelligence,*
232 *AAAI 2022, Thirty-Fourth Conference on Innovative Applications of Artificial Intelligence, IAAI*
233 *2022, The Twelfth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022*
234 *Virtual Event, February 22 - March 1, 2022*, pages 9504–9511. AAAI Press, 2022. URL <https://ojs.aaai.org/index.php/AAAI/article/view/21183>.
- 235
- 236 Lijie Chen, Badih Ghazi, Ravi Kumar, and Pasin Manurangsi. On distributed differential pri-
237 vacy and counting distinct elements. In James R. Lee, editor, *12th Innovations in Theoretical*
238 *Computer Science Conference, ITCS 2021, January 6-8, 2021, Virtual Conference*, volume 185
239 of *LIPICs*, pages 56:1–56:18. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2021. doi:
240 10.4230/LIPICs.ITCS.2021.56. URL <https://doi.org/10.4230/LIPICs.ITCS.2021.56>.
- 241 Damien Desfontaines, Andreas Lochbihler, and David A. Basin. Cardinality estimators do not
242 preserve privacy. *Proc. Priv. Enhancing Technol.*, 2019(2):26–46, 2019. doi: 10.2478/
243 popets-2019-0018. URL <https://doi.org/10.2478/popets-2019-0018>.
- 244 Damien Desfontaines, James Voss, Bryant Gipson, and Chinmoy Mandayam. Differentially
245 private partition selection. *Proc. Priv. Enhancing Technol.*, 2022(1):339–352, 2022. doi:
246 10.2478/popets-2022-0017. URL <https://doi.org/10.2478/popets-2022-0017>.
- 247 Charlie Dickens, Justin Thaler, and Daniel Ting. Order-invariant cardinality estimators are differen-
248 tially private. In *NeurIPS*, 2022. URL [http://papers.nips.cc/paper_files/paper/2022/](http://papers.nips.cc/paper_files/paper/2022/hash/623307df18da128262aaf394cdcfb235-Abstract-Conference.html)
249 [hash/623307df18da128262aaf394cdcfb235-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/623307df18da128262aaf394cdcfb235-Abstract-Conference.html).
- 250 Irit Dinur and Kobbi Nissim. Revealing information while preserving privacy. In *Proceedings of the*
251 *twenty-second ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*,
252 pages 202–210, 2003.
- 253 Cynthia Dwork and Aaron Roth. The algorithmic foundations of differential privacy. *Foundations*
254 *and Trends® in Theoretical Computer Science*, 9(3–4):211–407, 2014. URL [https://www.](https://www.cis.upenn.edu/~aaroht/Papers/privacybook.pdf)
255 [cis.upenn.edu/~aaroht/Papers/privacybook.pdf](https://www.cis.upenn.edu/~aaroht/Papers/privacybook.pdf).
- 256 Cynthia Dwork and Guy N Rothblum. Concentrated differential privacy. *arXiv preprint*
257 *arXiv:1603.01887*, 2016. URL <https://arxiv.org/abs/1603.01887>.
- 258 Cynthia Dwork, Krishnaram Kenthapadi, Frank McSherry, Ilya Mironov, and Moni Naor. Our data,
259 ourselves: Privacy via distributed noise generation. In *Advances in Cryptology-EUROCRYPT*
260 *2006: 24th Annual International Conference on the Theory and Applications of Cryptographic*
261 *Techniques, St. Petersburg, Russia, May 28-June 1, 2006. Proceedings 25*, pages 486–503.
262 Springer, 2006a.
- 263 Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity
264 in private data analysis. In *Theory of Cryptography: Third Theory of Cryptography Conference,*
265 *TCC 2006, New York, NY, USA, March 4-7, 2006. Proceedings 3*, pages 265–284. Springer, 2006b.
- 266 Cynthia Dwork, Moni Naor, Toniann Pitassi, Guy N. Rothblum, and Sergey Yekhanin. Pan-private
267 streaming algorithms. In Andrew Chi-Chih Yao, editor, *Innovations in Computer Science -*
268 *ICS 2010, Tsinghua University, Beijing, China, January 5-7, 2010. Proceedings*, pages 66–
269 80. Tsinghua University Press, 2010. URL [http://conference.iis.tsinghua.edu.cn/](http://conference.iis.tsinghua.edu.cn/ICS2010/content/papers/6.html)
270 [ICS2010/content/papers/6.html](http://conference.iis.tsinghua.edu.cn/ICS2010/content/papers/6.html).

