

## A APPENDICES

### A.1 LDP Mechanisms

**Optimal Local Hash.** The Optimal Local Hash (OLH) mechanism [57] is designed for randomizing private values in a large domain. It maps the value with a randomly selected hash function to a new data domain with the size of  $g \ll d$  before randomizing the value. The randomization method is the same as GRR with  $p'$  and  $q'$  as follows

$$\begin{cases} p' = \frac{e^\epsilon}{e^\epsilon + g - 1}, \\ q' = \frac{1}{e^\epsilon + g - 1}. \end{cases} \quad (3)$$

It can also use the generic method to estimate the count  $\tilde{c}_i$  with  $p = p'$  and  $q = \frac{1}{g}p' + \frac{g-1}{g}q' = \frac{1}{g}$ . When the size of the new data domain  $g = e^\epsilon + 1$ , the variance of  $\tilde{c}_i$  can be minimized as

$$\text{Var}[\tilde{c}_i] = n \cdot \frac{4e^\epsilon}{(e^\epsilon - 1)^2} \quad (4)$$

**Hadamard Response.** The Hadamard Response (HR) mechanism [6, 7] encodes private values with a  $K \times K$  Hadamard matrix, where  $K = 2^{\lceil \log_2(d+1) \rceil}$ . Expect for the first row of the matrix (all values are '1'), each other row corresponds to a value in the data domain. When encoding the  $i^{\text{th}}$  value in the data domain, the output value is randomly selected from column indices that have '1' in the  $(i+1)^{\text{th}}$  row with the probability of  $p$ , and selected from other indices (columns have '0') with the probability of  $q$ , where

$$\begin{cases} p = \frac{e^\epsilon}{1+e^\epsilon} \\ q = \frac{1}{1+e^\epsilon} \end{cases} \quad (5)$$

Then  $\tilde{c}_i$  can be calculated as follows

$$\tilde{c}_i = \frac{2(e^\epsilon + 1)}{e^\epsilon - 1} (\hat{c}_i - \frac{n}{2}) \quad (6)$$

The variance of the estimation result  $\tilde{c}_i$  is

$$\text{Var}[\tilde{c}_i] = n \cdot \frac{4(e^\epsilon + 1)^2}{(e^\epsilon - 1)^2} \quad (7)$$

### A.2 Proof of Theorem 3.2

**PROOF.** We assume that  $\hat{f}_i$  is the frequency of noisy data recorded in  $\mathcal{H}\mathcal{G}$  according to the *ED strategy* in BGR. Meanwhile, the *ED strategy* would introduce additional error when recording  $\hat{f}_i$ , while the frequency actually recorded and used for debiasing by the GRR mechanism before publishing is  $\tilde{f}_i$ . According to Lemma 3.1, we have the upper bound of  $\hat{f}_i - \tilde{f}_i$  is

$$\begin{aligned} \Pr[\hat{f}_i - \tilde{f}_i \geq \alpha t] &\leq \frac{1}{2\alpha t} (\hat{f}_i - \sqrt{\hat{f}_i^2 - \frac{4P_{\text{weak}}E(V)}{b-1}}) \\ \Rightarrow \Pr[\hat{f}_i \leq \tilde{f}_i + \alpha t] &\geq (1 - \frac{1}{2\alpha t} (\hat{f}_i - \sqrt{\hat{f}_i^2 - \frac{4P_{\text{weak}}E(V)}{b-1}})) \end{aligned}$$

The distribution of  $\hat{f}_i$  can be decomposed into  $\text{Bin}(f_i, p) + \text{Bin}(t - f_i, q)$ . Denote  $\text{Bin}(f_i, p)$  as  $P_1$  and  $\text{Bin}(t - f_i, q)$  as  $P_2$ , according to

the hoeffding inequality, we have

$$\begin{aligned} \Pr[f_i p - P_1 \leq \zeta] &\geq 1 - e^{-\frac{2\zeta^2}{f_i}}, \text{ and} \\ \Pr[(t - f_i)q - P_2 \leq \zeta] &\geq 1 - e^{-\frac{2\zeta^2}{t-f_i}} \end{aligned}$$

