

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

Jiaxiang Liu<sup>1</sup>, Karl Knopf<sup>1</sup>, Yiqing Tan<sup>1</sup>, Bolin Ding<sup>2</sup>, Xi He<sup>1</sup>

<sup>1</sup>University of Waterloo <sup>2</sup>Alibaba Group



## Overview

**Problem:** Customized privacy policies can incur performance cost for sensitive dataset over a large domain.

**Motivation:** Answer queries accurately and efficiently with customized and provable privacy guarantees.

**Contribution:**

- A new privacy framework, Dynamic Blowfish privacy, that adaptively generates privacy policies at query time.
- A privacy-preserving system, BlowfishDB, that allows data curators to set policy and data analysts to query data.

## Research Objective

Data curators wish to release statistics with privacy guarantees; Analysts wish to access sensitive information with accuracy and performance requirements.

	DP	Blowfish	Dynamic Blowfish
Accuracy	☆	☆☆☆	☆☆☆
Usability	☆☆☆	☆	☆☆
Performance	☆☆☆	☆	☆☆

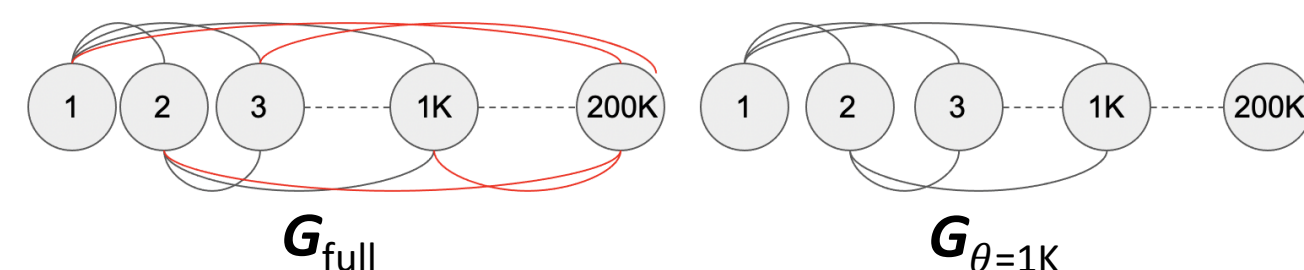
## Background

**Differential Privacy (DP) [DMN06]:** A randomized algorithm  $M$  satisfies  $\epsilon$ -DP if for all possible output sets, and any pair of neighboring databases ( $D_1, D_2$ ) that differ in one record, we have

$$S \in \text{Range}(M) : \Pr [M(D_1) \in S] \leq e^\epsilon \Pr [M(D_2) \in S]$$

**Policy Graph:**  $G = (V, E)$  is a secret discriminative graph over domain  $T$

- $V$ : a set of domain values  $V = T$   
e.g. Alice's salary is 10K
- $E$ : sensitive pairs of values  $V_{\text{pair}} \in V \times V$   
e.g. (Alice's salary is 10K, Alice's salary is 20K)



**Blowfish Privacy [HMD14]:** Let  $G = (V, E)$  be a policy graph. An algorithm  $M$  satisfies  $(\epsilon, G)$ -Blowfish Privacy if

$$S \in \text{Range}(M) : \Pr [M(D_1) \in S] \leq e^\epsilon \Pr [M(D_2) \in S]$$

for neighboring databases ( $D_1, D_2$ ) that differ in one record, such that  $(u, v) \in E$  where  $u$  is the record value in  $D_1$  and  $v$  is the corresponding record value in  $D_2$ .

**$(\alpha, \beta)$ -Accuracy:** The distance between the noisy output and true query answer is no more than  $\alpha$  with a high probability  $(1-\beta)$ .