- 271 Badih Ghazi, Ravi Kumar, Jelani Nelson, and Pasin Manurangsi. Private counting of distinct and
 272 k-occurring items in time windows. In Yael Tauman Kalai, editor, *14th Innovations in Theo-*
 273 *retical Computer Science Conference, ITCS 2023, January 10-13, 2023, MIT, Cambridge, Mas-*
 274 *sachusetts, USA*, volume 251 of *LIPICs*, pages 55:1–55:24. Schloss Dagstuhl - Leibniz-Zentrum
 275 für Informatik, 2023. doi: 10.4230/LIPICs.ITCS.2023.55. URL [https://doi.org/10.4230/](https://doi.org/10.4230/LIPICs.ITCS.2023.55)
 276 [LIPICs.ITCS.2023.55](https://doi.org/10.4230/LIPICs.ITCS.2023.55).
- 277 Arpita Ghosh, Tim Roughgarden, and Mukund Sundararajan. Universally utility-maximizing pri-
 278 vacy mechanisms. In *Proceedings of the forty-first annual ACM symposium on Theory of com-*
 279 *puting*, pages 351–360, 2009. URL <https://arxiv.org/abs/0811.2841>.
- 280 Sivakanth Gopi, Pankaj Gulhane, Janardhan Kulkarni, Judy Hanwen Shen, Milad Shokouhi, and
 281 Sergey Yekhanin. Differentially private set union. In *International Conference on Machine Learn-*
 282 *ing*, pages 3627–3636. PMLR, 2020. URL <https://arxiv.org/abs/2002.09745>.
- 283 John E. Hopcroft and Richard M. Karp. An $n^{5/2}$ algorithm for maximum matchings in bipartite
 284 graphs. *SIAM J. Comput.*, 2(4):225–231, 1973. doi: 10.1137/0202019. URL [https://doi.](https://doi.org/10.1137/0202019)
 285 [org/10.1137/0202019](https://doi.org/10.1137/0202019).
- 286 A.V. Karzanov. An exact estimate of an algorithm for finding a maximum flow, applied to the
 287 problem “on representatives”. *Issues of Cybernetics. Proc. of the Seminar on Combinatorial*
 288 *Mathematics*, pages 66–70, 1973.
- 289 Shiva Prasad Kasiviswanathan, Homin K Lee, Kobbi Nissim, Sofya Raskhodnikova, and Adam
 290 Smith. What can we learn privately? *SIAM Journal on Computing*, 40(3):793–826, 2011. URL
 291 <https://arxiv.org/abs/0803.0924>.
- 292 Aleksandra Korolova, Krishnaram Kenthapadi, Nina Mishra, and Alexandros Ntoulas. Releasing
 293 search queries and clicks privately. In Juan Quemada, Gonzalo León, Yoëlle S. Maarek, and
 294 Wolfgang Nejdl, editors, *Proceedings of the 18th International Conference on World Wide Web,*
 295 *WWW 2009, Madrid, Spain, April 20-24, 2009*, pages 171–180. ACM, 2009. doi: 10.1145/
 296 1526709.1526733. URL <https://doi.org/10.1145/1526709.1526733>.
- 297 Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher
 298 Potts. Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting*
 299 *of the Association for Computational Linguistics: Human Language Technologies*, pages 142–
 300 150, Portland, Oregon, USA, June 2011. Association for Computational Linguistics. URL [http:](http://www.aclweb.org/anthology/P11-1015)
 301 [//www.aclweb.org/anthology/P11-1015](http://www.aclweb.org/anthology/P11-1015).
- 302 Frank McSherry and Kunal Talwar. Mechanism design via differential privacy. In *48th Annual IEEE*
 303 *Symposium on Foundations of Computer Science (FOCS’07)*, pages 94–103. IEEE, 2007.
- 304 Lakshmi N. Imdb dataset of 50k movie reviews, 2020. URL [https://www.kaggle.com/](https://www.kaggle.com/datasets/lakshmi25npathi/imdb-dataset-of-50k-movie-reviews)
 305 [datasets/lakshmi25npathi/imdb-dataset-of-50k-movie-reviews](https://www.kaggle.com/datasets/lakshmi25npathi/imdb-dataset-of-50k-movie-reviews).
- 306 Jianmo Ni, Jiacheng Li, and Julian J. McAuley. Justifying recommendations using distantly-
 307 labeled reviews and fine-grained aspects. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiao-
 308 jun Wan, editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Lan-*
 309 *guage Processing and the 9th International Joint Conference on Natural Language Process-*
 310 *ing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, pages 188–197. As-
 311 sociation for Computational Linguistics, 2019. doi: 10.18653/v1/D19-1018. URL [https:](https://doi.org/10.18653/v1/D19-1018)
 312 [//doi.org/10.18653/v1/D19-1018](https://doi.org/10.18653/v1/D19-1018).
- 313 Sofya Raskhodnikova and Adam D. Smith. Efficient lipschitz extensions for high-dimensional graph
 314 statistics and node private degree distributions. *CoRR*, abs/1504.07912, 2015. URL [http://](http://arxiv.org/abs/1504.07912)
 315 arxiv.org/abs/1504.07912.
- 316 Adrian Rivera Cardoso and Ryan Rogers. Differentially private histograms under continual obser-
 317 vation: Streaming selection into the unknown. In Gustau Camps-Valls, Francisco J. R. Ruiz, and
 318 Isabel Valera, editors, *Proceedings of The 25th International Conference on Artificial Intelligence*
 319 *and Statistics*, volume 151 of *Proceedings of Machine Learning Research*, pages 2397–2419.
 320 PMLR, 28–30 Mar 2022. URL <https://arxiv.org/abs/2103.16787>.

