

Secure Data Retrieval on the Cloud

Homomorphic Encryption meets Coresets

Adi Akavia*, Dan Feldman[†] and Hayim Shaul[‡]

University of Haifa
akavia@cs.haifa.ac.il
dan.feldman@gmail.com
hayim.shaul@gmail.com

Abstract. *Secure report* is the problem of a client that retrieves all records matching specified attributes from a database table at the server (e.g. cloud), as in SQL SELECT queries, but where the query and the database are encrypted. Here, only the client has the secret key, but still the server is expected to compute and return the encrypted result. Secure report is theoretically possible with Fully Homomorphic Encryption (FHE). However, the current state-of-the-art solutions are realized by a polynomial of degree that is at least linear in the number m of records, which is too slow in practice even for very small databases.

We present the first solution that is realized by a polynomial that attains degree independent of the number of records m , as well as the first implementation of an FHE solution to Secure report. This is by suggesting a novel paradigm that forges a link between cryptography and modern data summarization techniques known as coresets (core-sets), and sketches in particular. The key idea is to compute only a coreset of the desired report. Since the coreset is small, the client can quickly decode the desired report that the server computes after decrypting the coreset.

We implemented our main reporting system in an open source library. This is the first implemented system that can answer such database queries when processing only FHE encrypted data and queries. As our analysis promises, the experimental results show that we can run Secure report queries on billions records in minutes on an Amazon EC2 server, compared to less than a hundred-thousands in previous FHE based solutions.

Keywords: Fully Homomorphic Encryption · Secure Computation · Secure Outsourcing of Computation · Secure Search · Coresets · Sketches · Group Testing

1 Introduction

Outsourcing storage and computation to large third-party systems often called “the cloud” (server) has become the norm for organizations and individuals (client). Typically this involves a trust-relationship where the client is required to reveal her private records to the server who executes the computations for her. To avoid this undesirable exposure of her personal, proprietary, or other sensitive information, secure computation techniques [Yao86, GMW87a] offer an appealing alternative: it is feasible for the client and server to jointly compute the desired functionality, while revealing nothing beyond the designated output.

*This work was supported in part by the Center for Cyber Law & Policy at the University of Haifa in conjunction with the Israel National Cyber Directorate in the Prime Minister’s Office.

[†]Supported by BSF grant no 2014627. This work was done in part while he was visiting the Simons Institute for the Theory of Computing.

[‡]Supported by grant no. 2014384 from the U.S.-Israeli Binational Science Foundation. Supported in part by the Bar Ilan Cyber Center.

A main challenge however in the design of secure computation is to attain low overhead over the un-secure version, in terms of both the communication and computation complexity. This translates to the goal of presenting the computation as a polynomial over a finite field with (a) low-degree and (b) low overall number of multiplications (see e.g. [BI05]).

Secure report. In this work we focus on secure outsourcing of *data storage and retrieval tasks* – a fundamental computational building block useful in numerous applications such as secure retrieval from a database, a corpus of tagged images, text documents, genomic data, and so forth. In the secure outsourcing context the client is the owner of data, queries, and retrieval outcome. The server is expected to provide the storage and desired data retrieval functionality while learning no (or minimal) information on the client’s data, queries and retrieved outcome.

We address the *report* data retrieval functionality, where the data is an unsorted array of elements $array = (x_1, \dots, x_m)$, the query ℓ corresponds to an agreed predicate $ISMATCH(x_i, \ell) \in \{0, 1\}$ (see below), and the retrieved outcome should be the list of (index, value) pairs for all array elements matching the query:

$$\mathcal{L} = \{(i, x_i) \mid ISMATCH(x_i, \ell) = 1\}.$$

Namely, the *report functionality* involves two parties, a client and a server, where the client’s input and output are $(array, \ell)$ and \mathcal{L} respectively, and the server has no input or output other than the shared parameters (the security parameter, and upper bounds on the number and size of records, queries and output \mathcal{L}).

isMatch Implementation. The predicate $ISMATCH(x_i, \ell)$ can be instantiated with any matching criterion, in a black-box fashion, to address versatile retrieval tasks. For example, *exact-match retrieval* is addressed by taking $ISMATCH$ to be an equality-test (i.e., $ISMATCH(x_i, \ell) = 1$ if-and-only-if $x_i = \ell$); *similarity-search* is addressed by taking $ISMATCH$ to be the predicate accepting 1 if-and-only-if $DIST(x_i, \ell) \leq \tau$ where τ is a threshold, and $DIST$ is a distance measure such as the Hamming, Euclidean, or Edit. Likewise, appropriate instantiations of $ISMATCH$ on encrypted data yield wild-card matching, range queries, Boolean queries, and sub-array queries; See [YSK⁺13, CKL15, LLN15, CKK16, KLL⁺16, TJW⁺16, KLL⁺17, AGHL18]. In this work we address the complementary problem of retrieving the set of matching elements, as defined by the provided $ISMATCH$ predicate.

