

# Catch a Blowfish Alive: A Demonstration of Policy-Aware Differential Privacy for Interactive Data Exploration

Jiaxiang Liu  
University of Waterloo  
j632liu@uwaterloo.ca

Yiqing Tan  
University of Waterloo  
y57tan@uwaterloo.ca

Karl Knopf  
University of Waterloo  
kknopf@uwaterloo.ca

Bolin Ding  
Alibaba Group  
bolin.ding@alibaba-inc.com

Xi He  
University of Waterloo  
xi.he@uwaterloo.ca

## ABSTRACT

Policy-aware differential privacy (DP) frameworks such as Blowfish privacy allow for more accurate query answers than standard DP. In this work, we build the first policy-aware DP system for interactive data exploration, BlowfishDB, that aims to (i) provide bounded and flexible privacy guarantees to the data curators of sensitive data and (ii) support accurate and efficient data exploration by data analysts. However, the specification and the processing of the customized privacy policies incur additional performance cost especially for datasets with a large domain. To address this challenge, we propose dynamic Blowfish privacy which allows for the dynamic generation of smaller privacy policy and their data representation at query time. BlowfishDB ensures the same levels of accuracy and privacy as one would get working on the static privacy policy. In this demonstration of BlowfishDB, we show how a data curator can fine-tune the privacy policies for a sensitive dataset and how a data analyst can retrieve accuracy-bounded query answers efficiently without being a privacy expert.

## 1. INTRODUCTION

Differential privacy (DP) [3] has arisen as the standard approach for protecting data contributors from privacy violations that could be caused through data exploration. Systems that implement DP mechanisms [12, 8, 9, 5, 13] have been built to support private data analysis. Under DP, the privacy loss is measured by a parameter,  $\epsilon$ , that guarantees that any neighbouring databases that differ in a record have similar output distributions. The smaller the value of  $\epsilon$ , the closer the output distributions, and hence the stronger the privacy guarantee, usually with a degradation to utility.

However, DP is often too strict to offer useful answers in many applications [11, 1], as relaxations of the privacy parameter  $\epsilon$  do not give meaningful utility improvements. A general approach that relaxes DP [7, 1, 14, 2], is to modify the notion of neighbouring inputs by applying a distance metric over the database domain. This relaxation gains a better privacy for utility trade-off. For example, Blowfish privacy [7] includes a privacy “policy graph” over the record domain that specifies which information must be kept secret about individuals to generalize the privacy guarantee of DP.

**Example 1:** Consider a database of individuals’ salary records. Each record takes one value in the domain  $\mathcal{T} = \{1, 2, \dots, 200k\}$ .

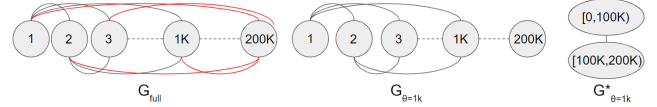


Figure 1:  $G_{full}$ : a fully connected policy graph for DP;  $G_{\theta=1k}$ : a  $\theta$ -distance policy graph for Blowfish privacy, e.g., no edges between (1k, 200k) or (1, 200k); and  $G_{\theta=1k}^*$ : a dynamic Blowfish graph based on a partition  $\{[1, 100k], [100k, 200k)\}$

A DP algorithm ensures that an adversary cannot distinguish whether an individual has a capital gain of  $x$  or  $y$ , for all  $x, y \in \mathcal{T}$ . The corresponding policy graph is a complete graph  $G_{full}$  as shown in Figure 1. We also show a weaker privacy policy graph, a  $\theta$ -distance graph  $G_{\theta=1k}$  that connects all domain values that differ at most  $\theta$ , where  $\theta = 1k$ . This specification intends to prevent the adversary from distinguishing close-by values, e.g. (1k vs. 1.2k), but not far-away values, e.g. (1k vs. 200k). By relaxing this privacy guarantee, the expected error per range query can be greatly improved from  $O(\log^3 |\mathcal{T}|)$  to  $O(\log^3 \theta)$ .  $\square$

Prior work [7, 6, 1, 14] only consider policy graphs over the full domain of a database record. General mechanisms for policy-aware DP [6] require the materialization of the policy graph in a matrix form and hence incur additional storage cost  $O(|\mathcal{T}|^3)$  and computation cost  $O(|\mathcal{T}|^3)$ , where  $|\mathcal{T}|$  is the domain size. These additional costs lead to challenges when extending policy-aware DP to exploration queries over high-dimensional relational data. First, it is unclear how to allow data curators the ability to specify privacy policies over the full domain of a record. There is no tool that helps data curators to understand the privacy-utility trade-off before setting the privacy policies. Second, though running a policy-aware mechanism offers more accurate answers than the standard DP, the additional storage and computation cost can hurt the user experience in an interactive setting.