- 321 Judy Hanwen Shen. Ask reddit, 2020. URL <https://github.com/heyjudes/differentially-private-set-union/tree/ea7b39285dace35cc9e9029692802759f3e1c8e8/data>.
- 324 Adam D. Smith, Shuang Song, and Abhradeep Thakurta. The flajolet-martin sketch itself preserves differential privacy: Private counting with minimal space. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/e3019767b1b23f82883c9850356b71d6-Abstract.html>.
- 330 Marika Swanberg, Damien Desfontaines, and Samuel Haney. DP-SIPS: A simpler, more scalable mechanism for differentially private partition selection. *CoRR*, abs/2301.01998, 2023. doi: 10.48550/arXiv.2301.01998. URL <https://doi.org/10.48550/arXiv.2301.01998>.
- 333 Salil Vadhan. The complexity of differential privacy. *Tutorials on the Foundations of Cryptography: Dedicated to Oded Goldreich*, pages 347–450, 2017. URL https://privacytools.seas.harvard.edu/files/privacytools/files/manuscript_2016.pdf.
- 336 Bing Zhang, Vadym Doroshenko, Peter Kairouz, Thomas Steinke, Abhradeep Thakurta, Ziyin Ma, Himani Apte, and Jodi Spacek. Differentially private stream processing at scale. *arXiv preprint arXiv:2303.18086*, 2023. URL <https://arxiv.org/abs/2303.18086>.

339 A Proofs

340 *Proof of Theorem 1.1.* First note that $q_\ell(D) = \text{DC}(D; \ell) - \frac{2\ell}{\varepsilon} \log(1/2\beta)$ has sensitivity ℓ . Since the generalized exponential mechanism is $\varepsilon/2$ -DP and adding Laplace noise is also $\varepsilon/2$ -DP, the overall algorithm is ε -DP by composition.