## Dynamic Blowfish Privacy

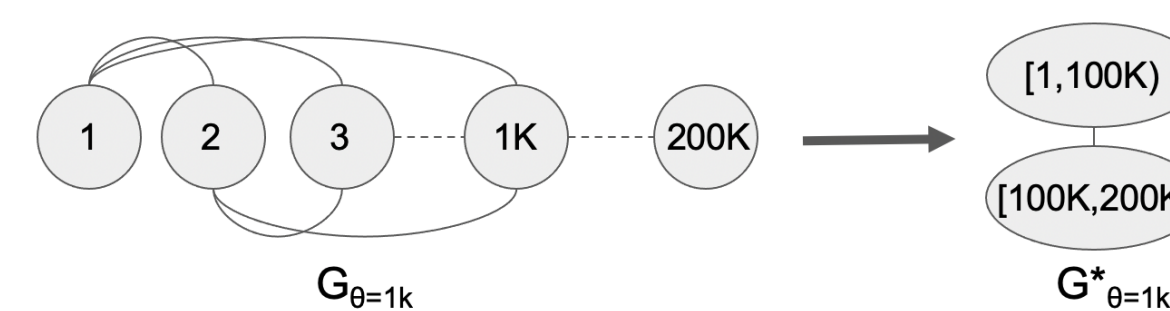
**Matrix Representation of Query:** Linear counting queries can be transformed into matrices:

- $x$ : a histogram over the full domain  $T$   
e.g. the count over each possible salary in  $\{1, \dots, 200K\}$
- $W$ : linear combinations of counts in  $x$   
e.g. the workload of range  $[1, 100K)$  and  $[1, 200K)$  is

$$W = \begin{pmatrix} 1 \dots 1 & 0 \dots 0 \\ 1 \dots 1 & 1 \dots 1 \end{pmatrix}$$

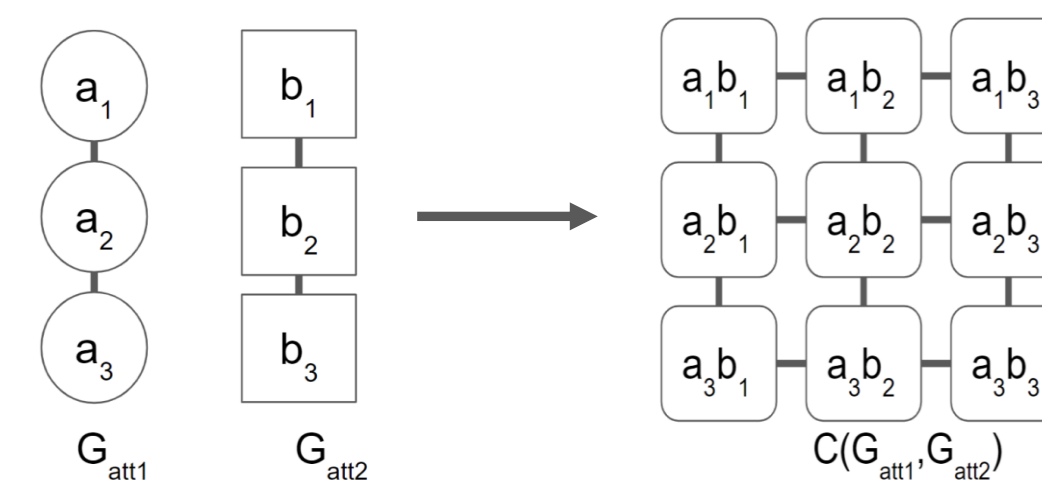
**Dynamic Blowfish Policy Projection:** Partition vertices in  $G$  based on the partition matrix  $T$  generated at query time.

- $x^*$ : the corresponding data vector partitioned by  $T$   
e.g.  $\{1, \dots, 200K\}$  into  $\{[1, 100K), [100K, 200K)\}$
- $W^*$ : the corresponding workload partitioned by  $T$
- $G^*$ : The policy graph  $G$  projected onto the partitioned domain