For *Secure report*, the server should provide the report functionality obliviously of data, query and outcome. Furthermore, we make the following requirements: the protocol should use a *single-server* in *single-round* with *low communication* proportional only to the size of encrypted input and output, and perform *no pre-processing* of the cleartext data.

Motivating use-cases scenarios for our above requirements arise in settings where either: (a) We cannot assume non-collusion between multiple servers, or do not have availability of multiple servers; (b) Communication is intermittent, unreliable or costly, e.g., when communicating with mostly off-line clients as with sensors networks, posing a limiting factor on the number of communication rounds; (c) Communication is slow or expensive leading to the low bandwidth requirement; (d) it is infeasible to pre-process the cleartext data, e.g., for sorting, indexing or insertion into search-oriented data-structures such as search trees or hash tables. We note that retrieval from an unsorted data necessitates a linear-scan of the data even on cleartext data; yet it naturally arises in many use-cases. Examples for settings disallowing pre-processing the cleartext data follow:

- A-priori unknown matching criterion, as in ad-hoc SQL queries and our generic $ISMATCH$ predicate;

- Range queries on high dimensional data, where index is infeasible due to size exponential in the number of attributes;
- Streaming data with each element discarded by the client immediately after being encrypted and uploaded to the server;
- Low capacity client too weak to store or sort the clear-text data, as in Internet-of-Things (IoT) devices;
- Fragmented data uploaded to the server from multiple distinct clients (e.g., agents, users, devices) with no single entity that holds entire clear-text data.

In all these examples, pre-processing the cleartext data is infeasible, whether because no single entity holds the entire cleartext data or has the capacity to process it, or because the required index is exponential large, or unknown in advance.

Threat Model. Our security requirement is that the adversary controlling the server learns no new information from participating in the protocol. For example, the adversary must not distinguish between two adversarially-chosen equal size queries ℓ (similarly, data arrays $array$) from participating in the protocol. We note that there is no need to consider adversaries controlling the client because the server has no input or output. We address computationally-bounded semi-honest adversaries, i.e., adversaries that follow the protocol but may try to learn additional information.

Secure report on FHE encrypted data. Protocols satisfying the above requirements are feasible using Fully Homomorphic Encryption (FHE) [RAD78, Gen09a, Gen09b]. FHE is an encryption scheme that enables homomorphically-evaluating polynomials f over encrypted input; that is, given a ciphertext $\llbracket x \rrbracket$ encrypting the input x , it is possible to compute the ciphertext $\llbracket y \rrbracket$ encrypting the output $y = f(x)$ (with no access to a decryption key that would compromise security).

Secure report protocol using FHE would easily follow by providing a polynomial f realizing the above report functionality, i.e., f s.t.

$$f(array, \ell) = \mathcal{L}.$$

The protocol using f would be as follows. The client first generates public and secret keys for the FHE, keeps the secret key to herself, and sends to the server the public key and encrypted data elements $\llbracket array \rrbracket = (\llbracket x_1 \rrbracket, \dots, \llbracket x_m \rrbracket)$. The client can then repeatedly issue encrypted retrieval queries $\llbracket \ell \rrbracket$ for the server to homomorphically evaluate f on $\llbracket array \rrbracket$ and $\llbracket \ell \rrbracket$, sends the encrypted outcome

$$\llbracket \mathcal{L} \rrbracket = \{(\llbracket i \rrbracket, \llbracket x_i \rrbracket) \mid \text{ISMATCH}(x_i, \ell) = 1\}$$

to the client which decrypts it and obtain the desired result \mathcal{L} .

Clearly this protocol has a single-server, single round, low bandwidth communication (proportional only to the encrypted input and output, while being independent of the complexity of the report functionality), and involves no preprocessing of the cleartext data. Moreover, security easily follows from the semantic security of the FHE scheme.