343 Since $\hat{\nu} \leftarrow q_{\hat{\ell}}(D) + \text{Lap}\left(2\hat{\ell}/\varepsilon\right)$, we have

$$\mathbb{P}_{\hat{\nu}} \left[\hat{\nu} \leq q_{\hat{\ell}}(D) + \frac{2\hat{\ell}}{\varepsilon} \log\left(\frac{1}{2\beta}\right) \right] = \mathbb{P}_{\hat{\nu}} \left[\hat{\nu} \geq q_{\hat{\ell}}(D) - \frac{2\hat{\ell}}{\varepsilon} \log\left(\frac{1}{2\beta}\right) \right] = 1 - \beta. \quad (8)$$

344 Substituting $q_\ell(D) = \text{DC}(D; \ell) - \frac{2\ell}{\varepsilon} \log(1/2\beta)$ into Equation (8) gives

$$\mathbb{P}_{\hat{\nu}} \left[\hat{\nu} \leq \text{DC}(D; \hat{\ell}) \right] = \quad (9)$$

$$\mathbb{P}_{\hat{\nu}} \left[\hat{\nu} \geq \text{DC}(D; \hat{\ell}) - \frac{4\hat{\ell}}{\varepsilon} \log\left(\frac{1}{2\beta}\right) \right] = 1 - \beta. \quad (10)$$

345 Combining Equation (9) with $\text{DC}(D; \hat{\ell}) \leq \text{DC}(D)$ yields the guarantee in Equation (3) that $\hat{\nu}$ is a lower bound on $\text{DC}(D)$ with probability $\geq 1 - \beta$.

347 The accuracy guarantee of the generalized exponential mechanism (Theorem 3.4) is

$$\mathbb{P}_{\hat{\ell}} \left[q_{\hat{\ell}}(D) \geq \max_{\ell \in [\ell_{\max}]} q_\ell(D) - \ell \cdot \frac{4}{\varepsilon/2} \log(\ell_{\max}/\beta) \right] \geq 1 - \beta$$

348 or, equivalently,

$$\mathbb{P}_{\hat{\ell}} \left[\text{DC}(D; \hat{\ell}) - \frac{2\hat{\ell}}{\varepsilon} \log\left(\frac{1}{2\beta}\right) \geq \max_{\ell \in [\ell_{\max}]} \text{DC}(D; \ell) - \frac{2\ell}{\varepsilon} \log\left(\frac{1}{2\beta}\right) - \frac{8\ell}{\varepsilon} \log\left(\frac{\ell_{\max}}{\beta}\right) \right] \geq 1 - \beta. \quad (11)$$

349 Combining Equations (10) and (11) with a union bound yields

$$\mathbb{P}_{(\hat{\ell}, \hat{\nu}) \leftarrow \mathcal{M}(D)} \left[\hat{\nu} \geq \max_{\ell \in [\ell_{\max}]} \text{DC}(D; \ell) - \frac{2\ell + 2\hat{\ell}}{\varepsilon} \log\left(\frac{1}{2\beta}\right) - \frac{8\ell}{\varepsilon} \log\left(\frac{\ell_{\max}}{\beta}\right) \right] \geq 1 - 2\beta. \quad (12)$$

350 To interpret Equation (12) we need a high-probability upper bound on $\hat{\ell}$. Let $A > 0$ be determined
 351 later and define

$$\ell_A^* := \arg \max_{\ell \in [\ell_{\max}]} \text{DC}(D; \ell) - \frac{A\ell}{\varepsilon}, \quad (13)$$

352 so that $\text{DC}(D; \ell) \leq \text{DC}(D; \ell_A^*) + (\ell - \ell_A^*) \frac{A}{\varepsilon}$ for all $\ell \in [\ell_{\max}]$. Assume the event in Equation (11)
 353 holds. We have

$$\begin{aligned} \text{DC}(D; \hat{\ell}) - \frac{2\hat{\ell}}{\varepsilon} \log \left(\frac{1}{2\beta} \right) &\leq \text{DC}(D; \ell_A^*) + (\hat{\ell} - \ell_A^*) \frac{A}{\varepsilon} - \frac{2\hat{\ell}}{\varepsilon} \log \left(\frac{1}{2\beta} \right), \quad (\text{by Equation (13)}) \\ \text{DC}(D; \hat{\ell}) - \frac{2\hat{\ell}}{\varepsilon} \log \left(\frac{1}{2\beta} \right) &\geq \max_{\ell \in [\ell_{\max}]} \text{DC}(D; \ell) - \frac{2\ell}{\varepsilon} \log \left(\frac{1}{2\beta} \right) - \frac{8\ell}{\varepsilon} \log \left(\frac{\ell_{\max}}{\beta} \right) \\ &\quad (\text{by assumption}) \\ &\geq \text{DC}(D; \ell_A^*) - \frac{2\ell_A^*}{\varepsilon} \log \left(\frac{1}{2\beta} \right) - \frac{8\ell_A^*}{\varepsilon} \log \left(\frac{\ell_{\max}}{\beta} \right). \end{aligned}$$