Then, according to Bonferroni inequality, we have

$$\begin{aligned} \Pr[P_1 + P_2 \geq f_i p + (t - f_i)q - 2\zeta] &\geq 1 - e^{-\frac{2\zeta^2}{f_i}} - e^{-\frac{2\zeta^2}{t-f_i}} \\ \Rightarrow \Pr[\hat{f}_i \geq f_i p + (t - f_i)q - 2\zeta] &\geq 1 - e^{-\frac{2\zeta^2}{f_i}} - e^{-\frac{2\zeta^2}{t-f_i}} \\ \Rightarrow \Pr[\frac{\hat{f}_i - tq}{p - q} - f_i \geq \frac{-2\zeta}{p - q}] &\geq 1 - e^{-\frac{2\zeta^2}{f_i}} - e^{-\frac{2\zeta^2}{t-f_i}} \\ \Rightarrow \Pr[f_i - \frac{\hat{f}_i + \alpha t - tq}{p - q} \leq f_i - \frac{\hat{f}_i - tq}{p - q} \leq \frac{2\zeta}{p - q}] \\ &\geq (1 - e^{-\frac{2\zeta^2}{f_i}} - e^{-\frac{2\zeta^2}{t-f_i}})(1 - \frac{f_i}{2\alpha t}(1 - \sqrt{1 - \frac{4P_{\text{weak}}E(V)}{\hat{f}_i^2(b-1)}})) \\ \Rightarrow \Pr[f_i - \tilde{f}_i \leq \frac{2\zeta + \alpha t}{p - q}] \\ &\geq (1 - 2e^{-\frac{2\zeta^2}{t}})(1 - \frac{1}{2\alpha}(1 - \sqrt{1 - \frac{4P_{\text{weak}}E(V)}{b-1}})) \\ \Rightarrow \Pr[f_i - \tilde{f}_i \leq (\sqrt{2t \log(2/\beta)} + \alpha t) \cdot \frac{e^{\epsilon+d-1}}{e^\epsilon - 1}] \\ &\geq (1 - \beta)(1 - \frac{1}{2\alpha}(1 - \sqrt{1 - \frac{4P_{\text{weak}}E(V)}{b-1}})) \end{aligned}$$

where  $\zeta = \sqrt{t \log(2/\beta)/2}$ ,  $P_{\text{weak}} = \frac{(i-1)!(d-k)!}{(d-1)!(i-k)!}$ , and  $E(V) = \sum_{j=i+1}^d f_j$ .  $\square$

### A.3 Proof of Theorem 4.2

**PROOF.** We assume that  $\hat{f}_i$  is the frequency of noisy data recorded in  $\mathcal{H}\mathcal{G}$  according to the *ED strategy* in BDR. Meanwhile, the *ED strategy* would introduce additional error when recording  $\hat{f}_i$ , while the frequency actually recorded and used for debiasing by the randomization mechanism before publishing is  $\tilde{f}_i$ . According to Lemma 3.1, we have the upper bound of  $\hat{f}_i - \tilde{f}_i$  is

$$\begin{aligned} \Pr[\hat{f}_i - \tilde{f}_i \geq \alpha t] &\leq \frac{1}{2\alpha t} (\hat{f}_i - \sqrt{\hat{f}_i^2 - \frac{4P_{\text{weak}}E(V)}{b-1}}) \\ \Rightarrow \Pr[\hat{f}_i \leq \tilde{f}_i + \alpha t] &\geq (1 - \frac{1}{2\alpha t} (\hat{f}_i - \sqrt{\hat{f}_i^2 - \frac{4P_{\text{weak}}E(V)}{b-1}})) \end{aligned}$$

The distribution of  $\hat{f}_i$  can be decomposed into  $\text{Bin}(f_i, p_1 p_2) + \text{Bin}(N_h - f_i, p_1 q_2) + \text{Bin}(t - N_h, q_1/k)$ , where  $N_h$  is the number of hot items,  $N_h \leq t$ . Denote  $\text{Bin}(f_i, p_1 p_2)$  as  $P_1$ ,  $\text{Bin}(N_h - f_i, p_1 q_2)$  as  $P_2$ , and  $\text{Bin}(t - N_h, q_1/k)$  as  $P_3$ , according to the hoeffding inequality, we

have

$$\begin{aligned} Pr[f_i p_1 p_2 - P_1 \leq \zeta] &\geq 1 - e^{-\frac{2\zeta^2}{f_i}}, \\ Pr[(N_h - f_i)p_1 q_2 - P_2 \leq \zeta] &\geq 1 - e^{-\frac{2\zeta^2}{N_h - f_i}}, \text{ and} \\ Pr[(t - N_h) \cdot \frac{q_1}{k} - P_3 \leq \zeta] &\geq 1 - e^{-\frac{2\zeta^2}{t - N_h}}. \end{aligned}$$

