

Authenticating Aggregate Range Queries over Multidimensional Dataset

Jia Xu
National University of Singapore
xujia@comp.nus.edu.sg

Ee-Chien Chang
National University of Singapore
changec@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 D of d -dimensional points, together with some authentication tag T , to an untrusted service provider Bob. Later, Alice issues some query over D to Bob, and Bob should produce a query result and a proof based on D and 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 T . We give a method that requires $O(d^2)$ communication bits to authenticate an aggregate query conditional on d -dimensional range selection. Besides counting, summing and finding of the minimum can also be supported. Furthermore, our scheme can be extended to authenticate d -dimensional usual (non-aggregate) range selection query with $O(d^2)$ bits communication overhead, improving known results that require $O(\log^{d-1} N)$ communication overhead, where N is the number of data points in the dataset.

Keywords

Authentication, Multidimensional Aggregate Query, Secure Outsourced Database, Provable Remote Computing

1. INTRODUCTION

Alice has a set D of d -dimensional points. She preprocesses the dataset D using her private key to generate some authentication tag T . She sends (outsources) D and T to an untrusted service provider Bob. Then Alice deletes the original copy of dataset D and tag T from her local storage. Later Alice may issue a query over 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 D and T . Alice wants to authenticate the query result, using only her private key.

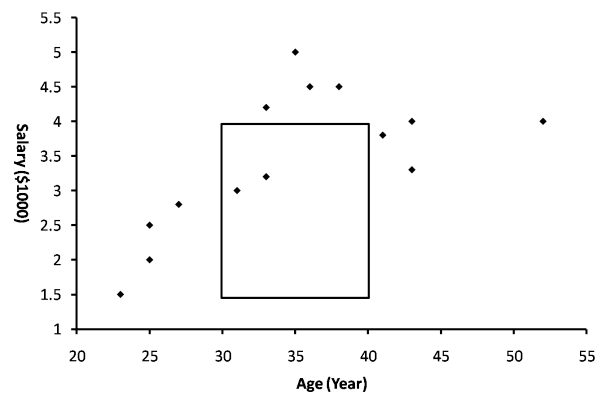
Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.
Copyright 20XX ACM X-XXXXX-XX-X/XX/XX ...\$10.00.

We are concerned about the communication cost and the storage overhead on Bob's side. Such requirements exclude the following two straightforward approaches: (1) Bob sends back the whole dataset D with its tag T ; (2) During preprocessing, Alice generates and signs answers to all possible queries.

The problem we study in this paper fits in the framework of the outsourced database applications [20, 51], which emerged in early 2000s as an example of "software-as-a-service" (SaaS). By outsourcing database management, backup services and other IT needs to a professional service provider, companies can reduce expensive cost in purchase of equipments and even more expensive cost in hiring or training qualified IT specialists to maintain the IT services [24].

Researches in secure outsourced database focus on two major aspects: privacy [51, 52, 53, 54] (i.e. protect the data confidentiality against both the service provider and any third party), and integrity [20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34] (i.e. authenticate the soundness and completeness of query results returned by the service provider). In the latter aspect, a lot of works are done for "identity query" [26], i.e. the query result is a subset of the database. Aggregate range query is arguably more challenging and only a few works (e.g. [7]) are devoted to the authentication of aggregate query.

Figure 1: An example of 2D aggregate range query: How many employees with age between 30 and 40 have salary between \$1500 and \$4000? The query range is $[30, 40] \times [1500, 4000]$ and the query result is 2.



In this paper, we are interested in the authentication of aggregate range query over static multidimensional dataset with sublinear (w.r.t. the number of data points within the query range) communication bits. A d -dimensional aggregate range query specifies a d -dimensional rectangular range (which can be represented by its two “end points”), and its outcome is the aggregated value over all points inside the range. Figure 1 shows a 2D aggregate range query. The aggregate operations we consider in this paper include counting, summing, and finding of the maximum and minimum. Besides aggregate query, we also discuss the usual range selection query as an extension.

Our results

We propose a scheme, which we call MAIA (Multidimensional Aggregate query Integrity Authentication). For a dataset \mathcal{D} with N d -dimensional points, the number of communication bits required is in $O(d^2)$ per query. The storage overhead on Bob’s side is $O(dN)$, which is linear w.r.t. the size of \mathcal{D} .

We now describe our main ideas in three steps: (1) We describe a plausible scheme, which is insecure and requires large amount of communication, to authenticate aggregate range query. (2) We point out problems in this scheme and difficulties to implement it. (3) We summarize the key techniques we used to conquer such problems or difficulties.

Plausible Scheme

For simplicity, let \mathcal{D} be a set of N 1D points. Let T be an additive homomorphic authentication tag function, such that $\mathsf{T}(\mathbf{x} + \mathbf{y}) = \mathsf{T}(\mathbf{x})\mathsf{T}(\mathbf{y})$. During setup, Alice preprocesses the dataset \mathcal{D} with a private key, to generate a tag value $\mathsf{T}(\mathbf{x})$ for each point $\mathbf{x} \in \mathcal{D}$. Next, Alice sends \mathcal{D} and tag values $\{\mathsf{T}(\mathbf{x}) : \mathbf{x} \in \mathcal{D}\}$ to Bob and deletes everything except the private key from her storage. Let us consider a summing query conditional on range \mathbf{B} which asks for the sum $\sum_{\mathbf{x} \in \mathcal{D} \cap \mathbf{B}} \mathbf{x}$. Bob is expected to send to Alice a number \mathbf{X} as the query result, and a proof to show that indeed $\mathbf{X} = \sum_{\mathbf{x} \in \mathcal{D} \cap \mathbf{B}} \mathbf{x}$.

To process this query, Alice sends an auxiliary message (called as *Help-Info*) $\{(\mathsf{T}(\mathbf{x}))^\alpha : \mathbf{x} \in \mathbf{B}\}$ to Bob, in order to help Bob to generate a proof for the query result. Here the secret random nonce α prevents Bob from abusing this *Help-Info* for other queries. Bob is expected to:

1. sum all $\mathbf{x} \in \mathcal{D} \cap \mathbf{B}$ to obtain \mathbf{X} : $\mathbf{X} = \sum_{\mathbf{x} \in \mathcal{D} \cap \mathbf{B}} \mathbf{x}$;
2. aggregate all $\mathsf{T}(\mathbf{x})$ ’s for $\mathbf{x} \in \mathcal{D} \cap \mathbf{B}$ to obtain Ψ_1 :

$$\Psi_1 = \prod_{\mathbf{x} \in \mathcal{D} \cap \mathbf{B}} \mathsf{T}(\mathbf{x}) = \mathsf{T}\left(\sum_{\mathbf{x} \in \mathcal{D} \cap \mathbf{B}} \mathbf{x}\right);$$

3. aggregate all $(\mathsf{T}(\mathbf{x}))^\alpha$ ’s for $\mathbf{x} \in \mathcal{D} \cap \mathbf{B}$ to obtain Ψ_2 :

$$\Psi_2 = \prod_{\mathbf{x} \in \mathcal{D} \cap \mathbf{B}} (\mathsf{T}(\mathbf{x}))^\alpha = \left(\mathsf{T}\left(\sum_{\mathbf{x} \in \mathcal{D} \cap \mathbf{B}} \mathbf{x}\right)\right)^\alpha.$$

Bob sends back \mathbf{X} as query result and (Ψ_1, Ψ_2) as proof. Alice can verify the correctness, by checking the following two equalities:

$$\begin{aligned} \Psi_1 &\stackrel{?}{=} \mathsf{T}(\mathbf{X}) \\ \Psi_1^\alpha &\stackrel{?}{=} \Psi_2 \end{aligned}$$

For data point $\mathbf{x} \in \mathcal{D} \cap \mathbf{B}$, Bob has both $\mathsf{T}(\mathbf{x})$ and $(\mathsf{T}(\mathbf{x}))^\alpha$. For any point $\mathbf{x} \notin \mathcal{D} \cap \mathbf{B}$, Bob cannot provide either $\mathsf{T}(\mathbf{x})$ or $(\mathsf{T}(\mathbf{x}))^\alpha$, assuming it is difficult for Bob to forge $\mathsf{T}(\mathbf{x})$ or $(\mathsf{T}(\mathbf{x}))^\alpha$ without the private key and the random nonce α . Hence, Bob is unable to include points outside $\mathcal{D} \cap \mathbf{B}$ in the computation of query result \mathbf{X} and proof (Ψ_1, Ψ_2) .

Problems and Difficulties

We identify 3 problems in the above scheme: (1) Bob still can cheat by double counting or undercounting some points in the set $\mathcal{D} \cap \mathbf{B}$. For example, Bob may send $2\mathbf{X}$ as query result and $(\Psi_1 \times \Psi_1, \Psi_2 \times \Psi_2)$ as proof. Due to homomorphism of T , Alice will consider this proof valid. (2) Additive homomorphism is unlikely secure [49], and the assumption of unforgeability of $\mathsf{T}(\mathbf{x})$ and $(\mathsf{T}(\mathbf{x}))^\alpha$ may not hold. (3) The size of *Help-Info* is linear w.r.t. the size of query range \mathbf{B} , which is huge and possibly comparable to the domain size of a data point. This implies large computation cost (to generate *Help-Info*) and communication cost (to send *Help-Info*).

Our solution

We employ or design several key techniques to solve these problems, including:

1. To prevent double counting or undercounting, Alice will check the complement query with range \mathbf{B}^c , which asks for the value of $\sum_{\mathbf{x} \in \mathcal{D} \cap \mathbf{B}^c} \mathbf{x}$. Nevertheless, a malicious Bob may undercount some points in $\mathcal{D} \cap \mathbf{B}$ and compensate it by double counting some points in $\mathcal{D} \cap \mathbf{B}^c$. To prevent such cheating, we introduce a secret and random attribute to each data point and aggregate it alongside with the requested value. Note such strategy that verifies completeness by checking the complement is also used in Chen *et al.* [19].
2. To decrease the size of query range \mathbf{B} , we “normalize” the dataset \mathcal{D} , so that the values of normalized data points are in $[1, N]$. Note N is the number of points in \mathcal{D} .
3. To further decrease the size of *Help-Info*, we separate each dimension of a query range, and are able to compress *Help-Info* for each dimension. For example, let $\mathbf{B} = [a_1, b_1] \times [a_2, b_2]$ be a 2D range. Alice compresses *Help-Info* for 1D range $[a_1, b_1]$ to obtain ϖ_1 , and compresses *Help-Info* for 1D range $[a_2, b_2]$ to obtain ϖ_2 . From (ϖ_1, ϖ_2) , Bob is able to reconstruct *Help-Info* for the 2D range \mathbf{B} . Such separation of dimensions might cause collusion attack [1].
4. We specially design an authentication tag function, which (1) allows Bob to do aggregation to generate a short proof; (2) is unforgeable; (3) allows *Help-Info* to be compressed and yet prevents collusion attack; (4) binds a secret and random attribute to each point and allows Alice to verify completeness by checking the complement.

Our tag function, denoted as Tag , consists of three components DTag , ITag and CTag . Let i be the normalized value of point $\mathbf{x} \in \mathcal{D}$.

1. $\mathsf{DTag}(\mathbf{x})$ binds different dimensions of a d -dimensional point \mathbf{x} , so Bob cannot forge a new point by mixing and combining different dimensions of different points (e.g. $(\mathbf{x}_2[1], \mathbf{x}_1[2], \mathbf{x}_1[3], \dots, \mathbf{x}_1[d])$);

Table 1: Performance of different authentication schemes for aggregate range query or range selection query. This table consists of two parts: the first two rows are for aggregate query; the rest four rows are for range selection query. Our scheme MAIA appears twice in this table, since it can authenticate both aggregate query and range selection query (with modifications in Section 6.2.4). Note: (1) The symbol “-” indicates that the authors do not provide such information in their paper; (2) Our scheme MAIA 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.

Scheme	Dimension d	Communication overhead (bits)	Key Size	Storage overhead	Computation (Verifier Alice)	Computation (Prover Bob)	Query
PDAS [7]	$d = 1$	$O(S \log N)$	$O(1)$	$O(N)$	$O(S \log N)$	$O(S + K^2)$	SUM, COUNT
MAIA (This paper)	$d \geq 1$	$O(d^2)$	$O(d)$	$O(dN)$	$O(d \log N)\dagger\star$	$O(dN \log N)\ddagger\star$	SUM, COUNT, MIN, MAX, MEDIAN, QUANTILE
Atallah <i>et al.</i> [30]	$d = 1, 2$	$O(1)$	$O(1)$	$O(N)$	$O(S)$	$O(1)$	Range Selection
Martel <i>et al.</i> [21]	$d \geq 1$	$O(\log^{d-1} N + S)$	-	-	-	-	Range Selection
Chen <i>et al.</i> [19]	$d \geq 1$	$O(\log^d M)$	-	$O(N \log^d M)$	$O(\log^d M)$	$O(\log^d M)$	Range Selection
MAIA (Section 6.2.4)	$d \geq 1$	$O(d^2)$	$O(d)$	$O(dN)$	$O(d \log N)\dagger\star$	$O(dN \log N)\ddagger\star$	Range Selection

N : The number of tuples in the dataset.

K : The number of servers in PDAS [7].

\dagger : $O(d \log N)$ modular multiplications.

\star : If the query range is 1D, the cost is $O(|S|)$.

S : The set of tuples satisfying the query condition.

M : The domain size of attributes in Chen *et al.* [19].

\ddagger : $O(dN \log N)$ modular exponentiations.

2. $\text{ITag}(i)$ is an authentication tag of i and allows Alice to compress the *Help-Info* (The compression is based on redactable signature scheme [49]);
3. $\text{CTag}(i)$ prevents double counting and undercounting.

With private key, Alice is able to “extract” $\text{DTag}(x)$, $\text{ITag}(i)$ or $\text{CTag}(i)$ from the value of $\text{Tag}(x, i)$. Without the private key, Bob cannot decompose $\text{Tag}(x, i)$. We construct Tag using three cyclic multiplicative subgroups of \mathbb{Z}_n^* , where n is the product of three safe primes [5, 55].

