

Authenticating Aggregate Range Queries over Dynamic Multidimensional Dataset

Jia Xu

National University of Singapore
Department of Computer Science
xujia@comp.nus.edu.sg

Abstract. We are interested in the integrity of the query results from an outsourced database service provider. Alice passes a set \mathbf{D} of d -dimensional points, together with some authentication tag \mathbf{T} , to an untrusted service provider Bob. Later, Alice issues some query over \mathbf{D} to Bob, and Bob should produce a query result and a proof based on \mathbf{D} and \mathbf{T} . Alice wants to verify the integrity of the query result with the help of the proof, using only the private key. In this paper, we consider aggregate query conditional on multidimensional range selection. In its basic form, a query asks for the total number of data points within a d -dimensional range. We are concerned about the number of communication bits required and the size of the tag \mathbf{T} . Xu and Chang [1] proposed a new method to authenticate aggregate count query conditional on d -dimensional range selection over static dataset, with $O(d^2 \log^2 N)$ communication bits, where N is the number of points in the dataset \mathbf{D} . We extend their method to support other types of queries, including summing, finding of the minimum/maximum/median and usual (non-aggregate) range selection, with similar complexity. Furthermore, dynamic operations, like insertion and deletion, over the outsourced dataset are also supported.

Keywords: Authentication, Multidimensional Aggregate Query, Secure Outsourced Database, Dynamic Database, Count, Sum, Average, Min, Max, Median, Range Selection

1 Introduction

Alice has a set \mathbf{D} of d -dimensional points. She preprocesses the dataset \mathbf{D} using her private key to generate some authentication tag \mathbf{T} . She sends (outsources) \mathbf{D} and \mathbf{T} to an untrusted service provider Bob. Then Alice deletes the original copy of dataset \mathbf{D} and tag \mathbf{T} from her local storage. Later Alice may issue a query over \mathbf{D} to Bob, for example, an aggregate query conditional on a multidimensional range selection, and Bob should produce the query result and a proof based on \mathbf{D} and \mathbf{T} . Alice wants to authenticate the query result, using only her private key. This problem fits in the framework of the outsourced database applications [2, 3], which emerged in early 2000s as an example of “software-as-a-service” (SaaS).

We are concerned about the communication cost and the storage overhead on Alice/Bob’s side. Such requirements exclude the following straightforward approaches: (1) Bob sends back the whole dataset \mathbf{D} with its tag \mathbf{T} ; (2) Alice keeps a local copy of the dataset; (3) During preprocessing, Alice generates and signs answers to all possible queries.

Very recently, Xu and Chang [1] proposed a new method to authenticate aggregate count query over a static d -dimensional outsourced dataset, with $O(d^2 \log^2 N)$ communication bits, where N is the number of points in the dataset. Their method combined two customer designed primitives: (1) a **GKEA** (Generalized Knowledge of Exponent Assumption [4]) based homomorphic authentication tag; (2) a functional encryption scheme supporting multidimensional range query. In this paper, we extend their method in two directions without sacrificing the communication complexity: (1) support other types of queries, including summing, finding of the minimum/maximum/median and usual (non-aggregate) range selection; (2) support dynamic operations like insertion and deletion over the outsourced dataset.

Table 1: Worst case performance of different authentication schemes for aggregate range query or range selection query. This table consists of two parts: the first three rows are for aggregate query; the rest four rows are for range selection query.

Note: (1) The symbol “-” indicates that the authors do not provide such information in their paper. (2) Our scheme is much more efficient in computation cost in 1D case, compared with high dimensional case (See annotation \star). (3) $dN \log N \leq \log^d N$, if $d > \log(dN)/\log \log N$. We point out that the high computation cost on prover can be mitigated with horizontal partition of the dataset and parallel execution on each partition. (4) We do not include [5, 6] in this table, since these works do not provide concise asymptotic bound on their schemes. However, their performances are limited by the underlying data structure they adopted, i.e. KD-tree [5] and R-Tree [6], which require exponential (in dimension) communication overhead in the worst case. (5) Our scheme supports private key verification, while the other works in this table support public key verification.

Scheme	Dimension d	Communication overhead (bits)	Storage overhead	Computation (Verifier Alice)	Computation (Prover Bob)	Query	Techniques
PDAS [7]	$d = 1$	$O(S \log N)$	$O(N)$	$O(S \log N)$	$O(S + K^2)$	SUM,COUNT	Aggregated commitment + Shamir’s Secret-Sharing Scheme
Li <i>et al.</i> [8]	$d \geq 1$	$O(dN + 2^d)$	$\Omega(dN)$	$O(dN + 2^d)$	$\Omega(N^{1-\frac{1}{d}})$	SUM or COUNT or MIN or MAX (One authentication data structure per query type)	MHT-like authentication structure for B-Tree/R-Tree
This paper and Xu <i>et al.</i> [1]	$d \geq 1$	$O(d^2 \log^2 \mathcal{Z})$	$O(dN)$	$O(d^2 \log^2 \mathcal{Z})\dagger\star$	$O(dN \log \mathcal{Z})\ddagger\star$	SUM,COUNT,MIN,MAX, MEDIAN	(customer designed) functional encryption + GKEA based homomorphic tag
Atallah <i>et al.</i> [9]	$d = 1, 2$	$O(1)$	$O(N)$	$O(S)$	$O(1)$	Range Selection	Precomputed prefix sum + BLS signature
Martel <i>et al.</i> [10]	$d \geq 1$	$O(\log^{d-1} N + S)$	-	-	-	Range Selection	Authentication Data Structure + Geometry Partition
Chen <i>et al.</i> [11]	$d \geq 1$	$O(\log^d \mathcal{Z})$	$O(N \log^d \mathcal{Z})$	$O(\log^d \mathcal{Z})$	$O(\log^d \mathcal{Z})$	Range Selection	Authentication Tree Structure + Access Control
This paper	$d \geq 1$	$O(d^2 \log^2 \mathcal{Z})$	$O(dN)$	$O(d^2 \log^2 \mathcal{Z} + S)\dagger\star$	$O(dN \log \mathcal{Z} + S)\ddagger\star$	Range Selection	(customer designed) functional encryption + GKEA based homomorphic tag
This paper	$d \geq 1$	$O(d^2 \log^2 \mathcal{Z})$	$O(dN \cdot 2^d)$	$O(d^2 \log^2 \mathcal{Z} + S)\dagger\star$	$O(dN \log \mathcal{Z} + S)\ddagger\star$	Range Selection with projection	(customer designed) functional encryption + GKEA based homomorphic tag

N : The number of tuples in the dataset.
 K : The number of servers in PDAS [7].
 \dagger : $O(d^2 \log^2 \mathcal{Z})$ group multiplications.
 \star : If the query range is 1D, the cost is $O(|S|)$.

S : The set of tuples satisfying the query condition.
 \mathcal{Z} : The domain size of attributes/points in one dimension.
 \ddagger : $O(dN \log \mathcal{Z})$ bilinear map operations.

1.1 Contribution

The main contribution of this paper can be summarized as below.

1. We propose a method to authenticate aggregate queries over static multidimensional dataset, including SUM, MIN, MAX, MEDIAN, with $O(d^2 \log^2 \mathcal{Z})$ communication bits, based on [1]. We prove that the new authentication method is secure.
2. We propose a method to authenticate range selection query over multidimensional static dataset, with $O(d^2 \log^2 \mathcal{Z})$ communication bits, based on [1]. We prove that the new authentication method is secure.
3. We propose a method to authenticate aggregate range query and non-aggregate range selection query over dynamic multidimensional dataset. We prove that the proposed method is secure.
4. We extend our method to support privacy protection and prevent frame attack.

The comparison between our result and previous work is given in Table 1.

2 Related work

Researches in secure outsourced database focus on two major aspects: (1) privacy (i.e. protect the data confidentiality against both the service provider and any third party) e.g. [3,12,13,14], and (2) integrity (i.e. authenticate the soundness and completeness of query results returned by the service provider) e.g. [2, 10, 15, 16, 17, 5, 18, 6, 19, 20, 9, 21, 22, 23, 24, 7, 8]. In the “integrity” track, a lot of works (e.g. [10, 16, 17, 5, 9, 6]) are done for “identity query” [18], i.e. the query result is a subset of the database. [16, 5] authenticated 1D range selection queries, with linear (in the number of tuples selected by the query condition) communication cost and storage overhead. [17] verified range selection queries using aggregated signatures (like RSA [25], BLS [26]). [6] proposed a linear (or superlinear) scheme, which uses chained signatures over a “verification R-Tree” built on a multidimensional data space, to authenticate windows query, range query, kNN query, and RNN query. To the best of our knowledge, the current most efficient authentication scheme for range selection queries is [9], which proposed an efficient authentication scheme for 1D and 2D range selection queries over a grid dataset (e.g. GIS or image data) with $O(1)$ communication cost and linear storage overhead. [18] claimed to authenticate arbitrary queries, but their security model is too weak: a playful adversary can easily break their scheme. Aggregate range query is arguably more challenging and only a few works (e.g. [5, 7, 8, 1]) are devoted to the authentication of aggregate query. We remark that our scheme can also protect privacy for aggregate attributes by using homomorphic encryption scheme like [7, 8]).

There are roughly four categories of approaches for outsourced database authentication in the literatures [2, 10, 15, 16, 17, 5, 18, 6, 19, 20, 9, 21, 22, 23, 24]. (1) (Homomorphic and/or aggregatable) Cryptographic primitives, like collision-resistant hash, digital signature, commitment [17, 27, 7]. (2) Geometry partition and authenticated data structure [10, 6, 9, 22, 19, 8]. For example, Merkle Hash Tree (typically for 1D case) and variants, KD-tree with chained signature [5], R-Tree with chained signature [6], and authenticated B-Tree/R-Tree [8]. (3) Authenticated precomputed partial result, e.g. authenticated prefix sum [9, 8] (the static case solution in [8]) (4) Inserting and auditing fake tuples [20]. Instead of leveraging on the standard or existing cryptographic primitives (e.g. digital signature scheme, cryptographic hash) like most of previous works, [1] designed a new functional encryption scheme and a new homomorphic authentication tag. Consequently, [1] achieves very good asymptotic performance, but their proof of security became much more challenging.

To the best of our knowledge, the existing few works (e.g. [5, 7, 8]) on authentication of aggregate query either only deal with 1D case, or have communication overhead¹ linear (or even superlinear) w.r.t. the number of data points in the query range, and/or exponential in dimension. Even for multidimensional (non-aggregate) range selection query, the communication overhead is still in $O(\log^{d-1} N + |S|)$ (Martel *et al.* [10], Chen *et al.* [11]), where S is the set of data points within the query range, N is the number of data points in the dataset, and d is the dimension.