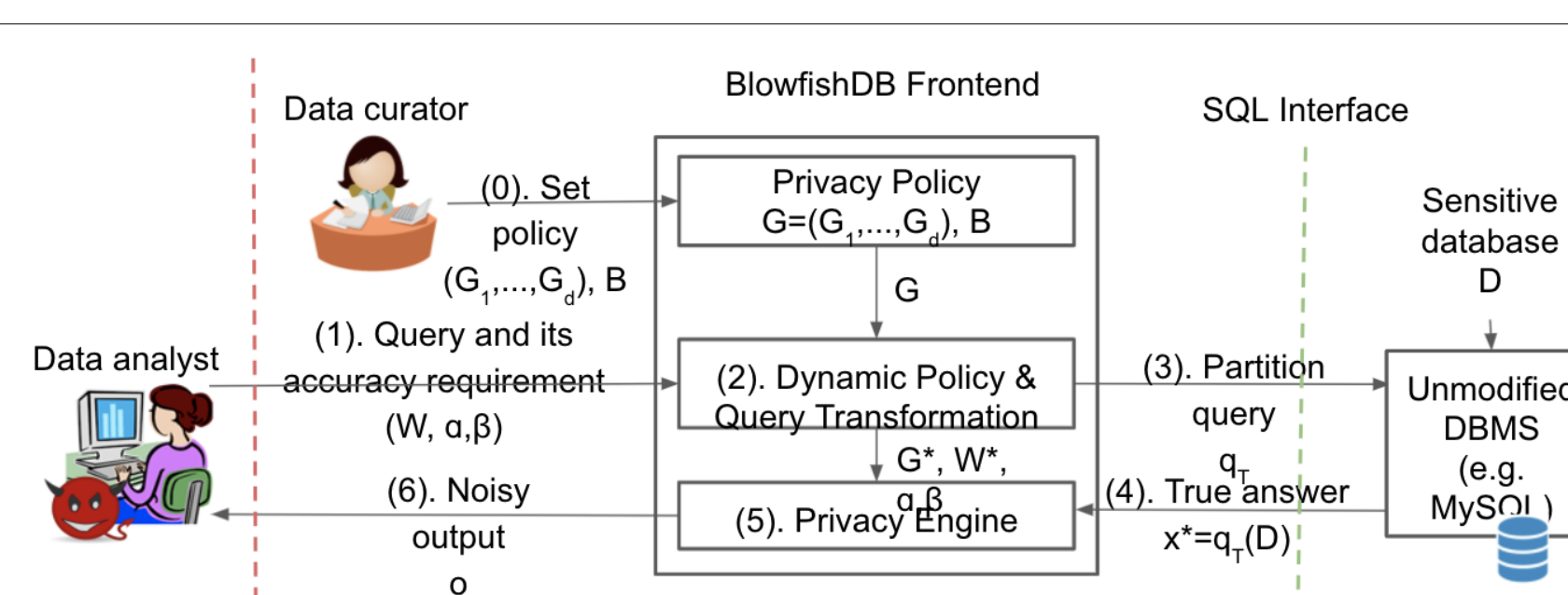


**Privacy and Accuracy:** Let  $M$  be a matrix mechanism that answers  $W$  under  $(\epsilon, G)$ -Blowfish Privacy, applying an optimal  $\epsilon$ -DP mechanism  $M^*$  that answers  $W^*$  also satisfies  $(\epsilon, G)$ -Blowfish Privacy on  $W$  with the same error.

**Attribute Policy Composition:** We create higher dimensional privacy policies by composing single attribute privacy definitions.



## BlowfishDB Overview



- Privacy policy exploration from data curator
- High-level language with accuracy requirements from data analysts
- Query-dependent partition matrix  $T$  for efficient query processing
- Overall privacy loss is bounded by privacy set by the data curator

## Evaluation

**Dataset:** Adult (Income) dataset

**Query:** Histogram/Cumulative histogram with varying granularities  $G_5/G_{10}$

**Privacy Policy:** DP, (Dynamic) Blowfish Privacy (distance threshold policies with varying thresholds, T1, T5, and T10)