Then, according to Bonferroni inequality, we have

$$\begin{aligned} &Pr[f_i p_1 p_2 - P_1 + (N_h - f_i)p_1 q_2 - P_2 + (t - N_h) \cdot \frac{q_1}{k} \\ &- P_3 \leq 3\zeta] \geq 1 - e^{-\frac{2\zeta^2}{f_i}} + 1 - e^{-\frac{2\zeta^2}{N_h - f_i}} + 1 - e^{-\frac{2\zeta^2}{t - N_h}} - 2 \\ \Rightarrow &Pr[P_1 + P_2 + P_3 \geq f_i p_1 p_2 + (N_h - f_i)p_1 q_2 \\ &+ (t - N_h) \cdot \frac{q_1}{k} - 3\zeta] \geq 1 - e^{-\frac{2\zeta^2}{f_i}} - e^{-\frac{2\zeta^2}{N_h - f_i}} - e^{-\frac{2\zeta^2}{t - N_h}} \\ \Rightarrow &Pr[\hat{f}_i \geq f_i p_1 p_2 + (N_h - f_i)p_1 q_2 + (t - N_h) \cdot \frac{q_1}{k} - 3\zeta] \\ \geq &1 - e^{-\frac{2\zeta^2}{f_i}} - e^{-\frac{2\zeta^2}{N_h - f_i}} - e^{-\frac{2\zeta^2}{t - N_h}} \\ \Rightarrow &Pr[\frac{\hat{f}_i - N_h p_1 q_2 - (t - N_h) \cdot \frac{q_1}{k}}{p_1(p_2 - q_2)} - f_i \geq \frac{-3\zeta}{p_1(p_2 - q_2)}] \\ \geq &1 - e^{-\frac{2\zeta^2}{f_i}} - e^{-\frac{2\zeta^2}{N_h - f_i}} - e^{-\frac{2\zeta^2}{t - N_h}} \\ \Rightarrow &Pr[f_i - \frac{\hat{f}_i + \alpha t - N_h p_1 q_2 - (t - N_h) \cdot \frac{q_1}{k}}{p_1(p_2 - q_2)} \\ \leq &f_i - \frac{\hat{f}_i - N_h p_1 q_2 - (t - N_h) \cdot \frac{q_1}{k}}{p_1(p_2 - q_2)} \leq \frac{3\zeta}{p_1(p_2 - q_2)}] \\ \geq &(1 - e^{-\frac{2\zeta^2}{f_i}} - e^{-\frac{2\zeta^2}{N_h - f_i}} - e^{-\frac{2\zeta^2}{t - N_h}}) \cdot (1 \\ &- \frac{f_i}{2\alpha t} (1 - \sqrt{1 - \frac{4P_{weak}E(V)}{\hat{f}_i^2(b-1)}})) \\ \Rightarrow &Pr[f_i - \tilde{f}_i \leq \frac{3\zeta + \alpha t}{p_1(p_2 - q_2)}] \geq (1 - 3e^{-\frac{2\zeta^2}{N_h}}) \cdot (1 \\ &- \frac{1}{2\alpha} (1 - \sqrt{1 - \frac{4P_{weak}E(V)}{b-1}})) \\ \Rightarrow &Pr[f_i - \tilde{f}_i \leq (3\sqrt{\frac{N_h \log(3/\beta)}{2}} + \alpha t) \\ &\cdot \frac{(e^{\epsilon_1} + 1)(e^{\epsilon_2} + k - 1)}{e^{\epsilon_1}(e^{\epsilon_2} - 1)}] \\ \geq &(1 - \beta)(1 - \frac{1}{2\alpha} (1 - \sqrt{1 - \frac{4P_{weak}E(V)}{b-1}})) \end{aligned}$$

□

#### A.4 Proof of Theorem 4.1

PROOF. Denote  $v$  and  $v'$  as two raw data,  $o_1$ ,  $o_2$  and  $o_3$  as the output of mechanisms,  $\Omega_h$  as the domain of hot items,  $\Omega_c$  as the

#### Algorithm 12 DSR\_Insert( $v, \mathcal{HG}, p_1, q_1, p_2, q_2$ )

---

```

1:  $\Omega_s = \{\mathcal{HG}.ID\} \cup \{\perp\}$ 
2: if the least count in  $\mathcal{HG}$  changed from  $> 1$  to  $\leq 1$  then
3:   for each  $\mathcal{HG}[j] \in \mathcal{HG}$  do
4:      $\mathcal{HG}[j] \leftarrow (\mathcal{HG}[j].C - q_1)(p_2 - q_2)/(p_1 - q_1)$ 
5:   Insert  $v$  into  $\mathcal{HG}$  following ED strategy;  $\triangleright r_i^t \in \Omega$ 
6:    $num_{entire} \leftarrow num_{entire} + 1$ 
7:   if the least count in  $\mathcal{HG} \leq 0$  then
8:     Replace the weakest KV pair with  $\langle r_i^t, 1 - p_1 \cdot num_{entire} / (p_1 - q_1) - p_2 \cdot num_{reduced} / (p_2 - q_2) \rangle$ 
9:   else if the least count in  $\mathcal{HG}$  changed from  $\leq 1$  to  $> 1$  then
10:    for each  $\mathcal{HG}[j] \in \mathcal{HG}$  do
11:       $\mathcal{HG}[j] \leftarrow (\mathcal{HG}[j].C - q_2)(p_1 - q_1)/(p_2 - q_2)$ 
12:    Insert  $v$  into  $\mathcal{HG}$  following ED strategy;  $\triangleright r_i^t \in \Omega_s$ 
13:     $num_{reduced} \leftarrow num_{reduced} + 1$ 
14:    else if the least count in  $\mathcal{HG} > 1$  then
15:      for each  $\mathcal{HG}[j] \in \mathcal{HG}$  do
16:         $\mathcal{HG}[j] \leftarrow \mathcal{HG}[j].C - q_1$ 
17:      Insert  $v$  into  $\mathcal{HG}$  following ED strategy;  $\triangleright r_i^t \in \Omega_s$ 
18:       $num_{reduced} \leftarrow num_{reduced} + 1$ 
19:    else
20:      for each  $\mathcal{HG}[j] \in \mathcal{HG}$  do
21:         $\mathcal{HG}[j] \leftarrow \mathcal{HG}[j].C - q_2$ 
22:      Insert  $v$  into  $\mathcal{HG}$  following ED strategy;  $\triangleright r_i^t \in \Omega$ 
23:       $num_{entire} \leftarrow num_{entire} + 1$ 
24:      if the least count in  $\mathcal{HG} \leq 0$  then
25:        Replace the weakest KV pair with  $\langle r_i^t, 1 - p_1 \cdot num_{entire} / (p_1 - q_1) - p_2 \cdot num_{reduced} / (p_2 - q_2) \rangle$ 
26:    return Updated  $\mathcal{HG}$ 