The main contributions in this paper can be summarized as below:

1. We propose a scheme called MAIA (Section 5), to authenticate multidimensional aggregate range queries, based on a specially designed tag function Tag (Section 4).
2. MAIA is efficient (See Table 1) and takes only $O(d^2)$ communication bits for a d -dimensional aggregate range query, independent of the query range size and dataset size.
3. We prove that MAIA is secure (Theorem 2) under reasonable assumptions (Assumption 1, Assumption 2, etc.).
4. We extend MAIA to authenticate multidimensional range selection query (Section 6.2.4). The performance is showed in Table 1.

2. RELATED WORK

Privacy-preserving computation and integrity verification are two major aspects of the security of outsourced computing. Many works in cryptography can be casted as privacy-preserving computation over ciphertexts, including homomorphic encryption [8, 9, 11], Attribute Based Encryption [2,

3], Predicate Encryption [4, 1, 5, 6], and Order Preserving Encryption [10].

Shi *et al.* [1] proposed MRQED (Multi-Dimensional Range Query over Encrypted Data), a public key encryption scheme supporting multidimensional range queries over ciphertexts. Both MRQED [1] and MAIA deal with multidimensional range selection and have to prevent collusion attack across different queries. But they are essentially different in at least these aspects: (1) MRQED dealt with privacy, and MAIA deals with integrity. (2) In MAIA, there is an aggregate operation after multidimensional range selection, and the verification of aggregated value is an additional requirement and not easy to handle when communication cost is concerned. (3) Besides collusion attack, MAIA also faces other challenges, like undercounting and double-counting attacks, which have no counterparts in researches of privacy-preserving computation, like MRQED. It is unlikely to solve our problem using MRQED as a black-box, and we doubt whether it is possible to design an alternative solution based on MRQED without significant modifications. However, MAIA and MRQED indeed can be benefited from each other, especially in how to represent a range compactly. We present an alternative compression method in Section 6.2.1 based on binary interval tree, which is employed by MRQED. On the other hand, it is possible to design an alternative solution for MRQED problem using redactable signature scheme [49], which we employ as our compression method.

Several works [12, 13, 14, 15, 16, 17] in verification of integrity of data stored in remote storage server also adopted some homomorphic and/or aggregatable verification tags to achieve efficient communication cost.

There are roughly four categories of approaches for outsourced database authentication in the literatures [20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34]. (1) (Homomorphic, or aggregatable) Cryptographic primitives, like collision-resistant

hash, digital signature, commitment [24, 35, 7]. (2) Merkle Hash Tree and variants [32, 28, 27]. (3) Computational geometry approach [21, 27, 30]. (4) Inserting and auditing fake tuples [29].

Thompson *et al.* [7] proposed a scheme called PDAS, with superlinear communication cost, more precisely $O(\min\{|S| \log N, N\})$ (S is the subset of tuples selected by the query condition), and linear storage overhead, to authenticate 1D aggregate (more precisely, SUM, COUNT) queries. PDAS is based on Shamir’s threshold secret sharing [36] and Pedersen commitment scheme [37], and protected the privacy of the aggregated attributes from the verifier. The authors briefly mentioned that their scheme could handle aggregate queries conditional on multidimensional range selections over insensitive attributes (stored in plaintext), but no details are provided. It seems that PDAS still requires superlinear communication cost for multidimensional aggregate query. It is worthy to point out that, the techniques in PDAS (like secret sharing, and how to compare two integers privately) can be integrated into MAIA to provide similar privacy protection, without much tradeoff of communication cost (it may increase by the factor of τ , where τ is the bit-length of the maximum attribute value in the dataset).

3. FORMULATION

In this section, we formalize the problem and security model. We introduce some key notations in Section 3.1, formally describe the dataset and query in Section 3.2, give a general security model in Section 3.3 for authentication of outsourced computing, and introduce our security assumptions in Section 3.4.

3.1 Notations

For the convenience of presentation, we introduce several operators \oplus , \otimes and \circledast on vectors as follows: Let $\mathbf{u} = (u_1, u_2, \dots, u_d)$ and $\mathbf{v} = (v_1, v_2, \dots, v_d)$ be two d -dimensional vectors, where $u_j, v_j \in \mathbb{N}$ for $1 \leq j \leq d$.

$$\oplus : \mathbf{u} \oplus \mathbf{v} = (u_1 + v_1, u_2 + v_2, \dots, u_d + v_d); \quad (1)$$

$$\otimes : \mathbf{u} \otimes \mathbf{v} = (u_1 v_1, u_2 v_2, \dots, u_d v_d); \quad (2)$$

$$\circledast : \mathbf{u} \circledast \mathbf{v} = (u_1^{v_1}, u_2^{v_2}, \dots, u_d^{v_d}). \quad (3)$$

We also introduce scalar multiplication $a\mathbf{u} = (au_1, au_2, \dots, au_d)$, where a is a scalar, and modulo operator $\mathbf{u} \bmod n$ on vectors:

$$\mathbf{u} \bmod n = (u_1 \bmod n, \dots, u_d \bmod n).$$

Additionally, we have vector $e = (\underbrace{1, 1, \dots, 1}_d)$.

Here in Table 2, we summarize the key notations in this paper.

3.2 Dataset and Query

We give a simplified (i.e. normalized [38]) data model, and the main construction in Section 5 will be based on this simplified data model. Later, in Section 6.2.2, we extend our scheme to general data model. We emphasize that this simplification does not lose generality, and serves two purposes: (1) reduce computation cost (precisely, reduce the size of *Help-Info*); (2) simplify the presentation.

Our data model consists of a set \mathbb{D} and a vector-valued attribute function Att , such that: (1) $\mathbb{D} = \{\mathbf{i}_1, \mathbf{i}_2, \dots, \mathbf{i}_N\} \subseteq [N]^d$, and for each $j \in [d]$, the set $\{\mathbf{i}_1[j], \mathbf{i}_2[j], \dots, \mathbf{i}_N[j]\} =$

Table 2: Summary of Key Notations

-
- \mathbb{N} : Set of all non-negative integers.
 - $x \stackrel{\$}{\leftarrow} S$: x is uniformly randomly chosen from a finite set S .
 - $[a, b]$: The set $\{a, a+1, a+2, \dots, b-1, b\}$, where $a \leq b$ are non-negative integers.
 - $[b]$: The set $\{1, 2, 3, \dots, b\}$, where b is a positive integer.
 - $|S|$: The size of the set (or multiset) S .
 - $\mathbf{i}[j]$: i_j , where $\mathbf{i} = (i_1, i_2, \dots, i_j, \dots)$ is a vector of dimension greater than or equal to j .
 - \mathbf{e} : The d -dimensional vector with each dimension equal to 1.
 - $\oplus, \otimes, \circledast$: Addition, multiplication and exponentiation operators over vectors.
 - **Prime**: Set of all odd primes.
 - $\text{Tag}(\mathbf{x}, \mathbf{i}; \xi)$: Randomized function Tag with input (\mathbf{x}, \mathbf{i}) and random coin ξ (Section 4).
 - **MAIA**: It stands for “Multidimensional Aggregate query Integrity Authentication” and it is our main scheme (Section 5).
 - **ProVer**: It stands for “Prover-Verifier” and it is an interactive algorithm (Section 5).
 - **CollRes**: It stands for “Collusion-Resistant” and it is a subroutine called by algorithm **ProVer** (Section 5).
-

$[N]$; (2) $\text{Att} : [N]^d \rightarrow \mathbb{X}^d$, where $\mathbb{X}^d \subset \mathbb{N}^d$ is the domain of attributes.

Let $\mathbf{R} = [a_1, b_1] \times [a_2, b_2] \times \dots \times [a_d, b_d] \subseteq [N]^d$ be a rectangular range. A range selection query with range \mathbf{R} over the dataset \mathbb{D} asks for the set of points \mathbf{i} ’s in $\mathbb{D} \cap \mathbf{R}$, and an aggregate sum query with range \mathbf{R} over the dataset \mathbb{D} , denoted as $\text{SUM}(\mathbf{R})$, asks for the sum $\bigoplus_{\mathbf{i} \in \mathbb{D} \cap \mathbf{R}} \text{Att}(\mathbf{i})$.

3.3 Security Model

We formalize the authentication problem described in Section 1.

Let us view the query on a database as the function $F : \mathbb{D} \times \mathbb{Q} \rightarrow \{0, 1\}^*$, where \mathbb{D} is the domain of databases, \mathbb{Q} is the domain of queries, and the output of F is a binary string. We define a remote computing protocol as follow:

DEFINITION 1 (RC). A \mathcal{RC} (Remote Computing) 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:

1. Given security parameter κ , Alice generates a key k : $k \leftarrow \text{KGen}(1^\kappa)$.
2. Alice encodes database $x \in \mathbb{D}$: $(p_x, s_x) \leftarrow \text{DEnc}(x, k)$, then sends p_x to Bob and keeps s_x .
3. Alice selects a query $r \in \mathbb{Q}$.
4. Algorithm $\text{Eval}(p_x)$ on Bob’s side, interacts with algorithm $\text{Ext}(s_x, r, k)$ on Alice’s side, to compute $y \leftarrow \langle \text{Eval}(p_x), \text{Ext}(s_x, r, k) \rangle$. If $y = \perp$, then Alice rejects. Otherwise, Alice believes that y is equal to $F(x, r)$.

In the setup phase, Alice executes Step (1) and (2). The query phase consists of multiple query sessions. In each query session, Alice and Bob execute Step (3) and (4).

We say a \mathcal{RC} protocol is provable, if the following conditions hold: (1) Alice accepts with o.h.p. (overwhelming high probability), when Bob follows the protocol honestly; (2) Alice rejects with o.h.p., when Bob returns a wrong result. Here we consider adversaries, i.e. malicious Bob, who are allowed to interact with Alice and learn for polynomial number of queries, before launching the attack. During the learning, the adversary may store whatever it has seen or learnt in a state variable.

DEFINITION 2 (*PRC*). A \mathcal{RC} protocol ($\text{KGen}, \text{DEnc}, \langle \text{Eval}, \text{Ext} \rangle$) w.r.t. function $F : \mathbb{D} \times \mathbb{Q} \rightarrow \{0, 1\}^*$, is called *PRC* (Provable Remote Computing) protocol, if the following two conditions hold: Let κ be the security parameter.

- *correctness*: for any $x \in \mathbb{D}$, $\mathbb{P}_{\text{Eval}}(x) \geq 1 - \text{negl}(\kappa)$ (asymptotically larger or equal);
- *soundness*: for any PPT adversary \mathcal{A} , for any $x \in \mathbb{D}$, $\mathbb{P}_{\mathcal{A}}(x) \geq 1 - \text{negl}(\kappa)$ (asymptotically larger or equal),

where \mathbb{P}_{Eval} and $\mathbb{P}_{\mathcal{A}}$ are defined as

$$\mathbb{P}_{\text{Eval}}(x) \stackrel{\text{def}}{=} \Pr \left[\begin{array}{l} k \leftarrow \text{KGen}(1^\kappa); \\ (p_x, s_x) \leftarrow \text{DEnc}(x, k); \\ r \stackrel{\$}{\leftarrow} \mathbb{Q}; \\ \zeta \leftarrow \langle \text{Eval}(p_x), \text{Ext}(s_x, r, k) \rangle; \\ \zeta = F(x, r) \end{array} \right],$$

$$\mathbb{P}_{\mathcal{A}}(x) \stackrel{\text{def}}{=} \Pr \left[\begin{array}{l} k \leftarrow \text{KGen}(1^\kappa); \\ (p_x, s_x) \leftarrow \text{DEnc}(x, k); \\ \text{for } 1 \leq i \leq \text{poly}(\kappa) \\ \quad r_i \stackrel{\$}{\leftarrow} \mathbb{Q}; \\ \quad \zeta_i \leftarrow \langle \mathcal{A}(p_x, \text{view}_{\mathcal{A}}^{\text{Ext}}), \text{Ext}(s_x, r_i, k) \rangle; \\ r \stackrel{\$}{\leftarrow} \mathbb{Q}; \\ \zeta \leftarrow \langle \mathcal{A}(p_x, \text{view}_{\mathcal{A}}^{\text{Ext}}), \text{Ext}(s_x, r, k) \rangle; \\ \zeta = \perp \vee \zeta = F(x, r) \end{array} \right].$$

The probability is taken over all random coins used by related algorithms, $\text{poly}(\cdot)$ is an arbitrary but fixed polynomial function, $\text{negl}(\cdot)$ is some negligible function, and $\text{view}_{\mathcal{A}}^{\text{Ext}}$ is a state variable¹ describing all random coins chosen by \mathcal{A} and all messages \mathcal{A} received from Ext during previous interactions.

We remark that the security model is similar to the formulation of \mathcal{POR} (Proof of Retrieability) [18], and it is not surprising that our scheme MAIA implies a (ρ, λ) -valid \mathcal{POR} system, with some parameters ρ and λ .

3.4 Assumptions

In this subsection, we state the assumptions that our proposed scheme relies on.

Let the group generator \mathcal{G} be a randomized algorithm, which takes as input a security parameter 1^κ and outputs a tuple $(n, p, q, r, G_p, G_q, G_r, g_p, g_q, g_r)$, where $n = (2p + 1)(2q + 1)(2r + 1)$ is a randomly chosen κ bits composite; all of $p, q, r, 2p + 1, 2q + 1$ and $2r + 1$ are distinct primes; p, q and r are of the same bit-length; G_p, G_q and G_r are three cyclic multiplicative subgroups of \mathbb{Z}_n^* , of order p, q and r respectively; and g_p, g_q and g_r are randomly chosen generators of G_p, G_q and G_r respectively.

¹The adversary \mathcal{A} may keep updating this state variable.

ASSUMPTION 1. Algorithm $\mathcal{G}(1^\kappa)$ is run to obtain $(n, p, q, r, G_p, G_q, G_r, g_p, g_q, g_r)$. Let group $G = G_p \times G_q \times G_r$. The following two distributions are computationally indistinguishable,

- $\mathcal{X} \stackrel{\text{def}}{=} \{\mathcal{X}_\kappa \stackrel{\$}{\leftarrow} G\}$, i.e. \mathcal{X}_κ is uniformly randomly distributed over G ;
- $\mathcal{Y} \stackrel{\text{def}}{=} \{\mathcal{Y}_\kappa \stackrel{\$}{\leftarrow} G_r\}$, i.e. \mathcal{Y}_κ is uniformly randomly distributed over G_r .