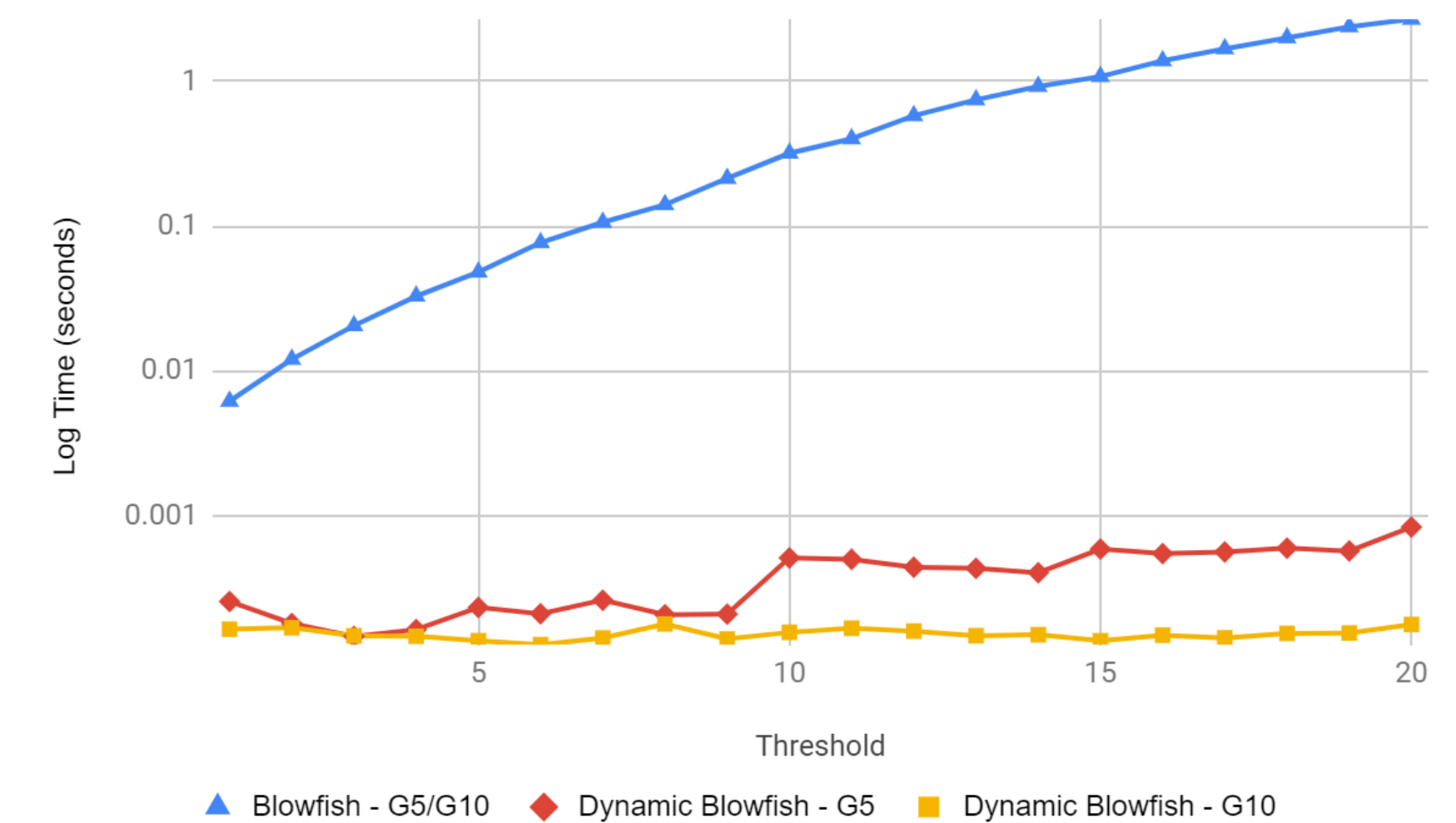


Figure 1: **Performance vs Privacy Guarantee.** For privacy guarantee characterized by each threshold, runtime performance is reported for significant thresholds.