Recently, Gennaro *et al.* [28] and Chung *et al.* [29] proposed methods to authenticate *any* outsourced (or delegated) polynomial time function, based on fully homomorphic encryption [30, 31, 32]. They [28, 29] also gave a good discussion on why previous techniques (e.g. interactive proofs, probabilistic checkable proof (PCP), and interactive arguments) are insufficient for authenticating outsourced function from the performance point of view. If a function has input size T_1 and output size T_2 , then both Gennaro *et al.* [28] and Chung *et al.* [29] have communication overhead in $\Omega(T_1 + T_2)$ to authenticate this function, where the hidden constant behind the big- Ω notation could be huge. The reason is two-fold: (1) First, using Gentry’s fully homomorphic encryption scheme, one bit plaintext will be expanded to $O(\kappa^3)$ bits ciphertext; (2) Second, in Gennaro *et al.* [28], before encrypting, each bit of plaintext will be replaced by a κ bits long message, which in turn will be encrypted by fully homomorphic encryption scheme; in Chung *et al.* [29], to authenticate a query, Alice has to generate $O(t)$ similar queries and issues all of these queries together and encrypts them using fully homomorphic encryption scheme, to achieve false positive probability 2^{-t} . The difference between their solutions and our work may become more clear when authenticating non-aggregate range selection query: Both Gennaro *et al.* [28] and Chung *et al.* [29] will require linear communication overhead (with huge constant factor), while our solution still requires $O(d^2 \log^2 \mathcal{Z})$ communication overhead.

¹ The original papers either do not provide a tight theoretical asymptotic bound, or do not relate the bound to generic parameters, including database size, domain size, dimension and security parameter.

3 Formulation

In this section, we restate the problem formulation and security model from [1], with modifications adapting our extension in this paper.

3.1 Dataset and Query

The dataset \mathbf{D} is a set of N d -dimensional points $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ from the domain $[\mathcal{Z}]^d$ where \mathcal{Z} is a big integer (e.g. 64 bits integer). Each point $\mathbf{x} \in \mathbf{D}$ is associated with a vector-valued attribute, denoted as $\text{Att}(\mathbf{x})$, where each component of the vector $\text{Att}(\mathbf{x})$ is an integer. Let $\mathbf{R} = [a_1, b_1] \times [a_2, b_2] \times \dots \times [a_d, b_d] \subseteq [\mathcal{Z}]^d$ be a d -dimensional rectangular range. Xu and Chang [1] focused on aggregate count query function COUNT :

$$\text{COUNT}(\mathbf{D}, \mathbf{R}) \stackrel{\text{def}}{=} \sum_{\mathbf{x} \in \mathbf{D} \cap \mathbf{R}} \text{Att}(\mathbf{x}) \pmod{p},$$

where the attribute $\text{Att}(\mathbf{x}) = 1$ for each point $\mathbf{x} \in \mathbf{D}$. Note that p is exponential in the security parameter κ and N is polynomial in κ .

In this paper, we are concerning the following queries together with multidimensional vector-valued attribute $\text{Att}(\mathbf{x})$, $\mathbf{x} \in \mathbf{D}$:

SUM: A sum query with range \mathbf{R} asks for the summation of attributes $\text{Att}(\mathbf{x})$ for all data points $\mathbf{x} \in \mathbf{D} \cap \mathbf{R}$.

$$\text{SUM}(\mathbf{D}, \mathbf{R}) = \bigoplus_{\mathbf{x} \in \mathbf{D} \cap \mathbf{R}} \text{Att}(\mathbf{x}) \pmod{p} \quad (1)$$

MIN: A min query with range \mathbf{R} and dimension $\iota \in [d]$ asks for the minimum attribute value along the ι -th dimension among all data points $\mathbf{x} \in \mathbf{D} \cap \mathbf{R}$.

$$\text{MIN}(\mathbf{D}, \mathbf{R}, \iota) = \min_{\mathbf{x} \in \mathbf{D} \cap \mathbf{R}} \text{Att}(\mathbf{x})[\iota] \quad (2)$$

MAX: A max query with range \mathbf{R} and dimension $\iota \in [d]$ asks for the maximum attribute value along the ι -th dimension among all data points $\mathbf{x} \in \mathbf{D} \cap \mathbf{R}$.

$$\text{MAX}(\mathbf{D}, \mathbf{R}, \iota) = \max_{\mathbf{x} \in \mathbf{D} \cap \mathbf{R}} \text{Att}(\mathbf{x})[\iota] \quad (3)$$

MEDIAN: A median query with range \mathbf{R} and dimension $\iota \in [d]$ asks for the median attribute value along the ι -th dimension among all data points $\mathbf{x} \in \mathbf{D} \cap \mathbf{R}$.

$$\text{MEDIAN}(\mathbf{D}, \mathbf{R}, \iota) = y, \text{ such that } y \in S = \{\text{Att}(\mathbf{x})[\iota] : \mathbf{x} \in \mathbf{D} \cap \mathbf{R}\} \text{ and } y \text{ is ranked } \lceil \frac{|S|}{2} \rceil\text{-th among the set } S \quad (4)$$

RANGESELECT: A range select query with range \mathbf{R} asks for all data points $\mathbf{x} \in \mathbf{D} \cap \mathbf{R}$.

$$\text{RANGESELECT}(\mathbf{D}, \mathbf{R}) = \{\mathbf{x} : \mathbf{x} \in \mathbf{D} \cap \mathbf{R}\} \quad (5)$$

3.2 Security Model

Xu and Chang [1] presented a formulation for the authentication problem over outsourced database, as a variant of *Verifiable Computation* [28]. Let us view a query on a dataset as the function $F : \mathbb{D} \times \mathbb{Q} \rightarrow \{0, 1\}^*$, where \mathbb{D} is the domain of datasets, \mathbb{Q} is the domain of queries, and the output of F is represented by a binary string. Note that a query $Q \in \mathbb{Q}$ is represented by combination of query type (like count, sum, etc), query range, and other parameters if any (e.g. a min query $\text{MIN}(\mathbf{R}, \iota)$). Xu and Chang [1] defined the remote computing protocol as follow:

Definition 1 (RC [1]) A Remote Computing (RC) protocol for a function $F : \mathbb{D} \times \mathbb{Q} \rightarrow \{0, 1\}^*$, between Alice and Bob, consists of a setup phase and a query phase. The setup phase consists of a key generating algorithm KGen and data encoding algorithm DEnc ; the query phase consists of a pair of interactive algorithms, namely the evaluator Eval and the extractor Ext . These four algorithms ($\text{KGen}, \text{DEnc}, \langle \text{Eval}, \text{Ext} \rangle$) run in the following way:

Setup Phase

1. Given security parameter κ , Alice generates a key $K : K \leftarrow \text{KGen}(1^\kappa)$.
2. Alice encodes dataset $\mathbf{D} \in \mathbb{D} : (\mathbf{D}_B, \mathbf{D}_A) \leftarrow \text{DEnc}(\mathbf{D}, K)$, then sends \mathbf{D}_B to Bob and keeps \mathbf{D}_A .

Query Phase The query phase consists of multiple query sessions. In each query session, Alice and Bob interact as below.

1. Alice selects a query $\mathbf{Q} \in \mathbb{Q}$.
2. Algorithm $\text{Ext}(\mathbf{D}_A, \mathbf{Q}, K)$ on Alice's side, interacts with algorithm $\text{Eval}(\mathbf{D}_B)$ on Bob's side to compute $(\zeta, X, \Psi) \leftarrow \langle \text{Eval}(\mathbf{D}_B), \text{Ext}(\mathbf{D}_A, \mathbf{Q}, K) \rangle$, where $\zeta \in \{\text{accept}, \text{reject}\}$ and Ψ is the proof of result X . If $\zeta = \text{accept}$, then Alice accepts that X is equal to $F(\mathbf{D}, \mathbf{Q})$. Otherwise, Alice rejects.

Definition 2 (Efficient RC [1]) A RC protocol is efficient, if

1. the size of K and \mathbf{D}_A are both in $O(\text{poly}(d))$ where d is the dimension of dataset \mathbf{D} ;
2. communication complexity is $O(\text{poly}(d, \log |\mathbf{D}|))$;
3. the size of \mathbf{D}_B is $O(\text{poly}(d, |\mathbf{D}|))$ (this implies the complexity of DEnc is in $O(\text{poly}(d, |\mathbf{D}|))$).
4. the algorithm Ext must be more efficient than computing F directly (This is similar with models in [28, 29]).

A RC protocol is *verifiable*, if the following conditions hold: (1) Alice always accepts, when Bob follows the protocol honestly; (2) Alice rejects with o.h.p. (overwhelming high probability), when Bob returns a wrong result. Here adversaries, i.e. malicious Bob, are allowed to interact with Alice and learn for polynomial number of query sessions, before launching the attack. During the learning, the adversary may store whatever it has seen or learnt in a state variable.

Definition 3 (VRC [1]) A RC protocol $\mathcal{E} = (\text{KGen}, \text{DEnc}, \langle \text{Eval}, \text{Ext} \rangle)$ w.r.t. function $F : \mathbb{D} \times \mathbb{Q} \rightarrow \{0, 1\}^*$, is called VRC (Verifiable Remote Computing) protocol, if the following two conditions hold: Let κ be the security parameter.

- *correctness*: for any $\mathbf{D} \in \mathbb{D}$, any $K \leftarrow \text{KGen}(1^\kappa)$ and any $\mathbf{Q} \in \mathbb{Q}$, it holds that $\langle \text{Eval}(\mathbf{D}_B), \text{Ext}(\mathbf{D}_A, \mathbf{Q}, K) \rangle = (\text{accept}, F(\mathbf{D}, \mathbf{Q}), \Psi)$ for some Ψ , where $(\mathbf{D}_B, \mathbf{D}_A) \leftarrow \text{DEnc}(\mathbf{D}, K)$.
- *soundness*: for any PPT (adaptive) adversary \mathcal{A} , the advantage $\text{Adv}_{\mathcal{E}, \mathcal{A}}(1^\kappa) \leq \text{negl}(\kappa)$ (asymptotically less or equal).

where $\text{Adv}_{\mathcal{E}, \mathcal{A}}(1^\kappa)$ is defined as

$$\text{Adv}_{\mathcal{E}, \mathcal{A}}(1^\kappa) \stackrel{\text{def}}{=} \Pr \left[(\zeta, X, \Psi, \text{view}_{\mathcal{A}}^{\mathcal{E}}, \mathbf{D}, \mathbf{Q}) \leftarrow \text{Exp}_{\mathcal{A}}^{\mathcal{E}}(1^\kappa) : \zeta = \text{accept} \wedge X \neq F(\mathbf{D}, \mathbf{Q}) \right];$$

Experiment $\text{Exp}_{\mathcal{A}}^{\mathcal{E}}(1^\kappa)$

$\mathbf{D} \leftarrow \mathcal{A}(\text{view}_{\mathcal{A}}^{\mathcal{E}});$
 $K \leftarrow \text{KGen}(1^\kappa);$
 $(\mathbf{D}_B, \mathbf{D}_A) \leftarrow \text{DEnc}(\mathbf{D}, K);$
loop until $\mathcal{A}(\text{view}_{\mathcal{A}}^{\mathcal{E}})$ decides to stop
 $\mathbf{Q}_i \leftarrow \mathcal{A}(\mathbf{D}_B, \text{view}_{\mathcal{A}}^{\mathcal{E}});$
 $(\zeta_i, X_i, \Psi_i) \leftarrow \langle \mathcal{A}(\mathbf{D}_B, \text{view}_{\mathcal{A}}^{\mathcal{E}}), \text{Ext}(\mathbf{D}_A, \mathbf{Q}_i, K) \rangle;$
 $\mathbf{Q} \leftarrow \mathcal{A}(\mathbf{D}_B, \text{view}_{\mathcal{A}}^{\mathcal{E}});$
 $(\zeta, X, \Psi) \leftarrow \langle \mathcal{A}(\mathbf{D}_B, \text{view}_{\mathcal{A}}^{\mathcal{E}}), \text{Ext}(\mathbf{D}_A, \mathbf{Q}, K) \rangle;$
Output $(\zeta, X, \Psi, \text{view}_{\mathcal{A}}^{\mathcal{E}}, \mathbf{D}, \mathbf{Q}).$

The probability is taken over all random coins used by related algorithms, $\text{negl}(\cdot)$ is some negligible function, and $\text{view}_{\mathcal{A}}^{\mathcal{E}}$ is a state variable² describing all random coins chosen by \mathcal{A} and all messages \mathcal{A} can access during previous interactions with \mathcal{E} .