To address these concerns, we introduce BlowfishDB, a practical data exploration system for policy-aware differential privacy over high dimensional data. BlowfishDB implements a novel form of Blowfish privacy, called *dynamic Blowfish privacy*, a class of privacy policies that can be dynamically generated based on a new partition of the full

domain (e.g.,  $G_{\theta=1k}^*$  in Figure 1). These privacy policies achieve the same privacy guarantee over the full domain while saving significant storage and computation cost. To use this, BlowfishDB includes a dynamic privacy policy generator that converts high-level semantic privacy policies specified by the data curator to lightweight representations for use in an  $\epsilon$ -DP mechanism. The policy generator also allows a data curator to efficiently explore the privacy-utility trade-off of their chosen policy, so they are able to make more informed choices. Data analysts can also efficiently query the sensitive data sets with accuracy requirements and even make more queries using BlowfishDB than under a standard DP system, such as APEX[5].

Two demo scenarios are presented for the attendees. They each describe one aspect of the trade-off between privacy and utility that are managed by the BlowfishDB system. The first allows the attendee to act as a data curator and explore the effects of policy selection on the accuracy of the queries. The second scenario allows the attendee to act as a data analyst, and explore how the privacy policies affect the utility of the query answers.

## 2. BACKGROUND

**Notation.** Consider a relation  $R = (att_1, \dots, att_d)$  with  $d$  attributes. Let  $dom(att_j)$  denote the domain of attribute  $att_j$ , and  $\mathcal{T}$  denote the full domain of a record  $dom(att_1) \times \dots \times dom(att_d)$ . Let the domain size  $|\mathcal{T}| = k$ . Let  $D$  be a database instance consisting of  $n$  records  $\{t_1, \dots, t_n\}$ . Each record  $t_i$  takes a value from  $\mathcal{T}$ . We can represent the database  $D$  by a histogram  $x \in \mathbb{R}^k$  over the full domain  $\mathcal{T}$ . A linear query can be expressed as a linear combination of the counts in  $x$ , i.e.,  $wx$ , where  $w$  is a  $k$ -dimensional row vector of real numbers. A workload of  $q$  linear queries can be then represented as  $W = [w_1, w_2, \dots, w_q]^T \in \mathbb{R}^{q \times k}$ .

**Differential Privacy (DP).** A randomized algorithm  $\mathcal{M}$  satisfies  $\epsilon$ -DP if for all possible output set, and all pairs of *neighboring databases*  $(D, D')$  that differ in a record, we have  $\mathcal{S} \in \text{Range}(\mathcal{M})$ :  $\Pr[\mathcal{M}(D) \in \mathcal{S}] \leq e^\epsilon \Pr[\mathcal{M}(D') \in \mathcal{S}]$ . This privacy guarantee has become a gold standard privacy notion and been used in practice [8].

A common approach to achieve DP is to perturb the query answers directly using such techniques as the Laplace and Gaussian mechanisms [4]. However, this approach can result in large errors in the final output. For linear queries, the error of the query answers can be optimized by using the matrix mechanism [10]. This approach first privately answers a representative strategy query workload  $A$  as  $Ax$  and then reconstruct a solution to the original workload  $W$ .

Formally, let  $A$  be a  $p \times k$  matrix that supports the workload  $W$  [10] such that  $Wx = WA^+Ax$ , where  $A^+$  be the Moore-Penrose pseudo-inverse of  $A$ . Let  $Lap(\sigma)^p$  represents a  $p$ -dimensional vector of independent samples, where each sample is drawn from  $\eta \sim \frac{1}{2\sigma} \exp(\frac{-|\eta|}{\sigma})$ . Then the matrix mechanism can be written as:

$$\mathcal{M}_A(W, x) = WA^+(Ax + Lap(\frac{\Delta_A}{\epsilon})^p) \quad (1)$$

where  $\Delta_A$  denotes the *sensitivity* of a workload  $A$ , that is the maximum difference in the answers to  $A$  between neighboring databases that differ in a record.