The problem is that known polynomials f for the report functionality have high degree $\Omega(m \cdot d)$, for m the number of data records and d the degree of the polynomial realizing the matching condition `ISMATCH`; See Appendix A. This is too slow with current FHE candidates and implementations.

The motivation for our work is to answer affirmatively the following question: **Is there an *efficient* protocol for Secure report on FHE encrypted data and query?**

1.1 Our Contribution

We propose a novel paradigm, named *coresets for homomorphic encryption* (CHE), that forges a link between data summarization techniques (known as “coresets” or “sketches”) and secure computation on FHE encrypted data: Whereas it is not clear if low-degree polynomials exist for many classic tasks, such as search and report, our paradigm suggests instead for the server to compute a coreset for these problems for the goal of gaining a dramatic reduction in server’s complexity, while essentially conserving the communication and client’s complexity; See Section 2.

We demonstrate the strength of our CHE paradigm by applying it on the Secure report problem to dramatically reduce the server’s complexity: reducing the degree of the polynomial evaluated by the server to be as low as the degree d of the polynomial realizing ISMATCH (improving over degree $\Omega(m \cdot d)$ in the direct polynomial). We stress that the degree in our solution is independent of the number of data records m . See Section 3.

Elaborating on the above, our Secure report solution improves over the direct polynomial in both the server’s degree and overall multiplications by multiplicative factors $\Omega(m)$ and $\Omega(|\mathcal{L}| \cdot m)$ respectively. Conversely, our Secure report exhibit a slight degradation in the communication and client’s complexity (by a factor $|\mathcal{L}|$ and $(|\mathcal{L}| + \log m)^{O(1)}$ respectively); this is a minor degradation when the number of matches $|\mathcal{L}|$ is moderate, say, logarithmic in m (whereas for large $|\mathcal{L}|$ transmitting all these matches would typically be undesirable, due to the high communication it entails even without the discussed overhead).

Experimental results. We implemented our Secure report solution in a C++ library building on HELib [Hal13] FHE implementation, and ran extensive experiments. Our experimental results show that we can run our Secure report on *billions* of database records in a minute, on a single Amazon AWS server (compared to less than a hundred-thousand records with the direct polynomial, on the same hardware); See Fig. 1 and Section 5.

Our code is provided for the community in an open source library of Coresets for Homomorphic Encryption (CHElib) [AFS18a], to reproduce our experiments, to extend our results for real-world applications, and for practitioners at industry or academy that wish to use these results for their future papers or products.

Future research. We expect our coresets for homomorphic encryption paradigm to be further employed for degree reduction in the context of other secure computation tasks, possibly together with classical degree reduction techniques such as low-degree approximation [Raz87, Smo87] and randomized polynomials [IK00]. Indeed, our paradigm was already employed in a follow-up work [AFS18c] to solve the secure search problem of returning the first-match (i.e., returning (i, x_i) for $i = \min \{i \in [m] \mid \text{ISMATCH}(x_i, \ell) = 1\}$), motivated by use-cases where the number of matches to report is too large.

1.2 Related Works

Secure search has been extensively studied employing a variety of cryptographic tools leading to solutions with versatile properties.

Most relevant scenario: Secure search on FHE encrypted data. The scenario most relevant to our settings is secure search on FHE encrypted data [Gen09b], when focusing on protocols with single-server, single-round, low-communication proportional only to the size of encrypted input and output, and with no pre-processing on cleartext data. We note that disallowing pre-processing necessitates a linear scan of the data, even if we were to compute on cleartext data and query.

Private Information Retrieval (PIR) [CGKS95] enables a client to retrieve a data item from a server holding a data array, while hiding from the server the client’s query, the

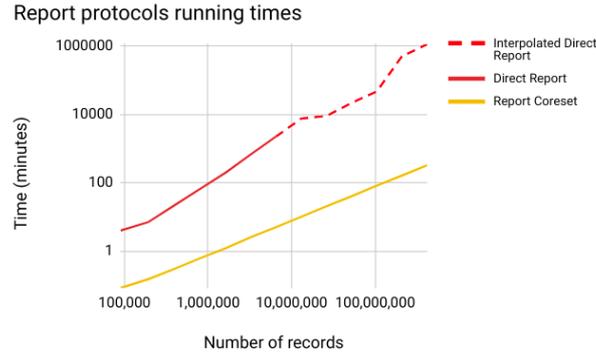


Figure 1: A log-scaled graph showing the server’s running time (y -axis) on a single machine of Amazon’s cloud, for different vector size (x -axis) of Secure Report (Protocol 3) over encrypted vector. In these experiments the vectors had 40 matching elements to be reported. The yellow line marks the running time of our report algorithm. The red line marks the running time of the direct report algorithm, where a solid line marks experiments that were made and a dashed line marks our interpolation based on our analysis.