```

---

domain of cold items. Firstly,  $\mathcal{M}_{judge}$  satisfies  $\epsilon_1$ -LDP, we have

$$\frac{Pr[o = \text{"Hot"} | v \in \Omega_h]}{Pr[o = \text{"Hot"} | v' \in \Omega_c]} \leq \frac{p_1}{q_1} = e^{\epsilon_1}$$

The same result can be proved when  $o = \text{"Cold"}$ .

Secondly, we prove that  $\mathcal{M}_{hot}$  mechanism satisfies  $\epsilon_2$ -LDP. The probability ratio of  $v$  and  $v'$  to get the same randomized item  $o_2 \in \Omega_h$  is

$$\frac{Pr[o_2 | v \in \Omega_h]}{Pr[o_2 | v' \in \Omega_c]} \leq \frac{p_2}{1/k} \leq \frac{p_2}{q_2} = e^{\epsilon_2}$$

Similarly, the  $\mathcal{M}_{cold}$  mechanism satisfies

$$\frac{Pr[o_3 | v \in \Omega_c]}{Pr[o_3 | v' \in \Omega_h]} \leq e^{\epsilon_2}$$

According to the composition theorem of DP [32], BDR satisfies  $(\epsilon_1 + \epsilon_2)$ -LDP at each timestamp where  $\epsilon_1 + \epsilon_2 = \epsilon$ . Therefore, BDR satisfies  $\epsilon$ -LDP. □

#### A.5 Algorithm of Function DSR\_Insert and DSR\_FinalDebias

#### A.6 Normalized Discounted Cumulative Gain (NDCG)

It measures the ordering quality of the heavy hitters captured by  $\mathcal{HG}$ , which is a common effectiveness in recommendation systems and other related applications. Specifically, let  $V = \{v_1, v_2, \dots, v_k\}$  as the Top- $k$  heavy hitters in  $\mathcal{HG}$ . If  $v_i$  is one of a true Top- $k$  heavy

**Algorithm 13** DSR\_FinalDebias( $\mathcal{HG}$ ,  $p_1$ ,  $q_1$ ,  $p_2$ ,  $q_2$ )

---

```

1: # Complete final debias based on the state of the previous  $\mathcal{HG}$ .
2: if the least count in  $\mathcal{HG}$  changed from  $> 1$  to  $\leq 1$  then
3:   for each  $\mathcal{HG}[j] \in \mathcal{HG}$  do
4:      $\mathcal{HG}[j] \leftarrow \mathcal{HG}[j].C / (p_2 - q_2)$ 
5: else if the least count in  $\mathcal{HG}$  changed from  $\leq 1$  to  $> 1$  then
6:   for each  $\mathcal{HG}[j] \in \mathcal{HG}$  do
7:      $\mathcal{HG}[j] \leftarrow \mathcal{HG}[j].C / (p_1 - q_1)$ 
8: else if the least count in  $\mathcal{HG}$   $> 1$  then
9:   for each  $\mathcal{HG}[j] \in \mathcal{HG}$  do
10:     $\mathcal{HG}[j] \leftarrow \mathcal{HG}[j].C / (p_1 - q_1)$ 
11: else
12:   for each  $\mathcal{HG}[j] \in \mathcal{HG}$  do
13:     $\mathcal{HG}[j] \leftarrow \mathcal{HG}[j].C / (p_2 - q_2)$ 
14: return Updated  $\mathcal{HG}$ 

```

---

hitter, the relevance score  $rel_i$  is

$$rel_{v_i} = |k - |rank_{actual}(v_i) - rank_{estimated}(v_i)||.$$

If  $v_i$  is not a true Top- $k$  heavy hitter, we directly set its  $rel_i$  as 0. Then, the Discounted Cumulative Gain (DCG) is

$$DCG_k = rel_{v_1} + \sum_{i=2}^k \frac{rel_{v_i}}{\log_2(i)}.$$

Finally, we normalize the DCG of  $\mathcal{HG}$  by comparing it with the Ideal DCG (IDCG), which is the DCG when  $\mathcal{HG}$  records an actual list of Top- $k$  heavy hitters.

$$NDCG_k = \frac{DCG_k}{IDCG_k}.$$

$NDCG_k$  is between 0 and 1 for all  $k$ , and the closer it is to 1 means the ordering quality of  $\mathcal{HG}$  is higher.