The error of a DP mechanism  $\mathcal{M}$  for a linear query  $w$  is commonly measured using mean squared error (MSE) per query [10], i.e.,  $\mathbb{E}((wx - \mathcal{M}(w, x))^2)$ ; or  $(\alpha, \beta)$ -accuracy [5],

i.e., the query error  $|wx - \mathcal{M}(w, x)|$  is no more than  $\alpha$  with a high probability  $(1 - \beta)$ .

**Blowfish Privacy.** Unlike DP that provides a blanket privacy guarantee for all values in the domain which can be too restrictive to offer sufficient utility, Blowfish privacy [7, 6] specifies pairs of domain values which the data curator wish to protect using a *policy graph*. A policy graph over the record domain  $\mathcal{T}$  is a graph  $G = (V, E)$ , where  $V = \mathcal{T}$  and  $E \subseteq V \times V$ . Based on a given policy graph, neighboring databases<sup>1</sup> and Blowfish privacy can be defined as follows.

**Definition 2.1** (Blowfish Privacy). Let  $G = (V, E)$  be a policy graph. An algorithm  $\mathcal{M}$  satisfies  $(\epsilon, G)$ -Blowfish privacy if  $\forall \mathcal{S} \in \text{Range}(\mathcal{M})$ :  $\Pr[\mathcal{M}(D) \in \mathcal{S}] \leq e^\epsilon \Pr[\mathcal{M}(D') \in \mathcal{S}]$  for neighboring databases  $(D, D')$  that differ in the value of exactly one record, such that  $(u, v) \in E$  where  $u$  is the record value in  $D$  and  $v$  is the corresponding record value in  $D'$ .

A broad class of policy graphs are known as  $\theta$ -distance-threshold policy graphs [7], where the edges are defined as  $E = \{(u, v) \mid d(u, v) \leq \theta\}$  for a given distance metric  $d(\cdot, \cdot)$ , (e.g. graphs in Figure 1). When  $\theta = 1$ , we call it *line-graph* policy. More types of policy graphs can be found in [7].

Given an arbitrary graph  $G$ , rather than designing a  $(\epsilon, G)$ -blowfish private mechanism from scratch, we can build such a mechanism from an  $\epsilon$ -DP mechanism [6]. This requires the materialization of the policy graph as a matrix form using  $P_G$  of  $(|V| - 1)$  rows and  $|E|$  columns. Applying an optimal DP mechanism (e.g., matrix mechanism Eqn. (1)) to the transformed workload  $W_G = WP_G$  on the transformed data vector  $x_G = P_G^{-1}x$ , will result in an optimal error for the given privacy policy graph  $G$ . Hence, the computation overhead directly relates to the graph size.

## 3. SYSTEM DESIGN

We design and build the first prototype, BlowfishDB, for policy-aware DP data exploration. This system ensures that (i) bounded and flexible privacy guarantees for the data curators; and (ii) accurate and efficient query processing for the data analyst. In this section, we first present the system overview and then introduce *dynamic Blowfish privacy*, a theoretical framework of Blowfish privacy that allows better performance in the interactive setting.

**System Overview.** BlowfishDB is designed as a flexible software proxy to operate over any standard relational database system. Figure 2 shows the components of BlowfishDB and how it interacts with both the database servers and its users. It connects to the relational database through a SQL interface. Privacy policies are specified by the data curator and are controlled through a policy manager module. A privacy engine, which creates accuracy aware private queries, allows data analysts the ability to explore the sensitive data without requiring extensive privacy knowledge. Through its interfaces, BlowfishDB provides its users with an complete private data exploration system.

**Frontend Interfaces.** When the data curator first saves sensitive data  $D$  at a standard DBMS server, the privacy policy will default to DP, the strongest privacy level. The data curator can then adjust the privacy policies at the attribute level to a desired policy graph  $G$  and set the total privacy budget  $B$  for the data exploration. BlowfishDB

<sup>1</sup>This assumes known data size and no other database constraints. Variants of the definition can be found in [6].