ASSUMPTION 2. Algorithm $\mathcal{G}(1^\kappa)$ is run to obtain $(n, p, q, r, G_p, G_q, G_r, g_p, g_q, g_r)$. Let group $G = G_p \times G_q \times G_r$. Define function $\mathcal{R} : G \rightarrow \mathbb{Z}_r$, such that for any $(x, y, z) \in \mathbb{Z}_p \times \mathbb{Z}_q \times \mathbb{Z}_r$, $\mathcal{R}(g_p^x g_q^y g_r^z) \stackrel{\text{def}}{=} z$. Let Δ be a set $\{g_p^x g_q^y g_r^z \in G : z \stackrel{\$}{\leftarrow} \mathbb{Z}_r \text{ and } (x, y) \text{ are chosen from } \mathbb{Z}_p \times \mathbb{Z}_q \text{ under any distributions}\}$. Given only (n, Δ) , it is hard to compute (A, B) , such that

$$A \not\equiv 1 \pmod{n}, B = \mathcal{R}(A) \pmod{r}.$$

Assumption 1 can be considered as a variant version of the p-Subgroup Assumption [39], and also a variant version of the Subgroup Membership Problem (\mathcal{SMP}) Assumption [9, 40, 41, 42, 43, 44, 45, 46], and Assumption 2 is closely related to the Projection Problem (\mathcal{PP}) Assumption [40, 42].

4. AUTHENTICATION TAG

In this section, we first construct an authentication tag function Tag (Section 4.1), which consists of three components. Next, we give an algorithm to decompose an output of the tag function into components (Section 4.2). At the end of this section, we show that the tag function is unforgeable (Section 4.3).

4.1 Construction of tag functions

We construct several tag functions dtag (stands for “tag of data”), itag (stands for “tag of index”²), and ctag (stands for “tag for completeness”). We combine these functions together to obtain tag function tag .

Let \mathcal{G} be as in Section 3.4.

DEFINITION 3. Let $(n, p, q, r, G_p, G_q, G_r, g_p, g_q, g_r)$ be an output of \mathcal{G} . Let $(s_1, s_2, s_3) \in \mathbb{Z}_p^* \times \mathbb{Z}_q^* \times \{0, 1\}^\kappa$. Let key $k = (s_1, s_2, s_3, n, p, q, r, g_p, g_q, g_r)$. Define functions dtag , itag , ctag and tag as follows:

$$\text{dtag}_k(x; \xi) = g_p^x g_q^{x+s_2\xi} \pmod{n},$$

$$\text{itag}_k(i) = g_p^{h(i)} \pmod{n},$$

$$\text{ctag}_k(i) = g_r^{f_{s_3}(i)} \pmod{n},$$

$$\text{tag}_k(x, i; \xi) = \text{dtag}_k(x; \xi) \times \text{itag}_k(i) \times \text{ctag}_k(i) \pmod{n},$$

where the random coin $\xi \stackrel{\$}{\leftarrow} \mathbb{Z}_q^*$, $h(\cdot) : \mathbb{N} \rightarrow \text{Prime}$ is a collision-resistant³ hash function, and $\{f_s\}_{s \in \{0, 1\}^\kappa}$ is PRF [48].

Note that itag is a variant of signature scheme [47, 49].

The data is in multidimensional vector form. For the convenience of presentation, we also define similar tag functions in vector forms.

²Here “index” means the normalized value i of data x . See the extended data model in Section 6.2.2.

³In other words, h should be *division intractable*. See Genaro *et al.* [47] for the definition of “division intractable”.

DEFINITION 4. Let $(n, p, q, r, G_p, G_q, G_r, g_p, g_q, g_r)$ be an output of \mathcal{G} . Let $(s_{\ell,1}, s_{\ell,2}, s_{\ell,3}) \in \mathbb{Z}_p^* \times \mathbb{Z}_q^* \times \{0, 1\}^k$ for each $\ell \in [d]$. Let $\mathbf{s}_1 = (s_{1,1}, s_{2,1}, \dots, s_{d,1})$, $\mathbf{s}_2 = (s_{1,2}, s_{2,2}, \dots, s_{d,2})$ and $\mathbf{s}_3 = (s_{1,3}, s_{2,3}, \dots, s_{d,3})$. Let $\mathcal{K} = (\mathbf{s}_1, \mathbf{s}_2, \mathbf{s}_3, n, p, q, r, g_p, g_q, g_r)$. Define vector-valued functions DTag , ITag , CTag and Tag as follows:

$$\begin{aligned} \text{DTag}_{\mathcal{K}}(\mathbf{x}; \xi) &= (\text{dtag}_{k_1}(\mathbf{x}[1]; \xi), \dots, \text{dtag}_{k_d}(\mathbf{x}[d]; \xi)), \\ \text{ITag}_{\mathcal{K}}(\mathbf{i}) &= (\text{itag}_{k_1}(\mathbf{i}[1]), \dots, \text{itag}_{k_d}(\mathbf{i}[d])), \\ \text{CTag}_{\mathcal{K}}(\mathbf{i}) &= (\text{ctag}_{k_1}(\mathbf{i}[1]), \dots, \text{ctag}_{k_d}(\mathbf{i}[d])), \\ \text{Tag}_{\mathcal{K}}(\mathbf{x}, \mathbf{i}; \xi) &= (\text{tag}_{k_1}(\mathbf{x}[1], \mathbf{i}[1]; \xi), \dots, \text{tag}_{k_d}(\mathbf{x}[d], \mathbf{i}[d]; \xi)), \end{aligned}$$

where for each $\ell \in [d]$, $k_\ell = (s_1[\ell], s_2[\ell], s_3[\ell], n, p, q, r, g_p, g_q, g_r)$, and the random coin $\xi \in \mathbb{Z}_q^*$.

In another equivalent expression,

$$\begin{aligned} \text{Tag}_{\mathcal{K}}(\mathbf{x}, \mathbf{i}; \xi) &= \text{DTag}_{\mathcal{K}}(\mathbf{x}; \xi) \otimes \text{ITag}_{\mathcal{K}}(\mathbf{i}) \otimes \text{CTag}_{\mathcal{K}}(\mathbf{i}) \\ &= ((g_p g_q e) \otimes \mathbf{x}) \otimes \text{ITag}_{\mathcal{K}}(\mathbf{i}) \otimes \text{DTag}_{\mathcal{K}}(\mathbf{0}; \xi) \otimes \text{CTag}_{\mathcal{K}}(\mathbf{i}). \end{aligned} \quad (4)$$

Note: (1) We call g_q^ξ binding factor among d dimensions of $\text{Tag}_{\mathcal{K}}(\mathbf{x}, \mathbf{i}; \xi)$. Given a candidate Tag value \mathbf{t} , we can extract the binding factor from $\mathbf{t}[j]$ for each dimension $j \in [d]$, and check whether they are equal. (2) We do not store the values of ξ 's when we use the tag function Tag in our scheme. Sometimes we may ignore ξ , and just write $\text{Tag}(\mathbf{x}, \mathbf{i})$.

4.2 Decompose a tag value

As mentioned in the introduction, Alice is able to extract an ITag (or CTag , DTag) value from a Tag value, with the private key. Such extraction or decomposition play an important role in the verification of our scheme (See Step A2 of CollRes in Section 5.3). We construct a pair of encoding and decoding algorithms $(\mathcal{E}, \mathcal{D})$ to reflect the idea of decomposition: (1) \mathcal{E} encodes outputs of DTag , ITag and CTag into an output of Tag ; (2) \mathcal{D} decodes an output of Tag into outputs of DTag , ITag and CTag . The rest of this subsection provides the details.

4.2.1 In Scalar Form

The output of function tag is an element from G . Recall that $n = (2p+1)(2q+1)(2r+1)$. With the factorization of n , i.e. p, q, r , one can extract out the G_p, G_q or G_r component from an element from G (or say, project a G element onto subgroup G_p, G_q and G_r [40]).

Let $\lambda_p = (qr)^{-1} \bmod p$, $\lambda_q = (pr)^{-1} \bmod q$, $\lambda_r = (pq)^{-1} \bmod r$. Let $t = g_p^x g_q^y g_r^z \bmod n$. We have

$$\begin{aligned} \Upsilon_p(t) &\stackrel{\text{def}}{=} (t^{qr})^{\lambda_p} \equiv g_p^x \pmod{n}; \\ \Upsilon_q(t) &\stackrel{\text{def}}{=} (t^{pr})^{\lambda_q} \equiv g_q^y \pmod{n}; \\ \Upsilon_r(t) &\stackrel{\text{def}}{=} (t^{pq})^{\lambda_r} \equiv g_r^z \pmod{n}. \end{aligned}$$

Based on these, one can also extract out dtag , itag , or ctag component from an output of function tag . Let $t = \text{tag}_k(\mathbf{x}, \mathbf{i}; \xi)$. Note $s_1, s_2, p, g, r, g_p, g_q$ and g_r are parts of the key k defined in Definition 3. We define

$$\begin{aligned} \gamma_1(t, g_p^x) &\stackrel{\text{def}}{=} \Upsilon_p(t) (g_p^x)^{-1} \equiv \text{itag}_k(\mathbf{i}) \pmod{n}; \\ \gamma_2(t, g_q^x) &\stackrel{\text{def}}{=} \Upsilon_q(t) (g_q^x)^{-1} \equiv \text{dtag}_k(\mathbf{0}; \xi) \pmod{n}; \\ \gamma_3(t) &\stackrel{\text{def}}{=} \Upsilon_r(t) \equiv \text{ctag}_k(\mathbf{i}) \pmod{n}. \end{aligned}$$

4.2.2 In Vector Form

Here we present the same decomposition in vector form. We may consider Tag as a homomorphic encoding algorithm \mathcal{E} (See equation (4)):

$$\begin{aligned} \text{Tag}_{\mathcal{K}}(\mathbf{x}, \mathbf{i}; \xi) &= \mathcal{E}_{\mathcal{K}}(\text{ITag}_{\mathcal{K}}(\mathbf{i}), \text{DTag}_{\mathcal{K}}(\mathbf{0}; \xi), \text{CTag}_{\mathcal{K}}(\mathbf{i}), \mathbf{x}); \\ \mathcal{E}_{\mathcal{K}}(\mathbf{u}_1, \mathbf{v}_1, \mathbf{w}_1, \mathbf{x}_1) \otimes \mathcal{E}_{\mathcal{K}}(\mathbf{u}_2, \mathbf{v}_2, \mathbf{w}_2, \mathbf{x}_2) &= \\ \mathcal{E}_{\mathcal{K}}(\mathbf{u}_1 \otimes \mathbf{u}_2, \mathbf{v}_1 \otimes \mathbf{v}_2, \mathbf{w}_1 \otimes \mathbf{w}_2, \mathbf{x}_1 \oplus \mathbf{x}_2) & \quad (5) \end{aligned}$$

The corresponding decoding algorithm \mathcal{D} is: Let $\mathbf{u}, \mathbf{v}, \mathbf{w}$ are the outputs of ITag , DTag , CTag , respectively, and $\boldsymbol{\psi} = \mathcal{E}_{\mathcal{K}}(\mathbf{u}, \mathbf{v}, \mathbf{w}, \mathbf{x})$.

$$\mathcal{D}_{\mathcal{K}}(\boldsymbol{\psi}; \mathbf{x}) \stackrel{\text{def}}{=} (\mathbf{\Gamma}_1(\boldsymbol{\psi}, \mathbf{x}), \mathbf{\Gamma}_2(\boldsymbol{\psi}, \mathbf{x}), \mathbf{\Gamma}_3(\boldsymbol{\psi})) = (\mathbf{u}, \mathbf{v}, \mathbf{w}) \quad (6)$$

where $\mathbf{\Gamma}_1, \mathbf{\Gamma}_2$ and $\mathbf{\Gamma}_3$ are defined below,

$$\begin{aligned} \mathbf{\Gamma}_1(\mathbf{t}, \mathbf{u}) &\stackrel{\text{def}}{=} (\gamma_1(\mathbf{t}[1], g_p^{\mathbf{u}[1]}), \dots, \gamma_1(\mathbf{t}[d], g_p^{\mathbf{u}[d]})); \\ \mathbf{\Gamma}_2(\mathbf{t}, \mathbf{u}) &\stackrel{\text{def}}{=} (\gamma_2(\mathbf{t}[1], g_q^{\mathbf{u}[1]}), \dots, \gamma_2(\mathbf{t}[d], g_q^{\mathbf{u}[d]})); \\ \mathbf{\Gamma}_3(\mathbf{t}) &\stackrel{\text{def}}{=} (\gamma_3(\mathbf{t}[1]), \dots, \gamma_3(\mathbf{t}[d])). \end{aligned}$$

4.3 Security

LEMMA 1. Let k be a key as in Definition 3. Given $S = \{(x_i, i, t_i) : \xi \stackrel{\$}{\leftarrow} \mathbb{Z}_q^*, t_i \leftarrow \text{tag}_k(x_i, i; \xi)\}_{1 \leq i \leq N}$, it is computationally hard to forge a tuple $(y, j, t) \notin S$, such that $\gamma_1(t, g_p^y) = \text{itag}_k(j)$ and $\gamma_3(t) = \text{ctag}_k(j)$, under Assumption 1 and Assumption 2.

The proof is given in Appendix B.1.

5. MAIA: MAIN CONSTRUCTION

Our scheme MAIA (Multidimensional Aggregate query Integrity Authentication) is a \mathcal{RC} protocol for aggregate range query, and consists of algorithms $(\text{KGen}, \text{DEnc}, \langle \text{Eval}, \text{Ext} \rangle)$. The protocol consists of setup phase (Section 5.2) and query phase (Section 5.3). The setup phase is associated with algorithms KGen and DEnc , and the query phase is associated with interactive algorithm $\langle \text{Eval}, \text{Ext} \rangle$, which calls subroutine CollRes .

Our scheme achieves efficient communication complexity, due to algorithms Compress and Uncompress . We defer the constructions of Compress and Uncompress to Section 5.4.