## A.7 Implementation Details

**A.7.1 Re-implementation for LDP Mechanisms.** We treat existing LDP frequency estimation approaches as privacy-preserving baselines. Specifically, we estimate the frequency of all items under LDP, and output the Top- $k$  counts as the heavy hitter. Cormode, Maddock, and Maple [24] placed various LDP frequency estimation approaches into a common framework, and performed an series of experiments in Python<sup>2</sup>. Their work offered a starting point of our implementations.

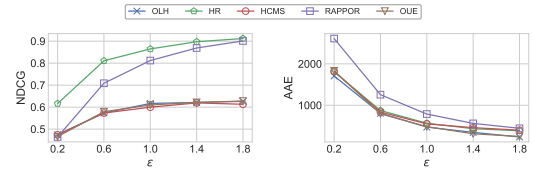
We carefully studied the source codes, and fully re-implemented all baseline LDP mechanisms with the following optimizations.

**Data serialization.** In [24], the client outputs the perturbed data as an object, which server takes as its input to do data aggregation. In practice, the server and the client would communicate via a network channel. This requires object serializations and introduces additional communication and computation costs. In our implementation, we manually serialize the perturbed data to 'byte[]' based on the underlying approaches. If the client outputs bit strings (e.g., OUE and RAPPOR), we compress the output bit string by representing each 8 bits into 1 byte. If the client outputs integers (e.g., OLH, HR), we represent the integer with the minimal byte length, i.e., 1-4 byte for integers in range  $[0, 2^8)$ ,  $[0, 2^{16})$ ,  $[0, 2^{24})$ , and  $[0, 2^{32})$ , respectively.

<sup>2</sup><https://github.com/Samuel-Maddock/pure-LDP>

**Choices of the hash.** Some frequency estimation approaches leverage (non-cryptographic) hash to map input to Boolean (BLH) or integer(s) (RAPPOR, OLH, HCMS). The performances of these approaches are greatly affected by the efficiency of the underlying hash. Meanwhile, HeavyGuardian also leverages (non-cryptographic) hash to partition data into buckets. Note that [24] and [62] respectively use xxHash and BobHash. We invoke BobHash in all schemes since our test shows that BobHash are more efficient<sup>3</sup>.

Besides, we find that debiasing the randomized data before storing it into *HeavyGuardian* can avoid the bias introduced by the ED strategy being amplified by the debiasing process. However, if a complete debiasing is performed every time randomized data comes, it can cause the previously accumulated count to be debiased repeatedly, e.g., divided by denominator  $p - q$  of the debiasing formula repeatedly in BGR. Therefore, we perform partial debiasing when collecting and dividing all counts by the denominator of the debiasing formula only before publishing the results.



**Figure 11: The comparison of different LDP mechanisms on a synthetic dataset with Normal distribution.**

**A.7.2 Re-implementation for HeavyGuardian.** We treat the original *HeavyGuardian* as the non-private baseline. We carefully studied existing open-source C/C++ codes provided by Yang et al.<sup>4</sup> and fully re-implemented the *HeavyGuardian* using Java. Our *HeavyGuardian* re-implementation has some improvements compared with the original implementation. First, the original implementation contains some hard-coded parameters for different tasks, while our re-implementation allows developers to dynamically config for different tasks. Second, the ED strategy in *HeavyGuardian* contains a Bernoulli sampling procedure, i.e., sampling a Boolean value with probability  $\mathcal{P} = b^{-C}$  being True, where  $b > 1$  is a predefined constant number, and  $C$  is a counting value. The naive method of sampling used in original *HeavyGuardian* implementation is to randomly sample  $r \in [0, 1)$  and test whether  $r < b^{-C}$ . However, since finite computers cannot faithfully represent real numbers, the naive method would not produce the Boolean value with the correct distribution. In our implementation, we parse  $b^{-C} = \exp(-C \cdot \ln(b))$  and leverage the method of Bernoulli( $\exp(-\gamma)$ ) proposed by Ganonne et al. [15] to do the sampling with no loss in accuracy.

Recall that the basic version of *HeavyGuardian* is a hash table with  $w \geq 1$  buckets storing KV pairs ( $\langle ID, count \rangle$ ). Each bucket is divided into the heavy part with size  $\lambda_h > 0$  and the light part with size  $\lambda_l \geq 0$ . Because heavy hitter detection focus on only hot items, Yang et al. [62] recommend setting  $\lambda_l = 0$  when using *HeavyGuardian* for heavy hitter detection tasks. Our experiments

<sup>3</sup>Our test shows that on MacBook M1, xxHash takes about 0.01us to provide an output while BobHash takes about 0.1us.

<sup>4</sup><https://github.com/Gavindeed/HeavyGuardian>

follow this recommendation and set  $\lambda_l = 0$  (except CNR). The basic version of *HeavyGuardian* also allows using different  $\lambda_h$  and different numbers of buckets  $w$  when counting the most frequent  $k$  items in the heavy hitter task. Although our implementation also allows setting  $w$  and  $\lambda_h$ , our experiments focus on the basic cases, i.e.,  $w = 1$  and  $\lambda_h = k$ , to better demonstrate the effectiveness of our schemes.