retrieved data item, and the query and access patterns. The PIR query is typically the item’s index $i \in [m]$. Nevertheless, constructions of PIR on FHE encrypted data [Gen09a, BV11, DSH14] are easily extended for retrieving encrypted keywords provided the keyword *uniquely identifies* at most a single data item (uniqueness constraint). In contrast, in our work we do not require the query to be a unique identifier to a single data item.

Secure search on encrypted data, eliminating the aforementioned uniqueness constraint, was recently achieved [AFS18c], providing secure search on FHE encrypted hiding the content of data items and queries, as well as the query and access patterns. In [AFS18c] the retrieved item is the *first matching item* in the array; whereas retrieving subsequent items requires further interaction, another round for each fetch-next item to be retrieved. In contrast, in our work we retrieve all matches in a single round.

Other scenarios: relaxed settings. Numerous other related works on secure search are relevant, if relaxing the above settings to allow, for example, leaking information on the access or search patterns, searching on cleartext data, $k > 1$ servers, $R > 1$ rounds, communication bandwidth growing with database size m or with complexity of the computed functionality, allowing pre-processing the cleartext data, etc.

When allowing information leakage, as in leaking access and query patterns or revealing order, searchable and structured encryption (SE and STE) [SWP00, CK10] allow for highly efficient construction approaching plaintext search time. Furthermore, these constructions typically allow pre-processing the cleartext data for achieving sub-linear search time. Yet, the inherent information leakage has been often exploited to obtain the content of query and data items [CGPR15, ZKP16, IKK12, AAG17, GAYBL17, GMN⁺16, GSB⁺17, PW16].

When searching on unencrypted data (in contrast to encrypted data in our settings), indexing can be used to enforce uniqueness to allow using the aforementioned PIR on FHE encryption data protocols [qCDD⁺16, RVVV17]; this indexing however incurs considerable time and memory overhead. Furthermore, if allowing pre-processing of cleartext data on top of searching on unencrypted data, then we could employ PIR-by-Keywords [CGN97] techniques to get sub-linear complexity.

When allowing $k > 2$ non-colluding servers, fast private queries on public (unencrypted) data systems were presented in [BGH⁺13, YYG⁺17].

When allowing $R > 1$ rounds, $\log m$ communication overhead, and a client maintaining and updating state – Oblivious RAM (ORAM) [Gol87] enable hiding data access patterns with sub-linear server time (breaking the linear scan barrier of our settings).

When allowing communication bandwidth to grow with the time to compute the search functionality, classical “pre-FHE” secure two-party computation (2PC*) techniques can be employed, such as garbled circuits and secret sharing [Yao86, GMW87b].

Likewise, when allowing communication bandwidth to grow with the number of records m , secure pattern matching (SPM) shows how to efficiently compute ISMATCH on FHE encrypted data item. Query [YSK⁺13, CKL15, LLN15, CKK16, TJW⁺16, KLL⁺17] can be employed by the server to compute the encrypted length m indicator vector indicating the desired matches, and transmitting this entire vector of length m to the client. This is clearly impractical for large databases where m is large.

Table 1: Comparison to related works for secure data retrieval on FHE encrypted data with single-server, single-round protocols, and no pre-processing of cleartext data. Rows correspond to approaches: Secure Pattern Matching [YSK⁺13, CKL15, LLN15, CKK16, TJW⁺16, KLL⁺17] (SPM), Private Information Retrieval [BV11, DSH14, qCDD⁺16, RVVV17](PIR), Secure Search [AFS18c], the Direct polynomial (see Appendix A), and our Report-Coreset approach (CHELib). Columns correspond to desired properties: sub-linear communication and client’s complexity, low degree for the polynomial homomorphically-evaluated by the server, handling query without uniqueness constraint, returning all matches (rather than a single match), and returning both index and record (rather than only the index). The parameter d is the degree of the specified matching predicate ISMATCH; and m is the number of records in the data array.