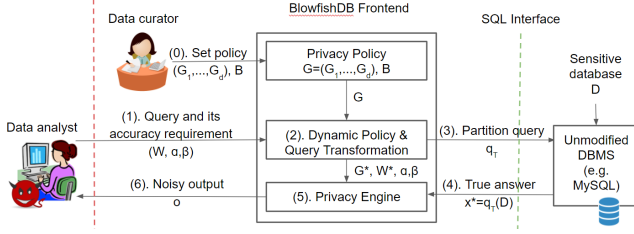


Figure 2: BlowfishDB overview. (0) set privacy policy  $G$  and budget  $B$ ; (1-2) generate dynamic policy and query; (3-4) query over DBMS; (5-6) run DP mechanism that outputs  $o$  with  $(\alpha, \beta)$ -accuracy guarantee.

will save this information to control the privacy loss when the data analyst queries the sensitive data. In addition, the data curator is allowed to query the databases and test the accuracy-privacy tradeoffs (Figure 3) at different privacy policies and privacy budgets, so that they can make informed decision on the policies based on the explorations.

A data analyst interface is also provided. The data analyst is provided with a set of exploration query templates to choose from. They then select a query  $W$  with its accuracy requirement  $(\alpha, \beta)$ . BlowfishDB will generate a projected workload, data vector query, and privacy policy  $(W^*, x^*, G^*)$  based on  $W$  and  $G$  which allows more efficient processing than the static version  $(W, x, G)$ . Then BlowfishDB will construct and send a SQL query  $q_T(D)$  that corresponds to answers for the partition on the database  $x^*$ . The accuracy translation engine then determines an  $(\epsilon, G^*)$ -blowfish private algorithm that runs on  $(W^*, x^*)$  with an  $(\alpha, \beta)$ -accuracy guarantee. In the prototype, BlowfishDB also offers noisy plots of the query answers and provides a log of the previously executed queries for comparison.

Next, we will explain how to construct suitable dynamic privacy policy  $G^*$  to improve the performance of BlowfishDB.

**Dynamic Blowfish Policy Projection.** To optimize performance, BlowfishDB dynamically generates a policy graph based on the query workload  $W$  and the static policy graph  $G$ . For example, given a policy graph  $G$  on a high-dimensional domain  $\mathcal{T} = \text{dom}(att_1) \times \dots \times \text{dom}(att_d)$ , if a workload  $W$  only conditions on a single attribute  $att_i$ , we can project the workload, the data, and the policy graph onto a partitioned domain based on  $att_i$ . The projected policy graph will have only  $|\text{dom}(att_i)|$  nodes instead of  $|\mathcal{T}|$  nodes.

**Definition 3.1.** [Projected Privacy Policy] Let  $\mathcal{T}^*$  be a partition over  $\mathcal{T}$ . Given a Blowfish privacy policy graph  $G = (V, E)$  over the full domain  $\mathcal{T}$ , the projected privacy policy  $G^*$  based on  $\mathcal{T}^*$  is defined as  $G^* = (V^*, E^*)$ , where the node set  $V^*$  is the partitioned domain  $\mathcal{T}^*$ , and the edge set  $E^*$  includes an edge  $(v_1^*, v_2^*)$  if exists an edge  $(v_{i_1}, v_{i_2}) \in E$  such that  $v_{i_1} \in v_1^*$  and  $v_{i_2} \in v_2^*$ .