We obtain memory consumption by measuring the deep sizes (i.e., the size of an object including the size of all referred objects, in addition to the size of the object itself) of Objects packaging the *HeavyGuardian* and our schemes. The tool we use is the JOL (Java Object Layout) library<sup>5</sup>. Although the error bounds we give in the theoretical analyses are debiasing after storing, the actual error in our implementation is still bounded by and even lower than the theoretical results.

## A.8 Supplementary Experiments

## A.9 Extended Related Work

**Differential Private Data Stream Collection** The earliest studies in differential privacy for streaming data collection originate from continuous observation of private data [9, 28, 31, 34, 35]. Long-term data collection from users can be regarded as the collection of data streams. These studies mainly consider the degradation of privacy guarantee due to the repeated appearance of private data when the user's state not changing for a period of time.

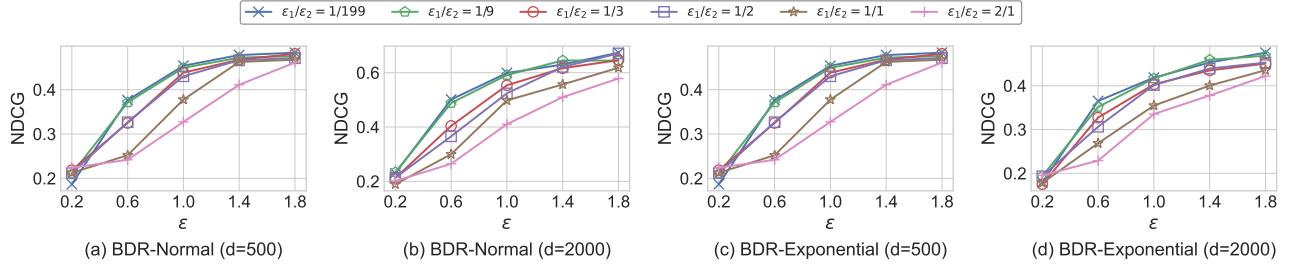
Recent works on differential private data stream collection mainly focus on Centralized Differential Privacy (CDP). Some works study how to publish the summation of the streaming data privately. To avoid the overestimation of sensitivity of the streaming data caused by outliers, Perrier et al. [50] propose truncating the data exceeding a threshold to reduce the sensitivity. Wang et al. [58] point out that the threshold should be dynamically adjusted for different distributions instead of using the fixed quantile value. Therefore, they propose to use an exponential mechanism to get a more reasonable sensitivity. Some works are devoted to solving the problem that the privacy guarantee is continuously degraded due to the repeated use of data in streams in consecutive periods. Kellaris et al. [43] propose  $w$ -event DP, regarding the statistical results published in a sliding window as an event. Farokhi et al. [36] propose to address the problem of the exploding privacy budget by reducing the privacy guarantee provided for data that appeared in the past over time. Some works study the release of correlated streaming data. Wang et al. [56] propose a correlated Laplace noise mechanism, and Bao et al. [10] propose a correlated Gaussian noise mechanism.

There are also some recent works on data stream collection with Local Differential Privacy (LDP) [41, 52, 58]. Joseph et al. [41] design a protocol to submit LDP-protected data only when the streaming data has changed and can have a greater impact on the statistical results. Ren et al. [52] propose a privacy budget segmentation framework that provides  $w$ -event LDP protection to prevent the degradation of privacy due to continuous data collection. Besides, Wang et al. [58] extend the proposed truncation-based CDP mechanism to LDP for the release of streaming data.

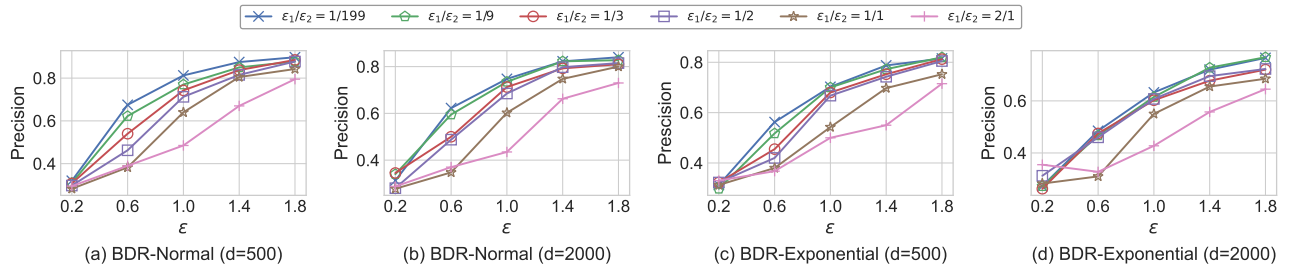
**Tracking Heavy Hitters in Data stream** Mining streaming data faces three principal challenges: *volume*, *velocity*, and *volatility*