We remark that this security model is also related to the formulation of \mathcal{POR} (Proof of Retrievability) [33] and it is not surprising that our scheme could imply a (ρ, λ) -valid \mathcal{POR} system, with some proper parameters ρ and λ .

4 Background

In this section, we summarize the authentication scheme for aggregate count query over static multidimensional dataset proposed by Xu and Chang [1], which serves as the base of this paper. For the sake of presentation of our extension in this paper, we make a very slight modification to the original scheme proposed by Xu and Chang [1].

Let \mathbf{D} be a set of N d -dimensional points in domain $[1, \mathcal{Z}]^d$, and each point $\mathbf{x} \in \mathbf{D}$ is associated with an attribute $\text{Att}(\mathbf{x})$. In [1], the attribute function is $\text{Att}(\mathbf{x}) = 1$, since it dealt with COUNT query.

Overview Xu and Chang [1] *implicitly* defined homomorphic authentication tag functions DTag and QTag . Their scheme is an interactive protocol between Alice and Bob, and contains a setup phase followed by a query phase. In the setup phase, Alice preprocesses the dataset by generating a tag $\text{DTag}_{\mathcal{K}}(\mathbf{x})$ for each point \mathbf{x} in the dataset \mathbf{D} with her private key \mathcal{K} . At the end of setup phase, Alice sends both the dataset \mathbf{D} and tags $\mathbf{T} = \{\text{DTag}_{\mathcal{K}}(\mathbf{x}) : \mathbf{x} \in \mathbf{D}\}$ to Bob and removes them from her storage. Later in the query phase, Alice may issue many queries over the dataset to Bob. For example, Alice may want to know how many points are within a range \mathbf{R} . Alice sends \mathbf{R} to Bob. Meanwhile, in order to help Bob to generate a proof, Alice chooses a random nonce ρ and sends $\Phi = \{\text{QTag}_{\mathcal{K}}(\mathbf{x}, \rho) : \mathbf{x} \in \mathbf{R}\}$ to Bob. After receiving \mathbf{R} and Φ , Bob is supposed to return $X = \sum_{\mathbf{x} \in \mathbf{D} \cap \mathbf{R}} \text{Att}(\mathbf{x}) = |\mathbf{D} \cap \mathbf{R}|$ as result and $\Psi_1 = \otimes_{\mathbf{x} \in \mathbf{D} \cap \mathbf{R}} \text{DTag}_{\mathcal{K}}(\mathbf{x})$ and $\Psi_2 = \otimes_{\mathbf{x} \in \mathbf{D} \cap \mathbf{R}} \text{QTag}_{\mathcal{K}}(\mathbf{x}, \rho)$ as proof. Since the tag functions DTag , QTag are homomorphic, Alice can verify the consistency between Ψ_1 and Ψ_2 using her private key \mathcal{K} . To ensure completeness, Alice has to interact with Bob and perform the above procedure for the complement query range \mathbf{R}^c .

However, the size of Φ is proportional to the size of range \mathbf{R} , which could be huge. One of main contributions of [1] is that the paper proposed a new functional encryption scheme and use it to reduce communication cost in the following way:

- In the setup phase, Alice produces a ciphertext $\text{CT}_{\mathbf{x}}$ for each data point $\mathbf{x} \in \mathbf{D}$ using the functional encryption scheme. Alice sends all ciphertexts $\text{CT}_{\mathbf{x}}$'s together with the dataset and tags to Bob at the end of setup.
- In a query session, for a count query with rectangular range \mathbf{R} , Alice chooses a random nonce ρ and generates a *short* delegation key δ w.r.t. the range \mathbf{R} and the random nonce ρ , using the functional encryption scheme. Alice sends the delegation key δ to Bob together with the query.
- For each data point $\mathbf{x} \in \mathbf{D} \cap \mathbf{R}$, Bob can decrypt ciphertext $\text{CT}_{\mathbf{x}}$ and obtain $\text{QTag}_{\mathcal{K}}(\mathbf{x}, \rho)$ as the decrypted value using the functional encryption scheme and the delegation key δ . For points $\mathbf{y} \notin \mathbf{R}$, Bob learns nothing about $\text{QTag}_{\mathcal{K}}(\mathbf{y}, \rho)$.

Since the size of delegation key δ is in $O(d \log^2 \mathcal{Z})$, the communication cost is reduced dramatically.

Formal Algorithm The homomorphic authentication tag (DTag , QTag , Verify) implied in Xu and Chang [1] is as below: Let key $\mathcal{K} = (\theta, \beta, \gamma) \in \tilde{\mathbb{G}} \times \mathbb{Z}_p^* \times \mathbb{Z}_p^*$, $v_{\mathbf{x}}, w_{\mathbf{x}} \in \tilde{\mathbb{G}}$ be random coins chosen for point \mathbf{x} , and

² The adaptive adversary \mathcal{A} may keep updating this state variable.

$$\Psi = (\Psi_1, \Psi_2, \Psi_3).$$

$$\text{DTag}_{\mathcal{K}}(\mathbf{x}) = \left(\theta^{\text{Att}(\mathbf{x})} v_{\mathbf{x}}, v_{\mathbf{x}}^{\beta}, w_{\mathbf{x}} \right) \quad (6)$$

$$\text{QTag}_{\mathcal{K}}(\mathbf{x}, \rho) = v_{\mathbf{x}}^{\gamma} w_{\mathbf{x}}^{\rho} \quad (7)$$

$$\text{Verify}_{\rho, \mathcal{K}}(Y, \Psi, \Psi_4) = \begin{cases} 1 & \left(\text{if } (\Psi_1 \theta^{-Y})^{\beta} = \Psi_2 \text{ and } (\Psi_1 \theta^{-Y})^{\gamma} \Psi_3^{\rho} = \Psi_4 \right) \\ 0 & \text{(otherwise)} \end{cases} \quad (8)$$

The authentication tag (DTag, QTag, Verify) is homomorphic and satisfies the following properties:

$$\text{Verify}_{\rho, \mathcal{K}}(\text{Att}(\mathbf{x}), \text{DTag}_{\mathcal{K}}(\mathbf{x}), \text{QTag}_{\mathcal{K}}(\mathbf{x}, \rho)) = 1 \quad (9)$$

$$\text{Verify}_{\rho, \mathcal{K}} \left(\sum_{\mathbf{x} \in \mathbf{R}} \text{Att}(\mathbf{x}), \bigotimes_{\mathbf{x} \in \mathbf{R}} \text{DTag}_{\mathcal{K}}(\mathbf{x}), \prod_{\mathbf{x} \in \mathbf{R}} \text{QTag}_{\mathcal{K}}(\mathbf{x}, \rho) \right) = 1 \quad (10)$$

We restate the scheme in [1] in Figure 1 with authentication tag (DTag, QTag, Verify) and hide details of the applications of the functional encryption scheme.

Security Since the authentication tag function is homomorphic, an adversary (i.e. a dishonest Bob) may attempt to cheat and convince Alice to accept a wrong result in this way: choose some integer $\mu_{\mathbf{x}}$ for each point \mathbf{x} , and in Step B1 of algorithm CollRes compute the proof (Ψ, Ψ_4) as below

$$\Psi \leftarrow \bigotimes_{\mathbf{x} \in \mathbf{D} \cap \mathbf{R}} t_{\mathbf{x}}^{\mu_{\mathbf{x}}} = \bigotimes_{\mathbf{x} \in \mathbf{D}} \text{DTag}(\mathbf{x})^{\mu_{\mathbf{x}}}, \quad \Psi_4 \leftarrow \prod_{\mathbf{x} \in \mathbf{D}} \text{QTag}(\mathbf{x}, \rho)^{\mu_{\mathbf{x}}} \quad (13)$$

It is easy to verify that the above forged proof passes the verification in Step A2 of CollRes in Figure 1, but may not pass the second equality test in Step 3 of Count in Figure 1. Such adversary looks “restricted” in its attack strategy. However, [1] showed that, under **GKEA** assumption, such adversary’s power is not restricted at all: If there exists an efficient (arbitrary) adversary that breaks their scheme, then there exists such “restricted” adversary that breaks their scheme.

[1] considered various types of PPT adversaries, which interacts with Alice by playing the role of Bob and intends to output a wrong query result and a forged but valid proof:

- Type I adversary: This adversary is not confined in any way in its attack strategy and produces a tuple $(X, \Psi = (\Psi_1, \Psi_2, \Psi_3), \Psi_4)$ on a query range \mathbf{R} .
- Type II adversary: A restricted adversary which can produce the same forgery³ from the same input as Type I adversary, additionally, it finds N integers⁴ μ_i ’s, $1 \leq i \leq N$, such that

$$\Psi \leftarrow \bigotimes_{\mathbf{x} \in \mathbf{D} \cap \mathbf{R}} t_{\mathbf{x}}^{\mu_{\mathbf{x}}} = \bigotimes_{\mathbf{x} \in \mathbf{D}} \text{DTag}(\mathbf{x})^{\mu_{\mathbf{x}}},$$

- Type III adversary: The same as Type II adversary, with additional constraint: $\mu_i = 0$ for $\mathbf{x}_i \in \mathbf{D} \cap \mathbf{R}^c$.
- Type IV adversary: The same as Type III adversary, with additional constraint: $\mu_i = 1$ for $\mathbf{x}_i \in \mathbf{D} \cap \mathbf{R}$.

It seems that from Type I to Type IV adversaries are more and more restricted, in the sense that

$$\{\text{Type I Adversary}\} \supseteq \{\text{Type II Adversary}\} \supseteq \{\text{Type III Adversary}\} \supseteq \{\text{Type IV Adversary}\} \quad (14)$$

However, [1] showed that, in the above formula (14), (*informally*) each inclusion relation \supseteq can be replaced by equality $=$, under related cryptographic assumptions (**GKEA**, computational diffie-hellman assumption etc). Furthermore, [1] proved that there exists no Type IV adversary under certain cryptographic assumptions.

³ This is possible, if the Type II adversary just invokes Type I adversary as a subroutine using the same random coin.

⁴ Note that μ_i can take negative integer value, and $\mu_i > 1$ ($\mu_i < 1$, respectively) corresponds to the case of double counting (undercounting, respectively) point \mathbf{x}_i .