Basically, the authentication consists of two parts. First, Alice has a mechanism, i.e. the algorithm CollRes , to ensure that Bob produces the result and proof *solely* using data points within the query range \mathbf{R} . Next, to verify completeness, Alice also runs CollRes on aggregate query with the complement range \mathbf{R}^c , and checks whether the two query results and proofs “sum” to the whole dataset. Since the range \mathbf{R}^c is not in the form that CollRes can handle, Alice has to divide it into $O(d)$ rectangular ranges, and runs CollRes on each of them.

First of all, let us review the authentication tag function Tag , which serves as the basis of our scheme.

5.1 Summary of tag function

Recall that, in section 4, we construct a key-ed tag function Tag , which consists of functions DTag , ITag and CTag as components. We also wrap Tag as a homomorphic encoding

algorithm \mathcal{E} : Let \mathcal{K} be the key as in Definition 4.

$$\begin{aligned} \text{Tag}_{\mathcal{K}}(\mathbf{x}, i; \xi) &= \mathcal{E}_{\mathcal{K}}(\text{ITag}_{\mathcal{K}}(i), \text{DTag}_{\mathcal{K}}(\mathbf{0}; \xi), \text{CTag}_{\mathcal{K}}(i), \mathbf{x}); \\ \mathcal{E}_{\mathcal{K}}(\mathbf{u}_1, \mathbf{v}_1, \mathbf{w}_1, \mathbf{x}_1) \otimes \mathcal{E}_{\mathcal{K}}(\mathbf{u}_2, \mathbf{v}_2, \mathbf{w}_2, \mathbf{x}_2) &= \\ \mathcal{E}_{\mathcal{K}}(\mathbf{u}_1 \otimes \mathbf{u}_2, \mathbf{v}_1 \otimes \mathbf{v}_2, \mathbf{w}_1 \otimes \mathbf{w}_2, \mathbf{x}_1 \oplus \mathbf{x}_2). \end{aligned}$$

The corresponding decoding algorithm \mathcal{D} is: Let $\mathbf{u}, \mathbf{v}, \mathbf{w}$ are the outputs of $\text{ITag}, \text{DTag}, \text{CTag}$, respectively, and $\psi = \mathcal{E}_{\mathcal{K}}(\mathbf{u}, \mathbf{v}, \mathbf{w}, \mathbf{x})$.

$$\mathcal{D}_{\mathcal{K}}(\psi, \mathbf{x}) = (\mathbf{u}, \mathbf{v}, \mathbf{w}).$$

5.2 Setup phase

Alice has a dataset \mathbf{D} associated with attribute function Att (See data model in Section 3.2). In the setup phase, Alice generates a private key \mathcal{K} by running key generating function KGen , and preprocesses the dataset \mathbf{D} using data encoding algorithm DEnc with key \mathcal{K} , in order to produce an authentication tag \mathbf{T} . At the end of setup phase, Alice removes \mathbf{D} and \mathbf{T} from her storage after sending them to Bob.

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

Alice

1. runs algorithm $\mathcal{G}(1^\kappa)$ to generate $(n, p, q, r, G_p, G_q, G_r, g_p, g_q, g_r)$, where \mathcal{G} is defined in Section 3.4;
2. chooses $(s_{\ell,1}, s_{\ell,2}, s_{\ell,3})$ at random from $\mathbb{Z}_p^* \times \mathbb{Z}_q^* \times \{0, 1\}^\kappa$ for each $\ell \in [d]$.

Let $\mathbf{s}_1 = (s_{1,1}, s_{2,1}, \dots, s_{d,1})$, $\mathbf{s}_2 = (s_{1,2}, s_{2,2}, \dots, s_{d,2})$ and $\mathbf{s}_3 = (s_{1,3}, s_{2,3}, \dots, s_{d,3})$. Let $\mathcal{K} = (\mathbf{s}_1, \mathbf{s}_2, \mathbf{s}_3, n, p, q, r, g_p, g_q, g_r)$. Output \mathcal{K} .

$\text{DEnc}(\mathbf{D}, \mathcal{K})$

For each data point $i \in \mathbf{D}$, Alice

1. chooses a number ξ_i at random from \mathbb{Z}_q^* ;
2. computes a tag t_i using the private key \mathcal{K} :

$$t_i \leftarrow \text{Tag}_{\mathcal{K}}(\text{Att}(i), i; \xi_i). \quad (7)$$

Next, Alice sends dataset \mathbf{D} , tag values $\mathbf{T} = \{t_i : i \in \mathbf{D}\}$ and modulus n to Bob, deletes the local copy of \mathbf{D} and \mathbf{T} from her storage, and keeps (\mathcal{K}, n, N) , where $N = |\mathbf{D}|$ is the size of the dataset.

5.3 Query phase

During the query phase, Alice has private key \mathcal{K} ; Bob has dataset \mathbf{D} and tag values \mathbf{T} . Alice and Bob share the modulus n and the size $N = |\mathbf{D}|$ of the dataset \mathbf{D} . The query phase consists of multiple query sessions. In a query session, Alice has a query $\text{SUM}(\mathbf{R})$, where $\mathbf{R} = [a_1, b_1] \times [a_2, b_2] \times \dots \times [a_d, b_d] \subseteq [N]^d$. The result to the query w.r.t. dataset \mathbf{D} is $\bigoplus_{i \in \mathbf{D} \cap \mathbf{R}} \text{Att}(i)$.

To process this query, Alice and Bob run an interactive algorithm ProVer (stands for ‘‘Prover-Verifier’’), i.e. Alice running Ext interacts with Bob running Eval , where $\text{ProVer} = (\text{Eval}, \text{Ext})$. The algorithm ProVer calls algorithm CollRes (stands for ‘‘Collusion-Resistant’’) as subroutine. Note that CollRes is a strengthened version of the plausible scheme presented in Section 1.

$\text{CollRes}(\text{SUM}(\mathbf{R}))$

Precondition. Range $\mathbf{R} = \mathbf{A}_1 \times \mathbf{A}_2 \times \dots \times \mathbf{A}_d$, where $\mathbf{A}_j \subseteq [N], j \in [d]$ (This is required by the subroutine Compress).

A1: (Alice’s first step) Alice

- (a) chooses $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_d)$ at random from $(\mathbb{Z}_p^*)^d$;
- (b) computes δ using $sk = (n, g_p, \phi(n), \mathbf{s}_1)$:

$$\delta \leftarrow \text{Compress}_{sk}(\mathbf{R}, \alpha). \quad (8)$$

Next, Alice sends \mathbf{R} and δ to Bob.

Note: (1) From δ , Bob is able to recover $\varpi = \{\mathcal{E}(\text{ITag}(i) \otimes \alpha, \mathbf{v}_i, \mathbf{w}_i, \mathbf{0}) : i \in \mathbf{R}\}$ for some $\mathbf{v}_i \in (G_q)^d$ and $\mathbf{w}_i \in (G_r)^d$. We call ϖ *Help-Info*, since it will help Bob to produce a proof. Without considering randomness \mathbf{v}_i ’s and \mathbf{w}_i ’s, one may view ϖ as $\{\text{ITag}(i) \otimes \alpha : i \in \mathbf{R}\}$. See Section 5.4 for details of function Compress . (2) The secret random nonce α prevents Bob from abusing this information for other queries, i.e. prevent collusion attack. (3) \mathbf{s}_1 is part of key \mathcal{K} .

B1: (Bob’s first step) Bob

- (a) computes the sum \mathbf{X} :

$$\mathbf{X} = \bigoplus_{i \in \mathbf{D} \cap \mathbf{R}} \text{Att}(i); \quad (9)$$

- (b) aggregates all Tag values for $i \in \mathbf{D} \cap \mathbf{R}$:

$$\Psi_1 \leftarrow \bigotimes_{i \in \mathbf{D} \cap \mathbf{R}} t_i \quad \text{mod } n; \quad (10)$$

- (c) uncompresses δ with $pk = n$: for each $i \in \mathbf{D} \cap \mathbf{R}$,

$$\mathbf{u}_i \leftarrow \text{Uncompress}_{pk}(\mathbf{R}, \delta, i); \quad (11)$$

- (d) aggregates all \mathbf{u}_i ’s for $i \in \mathbf{D} \cap \mathbf{R}$:

$$\Psi_2 \leftarrow \bigotimes_{i \in \mathbf{D} \cap \mathbf{R}} \mathbf{u}_i \quad \text{mod } n. \quad (12)$$

Bob sends back to Alice \mathbf{X} as query result, and (Ψ_1, Ψ_2) as the proof.

Note: In Step B1(c), $\mathbf{u}_i = \mathcal{E}(\text{ITag}(i) \otimes \alpha, \mathbf{v}_i, \mathbf{w}_i, \mathbf{0})$ for some $\mathbf{v}_i \in (G_q)^d$ and $\mathbf{w}_i \in (G_r)^d$.

A2: (Alice’s second step) Alice verifies the query result in the following way.

- (a) Decode Ψ_1 using $(\mathcal{K}, \mathbf{X})$, and decode Ψ_2 using $(\mathcal{K}, \mathbf{0})$:

$$(\mathbf{U}_1, \mathbf{V}, \mathbf{W}) \leftarrow \mathcal{D}_{\mathcal{K}}(\Psi_1, \mathbf{X})$$

$$(\mathbf{U}_2, \hat{\mathbf{V}}, \hat{\mathbf{W}}) \leftarrow \mathcal{D}_{\mathcal{K}}(\Psi_2, \mathbf{0})$$

- (b) Check whether the ITag component of Ψ_1 is consistent with Ψ_2 using the secret random nonce α :

$$\mathbf{U}_1 \otimes \alpha \stackrel{?}{=} \mathbf{U}_2 \quad \text{mod } n; \quad (13)$$

- (c) Check whether the binding factor $g_q^{\xi_i}$ ’s among different dimensions are consistent: Is there some $a \in \mathbb{Z}_n$, such that

$$\mathbf{V} \otimes \mathbf{s}_2^{-1} = a\mathbf{e} \quad \text{mod } n? \quad (14)$$

Note: (1) \mathbf{e} is the d -dimensional vector with each dimension equal to 1. (2) \mathbf{s}_2 is part of key \mathcal{K} and $\mathbf{s}_2^{-1} \otimes \mathbf{s}_2 = \mathbf{e} \quad \text{mod } q$.

(d) Alice extracts the **CTag** component from Ψ_1 :

$$\pi \leftarrow \mathbf{W}[1].$$

Note: This π is required by algorithm **ProVer** to prevent double counting or undercounting.

If all verifications succeed, Alice outputs (\mathbf{X}, π) . Otherwise, Alice output \perp .

Remark.

1. If Bob follows the protocol honestly, we have

$$\begin{aligned} \Psi_1 &= \mathcal{E} \left(\bigotimes_{i \in \mathbf{D} \cap \mathbf{R}} \text{ITag}(i), \bigotimes_{i \in \mathbf{D} \cap \mathbf{R}} \text{DTag}(\mathbf{0}; \xi_i), \right. \\ &\quad \left. \bigotimes_{i \in \mathbf{D} \cap \mathbf{R}} \text{CTag}(i), \bigoplus_{i \in \mathbf{D} \cap \mathbf{R}} \text{Att}(i)e \right); \end{aligned}$$

$$U_1 = \bigotimes_{i \in \mathbf{D} \cap \mathbf{R}} \text{ITag}(i) \pmod n;$$

$$U_2 = \left(\bigotimes_{i \in \mathbf{D} \cap \mathbf{R}} \text{ITag}(i) \right) \otimes \alpha \pmod n;$$

$$V = \bigotimes_{i \in \mathbf{D} \cap \mathbf{R}} \text{DTag}(\mathbf{0}; \xi_i) = \left(\left(\prod_{i \in \mathbf{D} \cap \mathbf{R}} g_q^{\xi_i} \right) e \right) \otimes s_2 \pmod n.$$

2. The equality test in equation (13) prevents Bob from bringing in points that do not satisfy the query condition into the computation of query result and its proof. For each point $i \in \mathbf{D} \cap \mathbf{R}$, Bob has both t_i (provided by Alice during setup) and $\text{ITag}(i) \otimes \alpha$ (provided by Alice during the current query session). For any point $i \notin \mathbf{D} \cap \mathbf{R}$, it is infeasible for Bob to forge t_i or $\text{ITag}(i) \otimes \alpha$, without the private key \mathcal{K} and secret random nonce α .
3. The equality test in equation (14) prevents Bob from forging a new point by mixing and combining different dimensions of different points. For example, without this verification, Bob may fool Alice with a faked point $(i_2[1], i_1[2], i_1[3], \dots, i_1[d])$ with tag $(t_2[1], t_1[2], t_1[3], \dots, t_1[d])$, where points $i_1, i_2 \in \mathbf{D}$ and t_1, t_2 are their corresponding **Tag** values. Furthermore, s_2 prevents permutation on dimensions of the same data point.
4. As mentioned in the introduction, the above proof system can not ensure the completeness, due to the homomorphism of \mathcal{E} . Bob could cheat by double counting or undercounting some points in $\mathbf{D} \cap \mathbf{R}$. For example, Bob may claim the result is $2 \sum_{i \in \mathbf{D} \cap \mathbf{R}} \text{Att}(i)$ by sending $(\Psi_1 \otimes \Psi_1, \Psi_2 \otimes \Psi_2)$ as proof. Due to homomorphism of \mathcal{E} , this forged proof can pass Alice's verification in equations (13) and (14). Such attacks can be detected by algorithm **ProVer**.
5. The specially designed tag function **Tag** has the following properties:

- It allows Bob to do aggregation, since \mathcal{E} is homomorphic w.r.t. its inputs $\text{ITag}(i), \text{DTag}(\mathbf{0}, \xi), \text{CTag}(i)$ and \mathbf{x} (See equation (5)).
- It is unforgeable (See Lemma 1), and $\text{Tag}(\mathbf{x}, i; \xi)$ itself is not homomorphic w.r.t. its inputs.
- It binds the point i and its attribute value \mathbf{x} together in a secure way.

⁴Without considering the randomization, one may view *Help-Info* ϖ as $\{\text{ITag}(i) \otimes \alpha : i \in \mathbf{R}\}$. See Section 5.4.