We can represent a partition  $\mathcal{T}^*$  by a partition matrix  $T$ , where  $T[i, j] = 1$  if  $v_j \in \mathcal{T}$  is a value in the  $i$ th bin of the partition  $\mathcal{T}^*$ . Then the projected data vector  $x^* = Tx$  represents the counts for the bins in the partition, where  $x$  is the histogram of the full domain  $\mathcal{T}$ . We say a partition  $\mathcal{T}^*$  supports a workload query  $W$  if each query in  $W$  can be expressed as a linear combination of the counts in  $Tx$ , in particular  $WT^+Tx = Wx$ . In practice, we can directly construct the projected workload  $W^* = WT^+$  and the

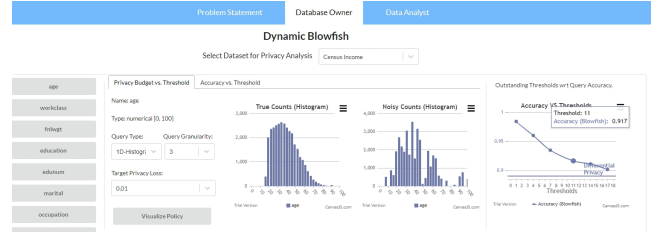


Figure 3: A snapshot of data owner interface that displays data schema, query template, true query answer plot, noisy answer plot, utility-privacy tradeoff plot, and policy visualization feature.

projected data vector  $x^* = Tx$  based on the new partition without materializing any of  $W, x, T$ .

**Example 2:** Continue from Example 1. A cumulative workload  $W$  at granularity of 100k, i.e.,  $\{[1, 100k], [1, 200k]\}$ , can be supported by a partition  $\{[1, 100k] \text{ and } [100k, 200k]\}$ . We can directly construct the projected data vector  $x^*$  using 2 selection queries and materialize the corresponding projected workload  $W^* = \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix}$ . The projected policy graph  $G_{\theta=1k}^*$  for  $G_{\theta=1k}$  will involve only 2 nodes as shown in Figure 1. It is not hard to verify that  $Wx = WT^+Tx = W^*x^*$ , where  $W = \begin{pmatrix} 1 \dots 1 & 0 \dots 0 \\ 1 \dots 1 & 1 \dots 1 \end{pmatrix}$  and  $T = \begin{pmatrix} 1 \dots 1 & 0 \dots 0 \\ 0 \dots 0 & 1 \dots 1 \end{pmatrix}$ .  $\square$

When the partition size is much smaller than the full domain size, the projected set of data vector, workload, and graph  $(x^*, W^*, G^*)$  will have a smaller representation than  $(x, W, G)$ , and results in a much faster computation time. In addition, directly applying an optimal  $\epsilon$ -DP matrix mechanism (e.g., Eqn. (1)) on  $(x^*, W^*, G^*)$  achieves  $(\epsilon, G)$ -Blowfish privacy with the same accuracy guarantee as on  $(x, W, G)$  (see Theorem 3 in the Appendix).

**Dynamic Blowfish Policy Composition.** To allow the data curator flexibly specify privacy policies, BlowfishDB consider a special class of Blowfish privacy policies that can be composed from a set of attribute-based privacy policies.

**Definition 3.2** (Attribute Composability). Given a policy  $G(V, E)$ , let  $G_{att_1}(V_1, E_1), \dots, G_{att_d}(V_d, E_d)$  denote the projected privacy policies from  $G$  onto the respective attribute partitions. We say  $G$  is *attribute composable* if  $G(V, E) = G_c(V_c, E_c)$ , where  $V_c = V_1 \times \dots \times V_d$  and  $E_c$  includes edge  $(u, v) \in V_c \times V_c$  if  $\sum_{j=1}^d \text{dist}_{G_{att_j}}(u.A_j, v.A_j) = 1$ , and the distance function  $\text{dist}_{G_{att_j}}(u.A_j, v.A_j)$  denotes the shortest distance between  $u$  and  $v$  in  $G_{att_j}$ .

An example for attribute composition is illustrated in Figure 4. Note this is not the only possible composition function. For this class of attribute composable privacy policies, the data curator can customize the privacy policy of each attribute one by one. For a given workload query  $W$  that can be supported by the partition over a subset of attributes  $\mathcal{T}^* = \{att_{i_1}, \dots, att_{i_j}\}$ , we can project the workload and data over  $\mathcal{T}^*$  and form the corresponding policy graph  $G_{\mathcal{T}^*}$  by composing the relevant attribute policies  $C(G_{att_{i_1}}, \dots, G_{att_{i_j}})$ . In an interactive setting, each query workload works with different dynamic privacy policies, but the total privacy loss can be bounded as follows.