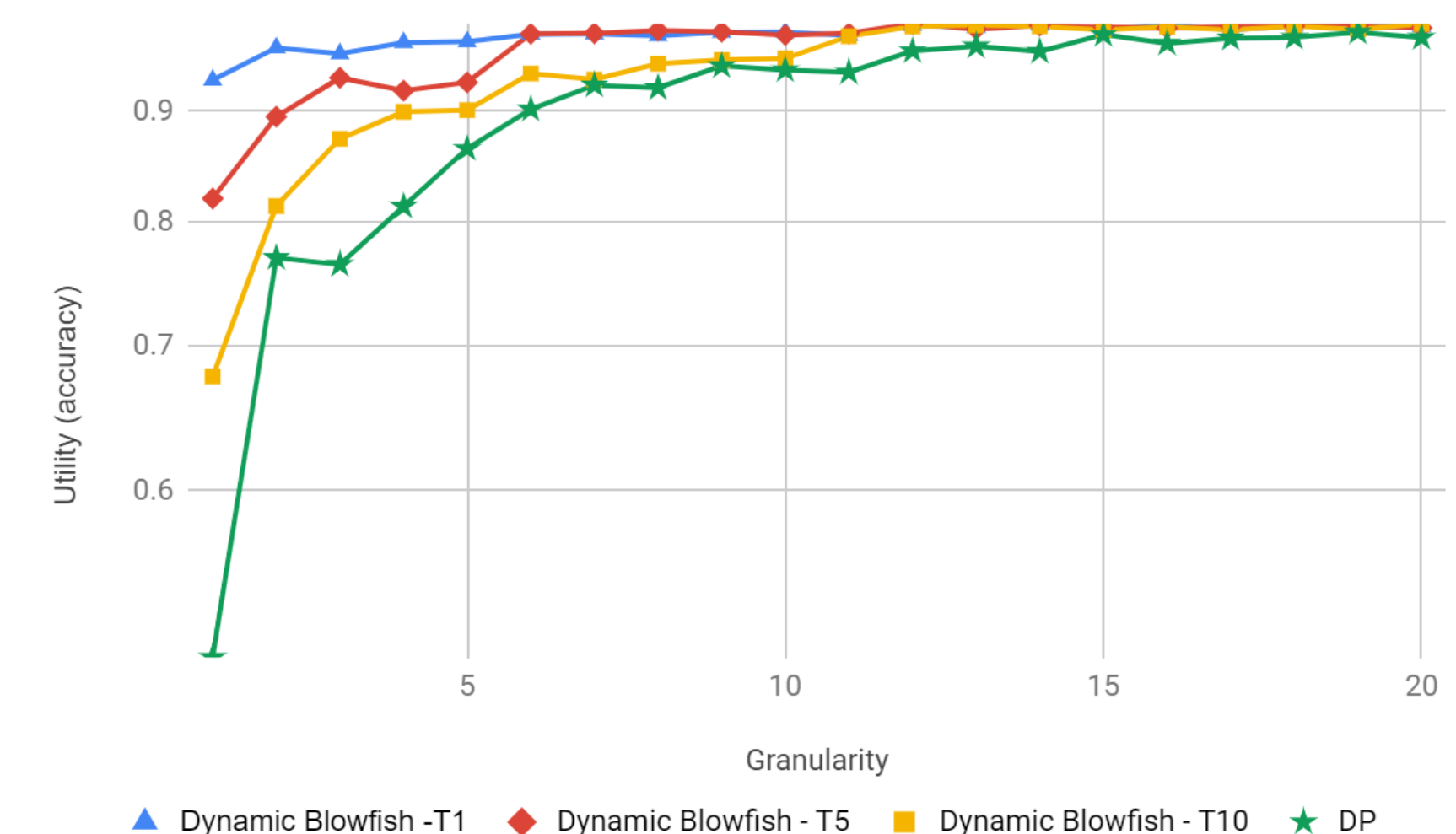


Figure 2: **Utility vs Query Granularity (Cumulative Histogram).** For partitioning matrix  $T$  characterized by each granularity, accuracy is reported for significant thresholds.

## Conclusions

**Conclusion:**

- Dynamic Blowfish achieves better utility versus DP and better performance than Blowfish.
- BlowfishDB provides a way for data curators to set better policies and data analysts to retrieve more accurate results.

**Future work:**

- Incorporate more complex privacy policies and queries.