(d) It allows compression of *Help-Info*. More precisely, **ITag** allows compression.

To execute an aggregate query with range $\mathbf{R} = [a_1, b_1] \times [a_2, b_2] \times \dots \times [a_d, b_d] \subseteq [N]^d$, Alice runs algorithm **Ext** and interacts with Bob who runs algorithm **Eval**. To verify the completeness, Alice will divide the complement range \mathbf{R}^c carefully into many subregions, and invokes **CollRes** on each subregion. The details is given below.

$$\text{ProVer}(\text{SUM}(\mathbf{R})) = \langle \text{Eval}(\mathbf{D}, \mathbf{T}, n), \text{Ext}(N, \text{SUM}(\mathbf{R}), \mathcal{K}) \rangle$$

Precondition: Alice knows the value⁵ $\pi^* = \prod_{i \in \mathbf{D}} \text{ctag}(i[1])$. For each $\ell \in [d]$, let $\mathbf{R}_\ell = [a_1, b_1] \times [a_2, b_2] \times \dots \times [a_\ell, b_\ell] \times [N]^{d-\ell} \subseteq [N]^d$ be an ℓ -dimensional range. Let \mathbf{R}_0 be the universal set $[N]^d$.

- Alice interacts with Bob to simulate **CollRes** on $(d+1)$ queries $\text{SUM}(\mathbf{R}_d)$, and $\text{SUM}(\mathbf{R}_{\ell-1} \setminus \mathbf{R}_\ell)$, $\ell \in [d]$.
- Bob interacts with Alice to simulate **CollRes** on $(d+1)$ queries $\text{SUM}(\mathbf{R}_d)$, and $\text{SUM}(\mathbf{R}_{\ell-1} \setminus \mathbf{R}_\ell)$, $\ell \in [d]$.
- Alice obtains output $(\mathbf{X}_{d+1}, \pi_{d+1})$ of **CollRes** for query $\text{SUM}(\mathbf{R}_d)$, and $(\mathbf{X}_\ell, \pi_\ell)$ for query $\text{SUM}(\mathbf{R}_{\ell-1} \setminus \mathbf{R}_\ell)$, $\ell \in [d]$. Alice accepts Bob's reply and output \mathbf{X} , if

$$\prod_{1 \leq \ell \leq d+1} \pi_\ell \stackrel{?}{=} \pi^* \pmod n. \quad (15)$$

Otherwise, Alice outputs \perp .

Note: The $(d+1)$ ranges $\mathbf{R}_{\ell-1} \setminus \mathbf{R}_\ell, \ell \in [d]$, and \mathbf{R}_d forms a partition of the universal set $[N]^d$.

Remark. Simply, Alice may precompute and store π^* . We will explain more in Section 6.1.1.

Such division of the complement range \mathbf{R}^c is essential, since our compression algorithm **Compress** only works on range of form $\mathbf{A}_1 \times \mathbf{A}_2 \times \dots \times \mathbf{A}_d$, where $\mathbf{A}_j \subseteq [N], j \in [d]$. Note that both \mathbf{R}_ℓ and $\mathbf{R}_{\ell-1} \setminus \mathbf{R}_\ell$ are of this form.

The main reasons that our scheme achieves efficient communication complexity are: (1) Alice can compress *Help-Info* sent from Alice to Bob, using an algorithm **Compress**, and Bob can uncompress and recover *Help-Info*, using an algorithm **Uncompress**; (2) The specially designed authentication tag functions **Tag** allows aggregation, so Bob just sends back some aggregated values as proof.

5.4 Compress and Uncompress

Now we construct the algorithms **Compress** and **Uncompress**. Johnson *et al.* [49] (in Section 5 "Set Homomorphic Signatures" of that paper) proposed a redactable signature scheme, based on the signature scheme in Gennaro *et al.* [47]. We wrap this redactable signature scheme as two algorithms **compress** and **uncompress**: Let (n, g_p, \mathbf{s}_1) be as in Definition 4 and hash function $h(\cdot)$ be as in Definition 3, Let $pk \leftarrow n; sk \leftarrow (n, g_p, \phi(n), \mathbf{s}_1)$, for a set \mathcal{I} of integers,

- compress** _{s_k} (\mathcal{I}): Output $\delta \leftarrow g_p^{\prod_{i \in \mathcal{I}} h(i)^{-1}} \pmod n$.
- uncompress** _{pk} ($\mathcal{I}, \delta, \hat{\mathcal{I}}$): If $\hat{\mathcal{I}} \not\subseteq \mathcal{I}$, output \perp . Otherwise, output $\hat{\delta}: \hat{\delta} \leftarrow \delta^{\prod_{i \in \mathcal{I} \setminus \hat{\mathcal{I}}} h(i)}$ $\pmod n$.

⁵The dimension 1 as in expression $i[1]$ is an arbitrary choice in $[d]$, and this choice will influence how to divide \mathbf{R}^c .

Then we define the compress/uncompress algorithms in vector form. Let $\mathbf{R} \subset [N]^d$. Let $\text{Proj}_j(\mathbf{R})$ denote the projection of \mathbf{R} on j -th dimension, i.e. $\text{Proj}_j(\mathbf{R}) = \{i[j] : i \in \mathbf{R}\}$, $j \in [d]$.

- **Compress_{sk}(\mathbf{R}, α)**: for each $1 \leq j \leq d$, let $\delta_j \leftarrow \text{compress}_{sk}(\text{Proj}_j(\mathbf{R}))$. Let $\mathbf{v} \leftarrow (\delta_1, \delta_2, \dots, \delta_d) \otimes \mathbf{s}_1$, $\mathbf{u}_1 \xleftarrow{\$} (G_q)^d$ and $\mathbf{u}_2 \xleftarrow{\$} (G_r)^d$. $\delta \leftarrow (\mathbf{v} \otimes \alpha) \otimes \mathbf{u}_1 \otimes \mathbf{u}_2$. Output δ .
- **Uncompress_{pk}(\mathbf{R}, δ, i)**: for each $1 \leq j \leq d$, let set $\hat{\mathcal{I}}_j = \{i[j]\}$ and let $\hat{\delta}_j \leftarrow \text{uncompress}_{pk}(\text{Proj}_j(\mathbf{R}), \delta[j], \hat{\mathcal{I}}_j)$. Output $\hat{\delta} \leftarrow (\hat{\delta}_1, \hat{\delta}_2, \dots, \hat{\delta}_d)$.

Remark.

1. The operation $\mathbf{v} \otimes \alpha$ aims to prevent reuse of δ , and permutation on dimensions of δ (So does \mathbf{s}_1).
2. The operation $(\mathbf{v} \otimes \alpha) \otimes \mathbf{u}_1 \otimes \mathbf{u}_2$ aims to randomize $\mathbf{v} \otimes \alpha$ with a random vector \mathbf{u}_1 from $(G_q)^d$ and a random vector \mathbf{u}_2 from $(G_r)^d$.
3. The output $\hat{\delta}$ of **Uncompress_{pk}(\mathbf{R}, δ, i)** equals to $\mathcal{E}(\text{ITag}(i) \otimes \alpha, \mathbf{v}, \mathbf{w}, \mathbf{0})$ for some $\mathbf{v} \in (G_q)^d$ and $\mathbf{w} \in (G_r)^d$. From $\hat{\delta}$, Alice can extract $\text{ITag}(i) \otimes \alpha$ with private key.
4. Without considering those randomization factors, one may view the output δ of **Compress_{sk}(\mathbf{R}, α)** as the compression of $\{\text{ITag}(i) \otimes \alpha : i \in \mathbf{R}\}$.
5. The range \mathbf{R} has to be of form $\mathbf{A}_1 \times \mathbf{A}_2 \times \dots \times \mathbf{A}_d$, where $\mathbf{A}_j \subseteq [N], j \in [d]$. Here \mathbf{A}_j is not necessarily continuous.

We also provide an alternative constructions of **compress** and **uncompress** in Section 6.2.1.

5.5 Security

Our scheme is secure, under Definition 2. The proof is in Appendix C.

THEOREM 2. MAIA is a PRC (Provable Remote Computing) protocol, w.r.t. d -dimensional aggregate SUM query, under Assumption 1 and Assumption 2.

We defer the complexity analysis to the next section.

6. PERFORMANCE AND EXTENSIONS

In Section 6.1, we discuss some speedup methods and analyze the complexity of our scheme MAIA. In Section 6.2, we propose an alternative compression method, generalize the dataset and query, and extend our scheme to authenticate COUNT, MIN queries, and range selection query.

6.1 Performance

6.1.1 1D Aggregate Range Query

Suppose the range \mathbf{R} is 1D, i.e. $\mathbf{R} = [a_1, b_1] \times [N]^{d-1}$. Due to normalization, the first dimension of points within 1D range \mathbf{R} is predictable and continuous: $\{i[1] : i \in \mathbf{D} \cap \mathbf{R}\} = [a_1, b_1]$. Hence, when running **ProVer** on query **SUM**(\mathbf{R}), Alice can set $\pi^* = \prod_{i \in [a_1, b_1]} \text{ctag}(i)$, which she can compute by herself, and invoke only one instance of **CollRes** on **SUM**(\mathbf{R}). Furthermore, Alice does not need to provide *Help-Info* to Bob in **CollRes**, since Alice can generate the right hand side of equation (13) by herself. As a result, **ProVer** is much more efficient when the query range is 1D, compared with higher dimensional case.

6.1.2 ℓ -Dimensional Aggregate Range Query

The algorithms **ProVer** and **CollRes** can also handle query like $\sum_{i \in \mathbf{D} \cap \mathbf{R}_\ell} \mathbf{x}[1]$, where $\mathbf{R}_\ell = [a_1, b_1] \times [a_2, b_2] \times \dots \times [a_\ell, b_\ell] \times [N]^{d-\ell}$ is an ℓ -dimensional range and $\ell \in [d]$. For such queries, we can do some optimizations: (1) Treat all vectors $\text{Att}(i)$, $\mathbf{t}_i, \delta, \mathbf{X}, \Psi_1$ and Ψ_2 etc. as ℓ -dimensional vectors, and ignore the remaining $(d-\ell)$ dimensions. So the communication cost of **CollRes** will be $O(\ell)$ instead of $O(d)$. (2) To ensure completeness, check the relative complement $\mathbf{R}_1 \setminus \mathbf{R}_d$, instead of the absolute complement. So we save one function call **CollRes**($\mathbf{R}_0 \setminus \mathbf{R}_1$). (3) Run all ℓ calls to **CollRes** simultaneously in parallel.

The parallel execution of ℓ instances of **CollRes** can save us a lot: (1) Alice only sends values of $a_1, b_1, \dots, a_\ell, b_\ell$ to represent all of ℓ ranges \mathbf{R}_ℓ and $\mathbf{R}_{j-1} \setminus \mathbf{R}_j$, $j \in [2, \ell]$. (2) Recall that, on input (\mathbf{R}, α) , **Compress** will separate each dimension of \mathbf{R} , and generate a *Help-Info* for each 1D range $[a_j, b_j]$. For the above ℓ ranges, Alice can only generate compressed *Help-Info* for 1D range $[a_j, b_j]$ and $[N] \setminus [a_j, b_j]$, $j \in [\ell]$, from which Bob is able to reconstruct *Help-Info* for each of ℓ ranges.

Note that the query only asks for the sum of the first dimension. For dimension $j \in [2, \ell]$, Bob may send back $g^{\mathbf{X}[j]} \bmod n$ instead of the sum $\mathbf{X}[j]$ for the j -th dimension, when $\mathbf{X}[j]$ is too large, so that the communication cost is always independent of the size of dataset. Here g is some generator of group G provided by Alice.

6.1.3 Other speedup techniques

There are some other tradeoff or speedup techniques, including the following. (1) Batch executions of multiple queries: Alice may find the smallest union set \mathbf{U} of all query ranges, and issue a new query with range \mathbf{U} at the very beginning. Next, Alice authenticates other queries by checking complement w.r.t \mathbf{U} . (2) Precomputation: Partition the space into regions, and store necessary information for proof of completeness inside each region, so that Alice only need check the complement w.r.t. some union of regions that fully cover the query range. (3) For any rectangle range \mathbf{R} , Alice can find the smallest rectangle range $\bar{\mathbf{R}}$ with help of Bob, such that $\bar{\mathbf{R}} \subseteq \mathbf{R}$ and $\mathbf{D} \cap \mathbf{R} = \mathbf{D} \cap \bar{\mathbf{R}}$. Then Alice can generate *Help-Info* for $\bar{\mathbf{R}}$, instead of \mathbf{R} . In many scenarios, $\bar{\mathbf{R}}$ is much smaller than \mathbf{R} . For example, in Figure 2, such smallest rectangle range for $\mathbf{R}_1 = [3, 8] \times [2, 7]$ is $[4, 6] \times [4, 7]$. (4) Offline verifications with multiple queriers: Charlie sends query Q to Bob, and Bob returns (\mathbf{X}, σ) to Charlie, where \mathbf{X} is the result to query Q and σ is the signature of (Q, \mathbf{X}) under Bob's private key. Charlie may decide when (or whether) to ask Alice to do the verification.

6.1.4 Complexity Analysis

Taking the modifications in Section 6.1.1 and Section 6.1.2 into account, we analyze the complexity of our scheme. Let us consider query with range $\mathbf{R} \subseteq [N]^d$, which asks for the sum $\sum_{i \in \mathbf{D} \cap \mathbf{R}} \text{Att}(i)[1]$, where \mathbf{R} is an ℓ -dimensional rectangular range, $\ell \in [d]$. To process this query, **ProVer** requires $\sum_{j \in [d]} O(j) = O(\ell^2)$ communication bits and $O(1)$ rounds.

When $\ell = 1$, the computation complexity per query is very efficient: $O(|S|)$ on both sides. Here S denotes the set of points satisfying the query condition. For $\ell \in [2, d]$, the computation complexity per query on Bob's side is $O(\ell N^2)$ (modular exponentiations). The dominant computation step

is **Uncompress** in Step B1(c) of **CollRes** in Section 5.3. The computation complexity per query on Alice’s side is $O(\ell N)$ (modular multiplications). The dominant computation step is **Compress** in Step A1(b) of **CollRes** in Section 5.3. We can reduce the computation cost on Alice’s side to $O(\ell \log N)$ (modular multiplications) at the cost of additional $O(N)$ storage on Bob’s side without sacrificing communication complexity, and reduce the computation cost on Bob’s side to $O(\ell N \log N)$ (modular exponentiations; the constant factor behind big- O notation is about 1 and the base of log is 2), with $O(N)$ temporary storage. We save the details.