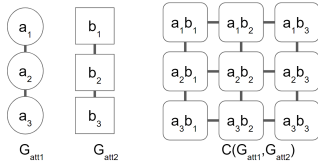


Figure 4: Compositing attribute-based policies  $G_{att_1}$  and  $G_{att_2}$  gives  $C(G_{att_1}, G_{att_2})$ .

**Theorem 1.** If  $G$  is attribute composable, let  $\mathcal{M}_1, \dots, \mathcal{M}_l$  be a sequence of algorithms, if each  $\mathcal{M}_i$  satisfies  $(\epsilon_i, G_{T_i^*})$ -Blowfish privacy for a partition  $T_i^*$ , then the overall privacy loss is  $(\sum_{i=1}^l \epsilon_i, G)$ -Blowfish privacy.

**Limitations and Future Work.** The prototype of BlowfishDB provided for the demonstration only supports (i) 1/2-D histogram queries and cumulative histogram queries and (ii) threshold policy graphs. The query templates supported by BlowfishDB also suggest a good partition choice. Besides expanding the scope of the queries and privacy policies, we will optimize the performance and accuracy over different choices of partitions and mechanisms as future work.

## 4. DEMO EXPERIENCE

An attendee can choose from two demo scenarios for the BlowfishDB system, which are presented as tabs on the interface. The first takes the perspective of a data curator, who wishes to explore the accuracy vs privacy policy trade-off. They will be able to explore the affect of private mechanisms on query outputs, and will be able to explore the dataset so they can set an appropriate privacy policy. The second scenario takes the perspective of a data analyst, who wishes to explore the privacy vs utility trade-off. They will be able to run a variety of queries to see how a fixed privacy policy is able to affect the quality of their data exploration.

**Demo 1: Data Curator Interface.** We will provide a set of datasets in .csv format that attendees can explore. The attendee, acting as a data curator will select one of these datasets, prompting a menu listing the data set attributes to appear. The attendee will then select one of these attributes, as well as a query type. The true answers to the chosen query will be displayed as a bar chart as shown in Figure 3.

The attendee can then begin the trade-off exploration by first choosing a privacy budget value  $\epsilon$ . This will cause a noisy answer under  $\epsilon$ -differential privacy to be calculated for this query. The result will be plotted besides the true answer plot, so a direct comparison can be made.

The data curator can then choose to display the current privacy policy, to allow for visual exploration. By default, the policy is equivalent to DP so a fully connect graph over the chosen attribute is displayed. It is then possible to relax this policy by specifying a threshold value, which will update the display. The attendee can then select a set of possible threshold values, for which the total privacy loss will be calculated. These will be displayed as a line chart, where the attendee can then select a node representing a specific threshold value to display a noisy answer under that threshold policy as a histogram plot. By carrying out these steps, the attendee will be able to interactively explore the affect of privacy policies on the queries accuracy, and will be able to make an informed decision on how they would set an appropriate policy to protect their data.

**Demo 2: Data Analyst Interface.** Similarly to the previous demo, the attendee, acting as a data analyst will select one of these datasets, prompting a menu listing the data set attributes to appear. The attendee will then select one of these attributes, as well as a query type and granularity of the query. The analyst will also be asked to specify a set of accuracy requirements,  $\alpha$  and  $\beta$ , for their queries. A plot of the noisy answer to the chosen query will be displayed.

The analyst can then see their remaining privacy budget and chose to make another query if there is sufficient budget to do so. This process can be repeated using different attributes or query types. The analyst is restricted to 1-D histogram, 1-D cumulative histogram and 2-D histogram queries in the prototype. All prior query results are stored in a table that is visible to the data analyst. They are able to display the results of up to two previous noisy queries, so they are able to compare the results.

If a policy was specified for an attribute using the data curator demo scenario, then it will be applied in this demo. Otherwise, or if the system was reset in-between scenarios, a default differential privacy policy will be used. An attendee can then run this scenario under different policies to be able to see the difference in privacy consumption, as well as the difference in the accuracy of the noisy answers.