Fig. 1: Construction of \mathcal{RC} protocol $\mathcal{E} = (\text{KGen}, \text{DEnc}, \langle \text{Eval}, \text{Ext} \rangle)$ where $\langle \text{Eval}, \text{Ext} \rangle$ (namely **Count**) invokes $\langle \widetilde{\text{Eval}}, \widetilde{\text{Ext}} \rangle$ (namely **CollRes**) as a subroutine. The attribute function is $\text{Att}(\mathbf{x}) = 1$ for each $\mathbf{x} \in \mathbf{D}$.

<p>(Alice) $\text{KGen}(1^\kappa)$: Output a private key \mathcal{K}.</p> <hr/> <p>(Alice) $\text{DEnc}(\mathbf{D}; \mathcal{K})$:</p> <ol style="list-style-type: none"> 1. For each point $\mathbf{x} \in \mathbf{D}$, generate a tag $t_{\mathbf{x}} = \text{DTag}(\mathbf{x}, \mathcal{K}).$ 2. For each point $\mathbf{x} \in \mathbf{D}$, generate a ciphertext $\text{CT}_{\mathbf{x}}$, using the functional encryption scheme with key \mathcal{K}. 3. Send $\mathbf{D}_B = (\mathbf{D}, \mathbf{T} = \{t_{\mathbf{x}} : \mathbf{x} \in \mathbf{D}\}, \mathbf{C} = \{\text{CT}_{\mathbf{x}} : \mathbf{x} \in \mathbf{D}\})$ to Bob, and keep <i>only</i> key \mathcal{K} and $\mathbf{D}_A = (N, d, \Delta = \bigotimes_{\mathbf{x} \in \mathbf{D}} t_{\mathbf{x}})$ in local storage. <hr/> <p>(Alice, Bob) Count = $\langle \text{Eval}(\mathbf{D}_B), \text{Ext}(\mathbf{D}_A, \mathbf{R}, \mathcal{K}) \rangle$: $\mathbf{D}_A = (N, d, \Delta), \mathbf{D}_B = (\mathbf{D}, \mathbf{T}, \mathbf{C})$ Precondition: The query range $\mathbf{R} \subset [\mathcal{Z}]^d$ is a rectangular range.</p> <p>Step 1: Alice partitions the complement range \mathbf{R}^c into $2d$ rectangular ranges $\{\mathbf{R}_\ell \subset [\mathcal{Z}]^d : \ell \in [1, 2d]\}$, and sets $\mathbf{R}_0 = \mathbf{R}$.</p> <p>Step 2—Reduction: For $0 \leq \ell \leq 2d$, Alice and Bob invokes CollRes on range \mathbf{R}_ℓ. Denote the output as $(\zeta_\ell, X_\ell, \Psi^{(\ell)})$.</p> <p>Step 3: Alice sets $\zeta = \text{accept}$, if the following equalities hold</p> $\forall 0 \leq \ell \leq 2d, \zeta_\ell \stackrel{?}{=} \text{accept}, \quad \bigotimes_{0 \leq \ell \leq 2d} \Psi^{(\ell)} \stackrel{?}{=} \Delta; \quad (11)$ <p>otherwise sets $\zeta = \text{reject}$. Alice outputs (ζ, X_0, Δ).</p> <hr/> <p>(Alice, Bob) CollRes = $\langle \widetilde{\text{Eval}}(\mathbf{D}_B), \widetilde{\text{Ext}}(\mathbf{D}_A, \mathbf{R}, \mathcal{K}) \rangle$: $\mathbf{D}_A = (N, d, \Delta), \mathbf{D}_B = (\mathbf{D}, \mathbf{T}, \mathbf{C})$ Precondition. The query range $\mathbf{R} \subset [\mathcal{Z}]^d$ is a rectangular range.</p> <p>Step A1: (Alice's first step) Alice chooses a random nonce ρ from \mathbb{Z}_p^* and produces a delegation key δ w.r.t. range \mathbf{R} and nonce ρ, using the functional encryption scheme. Alice sends (\mathbf{R}, δ) to Bob.</p> <p>Step B1: (Bob's first step) Bob computes the query result X and proof $(\Psi_1, \Psi_2, \Psi_3, \Psi_4)$ as follows</p> $X \leftarrow \sum_{\mathbf{x} \in \mathbf{D} \cap \mathbf{R}} \text{Att}(\mathbf{x}); \quad \Psi \leftarrow \bigotimes_{\mathbf{x} \in \mathbf{D} \cap \mathbf{R}} t_{\mathbf{x}} = \bigotimes_{\mathbf{x} \in \mathbf{D} \cap \mathbf{R}} \text{DTag}_{\mathcal{K}}(\mathbf{x}); \quad \Psi_4 \leftarrow \prod_{\mathbf{x} \in \mathbf{D} \cap \mathbf{R}} \text{QTag}_{\mathcal{K}}(\mathbf{x}, \rho) \quad (12)$ <p>where for each point $\mathbf{x} \in \mathbf{D} \cap \mathbf{R}$, $\text{QTag}_{\mathcal{K}}(\mathbf{x}, \rho)$ is obtained by decrypting $\text{CT}_{\mathbf{x}}$ using the functional encryption scheme with delegation key δ. Bob sends (X, Ψ, Ψ_4) to Alice.</p> <p>Step A2: (Alice's second step) Alice verifies whether (Ψ, Ψ_4) are valid tags for X under DTag and QTag respectively, using the private key \mathcal{K} and ρ. If the following equation holds,</p> $\text{Verify}_{\rho, \mathcal{K}}(X, \Psi, \Psi_4) \stackrel{?}{=} 1$ <p>then sets $\zeta = \text{accept}$. Otherwise sets $\zeta = \text{reject}$. Alice outputs (ζ, X, Ψ)</p>

5 Authenticating Aggregate queries beyond count: SUM, MIN, MAX

5.1 SUM

Suppose each data point $\mathbf{x} \in \mathbf{D}$ is associated with an attribute value $\text{Att}(\mathbf{x})$. The sum query with range \mathbf{R} asks for the summation $\sum_{\mathbf{x} \in \mathbf{D}} \text{Att}(\mathbf{x})$.

5.1.1 Summing 1D Attribute Suppose the attribute value is in 1D, i.e. $\text{Att}(\mathbf{x}) \in [\mathcal{Z}]$. The authentication scheme for SUM is identical to the scheme in Figure 1 for COUNT, except that

- the attribute function $\text{Att}(\mathbf{x}) = 1$ is redefined as $\text{Att}(\mathbf{x}) = y_{\mathbf{x}} \in [\mathcal{Z}]$.
- Alice sends $\{\text{Att}(\mathbf{x}) : \mathbf{x} \in \mathbf{D}\}$ to Bob in the setup phase, and Bob keeps it along with the dataset and tags.

Denote the modified scheme as $(\text{KGen}, \text{DEnc}, \text{SUM}^{(1)})$.

Lemma 1 *Suppose attribute value is in 1D. The extended scheme $(\text{KGen}, \text{DEnc}, \text{SUM}^{(1)})$ is a \mathcal{VRC} w.r.t. SUM, i.e. it is correct and sound to authenticate SUM.*

5.1.2 Summing Multi-Dimensional Attribute Suppose the attribute value $\text{Att}(\mathbf{x}) = \mathbf{x} \in [\mathcal{Z}]^n$ for some integer n . The authentication scheme for summing n -dimensional attribute is identical to the scheme in Section 5.1.1 for 1D case, except that

- The secret key \mathcal{K} generated by KGen contains an additional element $\mathbf{s} = (s_1, \dots, s_n)$ which is randomly chosen from the domain $(\mathbb{Z}_p^*)^n$;
- The attribute function $\text{Att}(\mathbf{x}) = y_{\mathbf{x}} \in [\mathcal{Z}]$ and the tag function

$$\text{DTag}(\mathbf{x}) = \left(\theta^{\text{Att}(\mathbf{x})} v_{\mathbf{x}}, v_{\mathbf{x}}^{\beta}, w_{\mathbf{x}} \right) = \left(\theta^{y_{\mathbf{x}}} v_{\mathbf{x}}, v_{\mathbf{x}}^{\beta}, w_{\mathbf{x}} \right)$$

are redefined as:

$$\text{Att}(\mathbf{x}) = \mathbf{x} \in [\mathcal{Z}]^n; \quad \text{DTag}(\mathbf{x}) = \left(\theta^{\langle \mathbf{x}, \mathbf{s} \rangle} v_{\mathbf{x}}, v_{\mathbf{x}}^{\beta}, w_{\mathbf{x}} \right)$$

- Step B1 and Step A2 in CollRes are modified accordingly (i.e. In Step A2, X is replaced by the inner product $\langle \mathbf{X}, \mathbf{s} \rangle$).

Denote the modified scheme as $(\text{KGen}, \text{DEnc}, \text{SUM}^{(n)})$.

Theorem 2 *The extended scheme $(\text{KGen}, \text{DEnc}, \text{SUM}^{(n)})$ is a \mathcal{VRC} w.r.t. SUM, i.e. it is correct and sound to authenticate SUM.*

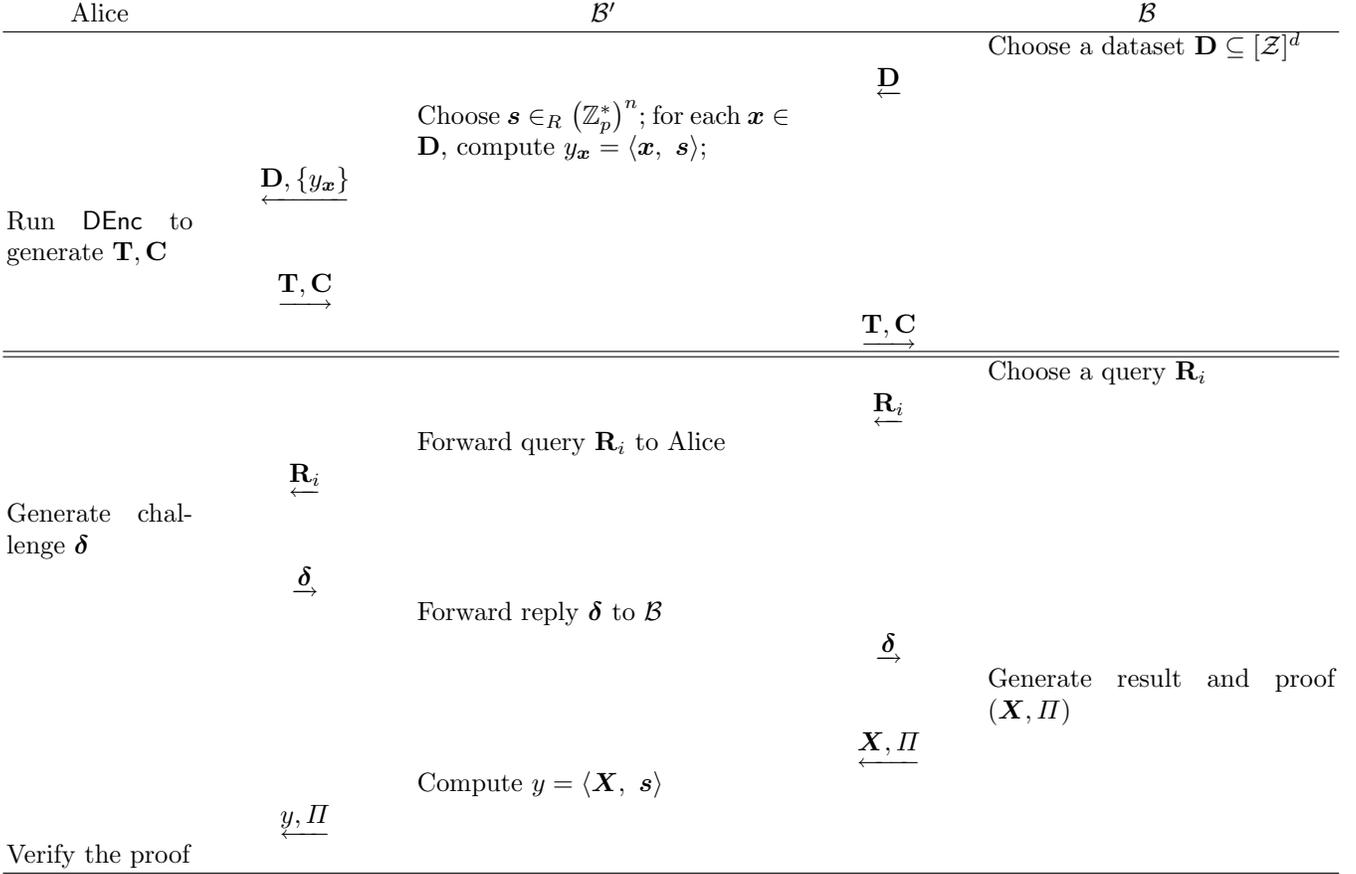
Proof (of Theorem 2). The correctness part is straightforward due to the homomorphic property (i.e. Equation (10)) of $(\text{DTag}, \text{QTag}, \text{Verify})$. We focus on the soundness part.