The storage overhead on Bob’s side, is $O(dN)$. The storage overhead on Alice’s side, i.e. the key size, is $O(d)$.

The summary of complexities of MAIA compared with related schemes is given in Figure 1 in Section 1. Note that when $d > \log(dN)/\log \log N$, we have $dN \log N \leq \log^d N$.

6.2 Extensions

6.2.1 Alternative Compression and Decompression

Here we present an alternative constructions of **ITag**, **Compress** and **Uncompress**, using a binary interval tree [1].

Assume $N = 2^H$ for some integer H . Build a complete binary tree with N leaves, such that: (1) Each tree node \mathbf{t} is associated with an interval, denoted as $\mathbf{t.int}$. For the j -th leaf \mathbf{t} (counting from the left to right), $\mathbf{t.int} = \{j\}$, and for an internal node \mathbf{t} with left child \mathbf{t}_l and right child \mathbf{t}_r , $\mathbf{t.int} = \mathbf{t}_l.int \cup \mathbf{t}_r.int$. (2) Each tree edge \mathbf{e} is associated with a public random prime, denoted as $\mathbf{e.p}$. (3) Each tree node \mathbf{t} is associated with an attribute, denoted as $\mathbf{t.val}$. Let $(\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_u)$ be the path from the root to a node \mathbf{t} . Then $\mathbf{t.val} = \mathbf{r.val}^{\prod_{1 \leq j \leq u} \mathbf{e}_j.p} \bmod n$, where \mathbf{r} denotes the root of the tree.

Note **ITag** and **Compress** are just vector form of **itag** and **compress**. We construct **itag**, **compress** as below. (1) **itag**(i): Set $\mathbf{r.val}$ to be some generator of a subgroup of \mathbb{Z}_n^* . **itag**(i) = $\mathbf{t.val}$, where \mathbf{t} is the i -th leaf. (2) **compress**($[a, b], \alpha$): Given an interval $[a, b] \subset [N]$, find the minimum set $\{\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_v\}$ of tree nodes, where $\mathbf{t}_j.int$ ’s, $j \in [v]$, form a partition of $[a, b]$. Output $\{\mathbf{t}_j.val^\alpha, j \in [v]\}$. Then a corresponding **uncompress** algorithm will recover $\{\ell_j.val^\alpha : \ell_j \text{ is the } j\text{-th leaf, } j \in [a, b]\}$, using the public primes associated to edges.

This alternative approach reduces the complexity of **uncompress** from $O(L)$ to $O(\log L)$, where L is the length of the interval $[a, b]$, at the cost of growth of output size (i.e. communication bandwidth of MAIA) by a factor of $O(\log N)$. Furthermore, the dynamic nature of tree provides us a clue to support dynamic operations on the dataset, for example, adding or deleting a point. We will look into this direction in the future.

6.2.2 Generalization of Dataset and Query

As mentioned previously, the dataset and query described in Section 3.2 are simplified for the convenience of presentation, and our scheme MAIA can handle more general case.

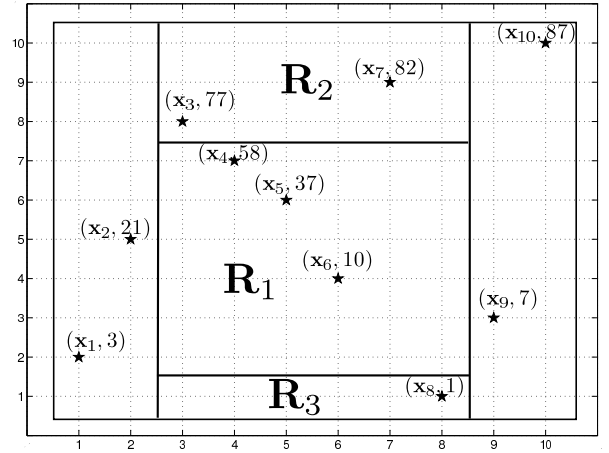
Let $\hat{\mathbf{D}} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\} \subset \mathbb{N}^d$ be a dataset. An aggregate **SUM** query with range $\mathbf{B} \subset \mathbb{N}^d$ over dataset $\hat{\mathbf{D}}$ asks for the sum $\bigoplus_{\mathbf{x} \in \hat{\mathbf{D}} \cap \mathbf{B}} \mathbf{x}$.

We can authenticate such query in the following way. (1) Normalize [38] the dataset $\hat{\mathbf{D}}$: for each dimension $\ell \in [d]$, sort all \mathbf{x}_j ’s along the ℓ -th dimension. Then associate a d -dimensional rank vector \mathbf{i}_j for each point \mathbf{x}_j , $j \in [N]$. Let $\mathbf{D} = \{\mathbf{i}_1, \mathbf{i}_2, \dots, \mathbf{i}_N\}$. Table 3 and Figure 2 show a normalized

Table 3: The dataset $\hat{\mathbf{D}} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{10}\}$. For each $1 \leq m \leq 10$, we write $\mathbf{x}_m = (x_{m,1}, x_{m,2})$, and $(i_{m,1}, i_{m,2})$ is the rank vector of \mathbf{x}_m , i.e. $i_{m,1}$ is the rank of $x_{m,1}$ among $\{x_{m',1} : 1 \leq m' \leq 10\}$. Similar for $i_{m,2}$ and $x_{m,2}$.

	\mathbf{x}_1	\mathbf{x}_2	\mathbf{x}_3	\mathbf{x}_4	\mathbf{x}_5	\mathbf{x}_6	\mathbf{x}_7	\mathbf{x}_8	\mathbf{x}_9	\mathbf{x}_{10}
$x_{m,1}$	5	7	11	24	31	45	58	61	83	97
$i_{m,1}$	1	2	3	4	5	6	7	8	9	10
$x_{m,2}$	3	21	77	58	37	10	82	1	7	87
$i_{m,2}$	2	5	8	7	6	4	9	1	3	10

Figure 2: The set $\hat{\mathbf{D}}$ of 10 points $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{10}$ are displayed on the $[1, 10] \times [1, 10]$ grid, where \mathbf{x}_m ’s, $1 \leq m \leq 10$, are as defined in Table 3. Each point \mathbf{x}_m is labeled with $(x_m, x_{m,2})$, and located at position $(i_{m,1}, i_{m,2})$. The 2D range $[3, 8] \times [2, 7]$ selects the region denoted by “ \mathbf{R}_1 ”.



2D dataset. (2) Define a vector-valued attribute function $\text{Att} : \mathbf{D} \rightarrow \hat{\mathbf{D}}$, such that $\text{Att}(i) = \mathbf{x}$, where $i \in \mathbf{D}$ is the rank vector of $\mathbf{x} \in \hat{\mathbf{D}}$. (3) Set up for dataset \mathbf{D} with attribute function Att , by running algorithms **KGen** and **DEnc**. (4) Translate query w.r.t. $\hat{\mathbf{D}}$ to query w.r.t. \mathbf{D} : find a range $\mathbf{R} \subseteq [N]^d$, such that $\bigoplus_{i \in \mathbf{D} \cap \mathbf{R}} \text{Att}(i) = \bigoplus_{\mathbf{x} \in \hat{\mathbf{D}} \cap \mathbf{B}} \mathbf{x}$. Alice can do this translation online with help of Bob and is able to verify the correctness of the translation by authenticating some tag values returned by Bob. (5) Run interactive algorithm **ProVer** to authenticate query **SUM**(\mathbf{R}) w.r.t. \mathbf{D} .

We remark that the translation of queries can be done with $O(d)$ communication bits and in $O(d)$ computation time, and this generalization will not change the complexity of our scheme asymptotically.

6.2.3 Authentication of COUNT, MIN Queries

Besides **SUM** query, we can generalize our scheme to authenticate aggregate **COUNT** query and **MIN** query. **COUNT** is just a special case of **SUM**: summing a constant attribute. **MIN** query can be converted to **COUNT** query: For any set S of numbers, $c = \min S \Leftrightarrow c \in S \wedge |S| = |\{x : x \in S \wedge x \geq c\}|$. To process a **MIN** query with range \mathbf{B} , which asks for $\min_{\mathbf{x} \in \hat{\mathbf{D}} \cap \mathbf{B}} \mathbf{x}[1]$, Alice and Bob interact in this way: (1) Bob sends back the query result c with the proof that

$c \in S$, where $S = \{\mathbf{x}[1] : \mathbf{x} \in \hat{\mathbf{D}} \cap \mathbf{B}\}$; (2) Alice generates two queries⁶ $\text{COUNT}(\mathbf{B})$ and $\text{COUNT}(\mathbf{B} \cap [c, +\infty) \times \mathbb{N}^{d-1})$ w.r.t. $\hat{\mathbf{D}}$, which ask for the sizes of set S and set $\{x : x \in S \wedge x \geq c\}$, respectively. Next, Alice issues the two queries to Bob and gets two count numbers with proofs. If all proofs are valid and the two count numbers are equal, Alice believes c is the minimum value.

Similarly, MAX, MEDIAN, QUANTILE queries can be authenticated by converting to COUNT query. So we conclude that: our scheme MAIA can authenticate all of SUM, COUNT, MIN, MAX, MEDIAN and QUANTILE aggregate queries.

6.2.4 Multidimensional Range Selection Query

It is straightforward to extend our scheme to authenticate range selection query with range \mathbf{R} , which asks for the set $\{\text{Att}(\mathbf{i}) : \mathbf{i} \in \mathbf{D} \cap \mathbf{R}\}$, with linear communication overhead: (1) Bob sends back $\{(\text{Att}(\mathbf{i}), \mathbf{i}, \mathbf{t}_\mathbf{i}) : \mathbf{i} \in \mathbf{D} \cap \mathbf{R}\}$; (2) Authenticate a COUNT query with range \mathbf{R} to ensure completeness.

Next, we brief a method to authenticate range selection query with $O(d^2)$ communication overhead. The key idea is that we can extend CollRes to authenticate weighted sum.

Suppose Alice wants to authenticate a range selection query with range \mathbf{R} , which asks for $S = \{\mathbf{x}[1] : \mathbf{i} \in \mathbf{D} \cap \mathbf{R}, \mathbf{x} = \text{Att}(\mathbf{i})\}$. She does it in three steps: (1) Issue the range selection query to Bob, and Bob returns the set S without any proof; (2) Issue a COUNT query with range \mathbf{R} and verify the result; (3) Simulate (modified) CollRes with Bob to authenticate a weighted sum query with range \mathbf{R} . More precisely, Alice asks for the weighted sum $\mathbf{X}[j] \leftarrow \sum_{\mathbf{i} \in \mathbf{D} \cap \mathbf{R}}^{\mathbf{x}=\text{Att}(\mathbf{i})} \mathbf{x}[j] \mathcal{H}(\mathbf{x}[1])$, $j \in [d]$, with $\mathcal{H}(\mathbf{x}[1])$ as weights, where \mathcal{H} is some collision-resistant hash function. The modification to CollRes is: (a) In equation (9), replace $\text{Att}(\mathbf{i})$ with $\text{Att}(\mathbf{i}) \otimes (\mathcal{H}(\text{Att}(\mathbf{i})[1])\mathbf{e})$, where \mathbf{e} is the d -dimensional vector with each dimension equal to 1; (b) Modify equations (10) and (12) accordingly, by “raising” $\mathbf{t}_\mathbf{x}$ and the output of Uncompress to power $\mathcal{H}(\text{Att}(\mathbf{i})[1])\mathbf{e}$ (See equation (3) for vector exponentiation). (c) Modify Alice’s verification steps accordingly, and add one more test⁷ $\mathbf{X}[1] \stackrel{?}{=} \sum_{\mathbf{x} \in S} \mathbf{x} \mathcal{H}(\mathbf{x})$. We save the details.

Assume all (legitimate) attribute values are distinct. From Step (3), Alice can verify the correctness of the result set S returned in Step (1). With the help of the count number from Step (2), Alice can ensure the completeness.

We remark that the generalization in Section 6.2.2 can be applied here, so that Alice can authenticate range selection query w.r.t. $\hat{\mathbf{D}} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, which asks for $\{\mathbf{x}[1] : \mathbf{x} \in \hat{\mathbf{D}} \cap \mathbf{B}\}$. The complexity of authentication of range selection query using our modified scheme is shown in Table 1 in Section 1.

6.2.5 Decrease the size of tag

In our scheme, the storage overhead, i.e. the size of tags, is $O(dN)$, which is linear w.r.t. the size of dataset, where d is the dimension and N is the number of points in the dataset. However, to be secure in practice [55], the security parameter κ is at least 1000. That means each point of size $32d$ bits⁸ has a tag of size $1000d$ bits. Here we present

⁶Note that with the generalization in Section 6.2.2, our scheme can authenticate aggregate such queries.

⁷If \mathbf{X} is too large, Bob can just send back $(ge) \otimes \mathbf{X}$ for some generator g .

⁸In real word applications, one dimension (of integer type)

an alternative construction of Tag which decreases the the storage overhead to $O(N)$, so that a point of size $32d$ bits has a tag of size 1000 bits.