## 5. REFERENCES

- [1] M. E. Andrés, N. E. Bordenabe, K. Chatzikokolakis, and C. Palamidessi. Geo-indistinguishability: Differential privacy for location-based systems. In *SIGSAC*, 2013.
- [2] K. Chatzikokolakis, M. E. Andrés, N. E. Bordenabe, and C. Palamidessi. Broadening the scope of differential privacy using metrics. In *PoPETS*, 2013.
- [3] C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private data analysis. In *TCC*, 2006.
- [4] C. Dwork, A. Roth, et al. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 2014.
- [5] C. Ge, X. He, I. F. Ilyas, and A. Machanavajjhala. APEX: Accuracy-aware differentially private data exploration. In *SIGMOD*, 2019.
- [6] S. Haney, A. Machanavajjhala, and B. Ding. Design of policy-aware differentially private algorithms. *VLDB*, 2015.
- [7] X. He, A. Machanavajjhala, and B. Ding. Blowfish privacy: Tuning privacy-utility trade-offs using policies. In *SIGMOD*, 2014.
- [8] N. Johnson, J. P. Near, and D. Song. Towards practical differential privacy for SQL queries. *VLDB*, 2018.
- [9] I. Kotsogiannis, Y. Tao, X. He, M. Fanaeepour, A. Machanavajjhala, M. Hay, and G. Miklau. Privatesql: a differentially private sql query engine. *VLDB*, 2019.
- [10] C. Li, G. Miklau, M. Hay, A. McGregor, and V. Rastogi. The matrix mechanism: optimizing linear counting queries under differential privacy. *The VLDB journal*, 2015.
- [11] A. Machanavajjhala, A. Korolova, and A. D. Sarma. Personalized social recommendations-accurate or private? *VLDB*, 2011.
- [12] F. D. McSherry. Privacy integrated queries: an extensible platform for privacy-preserving data analysis. In *SIGMOD*, 2009.
- [13] P. Mohan, A. Thakurta, E. Shi, D. Song, and D. Culler. GUPT: privacy preserving data analysis made easy. In *SIGMOD*, 2012.
- [14] Z. Xiang, B. Ding, X. He, and J. Zhou. Linear and range counting under metric-based local differential privacy. In *ISIT*, 2020.

## APPENDIX

**Blowfish Privacy Transformational Equivalence.** We will introduce a basic transformation equivalence result for the matrix mechanism. A policy graph  $G$  is represented using a matrix  $P_G$  of  $(|V| - 1)$  rows and  $|E|$  columns: first randomly pick a node  $v^* \in E$ ; for every  $(u, v) \in E \setminus \{v^*\}$ , add a column to  $P_G$  with a 1 in the row corresponding to value  $u$ , and  $-1$  in the row corresponding to value  $v$  (order of 1 and  $-1$  not important), and zeros in the rest of the rows; for every  $(u, v^*) \in E$ , add a column with a 1 in the row corresponding to  $u$  and zeros in the rest. We can transform the given workload  $W$  and  $x$  by  $P_G$  as  $W_G = WP_G$  on  $x_G = P_G^{-1}x$  and apply matrix mechanism.

**Theorem 2.** (Theorem 4.1 from [6]) If  $P_G$  has full row rank (and hence a right inverse  $P_G^{-1}$ ), then the matrix mechanism given in Eqn. (1) is both a  $(\epsilon, G)$ -Blowfish private mechanism for answering  $W$  and  $x$  and an  $\epsilon$ -DP algorithm for answering  $W_G$  and  $x_G$ . This mechanism also has the same error under both privacy definitions.

**Transformational Equivalence for Dynamic Blowfish Privacy.** Given the domain partition  $\mathcal{T}$ , the policy graph  $G$  can also be projected accordingly. We denote the projected workload and data vector on the new partition by  $W^* = WT^+$  and  $x^* = Tx$  respectively.

In practice, the matrix presentation  $P_{G^*}$  is directly materialized for the projected policy  $G^*$ , which is much smaller than  $P_G$ . It can also be obtained by first multiplying the partition matrix  $T$  with  $P_G$ , and then remove all resulted the zero columns in  $TP_G$ .