Approaches	Sub-linear Client	Server’s degree	Non-unique Identifiers	Returns All Matches	Returns Records
SPM	×	d	✓	✓	×
PIR	✓	d	×	×	✓
Secure Search	✓	$d \cdot \log^3 m$	✓	×	✓
Direct method	✓	$d \cdot m$	✓	✓	✓
CHELib	✓	d	✓	✓	✓

2 New Paradigm: Coresets for Homomorphic Encryption

In this section we give an overview of how we propose to employ coresets, sketches and group-testing for reducing the complexity of computing on FHE encrypted data.

Coreset is a data summarization C of a set P of items (e.g. points, vectors or database records) with respect to a set Q of queries (e.g. models, shapes, classifiers, points, lines) and a loss function f , such that $f(P, q)$ is approximately the same as $f(C, q)$ for every query $q \in Q$. The goal is to have provable bounds for (i) the size of C (say, $1/\varepsilon$), (ii) the approximation error (say, $\varepsilon \in (0, 1)$) and (iii) construction time of C given (P, Q, f) . We can then run (possibly inefficient) existing algorithms and heuristics on the small coreset C , to obtain provably approximated solution for the optimal query (with respect to f) of the original data.

Sketches are a special type of coresets, where given an $m \times d$ matrix P and a “fat”, $s \times m$ matrix $A_{s,m}$, with $s \ll m$, the result $C = A_{s,m}P$ is an $s \times d$ (“sketch”) vector. Many problems can be solved on the sketch C instead of the long vector P , by designing

a corresponding (sketch) matrix $A_{s,m}$. For example, if the entries of $A_{s,m}$ are random standard Gaussian variables, then we can approximate the k -means problem on C by the Johnson-Lindenstrauss Lemma; see [BZD10]. In this paper we use sketches in the context of compressed sensing, or more precisely, Group Testing [DHH00].

The coreset is a paradigm in the sense that its exact definition, structure and properties change from paper to paper. The term “coreset” was coined by Agarwal, Har-Peled, and Varadarajan [AHPV05] and originally used for optimization in computational geometry (e.g. [Phi16, Cla10, BMO⁺11]). Since then coresets were used in fields such as machine learning (e.g. [HCB16, BFL⁺17]), numerical algebra [W⁺14], graph theory (e.g. [CLMS13]), time series [IKM00, RLNZ09, GKMS01] streaming [CG05, GKS13, CK08, CG07] and pattern recognition [CZ17, CSZ05, GM12, CGHJ12, HK15].

Coresets for Homomorphic Encryption (CHE). In this paper we suggest a new paradigm that is inspired by the coreset paradigm, but for a very different application and context. Instead of running the complete secure algorithm on the server side and encrypted data, the server computes only (encrypted) coresets for the problem at hand and sends them back to the client. Since the coresets are small, communicating them to the client, decrypting them, and decoding the desired result from the coreset is relatively fast and require only little additional time on the client side.

In this paper we introduce the first two CHEs (cf. Fig. 2): one for reporting non-zero values in an array, and one for reducing the ring size from r to $O(\log r)$; see Section 4. These two coresets are then combined to give our secure solution to the report problem; see Theorem 1. The report coreset uses a sketch matrix to communicate the sparse binary vector of the reported indices to the client.

The key property that we show in this paper is that computing coresets on the server side securely (i.e., via polynomials) dramatically reduces its computation by reducing the degree of the polynomial from at least linear to logarithmic on the input size (CRT coreset) or even constant (report coreset).

Informally, in this paper a function S is an efficient coreset construction scheme for a problem if it has the following pair of properties.

- (i) $S(P) \in \mathbb{R}^k$ is a short vector that can be computed efficiently. Usually $k = (\log m)^{O(1)}$ for an input vector P of length m .
- (ii) $S(P)$ is a coreset or a “sketch” vector in the sense that the desired information, in a problem dependent sense, can be extracted efficiently and exactly from $S(P)$, without having access to (the large) original vector.

Another significant difference from traditional coreset and sketches papers is that the coresets in this paper are exact in the sense that the data reduction does not introduce any additional error ε .

Group Testing meets FHE for report coreset. Our report coreset in Protocol 3 uses a modern result from the field of Group Testing [DHH00]. To our knowledge, this is the first application of Group Testing to FHE. The motivation is to return to the client the very large (m bits) but sparse binary (indicator) vector χ of the desired s indices that were computed in Line 2a of our report coreset (Protocol 3) by communicating a number of bits that is only poly-logarithmic in m . Specifically, entries $\chi(i) = \text{ISMATCH}(x_i, \ell)$ indicate whether data elements x_i match the query ℓ .