[44]. The existing heavy hitters estimation algorithms in the data stream can be divided into three classes: Counter-based algorithms, Quantile algorithms, and Sketch algorithms [22]. Counter-based algorithms track the subset of items in the stream, and they quickly determine whether to record and how to record with each new arrival data. Manku et al. [48] propose two algorithms Sticky Sampling and Lossy Counting, which only record the item and its counts whose estimated counts exceed the threshold. Metwally et al. [49] design an algorithm called Space-Saving and record data with a data structure called Stream-Summary, which achieves rapid deletion, update, and insertion for each new arrival data. Subsequently, Yang et al. [62] propose a new algorithm called HeavyGuardian to improve Space-Saving. Zhou et al. [65] also propose a framework called Cold Filter (CF) to improve Space-Saving. The Quantile algorithms such as the GK algorithm [39] and QDigest algorithm [54], focus on finding the item which is the smallest item that dominates  $\phi n$  items from the data stream. Sketch algorithms record items with a data structure, which can be thought of as a linear projection of the input, hash functions are usually used to define the linear projection. Some existing works include Count Sketch[8], Count-Min Sketch[25], WavingSketch[46], and Moment[20], etc. However, the sketch algorithms may involve a large number of hash operations, which cannot meet the timeliness requirements of streaming data. Besides, the unimportant low-frequent items are all recorded, which leads to unnecessary memory consumption. Our design is based on Counter-based algorithms with an extended setting where streaming data is protected by LDP.

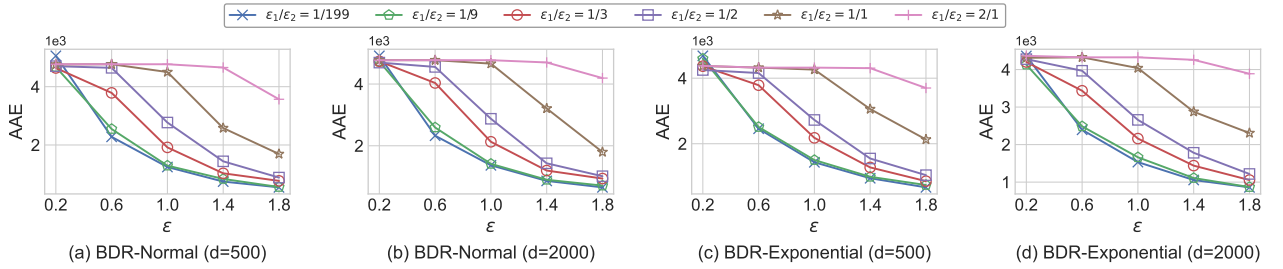
<sup>5</sup><http://hg.openjdk.java.net/code-tools/jol>



**Figure 12: The impact of different allocations of privacy budget on NDCG of BDR. The evaluation is conducted on the Synthetic dataset of Normal distribution and Exponential distribution with domain size 500 and 2000, while taking 1% data for warm-up stage.**



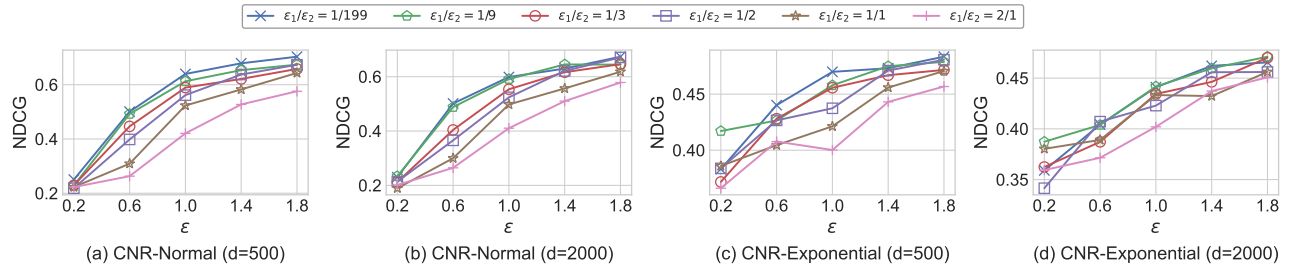
**Figure 13: The impact of different allocations of privacy budget on Precision of BDR. The evaluation is conducted on the Synthetic dataset of Normal distribution and Exponential distribution with domain size 500 and 2000, while taking 1% data for warm-up stage.**



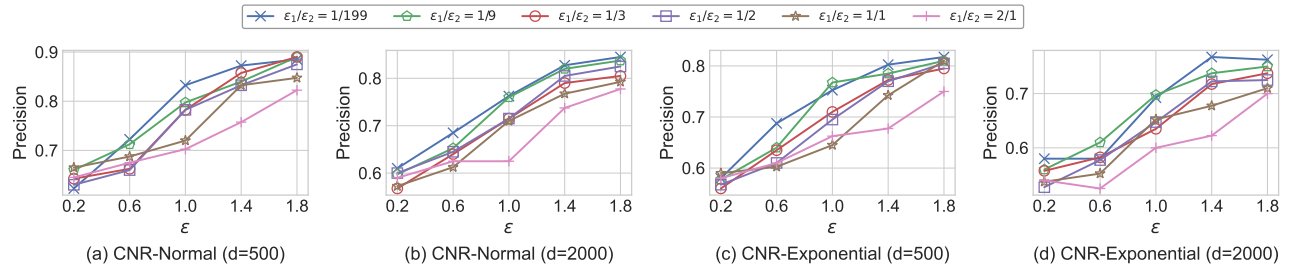
**Figure 14: The impact of different allocations of privacy budget on AAE of BDR. The evaluation is conducted on the Synthetic dataset of Normal distribution and Exponential distribution with domain size 500 and 2000, while taking 1% data for warm-up stage.**