Let n be a RSA modulus, such that $n = (2p + 1)(2q + 1)$, all of $p, q, 2p + 1, 2q + 1$ are primes and p and q has the same bit-length. Let G_p, G_q be cyclic multiplicative subgroups of \mathbb{Z}_n^* of order p and q , respectively, and $g_p \in G_p$ and $g_q \in G_q$ are generators of G_p and G_q , respectively. Let $\mathbf{i} \in \mathbf{D}$, $\mathbf{x} = \text{Att}(\mathbf{i})$, and let $\mathbf{s}_1, \mathbf{s}_2 \in (\mathbb{Z}_p^*)^d$ be two secret vectors. Let $\mathbf{h}(\mathbf{i}) = (\frac{1}{h(\mathbf{i}[1])}, \frac{1}{h(\mathbf{i}[2])}, \dots, \frac{1}{h(\mathbf{i}[d])})$. We define

$$\text{Tag}(\mathbf{x}, \mathbf{i}) \stackrel{\text{def}}{=} g_p^{\mathbf{h}(\mathbf{i}) \cdot \mathbf{s}_1 + \mathbf{x} \cdot \mathbf{s}_2} g_q^{f_{s_3}(\mathbf{i}[1])} \pmod{n}. \quad (16)$$

Note: (1) Since here the output of Tag is a scalar instead of vector, and the secret vectors \mathbf{s}_1 and \mathbf{s}_2 can prevent mixing and permutation on dimensions, we eliminate the need of binding factor g_q^{ξ} , whose role is to prevent mixing and permutation on dimensions. As a result, the modulus n becomes product of two primes instead of three primes. (2) We use $f_{s_3}(\mathbf{i}[1])$ instead of $f_{s_3}(\mathbf{i})$, so the optimization for 1D case given in Section 6.1.1 still works.

7. CONCLUSION

We proposed a scheme to authenticate aggregate range query over static multidimensional dataset, and the communication complexity (in term of bits) is $O(d^2)$ (d is the dimension) and independent of the query range size and the dataset size. Aggregate operations that our scheme can support include summing, counting, and finding of the minimum and maximum. Our scheme and techniques can be useful in various other applications. In particular, we showed that our scheme with slight modifications can authenticate multidimensional range selection query, and improved the communication overhead significantly compared with previous works. We will look into the possibility to support dynamic operations on the dataset, like adding or deleting a point, using a tree based compression method.

In practice, our scheme is very suitable for dataset with high dimensions but small number of tuples. For example, outsourcing a high dimensional database to a P2P system. The horizontal partition of the database by a P2P system ensures that each partition has a small number of tuples and our scheme allows horizontal partition of dataset. We will look into this direction.

8. REFERENCES

- [1] 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
- [2] Bethencourt, J., Sahai, A., Waters, B.: Ciphertext-policy attribute-based encryption. In: SP ’07: IEEE Symposium on Security and Privacy. (2007) 321–334
- [3] Goyal, V., Pandey, O., Sahai, A., Waters, B.: Attribute-based encryption for fine-grained access control of encrypted data. In: CCS ’06: ACM conference on Computer and communications security. (2006) 89–98
- [4] Boneh, D., Waters, B.: Conjunctive, subset, and range queries on encrypted data. In: TCC ’07: Theory of Cryptography. (2007) 535–554
- [5] Katz, J., Sahai, A., Waters, B.: Predicate encryption supporting disjunctions, polynomial equations, and inner

of a point is typically 32 bits.

- products. In: EUROCRYPT '08: International Conference on Advances in Cryptology. (2008) 146–162
- [6] Shi, E., Waters, B.: Delegating capabilities in predicate encryption systems. In: ICALP '08: International colloquium on Automata, Languages and Programming. (2008) 560–578
- [7] Brian Thompson, Danfeng Yao, S.H.W.G.H., Sander, T.: Privacy-preserving computation and verification of aggregate queries on outsourced databases. In: PETS '09: Privacy Enhancing Technologies Symposium (PETS). (2009)
- [8] Rivest, R.L., Shamir, A., Adleman, L.: A method for obtaining digital signatures and public-key cryptosystems. *Commun. ACM* **21**(2) (1978) 120–126
- [9] Paillier, P.: Public-key cryptosystems based on composite degree residuosity classes. In: EUROCRYPT '99: International Conference on Advances in Cryptology. (1999) 223–238
- [10] Boldyreva, A., Chenette, N., Lee, Y., O'Neill, A.: Order-preserving symmetric encryption. In: EUROCRYPT '09: International Conference on Advances in Cryptology. (2009) 224–241
- [11] Gentry, C.: Fully homomorphic encryption using ideal lattices. In: STOC '09: Annual ACM symposium on Theory of computing. (2009) 169–178
- [12] Ateniese, G., Burns, R., Curtmola, R., Herring, J., Kissner, L., Peterson, Z., Song, D.: Provable data possession at untrusted stores. In: CCS '07: ACM conference on Computer and communications security. (2007) 598–609
- [13] Chang, E.C., Xu, J.: Remote integrity check with dishonest storage server. In: ESORICS '08: European Symposium on Research in Computer Security. (2008) 223–237
- [14] Shacham, H., Waters, B.: Compact proofs of retrievability. In: ASIACRYPT '08: International Conference on the Theory and Application of Cryptology and Information Security. (2008) 90–107
- [15] Bowers, K.D., Juels, A., Oprea, A.: Hail: a high-availability and integrity layer for cloud storage. In: CCS '09: ACM Conference on Computer and Communications Security. (2009) 187–198
- [16] Dodis, Y., Vadhan, S., Wichs, D.: Proofs of retrievability via hardness amplification. In: TCC '09: Theory of Cryptography. (2009) 109–127
- [17] Ateniese, G., Kamara, S., Katz, J.: Proofs of storage from homomorphic identification protocols. In: ASIACRYPT '09: Annual International Conference on the Theory and Application of Cryptology and Information Security. (2009) 319–333
- [18] 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
- [19] 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
- [20] Devanbu, P.T., Gertz, M., Martel, C.U., Stubblebine, S.G.: Authentic third-party data publication. In: IFIP '00: IFIP TC11/ WG11.3 Working Conference on Database Security. (2000) 101–112
- [21] 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
- [22] Devanbu, P., Gertz, M., Martel, C., Stubblebine, S.G.: Authentic data publication over the internet. *J. Comput. Secur.* **11**(3) (2003) 291–314
- [23] 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
- [24] Mykletun, E., Narasimha, M., Tsudik, G.: Authentication and integrity in outsourced databases. *Trans. Storage* **2**(2) (2006) 107–138
- [25] Pang, H.H., Tan, K.L.: Verifying completeness of relational query answers from online servers. *ACM Trans. Inf. Syst. Secur.* **11**(2) (2008) 1–50
- [26] Sion, R.: Query execution assurance for outsourced databases. In: VLDB '05: International conference on Very large data bases. (2005) 601–612
- [27] Cheng, W., Tan, K.L.: Query assurance verification for outsourced multi-dimensional databases. *J. Comput. Secur.* **17**(1) (2009) 101–126
- [28] 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
- [29] 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
- [30] 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
- [31] 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
- [32] 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
- [33] Pang, H.H., Zhang, K.M.: Scalable verification for outsourced dynamic databases. In: VLDB '09: International Conference on Very Large Data Bases. (2009)
- [34] Goodrich, M.T., Tamassia, R., Triandopoulos, N.: Super-efficient verification of dynamic outsourced databases. In: CT-RSA '08: The Cryptographers' Track at the RSA Conference. (2008) 407–424
- [35] 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.
- [36] Shamir, A.: How to share a secret. *Commun. ACM* **22**(11) (1979) 612–613
- [37] Pedersen, T.P.: Non-interactive and information-theoretic secure verifiable secret sharing. In: CRYPTO '91: Advances in Cryptology. (1991) 129–140
- [38] Preparata, F.P., Shamos, M.I.: *Computational geometry: an introduction*. (1985)
- [39] Okamoto, T., Uchiyama, S.: A new public-key cryptosystem as secure as factoring. In: EUROCRYPT '98: International Conference on Advances in Cryptology. (1998) 308–318
- [40] González Nieto, J.M., Boyd, C., Dawson, E.: A public key cryptosystem based on a subgroup membership problem. *Des. Codes Cryptography* **36**(3) (2005) 301–316
- [41] Yamamura, A., Saito, T.: Private information retrieval based on the subgroup membership problem. In: ACISP '01: Australasian Conference on Information Security and Privacy. (2001) 206–220
- [42] Gjosteen, K.: Subgroup membership problems and public key cryptosystems. PhD thesis, NTNU (2004)
- [43] Boneh, D., Goh, E.-J., Nissim, K.: Evaluating 2-DNF Formulas on Ciphertexts. In: TCC '05: Theory of Cryptography. (2005) 325–341
- [44] Brown, J., Nieto, J.M.G., Boyd, C.: Concrete chosen-ciphertext secure encryption from subgroup membership problems. In: CANS '06: Cryptology and Network Security. (2006) 1–18
- [45] Gjosteen, K.: Symmetric subgroup membership problems. In: PKC '05: Public Key Cryptography. (2005) 104–119
- [46] Hofheinz, D., Kiltz, E.: The group of signed quadratic residues and applications. In: CRYPTO '09: Advances in

Cryptology. (2009) 637–653

- [47] Gennaro, R., Halevi, S., Rabin, T.: Secure hash-and-sign signatures without the random oracle. In: EUROCRYPT '99: International Conference on Advances in Cryptology. (1999) 123–139
- [48] Goldreich, O.: Foundations of Cryptography: Volume 1. (2006)
- [49] Johnson, R., Molnar, D., Song, D.X., Wagner, D.: Homomorphic signature schemes. In: CT-RSA '02: The Cryptographer's Track at the RSA Conference on Topics in Cryptology. (2002) 244–262
- [50] Boneh, D., Lynn, B., Shacham, H.: Short signatures from the weil pairing. *J. Cryptol.* **17**(4) (2004) 297–319
- [51] 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
- [52] Hacigümüş, H., Iyer, B.R., Mehrotra, S.: Efficient execution of aggregation queries over encrypted relational databases. In: DASFAA. (2004) 125–136
- [53] Mykletun, E., Tsudik, G.: Aggregation queries in the database-as-a-service model. In: IFIP '06: IFIP WG 11.3 Working Conference on Data and Applications Security. (2006) 89–103
- [54] 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
- [55] Hinek, M.J.: On the security of multi-prime rsa. *Journal of Mathematical Cryptology* **2**(2) (2008) 117–147
- [56] Miller, G.L.: Riemann's hypothesis and tests for primality. In: STOC '75: ACM Symposium on the Theory of Computing. (1975) 234–239
- [57] Xu, J., Chang, E.C.: Authenticating Aggregate Range Queries over Multidimensional Dataset. *Cryptology ePrint Archive*, Report 2010/050 <http://eprint.iacr.org/>.

APPENDIX

We remark that the design of Tag function is conservative and there may be some redundancy in our construction.

A. IDEA OF PROOF OF THEOREM 2

Here we brief the idea to prove Theorem 2. The security model is given in Definition 2. The correctness part is straightforward. We focus on the soundness part.

Let \mathcal{R} be the variant logarithm function as defined in Assumption 2. We define \mathcal{P} and \mathcal{Q} similar to \mathcal{R} . Define function $\mathcal{P} : G \rightarrow \mathbb{Z}_p$, such that for any $(x, y, z) \in \mathbb{Z}_p \times \mathbb{Z}_q \times \mathbb{Z}_r$, $\mathcal{P}(g_p^x g_q^y g_r^z) \stackrel{\text{def}}{=} x$. Define function $\mathcal{Q} : G \rightarrow \mathbb{Z}_q$, such that for any $(x, y, z) \in \mathbb{Z}_p \times \mathbb{Z}_q \times \mathbb{Z}_r$, $\mathcal{Q}(g_p^x g_q^y g_r^z) \stackrel{\text{def}}{=} y$.

CLAIM 1. *Assumption 2 holds for \mathcal{Q} and \mathcal{R}*

During all interactions with Alice (i.e. adversary's learning phase), Bob receives i 's, $\text{Att}(i)$, values a_j 's, b_j 's which specify query range, and a lot of elements from group G , including tag values t_i and compressed *Help-Info* δ . For any element $g_p^x g_q^y g_r^z \in G$ Bob receives, both v and w are (cryptographically secure) pseudo-random: (1) In a tag value t_i , ξ and $s_{\ell,2}$, $\ell \in [d]$, are true random and $f_s(\cdot)$ is PRF; (2) In each compressed *Help-Info* $\delta = (v \otimes \alpha) \otimes \mathbf{u}_1 \otimes \mathbf{u}_2$, where $v \in (G_p)^d$, $\mathbf{u}_1 \in (G_q)^d$ and $\mathbf{u}_2 \in (G_r)^d$, both \mathbf{u}_1 and \mathbf{u}_2 are true random. See Section 5.4 for *Compress*. Hence, the prerequisite of Assumption 2 is satisfied for both \mathcal{Q} and \mathcal{R} .

CLAIM 2. *It is hard for a dishonest Bob to pass tests in equations (13)(14), beyond homomorphism of \mathcal{E} .*

1. **tag** and **itag** are unforgeable, by Lemma 1 and [47]. The details of the proof will require Assumption 1.
2. One may view $\text{itag}(i)^\alpha$ as a variant version of multi-prime RSA signature⁹ [55] on message $\text{itag}(i)$, with private key α and public key $\alpha^{-1} \bmod p$, where both public key and private key are kept securely from adversaries. The definition of security of homomorphic signature scheme (intuitively, unforgeability except for homomorphism) is given in [49]. For simplicity, we ignore the requirement of security of multi-prime RSA signature in the statement of the main theorem.
3. Due to the binding factor g_q^ξ among different dimensions (Equation (14) checks the binding factor), $\text{ITag}(i) \otimes \alpha$ can not be forged beyond homomorphism. This can be proved by applying Assumption 2 on \mathcal{Q} . Note this homomorphism corresponds to double counting or undercounting points that satisfy the query condition.

CLAIM 3. *If an adversary cheats by double counting and undercounting, and passes the verification in ProVer, then he can break Assumption 2.*

Suppose the adversary (i.e. malicious Bob) double counts points in the multiset $S_1 = \{i_{u,j} : i_{u,j} \in \mathcal{D}, j \in [m]\}$ and undercounts points in set $S_2 \subseteq \mathcal{D}$, where $S_1 \neq S_2$, and passes the verification. That means, $\eta_1 = \eta_2$, where

$$\eta_1 = \prod_{i \in S_1} \text{ctag}(i[1]) \bmod n; \quad \eta_2 = \prod_{i \in S_2} \text{ctag}(i[1]) \bmod n.$$

According to the definition of **tag**, **itag** and \mathcal{R} , $\log_{g_r} \eta_1 = \mathcal{R}(\prod_{i \in S_1} t_i[1])$ and $\log_{g_r} \eta_2 = \mathcal{R}(\prod_{i \in S_2} t_i[1])$. Hence,

$$\mathcal{R}\left(\frac{\prod_{i \in S_1} t_i[1]}{\prod_{i \in S_2} t_i[1]}\right) = \mathcal{R}\left(\prod_{i \in S_1} t_i[1]\right) - \mathcal{R}\left(\prod_{i \in S_2} t_i[1]\right) = 0.$$

Let $A = \frac{\prod_{i \in S_1} t_i[1]}{\prod_{i \in S_2} t_i[1]} \bmod n$. If $A = 1$, this will lead to an efficient method to factorize n . See Lemma 6 in Appendix C. If $A \neq 1$, then $(A, 0)$ is a contradiction to Assumption 2. Note that the Assumption 2 implies that it is hard to factorize n .