354 Combining inequalities and simplifying yields

$$\hat{\ell} \cdot \left(2 \log \left(\frac{1}{2\beta} \right) - A \right) \leq \ell_A^* \cdot \left(2 \log \left(\frac{1}{2\beta} \right) + 8 \log \left(\frac{\ell_{\max}}{\beta} \right) - A \right). \quad (14)$$

355 Now we set $A = \log \left(\frac{1}{2\beta} \right)$ to obtain

$$\hat{\ell} \cdot \log \left(\frac{1}{2\beta} \right) \leq \ell_A^* \cdot \left(\log \left(\frac{1}{2\beta} \right) + 8 \log \left(\frac{\ell_{\max}}{\beta} \right) \right). \quad (15)$$

356 Substituting Equation (15) into Equation (12) gives

$$\mathbb{P}_{(\hat{\ell}, \hat{\nu}) \leftarrow \mathcal{M}(D)} \left[\hat{\nu} \geq \max_{\ell \in [\ell_{\max}]} \text{DC}(D; \ell) - \frac{2\ell + 2\ell_A^*}{\varepsilon} \log \left(\frac{1}{2\beta} \right) - \frac{8\ell + 16\ell_A^*}{\varepsilon} \log \left(\frac{\ell_{\max}}{\beta} \right) \right] \geq 1 - 2\beta. \quad (16)$$

357 We simplify Equation (16) using $\log \left(\frac{1}{2\beta} \right) \leq \log \left(\frac{\ell_{\max}}{\beta} \right)$ to obtain Equation (4).

358 Finally, the runtime of $\text{DPDISTINCTCOUNT}(D)$ is dominated by ℓ_{\max} calls to the
 359 $\text{SENSITIVEDISTINCTCOUNT}(D)$ subroutine, which computes the maximum size of a bipartite
 360 matching on a graph with $|E| = \sum_{i \in [n]} |u_i| \cdot \min\{\ell, |u_i|\} \leq |D| \cdot \ell_{\max}$ edges and $|V| + |U| =$
 361 $\text{DC}(D) + \sum_{i \in [n]} \min\{\ell, |u_i|\} \leq 2|D|$ vertices. The Hopcroft-Karp-Karzanov algorithm runs in
 362 time $O(|E| \cdot \sqrt{|V| + |U|}) \leq O(|D|^{1.5} \cdot \ell_{\max})$ time. \square

363 *Proof of Theorem 1.2.* We start from proving privacy guarantees. Note that Algo-
 364 rithm 2 produces the same result as Algorithm 4. Hence, it is enough to prove that
 365 $\text{SENSITIVEAPPROXDISTINCTCOUNT}(\cdot, \ell)$ has sensitivity ℓ . In addition, note that

$$\begin{aligned} \text{SENSITIVEAPPROXDISTINCTCOUNT}((u_1, \dots, u_n), \ell) = \\ \text{SENSITIVEAPPROXDISTINCTCOUNT}(\underbrace{(u_1, \dots, u_n, \dots, u_1, \dots, u_n)}_{\ell \text{ times}}, 1); \end{aligned}$$

366 therefore, it is enough to prove that $\text{SENSITIVEAPPROXDISTINCTCOUNT}(\cdot, 1)$ has sensitivity 1.