Using proof by contradiction, suppose that there exists a PPT adversary \mathcal{B} which can output a wrong query result $\mathbf{Y} \neq \sum_{\mathbf{x} \in \mathbf{D} \cap \mathbf{R}} \text{Att}(\mathbf{x}) \pmod{p}$ for sum query with range \mathbf{R} and passes all verifications with non-negligible probability ϵ .

Part I: *We will show that $\langle \mathbf{Y}, \mathbf{s} \rangle = \langle \mathbf{X}, \mathbf{s} \rangle \pmod{p}$, where \mathbf{s} is a part of private key and $\mathbf{X} = \sum_{\mathbf{x} \in \mathbf{D} \cap \mathbf{R}} \text{Att}(\mathbf{x}) \pmod{p}$ is the correct query result for the corresponding query \mathbf{R} .*

We construct an adversary \mathcal{B}' to against the scheme $(\text{KGen}, \text{DEnc}, \text{SUM}^{(1)})$ for 1D case, based on \mathcal{B} . Adversary \mathcal{B}' will simulate two instances of experiments:

- Experiment $\text{Exp}_{\mathcal{B}'}^{\mathcal{E}_{1D}}$ for scheme $(\text{KGen}, \text{DEnc}, \text{SUM}^{(1)})$: \mathcal{B}' takes the role of Bob and interacts with Alice.
- Experiment $\text{Exp}_{\mathcal{B}}^{\mathcal{E}_{dD}}$ for scheme $(\text{KGen}, \text{DEnc}, \text{SUM}^{(d)})$: \mathcal{B}' takes the role of Alice and \mathcal{B} takes the role of Bob.



By our hypothesis, with non-negligible probability ϵ , we have $\mathbf{Y} \neq \sum_{\mathbf{x} \in \mathbf{D} \cap \mathbf{R}} \mathbf{x}$ and the proof $\Pi_{\mathbf{Y}}$ passes all verification. As a result, with probability at least ϵ , \mathcal{B}' 's reply is $(\langle \mathbf{Y}, \mathbf{s} \rangle, \Pi_{\mathbf{Y}})$ passes all verification. On the other hand, Lemma 1 says, for any PPT adversary which can output a query result with a proof, if the proof is valid then the query result is correct with overwhelming high probability $1 - \epsilon'$ (ϵ' is some negligible function). Hence, with probability at least $\epsilon(1 - \epsilon')$, we have

$$\langle \mathbf{Y}, \mathbf{s} \rangle = \sum_{\mathbf{x} \in \mathbf{D} \cap \mathbf{R}} y_{\mathbf{x}} = \sum_{\mathbf{x} \in \mathbf{D} \cap \mathbf{R}} \langle \mathbf{x}, \mathbf{s} \rangle = \langle \mathbf{X}, \mathbf{s} \rangle, \text{ where } \mathbf{X} = \sum_{\mathbf{x} \in \mathbf{D} \cap \mathbf{R}} \mathbf{x} \neq \mathbf{Y}.$$

Therefore, with non-negligible probability $\epsilon(1 - \epsilon')$, the adversary \mathcal{B} can output two distinct values \mathbf{X} and \mathbf{Y} such that $\langle \mathbf{Y}, \mathbf{s} \rangle = \langle \mathbf{X}, \mathbf{s} \rangle \pmod{p}$.

Part II: We construct an algorithm to solve Discrete Log Problem (**DLP**) based on the adversary \mathcal{B} .

Solve Discrete Log Problem (**DLP**) based on Adversary \mathcal{B}

1. Input is $(u, u^a) \in \tilde{\mathbb{G}}^2$. The goal is to find $a \in \mathbb{Z}_p^*$.
2. For each $i \in [d]$, choose y_i, z_i from \mathbb{Z}_p^* at random and set $\theta_i = (u^a)^{y_i} u^{z_i} \in \tilde{\mathbb{G}}$.
3. Simulate the authentication scheme ($\text{KGen}, \text{DEnc}, \text{SUM}^{(n)}$) with the following modifications:
 - The secret key θ and $\mathbf{s} = (s_1, \dots, s_n)$ are implicitly defined by $\theta_i = \theta^{s_i}$.
 - The simulator does not know values of (θ, \mathbf{s}) , but still can compute $\theta^{\langle \text{Att}(\mathbf{x}), \mathbf{s} \rangle}$:

$$\theta^{\langle \text{Att}(\mathbf{x}), \mathbf{s} \rangle} = \prod_{i \in [n]} \theta_i^{\lambda_i}, \text{ where } \text{Att}(\mathbf{x}) = (\lambda_1, \dots, \lambda_n)$$

4. Invoke the adversary \mathcal{B} and obtains output (\mathbf{X}, \mathbf{Y}) . Let $\mathbf{w} = (w_1, \dots, w_n) = \mathbf{X} - \mathbf{Y} \pmod p$. Applying the result in Part I, with non-negligible probability, we have

$$\theta^{\langle \mathbf{w}, \mathbf{s} \rangle} = 1 \text{ and } \mathbf{w} \neq \mathbf{0} \pmod p$$

5. A univariable equation on unknown a can be formed: Let $\mathbf{y} = (y_1, \dots, y_n)$ and $\mathbf{z} = (z_1, \dots, z_n)$.

$$\prod_{i \in [n]} \theta_i^{w_i} = \prod_{i \in [n]} u^{w_i(\alpha y_i + z_i)} = 1.$$

$$a \langle \mathbf{y}, \mathbf{w} \rangle + \langle \mathbf{z}, \mathbf{w} \rangle = 0 \pmod p$$

Note: Given $\theta_i, i \in [n]$, y_i 's are truly random. Hence, the probability that $\langle \mathbf{y}, \mathbf{w} \rangle = 0$ is negligible.

6. Solve the equation and get the root a^* . Output a^* .

The constructed PPT algorithm solves DLP with non-negligible probability. The contradiction with DL Assumption, implies that our hypothesis is wrong and such adversary \mathcal{B} does not exist. Consequently, the soundness part of Theorem 2 is proved. \square

For the rest of remaining part of this paper, we assume the attribute function is

$$\text{Att}(\mathbf{x}) = (x_1, \dots, x_d, 1), \text{ where } \mathbf{x} = (x_1, \dots, x_d).$$

By Theorem 2, our scheme can support both summing over the first d dimensions of attribute values and counting over the last dimension of attribute values. The AVG query can be authenticated by combination of SUM and COUNT.

5.2 MIN

A min query $\text{MIN}(\mathbf{R}, \iota)$ with query range \mathbf{R} and dimension $\iota \in [d]$, asks for the minimum attribute values along the ι -th dimension among all data points within $\mathbf{D} \cap \mathbf{R}$. We find that MIN query can be converted to COUNT query. The conversion is based on this proposition:

Proposition 1 *For any finite set S of numbers,*

$$c = \min S \quad \Leftrightarrow \quad c \in S \wedge |S| = |\{x : x \in S \wedge x \geq c\}|. \quad (15)$$

Suppose Alice asks Bob for the minimum attribute value along the ι -th dimension of points within range \mathbf{R} . Bob returns a data point \mathbf{x} , such that $\text{Att}(\mathbf{x})[\iota]$ is minimum in the set S of attribute values along the ι -th dimension of all points within range \mathbf{R} (i.e. $S = \{\mathbf{x}[\iota] : \mathbf{x} \in \mathbf{R} \cap \mathbf{D}\}$). Meanwhile, Bob also sends a proof to show that $\mathbf{x} \in \mathbf{D}$. Then Alice issues two COUNT queries to Bob: (1) $\text{COUNT}(\mathbf{R})$, i.e. the size of set S ; (2) $\text{COUNT}(\mathbf{R} \cap ([\mathcal{Z}]^{\iota-1} \times [c, \mathcal{Z}] \times [\mathcal{Z}]^{d-\iota}))$ where $c = \text{Att}(\mathbf{x})[\iota]$, i.e. the size of set $\{x : x \in S \wedge x \geq c\}$. Bob is expected to return the two count numbers with proofs following the scheme in [1]. Alice believes c is the minimum value if all proofs are valid and the two count numbers are equal. The algorithm is showed in Figure 2.

Theorem 3 *The extended scheme is a VRC w.r.t. MIN, i.e. it is correct and sound to authenticate MIN.*

Similarly, MAX query can be authenticated.

5.3 Median

MEDIAN can also be converted into COUNT. Quartile or percentile queries can be handled in a similar way.

Proposition 2 *Let S be a finite set of numbers.*

$$c \text{ is the median in set } S \iff c \in S \wedge |\{x : x \in S \wedge x \leq c\}| \geq \lceil \frac{|S|}{2} \rceil \wedge |\{x : x \in S \wedge x \geq c\}| \geq \lceil \frac{|S|}{2} \rceil \quad (16)$$

Suppose Alice asks Bob for the median attribute value along the ι -th dimension of points within range \mathbf{R} . Bob returns a data point \mathbf{x} , such that $\text{Att}(\mathbf{x})[\iota]$ is the median in the set S of attribute values along the ι -th dimension of all points within range \mathbf{R} (i.e. $S = \{\mathbf{x}[\iota] : \mathbf{x} \in \mathbf{D} \cap \mathbf{R}\}$). Meanwhile, Bob also sends a proof to show that $\mathbf{x} \in \mathbf{D}$. Then Alice issues three COUNT queries to Bob: (1) $\text{COUNT}(\mathbf{R})$, i.e. the size of set S ; (2) $\text{COUNT}(\mathbf{R} \cap ([\mathcal{Z}]^{\iota-1} \times [c, \mathcal{Z}] \times [\mathcal{Z}]^{d-\iota}))$ where $c = \text{Att}(\mathbf{x})[\iota]$, i.e. the size of set $\{x : x \in S \wedge x \geq c\}$; (3) $\text{COUNT}(\mathbf{R} \cap ([\mathcal{Z}]^{\iota-1} \times [1, c] \times [\mathcal{Z}]^{d-\iota}))$ i.e. the size of set $\{x : x \in S \wedge x \leq c\}$. Bob is expected to return the three count numbers N_1, N_2 and N_3 with proofs following the scheme in [1]. Alice believes c is the median value if all proofs are valid and $N_2 \geq \lceil \frac{N_1}{2} \rceil$ and $N_3 \geq \lceil \frac{N_1}{2} \rceil$.