B. SECURITY OF AUTHENTICATION TAG FUNCTION

B.1 Proof of Lemma 1

LEMMA 1. *Let k be a key as in Definition 3. Given $S = \{(x_i, i, t_i) : \xi \stackrel{\$}{\leftarrow} \mathbb{Z}_q^*, t_i \leftarrow \text{tag}_k(x_i, i; \xi)\}_{1 \leq i \leq N}$ and the modulus n , it is computationally hard to forge a tuple $(y, j, t) \notin S$, such that $\gamma_1(t, g_p^y) = \text{itag}_k(j)$ and $\gamma_3(t) = \text{ctag}_k(j)$, under Assumption 1 and Assumption 2.*

Let $X \sim Y$ denote that X and Y are computationally indistinguishable.

PROOF. Suppose that some PPT adversary \mathcal{A} takes as input a set (S, n) , and outputs $(y, j, t) \notin S$, such that $\gamma_1(t, g_p^y) = \text{itag}_k(j)$ and $\gamma_3(t) = \text{ctag}_k(j)$ with non-negligible probability. Note Assumption 2 is applicable for \mathcal{R} .

⁹Or a symmetric key version of BLS signature [50].

Case 1: $j \in [N]$. Let $(x_j, j, t_j) \in S$ be a tuple in the set S . Since $(y, j, t) \notin S$, we have $(y, t) \neq (x_j, t_j)$.

Case 1.1: $t_j \neq t$. From $\gamma_3(t) = \text{ctag}_k(j) = \gamma_3(t_j)$, we have

$$\mathcal{R}(t) = \log_{g_r} \gamma_3(t) = \log_{g_r} \gamma_3(t_j) = \mathcal{R}(t_j).$$

Let $A = t_j t^{-1} \pmod n$. Let $A \not\equiv 1 \pmod n$, and $B = \mathcal{R}(A) = 0$. The tupe (A, B) is a contradiction with Assumption 2.

Case 1.2: $t_j = t$ and $x_j \neq y$. From $t_j = t$, we have,

$$\begin{aligned} \mathcal{P}(t_j) &= \mathcal{P}(t) \\ \Rightarrow x_j + \frac{s_1}{h(j)} &\equiv y + \frac{s_1}{h(j)} \pmod p \\ \Rightarrow x_j &\equiv y \pmod p. \end{aligned}$$

Hence, $p | (x_j - y)$. Simulate \mathcal{A} again, we can get another multiple of p , then apply the GCD algorithm to obtain p . Next, we show that with the knowledge of p , the adversary can break Assumption 1.

Assumption 1 can also be applied on group G_p or G_q . Let $Y_p = \{y_i : y_i \xleftarrow{\$} G_p\}_{1 \leq i \leq N}$, $Y_q = \{y_i : y_i \xleftarrow{\$} G_q\}_{1 \leq i \leq N}$, and $\mathcal{X} = \{y_i : y_i \xleftarrow{\$} G\}_{1 \leq i \leq N}$. Let $S_t = \{t_i : \exists x, i, (x, i, t_i) \in S\}$. Since ξ 's are true random, by Assumption 1 with multiple sampling (**Theorem 3.2.6** [48]), $S_t \sim Y_q \sim \mathcal{X} \sim Y_p$. Let $\tilde{Y}_q = \{(x_i, i, y) : y \in Y_q\}$, $\tilde{Y}_p = \{(x_i, i, y) : y \in Y_p\}$, and $\tilde{\mathcal{X}} = \{(x_i, i, y) : y \in \mathcal{X}\}$. Note that $S = \{(x_i, i, t_i) : t_i \in S_t\}$. We have $S \sim \tilde{Y}_q \sim \tilde{\mathcal{X}} \sim \tilde{Y}_p$, since ξ 's in t_i 's are true random independent of x_i 's and i 's.

However, the adversary can run \mathcal{A} with input (S, n) and apply GCD algorithm to find p . With the value of p , the adversary can distinguish S and \tilde{Y}_p . Contradiction!

Case 2: $j \notin [N]$. Note that in our use of **tag** in this paper, this case will not be considered as a valid forgery, since we are only interested in normalized value $i \in [N]$, and (y, j, t) with $j \notin [N]$ will be rejected by Alice immediately.

Note that in the above proof, we do not utilize the condition $\gamma_1(t, g_p^y) = \text{itag}_k(j)$. As mentioned previously, our design is conservative and may contain some redundancy. \square

B.2 Security of tag function Tag

Let $\mathbf{D} = \{i_1, i_1, \dots, i_N\} \subset [N]^d$ be a dataset as in Section 3.2. That is, \mathbf{D} satisfies property: For any $j \in [d]$, $\{i[j] : i \in \mathbf{D}\} = [1, N]$.

Let \mathcal{K} be a key as in Definition 4. We construct a set S as below: For any $i \in \mathbf{D}$,

1. let \mathbf{x}_i be any element in \mathbb{N}^d ;
2. choose ξ_i at random from \mathbb{Z}_q^* ;
3. compute $\mathbf{t}_i \leftarrow \text{Tag}_{\mathcal{K}}(\mathbf{x}_i, \mathbf{i}; \xi)$.

Let $S = \{(\mathbf{x}_i, \mathbf{i}, \mathbf{t}_i) : \mathbf{i} \in \mathbf{D}\}$.

LEMMA 3. *Given (S, n) , it is hard to forge $(\mathbf{y}, \mathbf{j}, \mathbf{t}) \notin S$, such that (1) $\mathbf{\Gamma}_1(\mathbf{t}, \mathbf{y}) = \text{lTag}_{\mathcal{K}}(\mathbf{j})$, (2) there exists $a \in \mathbb{Z}_n$, $\mathbf{\Gamma}_2(\mathbf{t}, \mathbf{y}) \otimes \mathbf{s}_2 = a\mathbf{e} \pmod n$, (3) $\mathbf{\Gamma}_3(\mathbf{t}) = \text{CTag}_{\mathcal{K}}(\mathbf{j})$, (4) $\mathbf{j} \in [N]^d$, under Assumption 1 and Assumption 2.*

PROOF. (Sketch) Without considering the binding factor g_q^{ξ} 's, this problem can be considered as d instances of the problem in Lemma 1.

The adversary either (1) forges a tuple (y, j, t) as in Lemma 1, or (2) permute on dimensions of a valid tuple in S ; (3) mix

and combine different dimensions of different valid tuples in S .

Case (1) is prevented by Lemma 1; Case (2) is prevented by $\mathbf{s}_1, \mathbf{s}_2, \mathbf{s}_3$, which are parts of key \mathcal{K} ; Case (3) is prevented by the binding factor. \square

C. PROOF OF THEOREM 2

The framework of proof is given in Appendix A. Here we fill in the missing parts.

C.1 For Claim 2

We extract out the verification in equation (13) and treat as a private key signature scheme.

LEMMA 4. *Let $S = \{(X_i Y_{i,1} Z_{i,1} \pmod n, X_i^\alpha Y_{i,2} Z_{i,2} \pmod n) : X_i \in G_p, Y_{i,j} \in G_q, Z_{i,j} \in G_r, j \in \{1, 2\}\}_{1 \leq i \leq m}$, where m is some polynomial of the security parameter κ , i.e. the bit length of the modulus n . Given (S, n) , it is hard to forge a valid pair $(XY_1 Z_1 \pmod n, X^\alpha Y_2 Z_2 \pmod n) \notin S$, $X \in G_p, Y_1, Y_2 \in G_q, Z_1, Z_2 \in G_r$, beyond multiplicative homomorphism, if Multiple-Prime RSA [55] is unforgeable beyond multiplicative homomorphism.*

PROOF. Let $(X, Y, Z, U) \xleftarrow{\$} (G_p, G_q, G_r, G)$. By Assumption 1, $X \sim U, Y \sim U, Z \sim U$. Then $XYZ \sim XU \sim U \sim X$ (Here XYZ means $X \times Y \times Z \pmod n$).

By multiple sampling (**Theorem 3.2.6** [48]), we have

$$(X, X^\alpha) \sim (XY_1 Z_1, X^\alpha Y_2 Z_2) \text{ and } (X, X^\alpha) \sim (U, U^\alpha).$$

Hence,

$$(XY_1 Z_1, X^\alpha Y_2 Z_2) \sim (U, U^\alpha).$$

Note that $(X, X^\alpha) \sim (U, U^\alpha)$. (U, U^α) is the multiple-Prime RSA signature scheme without revealing the public key. \square

Let \mathcal{E} and \mathcal{K} be as in Section 4. Let \mathbf{D} and S be the as in Section B.2. Note that $\mathbf{t}_i = \mathcal{E}_{\mathcal{K}}(\text{lTag}_{\mathcal{K}}(\mathbf{i}), \mathbf{v}_i, \mathbf{w}_i, \mathbf{x}_i)$, where $\mathbf{v}_i \in (G_q)^d$ and $\mathbf{w}_i \in (G_r)^d$.

Let $\mathbf{R} \subseteq [N]^d$ be range. For each $\mathbf{i} \in \mathbf{R}$,

1. choose α at random from $(\mathbb{Z}_p^*)^d$;
2. let $\mathbf{u}_i \leftarrow \mathcal{E}_{\mathcal{K}}(\text{lTag}_{\mathcal{K}}(\mathbf{i}), \mathbf{v}_i, \mathbf{w}_i, \mathbf{0})$, where $\mathbf{v}_i \in (G_q)^d$ and $\mathbf{w}_i \in (G_r)^d$.

Let $\bar{S} = \{\mathbf{u}_i : \mathbf{i} \in \mathbf{R}\}$. The following lemma is for the CLAIM 2.

LEMMA 5. *Given (S, \bar{S}, n) , it is hard to forge a tuple $(\Psi_1, \Psi_2, \mathbf{X})$ beyond homomorphism of \mathcal{E} , where $\Psi_1 = \mathcal{E}_{\mathcal{K}}(U, \mathbf{V}_1, \mathbf{W}_1, \mathbf{X})$, $\Psi_2 = \mathcal{E}_{\mathcal{K}}(U \otimes \alpha, \mathbf{V}_2, \mathbf{W}_2, \mathbf{0})$, and $U \in (G_p)^d$, $\mathbf{V}_1, \mathbf{V}_2 \in (G_q)^d$, $\mathbf{W}_1, \mathbf{W}_2 \in (G_r)^d$, such that $(\Psi_1, \Psi_2, \mathbf{X})$ passes verification in equations (13) (14).*

PROOF. (Sketch) Without considering the binding factor g_q^{ξ} 's, this problem can be considered as d instances of the problem in Lemma 4.

The adversary either (1) forges a **Tag** value for some $\mathbf{i} \in \mathbf{R} \setminus \mathbf{D}$ or $\mathcal{E}_{\mathcal{K}}(\text{lTag}_{\mathcal{K}}(\mathbf{i}) \otimes \alpha, \mathbf{v}, \mathbf{w}, \mathbf{0})$ for some $\mathbf{i} \in \mathbf{D} \setminus \mathbf{R}$, or (2) forges a tuple as in Lemma 4, or (3) permute on dimensions of a valid tuple, or; (4) mix and combine different dimensions of different valid tuples.

Case (1) is prevented since **Tag** and (loosely speaking) **lTag** are unforgeable; Case (1) is prevented by Lemma 4;

Case (2) is prevented by s_1, s_2 (Note Alice does not have enough information to check CTag with s_3 , so there is no s_3 here), which are parts of key \mathcal{K} ; Case (3) is prevented by the binding factor, which is checked in equation (14). \square

C.2 For Claim 3

Let group G be as in Assumption 1.

LEMMA 6. *Given a set $S = \{y_i : y_i \xleftarrow{\$} G\}_{1 \leq i \leq m}$ and the modulus n , it is hard to find a set $S_1 \subset S$ and a multiset $S_2 \subset S$, such that (1) $S_1 \neq S_2$ and (2) $\prod_{y \in S_1} y = \prod_{y \in S_2} y \pmod n$, assuming factorization of n is hard.*

Here multiset $S_2 \subset S$, means any element of S_2 is in S , without considering the number of duplications of elements in S_2 .

Note that by Assumption 1, $\{\mathbf{t}_i[j] : \mathbf{t}_i \in \mathbf{T}, j \in [d]\}$ is indistinguishable from uniform distribution over group G , where \mathbf{T} is the set of tag values generated by DEnc (See Section 5.2).

PROOF. Assume that an algorithm \mathcal{A} with inputs (S, n) can output such two sets with non-negligible probability.

1. Choose g at random from \mathbb{Z}_n . With high chance (non-negligible, more precisely, about $1/8$), g is a generator of G .
2. For $1 \leq i \leq m$, choose x_i at random from $[1, n^2]$ and compute $y_i \leftarrow g^{x_i} \pmod n$ (Here we purposely choose the domain with size much larger than $\lambda(n)$. So y_i only carries partial information of x_i).
3. Let $S = \{y_i : 1 \leq i \leq m\}$. Run $\mathcal{A}(S, n)$ to get sets S_1 and S_2 .
4. Compute $W_1 \leftarrow \sum_{y_i \in S_1} x_i$ and $W_2 \leftarrow \sum_{y_i \in S_2} x_i$.
5. We have $\lambda(n) | (W_1 - W_2)$ and with high chance $W_1 - W_2$ is not 0 (This is because we choose a domain of size larger than n , and \mathcal{A} has access to only a partial information of x_i 's). From $\lambda(n)$, we can factorize n [56].

\square