367 Assume $D' = (u_1, \dots, u_{j-1}, u_{j+1}, \dots, u_n)$ and let $S_1, \dots, S_n, v_1, \dots, v_n$ be states of S and
 368 v ($v_i = \perp$ if i is skipped), respectively, when run $\text{APPROXSENSITIVEDISTINCTCOUNT}(D)$ and
 369 $S'_1, \dots, S'_n, v'_1, \dots, v'_n$ be states of S and v ($v_i = \perp$ if i is skipped), respectively, when run
 370 $\text{APPROXSENSITIVEDISTINCTCOUNT}(D')$. Let $\{i_1, \dots, i_k\} = \{i : S_i \neq S'_i\}$. It is clear that
 371 $i_1 \geq j$ and $v'_{i_1} = v_j$; similarly v'_{i_2} is either \perp or $v'_{i_2} = v_{i_1}$ etc. As a result $|S'_n| \leq |S_n| \leq |S'_n| + 1$.

372 The sensitivity bound implies that $q_\ell(D) = |S| - \frac{2\ell}{\varepsilon} \log(1/2\beta)$ has sensitivity ℓ . Since the gen-
 373 eralized exponential mechanism is $\varepsilon/2$ -DP and adding Laplace noise is also $\varepsilon/2$ -DP, the overall
 374 algorithm is ε -DP by composition.

Algorithm 4 Approximate Distinct Count Algorithm

```

1: procedure SENSITIVEAPPROXDISTINCTCOUNT( $D = (u_1, \dots, u_n) \in (\Omega^*)^n; \ell \in \mathbb{N}$ )
2:    $S \leftarrow \emptyset$ .
3:   for  $\ell' \in [\ell]$  do
4:     for  $i \in [n]$  with  $u_i \setminus S \neq \emptyset$  do
5:       Choose lexicographically first  $v \in u_i \setminus S$ . ▷ Match  $(i, \ell)$  to  $v$ .
6:       Update  $S \leftarrow S \cup \{v\}$ .
7:     end for
8:   end for
9:   return  $|S|$ 
10: end procedure
11: procedure DPAPPROXDISTINCTCOUNT( $D = (u_1, \dots, u_n) \in (\Omega^*)^n; \ell_{\max} \in \mathbb{N}, \varepsilon > 0, \beta \in (0, \frac{1}{2})$ )
12:   for  $\ell \in [\ell_{\max}]$  do
13:     Define  $q_\ell(D) := \text{SENSITIVEAPPROXDISTINCTCOUNT}(D; \ell) - \frac{2\ell}{\varepsilon} \cdot \log\left(\frac{1}{2\beta}\right)$ .
14:   end for
15:    $\hat{\ell} \leftarrow \text{GEM}(D; \{q_\ell\}_{\ell \in [\ell_{\max}]}, \{\ell\}_{\ell \in [\ell_{\max}]}, \varepsilon/2, \beta)$ . ▷ Algorithm 3
16:    $\hat{v} \leftarrow q_{\hat{\ell}}(D) + \text{Lap}\left(2\hat{\ell}/\varepsilon\right)$ .
17:   return  $(\hat{\ell}, \hat{v}) \in [\ell_{\max}] \times \mathbb{R}$ .
18: end procedure

```

375 Let us denote by $\widehat{\text{DC}}(D; \ell)$ the value of $|S|$ we obtain on Line 8 in Algorithm 2. Note that $\widehat{\text{DC}}(D; \ell)$
 376 is size of a maximal matching in G_ℓ , where G_ℓ is the bipartite graph corresponding to the input D
 377 with ℓ copies of each person (see Algorithm 1 for a formal description of the graph). Since a maximal
 378 matching is a 2-approximation to a maximum matching, we have $\widehat{\text{DC}}(D; \ell) \geq \frac{1}{2}\text{DC}(D; \ell)$. Also
 379 $\widehat{\text{DC}}(D; \ell) \leq \text{DC}(D; \ell)$.