**Table 3: Total communication overheads (MB) consumed by the baseline and proposed three schemes for deriving the Top-20 on each entire dataset. The left side of “/” is the uploading overheads for users, while the right side is the downloading overheads for users.**

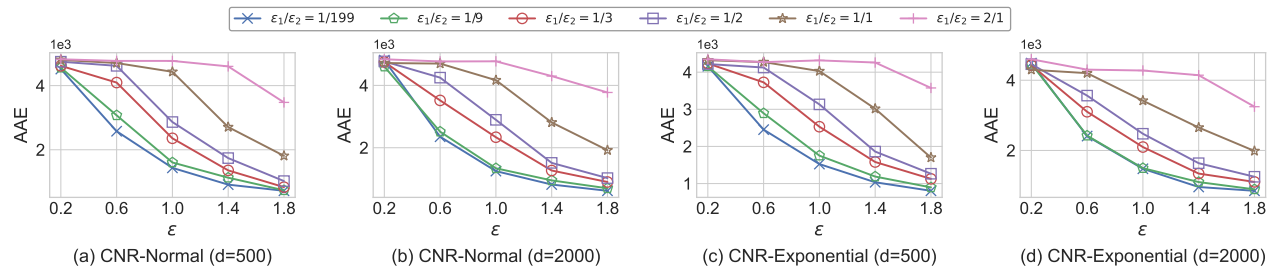
Dataset Scheme	Kosarak (30.5M)	Webdocs (1413.24M)
BGR	35.09/0	1883.36/0
DSR	35.07/1062.79	1875.19/49891.47
BDR	26.67/709.04	1668.71/26295.89
CNR	26.49/698.95	1680.25/26861.62



**Figure 15: The impact of different allocations of privacy budget on NDCG of CNR. The evaluation is conducted on the Synthetic dataset of Normal distribution and Exponential distribution with domain size 500 and 2000, while taking 1% data for warm-up stage.**



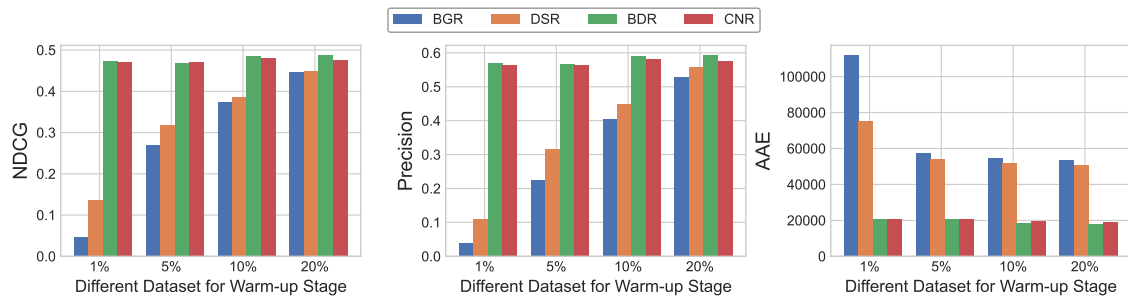
**Figure 16: The impact of different allocations of privacy budget on Precision of CNR. The evaluation is conducted on the Synthetic dataset of Normal distribution and Exponential distribution with domain size 500 and 2000, while taking 1% data for warm-up stage.**



**Figure 17: The impact of different allocations of privacy budget on AAE of CNR. The evaluation is conducted on the Synthetic dataset of Normal distribution and Exponential distribution with domain size 500 and 2000, while taking 1% data for warm-up stage.**

**Table 4: Total server runtime (second) of the baseline and three proposed schemes for deriving the Top-20 on each entire dataset.**

Scheme \ Dataset	Kosarak	Webdocs
BGR	5.375	215.814
DSR	5.586	245.314
BDR	5.224	257.32
CNR	7.563	417.184



**Figure 18: The impact of the warm-up stage on the accuracy of the proposed schemes for tracking Top-20 items on Kosarak dataset, where  $\epsilon = 2$ .**

**Table 5: Total client runtime (second) of the baseline and three proposed schemes for deriving the Top-20 on each entire dataset.**

Scheme \ Dataset	Kosarak	Webdocs
<b>BGR</b>	<b>1.88</b>	<b>184.481</b>
<b>DSR</b>	<b>2.929</b>	<b>226.501</b>
<b>BDR</b>	<b>3.549</b>	<b>197.065</b>
<b>CNR</b>	<b>3.653</b>	<b>196.88</b>