The algorithm is shown in Figure 2. Note that when the size of S is even, there are two medians. For simplicity of presentation of the algorithm, we request Bob to return either one of the two medians, instead of both.

Fig. 2: Authenticating MIN query and MEDIAN query.

(Alice, Bob) $\text{MIN}(\mathbf{R}, \iota)$:

1. Alice sends (\mathbf{R}, ι) to Bob.
2. Bob finds $\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbf{D} \cap \mathbf{R}} \mathbf{x}[\iota]$ and sends \mathbf{x}^* to Alice.
3. Alice issues a count query $\text{COUNT}(\{\mathbf{x}^*\})$ with range $\{\mathbf{x}^*\}$ to Bob and gets authenticated query result N_0 .
4. Alice sets $c = \mathbf{x}^*[\iota]$ and finds the range $\mathbf{R}_c = \mathbf{R} \cap ([\mathcal{Z}]^{\iota-1} \times [c, \mathcal{Z}] \times [\mathcal{Z}]^{d-\iota})$.
5. Alice issues two count queries $\text{COUNT}(\mathbf{R})$ and $\text{COUNT}(\mathbf{R}_c)$ to Bob and gets authenticated results N_1 and N_2 .
6. Alice accepts c as the minimum, if all verifications succeed and $N_0 \geq 1$ and $N_1 = N_2$.

(Alice, Bob) $\text{MEDIAN}(\mathbf{R}, \iota)$:

1. Alice sends (\mathbf{R}, ι) to Bob.
2. Bob finds \mathbf{x}^* such that $\mathbf{x}^*[\iota]$ is a median among $\{\mathbf{x}[\iota] : \mathbf{x} \in \mathbf{D}\}$ and sends \mathbf{x}^* to Alice.
3. Alice issues a count query $\text{COUNT}(\{\mathbf{x}^*\})$ with range $\{\mathbf{x}^*\}$ to Bob and gets authenticated query result N_0 .
4. Alice sets $c = \mathbf{x}^*[\iota]$ and finds the range $\mathbf{R}_c^+ = \mathbf{R} \cap ([\mathcal{Z}]^{\iota-1} \times [c, \mathcal{Z}] \times [\mathcal{Z}]^{d-\iota})$ and range $\mathbf{R}_c^- = \mathbf{R} \cap ([\mathcal{Z}]^{\iota-1} \times [1, c] \times [\mathcal{Z}]^{d-\iota})$.
5. Alice issues three count queries $\text{COUNT}(\mathbf{R})$, $\text{COUNT}(\mathbf{R}_c^+)$ and $\text{COUNT}(\mathbf{R}_c^-)$ to Bob and gets authenticated results N_1, N_2 and N_3 .
6. Alice accepts c as the median, if all verifications succeed and $N_0 \geq 1$ and $N_2 \geq \lceil \frac{N_1}{2} \rceil$ and $N_3 \geq \lceil \frac{N_1}{2} \rceil$.

5.4 Beyond Aggregate queries: Range Selection

In this section, we extend our method to support range selection and range selection with projection.

5.5 Range Selection

A range selection query with range \mathbf{R} asks for all data points within the range \mathbf{R} :

$$\text{RANGESELECT}(\mathbf{R}) = \{\mathbf{x} : \mathbf{x} \in \mathbf{D} \cap \mathbf{R}\}.$$

We assume the dataset \mathbf{D} is a set of *distinct* points. The authentication scheme for range selection query is as follows:

Authenticating Multidimensional Range Selection Query

1. In the setup, Alice generates a signature $\text{Sig}(\mathbf{x})$ for each data point $\mathbf{x} \in \mathbf{D}$ using an aggregate signature scheme, and sends all signatures to Bob.
 2. To answer a range selection query with range \mathbf{R} , Bob finds the set $S = \{\mathbf{x} : \mathbf{x} \in \mathbf{D} \cap \mathbf{R}\}$ and computes an aggregated signature $\text{Sig}(S)$ for set S from signatures $\text{Sig}(\mathbf{x})$'s for point $\mathbf{x} \in \mathbf{D} \cap \mathbf{R}$, using the aggregate signature scheme. Bob sends $(S, \text{Sig}(S))$ to Alice.
 3. Alice verifies: (1) Is S a set of distinct points? (2) Is S a subset of query range \mathbf{R} ? (3) Is $\text{Sig}(S)$ a valid signature for S ?
 4. Alice issues a count query with range \mathbf{R} to Bob and gets authenticated result N_0 .
 5. Alice verifies whether $|S| = N_0$.
 6. Alice accepts S as the query result, if all verifications succeed.
-

The above method has communication overhead equal to that of COUNT query: $O(d^2 \log^2 \mathcal{Z})$. To the best of our knowledge, this is the first efficient \mathcal{VRC} (See Definition 2) to authenticate multidimensional range selection query.

5.6 Range Selection with Projection

A range selection query with projection on the 1st dimension asks for the 1st dimension of all data points within the query range

$$\text{RANGESELECT}(\mathbf{R}) = \{\mathbf{x}[1] : \mathbf{x} \in \mathbf{D} \cap \mathbf{R}\}.$$

Authenticating this query with sublinear communication overhead is more challenging than range selection without projection. If we just apply the method in Section 5.5, the communication overhead will be linear: Since only the 1st dimension of data points within the query range is asked for, but all dimensions of such data points are returned as the query result. The requirement of sublinear communication overhead implies that Alice has to verify whether $\mathbf{x} \in \mathbf{R}$, with only the knowledge of the first dimension $\mathbf{x}[1]$ of point \mathbf{x} .

Our idea is that: Alice derives the randomness $v_{\mathbf{x}}$ from $\mathbf{x}[1]$ only using a pseudorandom function $F_{\varpi}(\cdot)$ when generating the authentication tag during the setup. Then Alice issues a count query with range \mathbf{R} to Bob and receives from Bob the query result N_0 and its proof Ψ . Meanwhile, Alice also receives $S = \{\mathbf{x}[1] : \mathbf{x} \in \mathbf{D} \cap \mathbf{R}\}$. Alice verifies whether the proof Ψ is consistent with $\prod_{\mathbf{x} \in S} F_{\varpi}(x)$ and whether $|S| = N_0$. The detailed algorithm is given in Figure 3.

This solution has a limitation: When generating authentication tag during the setup, if Alice derives the randomness $v_{\mathbf{x}}$ from $\mathbf{x}[1]$, then the resulting scheme only supports projection on the 1st dimension. To support projection on any combination of dimensions, Alice has to generate 2^d authentication tags for each data point, and one tag for one subset of $[d]$. As a result, we can authenticate range selection with projection, at the cost of $O(d^2 \log^2 \mathcal{Z})$ communication overhead per query and $O(dN \cdot 2^d)$ storage on Bob's side.

Theorem 4 *The extended scheme is a \mathcal{VRC} w.r.t. range selection, i.e. it is correct and sound to authenticate d -dimensional range selection.*

Fig. 3: Construction of \mathcal{RC} protocol $\mathcal{E} = (\text{KGen}, \text{DEnc}, \langle \text{Eval}, \text{Ext} \rangle)$ to authenticate multidimensional range selection query with projection on the 1st dimension. The attribute function is $\text{Att}(\mathbf{x}) = (x_1, \dots, x_d, 1)$ for each $\mathbf{x} = (x_1, \dots, x_d) \in \mathbf{D}$.

(Alice) $\text{KGen}(1^\kappa)$:

1. Generate a private key \mathcal{K} as in Figure 2 in MAIA.
2. Let $\{F_\varpi : [\mathcal{Z}]^d \rightarrow \widetilde{\mathbb{G}}\}_{\varpi \in \{0,1\}^\kappa}$ be a pseudorandom function. Choose a random seed $\varpi \in \{0,1\}^\kappa$.
3. Set $\mathcal{K} \leftarrow (\mathcal{K}, \varpi)$.
4. Output the private key \mathcal{K} .

(Alice) $\text{DEnc}(\mathbf{D}; \mathcal{K})$: The same as in DEnc in Figure 1 [1], except that the randomness $v_x \in \widetilde{\mathbb{G}}$ is generated in this way: Let $\mathbf{x}[1]$ denote the first component of vector value \mathbf{x} .

$$\forall \mathbf{x} \in \mathbf{D}, v_x = F_\varpi(\mathbf{x}[1]).$$

(Alice, Bob) $\text{RangeSelect} = \langle \text{Eval}(\mathbf{D}_B), \text{Ext}(\mathbf{D}_A, \mathbf{R}, \mathcal{K}) \rangle$: $\mathbf{D}_A = (N, d, \Delta)$, $\mathbf{D}_B = (\mathbf{D}, \mathbf{T}, \mathbf{C})$

Precondition: The query range $\mathbf{R} \subset [\mathcal{Z}]^d$ is a rectangular range.

Step 1: Alice partitions the complement range \mathbf{R}^c into $2d$ rectangular ranges $\{\mathbf{R}_\ell \subset [\mathcal{Z}]^d : \ell \in [1, 2d]\}$, and sets $\mathbf{R}_0 = \mathbf{R}$.

Step 2—Reduction: For $0 \leq \ell \leq 2d$, Alice and Bob invokes CollRes on range \mathbf{R}_ℓ . Denote the output as $(\zeta_\ell, X_\ell, \Psi^{(\ell)})$.

Step 3: Alice verifies whether the following equalities hold:

$$\forall 0 \leq \ell \leq 2d, \zeta_\ell \stackrel{?}{=} \text{accept}, \quad \bigotimes_{0 \leq \ell \leq 2d} \Psi^{(\ell)} \stackrel{?}{=} \Delta. \quad (17)$$

Note: Until this point, all are identical to the Count algorithm.

Step 4: Bob sends back $S = \{\mathbf{x}[1] : \mathbf{x} \in \mathbf{D} \cap \mathbf{R}\}$ to Alice.

Step 5: Alice verifies whether the following equalities hold: Let $\Psi^{(0)}[2]$ denote the 2nd component of vector value $\Psi^{(0)}$.

$$\Psi^{(0)}[2] \stackrel{?}{=} \prod_{\mathbf{x} \in S} F_\varpi(\mathbf{x}[1])^\beta; \quad |S| = X_0 \quad (18)$$

Alice sets $\zeta = \text{accept}$ if all verifications in equation (17) and equation (18) succeed; and sets $\zeta = \text{reject}$ otherwise. Alice outputs (ζ, S, Δ) .

(Alice, Bob) $\text{CollRes} = \langle \widetilde{\text{Eval}}(\mathbf{D}_B), \widetilde{\text{Ext}}(\mathbf{D}_A, \mathbf{R}, \mathcal{K}) \rangle$: $\mathbf{D}_A = (N, d, \Delta)$, $\mathbf{D}_B = (\mathbf{D}, \mathbf{T}, \mathbf{C})$

Identical with CollRes in Figure 1. Save the details.