380 Since $\hat{v} \leftarrow q_{\hat{\ell}}(D) + \text{Lap}\left(2\hat{\ell}/\varepsilon\right)$, we have

$$\mathbb{P}_{\hat{v}} \left[\hat{v} \leq q_{\hat{\ell}}(D) + \frac{2\hat{\ell}}{\varepsilon} \log\left(\frac{1}{2\beta}\right) \right] = \mathbb{P}_{\hat{v}} \left[\hat{v} \geq q_{\hat{\ell}}(D) - \frac{2\hat{\ell}}{\varepsilon} \log\left(\frac{1}{2\beta}\right) \right] = 1 - \beta. \quad (17)$$

381 Substituting $q_\ell(D) = \widehat{\text{DC}}(D; \ell) - \frac{2\ell}{\varepsilon} \log(1/2\beta)$ into Equation (17) gives

$$\mathbb{P}_{\hat{v}} \left[\hat{v} \leq \widehat{\text{DC}}(D; \hat{\ell}) \right] = \quad (18)$$

$$\mathbb{P}_{\hat{v}} \left[\hat{v} \geq \widehat{\text{DC}}(D; \hat{\ell}) - \frac{4\hat{\ell}}{\varepsilon} \log\left(\frac{1}{2\beta}\right) \right] = 1 - \beta. \quad (19)$$

382 Combining Equation (18) with $\widehat{\text{DC}}(D; \hat{\ell}) \leq \text{DC}(D)$ yields the guarantee in Equation (6) that \hat{v} is a
 383 lower bound on $\text{DC}(D)$ with probability $\geq 1 - \beta$.

384 The accuracy guarantee of the generalized exponential mechanism (Theorem 3.4) is

$$\mathbb{P}_{\hat{\ell}} \left[q_{\hat{\ell}}(D) \geq \max_{\ell \in [\ell_{\max}]} q_\ell(D) - \ell \cdot \frac{4}{\varepsilon/2} \log(\ell_{\max}/\beta) \right] \geq 1 - \beta$$

385 or, equivalently,

$$\mathbb{P}_{\hat{\ell}} \left[\widehat{\text{DC}}(D; \hat{\ell}) - \frac{2\hat{\ell}}{\varepsilon} \log\left(\frac{1}{2\beta}\right) \geq \max_{\ell \in [\ell_{\max}]} \widehat{\text{DC}}(D; \ell) - \frac{2\ell}{\varepsilon} \log\left(\frac{1}{2\beta}\right) - \frac{8\ell}{\varepsilon} \log\left(\frac{\ell_{\max}}{\beta}\right) \right] \geq 1 - \beta. \quad (20)$$

386 Combining Equations (19) and (20) with a union bound yields

$$\mathbb{P}_{(\hat{\ell}, \hat{v}) \leftarrow \mathcal{M}(D)} \left[\hat{v} \geq \max_{\ell \in [\ell_{\max}]} \widehat{\text{DC}}(D; \ell) - \frac{2\ell + 2\hat{\ell}}{\varepsilon} \log\left(\frac{1}{2\beta}\right) - \frac{8\ell}{\varepsilon} \log\left(\frac{\ell_{\max}}{\beta}\right) \right] \geq 1 - 2\beta. \quad (21)$$

387 To interpret Equation (21) we need a high-probability upper bound on $\hat{\ell}$. Let $A > 0$ be determined
 388 later and define

$$\ell_A^* := \arg \max_{\ell \in [\ell_{\max}]} \widehat{\text{DC}}(D; \ell) - \frac{A\ell}{\varepsilon}, \quad (22)$$

389 so that $\widehat{\text{DC}}(D; \ell) \leq \widehat{\text{DC}}(D; \ell_A^*) + (\ell - \ell_A^*) \frac{A}{\varepsilon}$ for all $\ell \in [\ell_{\max}]$. Assume the event in Equation (20)
 390 holds. We have

$$\begin{aligned} \widehat{\text{DC}}(D; \hat{\ell}) - \frac{2\hat{\ell}}{\varepsilon} \log\left(\frac{1}{2\beta}\right) &\leq \widehat{\text{DC}}(D; \ell_A^*) + (\hat{\ell} - \ell_A^*) \frac{A}{\varepsilon} - \frac{2\hat{\ell}}{\varepsilon} \log\left(\frac{1}{2\beta}\right), \quad (\text{by Equation (22)}) \\ \widehat{\text{DC}}(D; \hat{\ell}) - \frac{2\hat{\ell}}{\varepsilon} \log\left(\frac{1}{2\beta}\right) &\geq \max_{\ell \in [\ell_{\max}]} \widehat{\text{DC}}(D; \ell) - \frac{2\ell}{\varepsilon} \log\left(\frac{1}{2\beta}\right) - \frac{8\ell}{\varepsilon} \log\left(\frac{\ell_{\max}}{\beta}\right) \\ &\geq \widehat{\text{DC}}(D; \ell_A^*) - \frac{2\ell_A^*}{\varepsilon} \log\left(\frac{1}{2\beta}\right) - \frac{8\ell_A^*}{\varepsilon} \log\left(\frac{\ell_{\max}}{\beta}\right). \end{aligned} \quad (\text{by assumption})$$