**Theorem 3.** Given a workload  $W$  over the full domain  $\mathcal{T}$  and a privacy policy graph  $G$  which has a full row rank matrix representation  $P_G$ . Let  $T$  be a partition matrix that supports  $W$ . Consider the matrix mechanism  $\mathcal{M}_A$  defined in Eqn. (1), this mechanism is  $(\epsilon, G^*)$ -Blowfish private for answering  $W^* = WT^+$  on  $x^* = Tx$ , an  $(\epsilon, G)$ -Blowfish private for answering  $W$  on  $x$ , an  $\epsilon$ -differentially private for answering  $W_G = WP_G$  on  $x_G = P_G^{-1}x$ , and  $\epsilon$ -differentially private for answering  $W_{G^*} = W^*P_{G^*}$  on  $x_{G^*} = P_{G^*}^{-1}x^*$ . The mechanism has the same error for all four cases.

*Proof.* We first construct a policy graph  $G'(V', E')$  based on  $G^*(V^*, E^*)$ , let  $V' = \mathcal{T}$  where  $\mathcal{T}$  is the database domain before partitioning, and  $E' = \{(v'_1, v'_2) \in \mathcal{T} \times \mathcal{T} \mid \exists (v_1^*, v_2^*) \in E^* \text{ s.t. } v'_1 \in v_1^* \wedge v'_2 \in v_2^*\}$

For a  $(\epsilon, G^*)$ -Blowfish private mechanism  $\mathcal{M}_A$ , we have:

$$Pr[\mathcal{M}(D) \in \mathcal{S}] \leq e^\epsilon Pr[\mathcal{M}(D') \in \mathcal{S}]$$

for  $(D, D') \in \mathcal{N}(G^*)$

Based on the definition of Blowfish neighbors under bounded DP, and the construction of  $G'$ , it follows that if  $(D, D') \in \mathcal{N}(G^*)$ , then  $(D, D') \in \mathcal{N}(G')$ . Thus, it follows that  $\mathcal{M}_A$  is  $(\epsilon, G')$  Blowfish private.

By Definition 3.1,  $V' = V$  and  $E \subseteq E'$ . Using the same argument from Lemma C.2 in [6], we can show that  $\mathcal{M}_A$  is  $(\epsilon, G)$  Blowfish private.

We will now show that the accuracy is the same for all three cases. First, for a  $(\epsilon, G)$ -Blowfish private mechanism  $\mathcal{M}_A$  for answering  $W$  on  $x$ , and a  $(\epsilon, G^*)$ -Blowfish private mechanism  $\mathcal{M}_A$  for answering  $W^*$  on  $x^*$ :

$$Wx + WA^+ Lap\left(\frac{\Delta_A(G)}{\epsilon}\right) = W^*x^* + W^*(A^*)^+ Lap\left(\frac{\Delta_{A^*}(G^*)}{\epsilon}\right)$$

In addition,

$$W^*x^* = WT^+Tx = Wx$$

Next,  $\Delta_{A^*}(G^*) = \Delta_A(G)$  by construction of  $G^*$  from  $G$ .

$$\begin{aligned} W^*(A^*)^+ Lap\left(\frac{\Delta_{A^*}(G^*)}{\epsilon}\right) &= WT^+(AT^+)^T Lap\left(\frac{\Delta_{A^*}(G^*)}{\epsilon}\right) \\ &= WT^+TA^+ Lap\left(\frac{\Delta_{A^*}(G^*)}{\epsilon}\right) \\ &= WA^+ Lap\left(\frac{\Delta_{A^*}(G^*)}{\epsilon}\right) \end{aligned}$$

Thus, following the same argument as Theorem 4.1 in [6], we have:

$$Wx + WA^+ Lap\left(\frac{\Delta_A(G)}{\epsilon}\right) = W_Gx_G + W_GA_G^+ Lap\left(\frac{\Delta_{A_G}}{\epsilon}\right)^p \quad \square$$

**Theorem 1.** If  $G$  is attribute composable, let  $\mathcal{M}_1, \dots, \mathcal{M}_l$  be a sequence of algorithms, if each  $\mathcal{M}_i$  satisfies  $(\epsilon_i, G_{T_i^*})$ -Blowfish privacy for a partition  $T_i^*$ , then the overall privacy loss is  $(\sum_{i=1}^l \epsilon_i, G)$ -Blowfish privacy.

*Proof.* This follows by Theorem 3 and Composition rule of Blowfish privacy (Theorem 4.1 in [7]).  $\square$