6 Dynamic Dataset

6.1 Insertion

Insert($\hat{\mathbf{D}}, \mathcal{K}$):

Precondition: Alice has $\mathbf{D}_s = (\mathcal{K}, N, d, \Delta)$; Bob has $\mathbf{D}_p = (\mathbf{D}, \mathbf{T}, \mathbf{C})$.

Alice runs the algorithm $\text{DEnc}(\hat{\mathbf{D}}, \mathcal{K})$ to generate $(\hat{\mathbf{D}}_s = (\hat{N}, d, \hat{\Delta}), \hat{\mathbf{D}}_p = (\hat{\mathbf{D}}, \hat{\mathbf{T}}, \hat{\mathbf{C}}))$. Alice updates $\Delta \leftarrow \Delta \cdot \hat{\Delta}, N \leftarrow N + \hat{N}$. Alice sends $(\hat{\mathbf{D}}, \hat{\mathbf{T}}, \hat{\mathbf{C}})$ to Bob. Bob sets $\mathbf{D} = \mathbf{D} \cup \hat{\mathbf{D}}, \mathbf{T} = \mathbf{T} \cup \hat{\mathbf{T}}, \mathbf{C} = \mathbf{C} \cup \hat{\mathbf{C}}$.

We can prove the security if insertion is non-adaptive, i.e. the inserted items are sampled from a particular distribution.

Theorem 5 *The extended scheme is correct and sound to authenticate d -dimensional COUNT, SUM, AVG, MIN, MAX, MEDIAN and range selection queries over dynamic dataset that supports insertion.*

6.2 Deletion

Deletion is equivalent to insertion into another dataset.

Let $\mathcal{E} = (\text{KGen}, \text{DEnc}, \text{ProVer})$.

$\text{KGen}(1^\kappa)$: Run $\mathcal{E}.\text{KGen}(1^\kappa)$ twice independently and output two keys \mathcal{K} and $\bar{\mathcal{K}}$.

$\text{DEnc}(\mathbf{D}; \mathcal{K}, \bar{\mathcal{K}})$: Run $\mathcal{E}.\text{DEnc}(\mathbf{D}; \mathcal{K})$ to generate $(\mathbf{D}_A, \mathbf{D}_B)$. Set $\bar{\mathbf{D}} = \bar{\mathbf{T}} = \bar{\mathbf{C}} = \emptyset, \bar{\mathbf{D}}_A = (\bar{N} = 0, d, \bar{\Delta} = 1)$, and $\bar{\mathbf{D}}_B = (\bar{\mathbf{D}}, \bar{\mathbf{T}}, \bar{\mathbf{C}}, \bar{pk})$, where \bar{pk} is part of $\bar{\mathcal{K}}$.

Insert(\mathbf{D}', \mathcal{K}): Run $\mathcal{E}.\text{Insert}(\mathbf{D}', \mathcal{K})$.

Delete(\mathbf{D}'): Set $\mathbf{D}' \leftarrow \mathbf{D}' \cap \mathbf{D}$. Run $\mathcal{E}.\text{Insert}(\mathbf{D}', \bar{\mathcal{K}})$ to insert points in \mathbf{D}' into the complement dataset $\bar{\mathbf{D}}$.

Sum($\mathbf{R}; \mathcal{K}, \bar{\mathcal{K}}$): Run $\mathcal{E}.\text{Sum}(\mathbf{R}; \mathcal{K})$ over dataset \mathbf{D} to obtain (ζ, X, Δ) ; run $\mathcal{E}.\text{Sum}(\mathbf{R}; \bar{\mathcal{K}})$ over the complement dataset $\bar{\mathbf{D}}$ to obtain $(\bar{\zeta}, \bar{X}, \bar{\Delta})$. If $\zeta = \bar{\zeta} = \text{accept}$, then set $\varsigma = \text{accept}$; otherwise set $\varsigma = \text{reject}$. Output $(\varsigma, X - \bar{X}, \Delta/\bar{\Delta})$.

Min($\mathbf{R}, \iota; \mathcal{K}, \bar{\mathcal{K}}$): Assume $\sum_{\mathbf{x} \in \mathbf{D}} \mathbf{x} < (p, \dots, p)$

1. Alice sends range \mathbf{R} to Bob.
 2. Bob finds $\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbf{D} \cap \mathbf{R}} \mathbf{x}[\iota]$. Bob sends \mathbf{x}^* back to Alice.
 3. Alice issues SUM query with range $\mathbf{R} = \{\mathbf{x}^*\}$ to Bob over dataset \mathcal{D} and $\bar{\mathbf{D}}$, and obtains output $(\zeta, X - \bar{X}, \Delta/\bar{\Delta})$. If $\zeta = \text{accept}$ and $X - \bar{X} > 0$, then believes that \mathbf{x}^* is a valid data point.
 4. Alice issues a COUNT query.
-

Corollary 6 *The extended scheme is correct and sound to authenticate d -dimensional COUNT, SUM, AVG, MIN, MAX, MEDIAN and range selection queries over dynamic dataset that supports both insertion and deletion.*

7 Security beyond authentication

7.1 Privacy

At first, let us distinguish aggregate attributes and selection attributes: (1) Aggregate attributes are dimensions along which a query apply the aggregate operation (like sum, min, max); (2) Selection attributes are dimensions on which a query applys range constraint. Take an example, a sum query which asks for the sum of the 3rd dimension of data points with selection on the 1st and th 2nd dimensions: Let $\mathbf{R} = [a_1, b_1] \times [a_2, b_2] \times [1, \mathcal{Z}]^{d-2}$ be the query range.

$$\sum_{\mathbf{x} \in \mathbf{R}} \mathbf{x}[3]$$

In this example, the 3rd dimension is the aggregate attribute, and the 1st and the 2nd are selection attributes. Note that in some query, a dimension can be both aggregate attribute and selection attribute.

We found previous works on privacy preserving aggregate range query over outsourced dataset can be divided into two categories:

- Protect the privacy of aggregate attributes, where any aggregate attribute is not a selection attribute, e.g. [7]. These works typically employ homomorphic encryption scheme to hide the aggregate attribute values and reserve the capability of doing aggregation. Particularly, additive homomorphic encryption scheme (e.g. Paillier system [34]) for aggregate SUM,AVG query and order preserving encryption scheme (e.g. [35]) for aggregate min/max query.
- Protect the privacy of selection attributes. fully homomorphic encryption scheme [30]. It is worthy to point out that order preserving encryption scheme (e.g. [35]) and MRQED scheme [36] based approaches are not secure against adaptive adversary.

Privacy of Aggregate Attributes against Adaptive Adversary Our solution can achieve similar privacy protection as Yao *et al.* [7].

W.L.O.G, we assume only the 1st dimension is aggregate attribute, and in every query range \mathbf{R} has the form $\mathbf{R} = [1, \mathcal{Z}] \times [a_2, b_2] \times \dots \times [a_d, b_d] \subseteq [1, \mathcal{Z}]^d$, i.e. the query has no constraint on the 1st dimension. Let $(\mathbf{G}, \mathbf{E}, \mathbf{D}, \mathbf{H})$ be an additive homomorphic encryption scheme (e.g. Paillier system [34]).

The new scheme is identical to the solution for aggregate sum query in Section 5.1, except that

- Additionally, in the setup, Alice generates a key pair (K_E, K_D) by running the key generating algorithm \mathbf{G} , and for each point $\mathbf{x} \in \mathbf{D}$ replaces the first dimension $\mathbf{x}[1]$ with the ciphertext $E_{K_E}(\mathbf{x}[1])$. Next, apply \mathbf{DEnc} on the $(d-1)$ -dimensional dataset $\mathbf{D}' = \{(x_2, \dots, x_d) : (x_1, x_2, \dots, x_d) \in \mathbf{D}\}$.
- To answer a sum query over 1st dimension, Bob “sums” all ciphertexts $E_{K_E}(\mathbf{x}[1])$ for points $\mathbf{x} \in \mathbf{D} \cap \mathbf{R}$ using the homomorphic property of the encryption scheme, i.e. using algorithm \mathbf{H} , and sends the resulting ciphertext $\mathbf{CT} = E_{K_E}(\sum_{\mathbf{x} \in \mathbf{D} \cap \mathbf{R}} \mathbf{x}[1])$ of sum to Alice as the query result.
- Alice can decrypt ciphertext \mathbf{CT} with decryption key K_D to recover the sum $\sum_{\mathbf{x} \in \mathbf{D} \cap \mathbf{R}} \mathbf{x}[1]$.

Privacy of Selection Attributes against Memoryless Adversary [28,36]

To the best of our knowledge, Gennaro *et al.* [28] is the only available solution which can protect privacy against adaptive adversaries. Here, we present two methods based on homomorphic encryption scheme, which is secure against non-adaptive adversary who has no memories.

A straightforward approach is to apply order preserving encryption scheme [37,35]. During the setup, Alice can encrypt each selection attribute with an order preserving encryption scheme. Since the order between any two values are preserved, Bob can do comparison directly.

Alternatively, we may apply MRQED, which is predicate encryption scheme supporting multidimensional range query. Under MRQED, a message \mathbf{Msg} can be encrypted under an identity \mathbf{x} , which is a point in a d -dimensional space $[1, \mathcal{Z}]^d$. From the master secret key, a delegation key δ w.r.t. a d -dimensional rectangular range \mathbf{R} can be derived. With the delegation key δ , the ciphertext for the message \mathbf{Msg} under identity point \mathbf{x} can be decrypted to recover \mathbf{Msg} , iff the identity point \mathbf{x} is within the range \mathbf{R} .

Let \mathbb{M} be the domain of the messages to be encrypted under MRQED, and $\hat{\mathbb{M}}$ be a subset of \mathbb{M} such that (1) the size of $\hat{\mathbb{M}}$ is superpolynomial; (2) the ratio $\frac{|\hat{\mathbb{M}}|}{|\mathbb{M}|}$ is negligible. During the setup, for each data point $\mathbf{x} \in \mathbf{D}$, Alice choose a random message $\mathbf{Msg}_{\mathbf{x}} \in \hat{\mathbb{M}}$, and encrypts the message $\mathbf{Msg}_{\mathbf{x}}$ under identity point \mathbf{x} using MRQED encryption scheme. Alice replaces each data point with its corresponding ciphertext and sends all of N resulting ciphertexts to Bob. Later, Alice wants to query range \mathbf{R} , then she can derive the delegation key δ w.r.t \mathbf{R} and sends δ to Bob. With the delegation key δ , Bob can decide whether a ciphertext is corresponding to a data point \mathbf{x} within the query range \mathbf{R} , without knowing the value of \mathbf{x} : Bob decrypts each ciphertext $\mathbf{C}_{\mathbf{x}}$ with the delegation key, and gets the decrypted value $\mathcal{M}_{\mathbf{x}}$. If $\mathcal{M}_{\mathbf{x}} \in \hat{\mathbb{M}}$, then $\mathbf{x} \in \mathbf{R}$ with o.h.p. Otherwise, $\mathbf{x} \notin \mathbf{R}$ definitely.