391 Combining inequalities and simplifying yields

$$\hat{\ell} \cdot \left(2 \log\left(\frac{1}{2\beta}\right) - A\right) \leq \ell_A^* \cdot \left(2 \log\left(\frac{1}{2\beta}\right) + 8 \log\left(\frac{\ell_{\max}}{\beta}\right) - A\right). \quad (23)$$

392 Now we set $A = \log\left(\frac{1}{2\beta}\right)$ to obtain

$$\hat{\ell} \cdot \log\left(\frac{1}{2\beta}\right) \leq \ell_A^* \cdot \left(\log\left(\frac{1}{2\beta}\right) + 8 \log\left(\frac{\ell_{\max}}{\beta}\right)\right). \quad (24)$$

393 Substituting Equation (15) into Equation (21) gives

$$\mathbb{P}_{(\hat{\ell}, \hat{\nu}) \leftarrow \mathcal{M}(D)} \left[\hat{\nu} \geq \max_{\ell \in [\ell_{\max}]} \widehat{\text{DC}}(D; \ell) - \frac{2\ell + 2\ell_A^*}{\varepsilon} \log\left(\frac{1}{2\beta}\right) - \frac{8\ell + 16\ell_A^*}{\varepsilon} \log\left(\frac{\ell_{\max}}{\beta}\right) \right] \geq 1 - 2\beta. \quad (25)$$

394 We simplify Equation (25) using $\log\left(\frac{1}{2\beta}\right) \leq \log\left(\frac{\ell_{\max}}{\beta}\right)$ to obtain

$$\mathbb{P}_{(\hat{\ell}, \hat{\nu}) \leftarrow \widehat{\mathcal{M}}(D)} \left[\hat{\nu} \geq \max_{\ell \in [\ell_{\max}]} \widehat{\text{DC}}(D; \ell) - \frac{10\ell + 18\ell_A^*}{\varepsilon} \log\left(\frac{\ell_{\max}}{\beta}\right) \right] \geq 1 - 2\beta. \quad (26)$$

395 Note that $\widehat{\text{DC}}(D; \ell) \geq \frac{1}{2}\text{DC}(D; \ell)$; hence,

$$\mathbb{P}_{(\hat{\ell}, \hat{\nu}) \leftarrow \widehat{\mathcal{M}}(D)} \left[\hat{\nu} \geq \max_{\ell \in [\ell_{\max}]} \frac{1}{2}\text{DC}(D; \ell) - \frac{10\ell + 18\ell_A^*}{\varepsilon} \log\left(\frac{\ell_{\max}}{\beta}\right) \right] \geq 1 - 2\beta.$$

396 Finally, note that $\ell_A^* \leq \ell_*$; therefore, we proved Equation (7).

397 It only remains to verify that Algorithm 2 can be implemented in $O(|D|)$ time. We can implement
 398 S using a hash table to ensure that we can add an element or query membership of an element in
 399 constant time. (We can easily maintain a counter for the size of S .) We assume D is presented as
 400 a linked list of linked lists representing each u_i and furthermore that the linked lists u_i are sorted
 401 in lexicographic order. The outer loop proceeds through the linked list for $D = (u_1, \dots, u_n)$. For
 402 each u_i , we simply pop elements from the linked list and check if they are in S until either we find
 403 $v \in u_i \setminus S$ (and add v to S) or u_i becomes empty (in which case we remove it from the linked list
 404 for D .) Since each iteration decrements $|D|$, the runtime is $O(|D|)$. \square