It is worthy to point out that, the above two approaches using order preserving encryption and MRQED is secure in privacy protection if only one query is allowed. After multiple queries, Bob may be able to infer the value of some data point, using the information that whether the point is inside or outside previous query ranges.

7.2 Frame Attack

Let $(\mathbf{KG}, \mathbf{Sign}, \mathbf{Verify})$ be a secure digital signature scheme. Suppose Alice has signing key (PK_A, SK_A) and Bob has signing key (PK_B, SK_B) , where only Alice knows the private key SK_A , only Bob knows the private key SK_B , and the two public keys PK_A and PK_B are known to public. We assume there is no *Denial of Service* (DOS) attack.

7.3 Dynamic Dataset

In addition to our original scheme, Alice and Bob are required to do the following steps to prevent frame attack.

1. During the setup, Alice signs the dataset \mathbf{D} with her private key and sends the signature $\text{Sign}_{SK_A}(\mathbf{D})$ to Bob.
2. For each update command \mathcal{U} (e.g. insertion x or deletion y) that Alice issues to Bob, Alice has to sign it and sends the signature $\text{Sign}_{SK_A}(\mathcal{U})$ together with the update command to Bob. Next, Bob sends an ACK message with signature, i.e. $\text{Sign}_{SK_B}(\mathcal{U})$ to Alice.
3. For each query \mathcal{Q} Alice issues to Bob, Bob generates the query result \mathcal{X} and proof Π . Bob sends back to Alice the signed result with proof, i.e. $(\mathcal{X}, \Pi, \text{Sign}_{SK_B}(\mathcal{Q}, \mathcal{X}, \Pi))$.
4. In all above cases, the receiver of a signature will verify the validity of the signature using the corresponding public key, and rejects that reply message and asks the sender to resend the corrupted message if the signature is not valid.

To claim that Bob returned a wrong result for query \mathcal{Q} , Alice has to present to the third trusted party two pieces of information: (1) the query \mathcal{Q} and the corresponding signed result: $(\mathcal{X}, \Pi, \text{Sign}_B(\mathcal{Q}, \mathcal{X}, \Pi))$; (2) the set S_{ACK} of all ACKs with Bob's signatures. On the other hand, to prove his innocence, Bob has to present to the third trusted party the original dataset \mathbf{D} together with Alice's signature, and the set $S_{\mathcal{U}}$ of all update commands signed by Alice.

The third trusted party verifies all signatures using corresponding public keys. If Alice presents a message which is wrongly signed by Bob, then decides Alice cheats. Similarly, if Bob presents a message which is wrongly signed by Alice, then decides Bob cheats. The third trusted party then checks the authenticated messages in the following way:

- If $S_{ACK} \subsetneq S_{\mathcal{U}}$, then judges that Bob cheats.
- If $S_{\mathcal{U}} \subsetneq S_{ACK}$, then judges that Alice cheats.
- Until this point, All parties have a consensus on the current status of dataset: Let the dataset \mathbf{D}^* be the resulting dataset after applying the authenticated update commands in $S_{\mathcal{U}}$ to the original dataset \mathbf{D} . Then, compute the query \mathcal{Q} over the authenticated dataset \mathbf{D}^* , and gets the result \mathcal{Y} . If $\mathcal{X} \neq \mathcal{Y}$, then judges that Bob cheats.

For the static case, where the update commands are not allowed, all actions/signatures on updates can be saved and Bob only needs to provide the signed original dataset.

8 Conclusion

We propose efficient schemes to authenticate queries over static/dynamic outsourced dataset with $O(d^2 \log^2 \mathcal{Z})$ communication overhead, which conquer the “curse of dimensionality”. The supported queries include aggregate range query, i.e. COUNT, SUM, AVG, MIN, MAX, MEDIAN and range selection.

References

1. Xu, J., Chang, E.C.: Authenticating aggregate range queries over multidimensional dataset. Cryptology ePrint Archive, Report 2010/050 (2010) <http://eprint.iacr.org/>.
2. Devanbu, P.T., Gertz, M., Martel, C.U., Stubblebine, S.G.: Authentic Third-party Data Publication. In: Proceedings of the IFIP TC11/ WG11.3 Fourteenth Annual Working Conference on Database Security. (2001) 101–112
3. Hacigümüş, H., Iyer, B., Li, C., Mehrotra, S.: Executing SQL over encrypted data in the database-service-provider model. In: SIGMOD '02: ACM SIGMOD International conference on Management of data. (2002) 216–227
4. Wu, J., Stinson, D.: An Efficient Identification Protocol and the Knowledge-of-Exponent Assumption. Cryptology ePrint Archive, Report 2007/479 (2007) <http://eprint.iacr.org/>.
5. Pang, H., Tan, K.L.: Verifying Completeness of Relational Query Answers from Online Servers. ACM Trans. Inf. Syst. Secur. **11**(2) (2008) 1–50
6. Cheng, W., Tan, K.L.: Query assurance verification for outsourced multi-dimensional databases. J. Comput. Secur. **17**(1) (2009) 101–126
7. Thompson, B., Yao, D., Haber, S., Horne, W.G., Sander, T.: Privacy-Preserving Computation and Verification of Aggregate Queries on Outsourced Databases. In: PETS '09: Privacy Enhancing Technologies Symposium. (2009)
8. Li, F., Hadjieleftheriou, M., Kollios, G., Reyzin, L.: Authenticated index structures for aggregation queries. ACM Trans. Inf. Syst. Secur. **13** (2010) 32:1–32:35

9. Atallah, M.J., Cho, Y., Kundu, A.: Efficient Data Authentication in an Environment of Untrusted Third-Party Distributors. In: ICDE '08: IEEE International Conference on Data Engineering. (2008) 696–704
10. Martel, C., Nuckolls, G., Devanbu, P., Gertz, M., Kwong, A., Stubblebine, S.G.: A General Model for Authenticated Data Structures. *Algorithmica* **39**(1) (2004) 21–41
11. Chen, H., Ma, X., Hsu, W.W., Li, N., Wang, Q.: Access Control Friendly Query Verification for Outsourced Data Publishing. In: ESORICS '08: European Symposium on Research in Computer Security. (2008) 177–191
12. Hacigümüs, H., Iyer, B.R., Mehrotra, S.: Efficient Execution of Aggregation Queries over Encrypted Relational Databases. In: DASFAA. (2004) 125–136
13. Mykletun, E., Tsudik, G.: Aggregation Queries in the Database-As-a-Service Model. In: IFIP WG 11.3 Working Conference on Data and Applications Security. (2006) 89–103
14. Ge, T., Zdonik, S.B.: Answering Aggregation Queries in a Secure System Model. In: VLDB '07: International Conference on Very Large Data Bases. (2007) 519–530
15. Devanbu, P., Gertz, M., Martel, C., Stubblebine, S.G.: Authentic data publication over the internet. *J. Comput. Secur.* **11**(3) (2003) 291–314
16. Pang, H., Jain, A., Ramamritham, K., Tan, K.L.: Verifying completeness of relational query results in data publishing. In: SIGMOD '05: ACM SIGMOD International conference on Management of data. (2005) 407–418
17. Mykletun, E., Narasimha, M., Tsudik, G.: Authentication and Integrity in Outsourced Databases. *Trans. Storage* **2**(2) (2006) 107–138
18. Sion, R.: Query Execution Assurance for Outsourced Databases. In: VLDB '05: International Conference on Very Large Data Bases. (2005) 601–612
19. Li, F., Hadjieleftheriou, M., Kollios, G., Reyzin, L.: Dynamic authenticated index structures for outsourced databases. In: SIGMOD '06: ACM SIGMOD International conference on Management of data. (2006) 121–132
20. Xie, M., Wang, H., Yin, J., Meng, X.: Integrity auditing of outsourced data. In: VLDB '07: International conference on Very large data bases. (2007) 782–793
21. Yang, Y., Papadias, D., Papadopoulos, S., Kalnis, P.: Authenticated join processing in outsourced databases. In: SIGMOD '09: ACM SIGMOD International conference on Management of data. (2009) 5–18
22. Mouratidis, K., Sacharidis, D., Pang, H.: Partially materialized digest scheme: an efficient verification method for outsourced databases. *The VLDB Journal* **18**(1) (2009) 363–381
23. Pang, H., Zhang, J., Mouratidis, K.: Scalable Verification for Outsourced Dynamic Databases. *Proc. VLDB Endow.* **2** (2009) 802–813
24. Goodrich, M.T., Tamassia, R., Triandopoulos, N.: Super-Efficient Verification of Dynamic Outsourced Databases. In: CT-RSA '08: The Cryptographer's Track at the RSA Conference on Topics in Cryptology. (2008) 407–424
25. Rivest, R.L., Shamir, A., Adleman, L.: A method for obtaining digital signatures and public-key cryptosystems. *Commun. ACM* **21**(2) (1978) 120–126
26. Boneh, D., Lynn, B., Shacham, H.: Short Signatures from the Weil Pairing. *J. Cryptol.* **17**(4) (2004) 297–319
27. Haber, S., Horne, W., Sander, T., Yao, D.: Privacy-Preserving Verification of Aggregate Queries on Outsourced Databases. Technical report, HP Laboratories (2006) HPL-2006-128.
28. Gennaro, R., Gentry, C., Parno, B.: Non-interactive Verifiable Computing: Outsourcing Computation to Untrusted Workers. In: CRYPTO '10: Annual International Cryptology Conference on Advances in Cryptology. (2010) 465–482
29. Chung, K.M., Kalai, Y., Vadhan, S.P.: Improved Delegation of Computation Using Fully Homomorphic Encryption. In: CRYPTO '10: Annual International Cryptology Conference on Advances in Cryptology. (2010) 483–501
30. Gentry, C.: Fully Homomorphic Encryption using Ideal Lattices. In: STOC '09: ACM symposium on Theory of computing. (2009) 169–178
31. van Dijk, M., Gentry, C., Halevi, S., Vaikuntanathan, V.: Fully Homomorphic Encryption over the Integers. In: EUROCRYPT '10: Annual International Conference on Advances in Cryptology. (2010) 24–43
32. Gentry, C.: Toward Basing Fully Homomorphic Encryption on Worst-Case Hardness. In: CRYPTO '10: Annual International Cryptology Conference on Advances in Cryptology. (2010) 116–137
33. Juels, A., Kaliski, Jr., B.S.: Pors: proofs of retrievability for large files. In: CCS '07: ACM conference on Computer and communications security. (2007) 584–597
34. Paillier, P.: Public-Key Cryptosystems Based on Composite Degree Residuosity Classes. In: EUROCRYPT '99: Annual International Conference on Advances in Cryptology. (1999) 223–238
35. Boldyreva, A., Chenette, N., Lee, Y., O'Neill, A.: Order-Preserving Symmetric Encryption. In: EUROCRYPT '09: Annual International Conference on Advances in Cryptology. (2009) 224–241

36. Shi, E., Bethencourt, J., Chan, T.H.H., Song, D., Perrig, A.: Multi-Dimensional Range Query over Encrypted Data. In: SP '07: IEEE Symposium on Security and Privacy. (2007) 350–364
37. Agrawal, R., Kiernan, J., Srikant, R., Xu, Y.: Order preserving encryption for numeric data. In: SIGMOD '04: ACM SIGMOD international conference on Management of data. (2004) 563–574