The server however sees only the ciphertexts $\llbracket \chi(i) \rrbracket$, encrypting the indicators $\chi(i)$, that the server obtains by homomorphically-evaluating the ISMATCH polynomial on encrypted elements $\llbracket x_i \rrbracket$ and query $\llbracket \ell \rrbracket$. Since the vector χ is encrypted, an efficient algorithm must be realized by a low-degree polynomial, so simple comparison of each entry to 1 is impossible

by such a polynomial. Instead, Group Testing suggests to extract the set of 1's entries in χ from a small number of predefined sums on subsets entries in χ .

Our Report Coreset multiplies the (encrypted) indicator vector χ by a (sketch) matrix $A_{s,m}$ of size $k \times m$ where k is polynomial in s and $\log m$, i.e., we compute few *linear* combinations $A_{s,m} \cdot \chi$ of the entries of χ . In fact, since the sketch matrix $A_{s,m}$ is public, computing these linear combinations requires only an additively homomorphic encryption. Indyk, Ngo and Rudra [INR10] suggested such a binary matrix $A_{s,m}$ with $k = O(s^2 \log m)$ rows, so that our χ vector can be decoded from $A_{s,m} \cdot \chi$ in time that is also polynomial in s and $\log m$ on the client side.

Unfortunately, unlike χ , the vector $A_{s,m} \cdot \chi$ is not a binary vector; moreover, correct decoding of the sketch assumes this product is computed over the reals. In contrast, in the context of FHE, computation is typically within a finite ring, e.g., computation modulo p for a prime p , and not over the reals.¹ Taking a sufficiently large p (e.g., $p > m$ if $A_{m,s}$ is binary), would have resolved the correctness issue, albeit introducing an undesirable complexity increase, as larger p leads to a general slowdown in the homomorphic operations as well as size inflation of the keys and ciphertexts. To aggravating the problem further, if $\text{ISMATCH}(x_i, \ell)$ uses Fermat's Little Theorem (e.g., for testing equality), then the degree of the polynomial for computing χ in Line 2a of Protocol 3 grows at least linearly with p , resulting in a high-degree polynomial.

Our CRT Coreset offers a solution that guarantees to resolve both the correctness and complexity issues by executing the homomorphic evaluations in multiple rings (in parallel), where the CRT coreset is the tuple of results in all the rings. Our analysis follows from the Chinese Remainder Theorem (CRT) showing that computing modulo $O(\log m)$ small primes $p = \tilde{O}(\log m)$, not only guarantees the efficiency and low-degree of the polynomials we evaluate, but also ensures that the correct outcome can be efficiently decoded from the CRT coreset.

In addition, the value s of the sparsity of χ is required to compute the sketch matrix $A_{s,m}$ but unknown to the server. These and other challenges are handled by the algorithms of our system.

3 Problem Statement and Main Result

In this section we formally define the Secure report problem and state our main theorem.

3.1 The Secure Report Problem

Our main example application for using coresets in the service of FHE is to solve the following report queries problem efficiently on a secure database.

For simplicity of the presentation, we first address the problem of reporting the list \mathcal{I} consisting of all *indices* $i \in \{1, \dots, m\}$ for the matching elements, i.e., all x_i so that $\text{ISMATCH}(x_i, \ell) = 1$. The extension for reporting the list \mathcal{L} of index and value pairs (i, x_i) , for $i \in \mathcal{I}$, is simple and requires no increase in the round complexity; see Section 4.2.2.

We focus here on key components of the definition, deferring to Section 3.3 some implementation details regarding the input representation and compatibility requirements for the ISMATCH polynomial and the FHE.

Definition 1 (Report problem). Let m be a positive integer, let \mathcal{M} and \mathcal{Q} the domains of data-elements and queries respectively, $\text{ISMATCH}: \mathcal{M} \times \mathcal{Q} \rightarrow \{0, 1\}$ a polynomial specifying the matching condition. In the *Report problem*, on input *array* $= (x_1, \dots, x_m) \in \mathcal{M}^m$ and

¹An exception is a recent FHE scheme computing over the reals, albeit inherently introducing rounding errors in each computation [CKKS17].

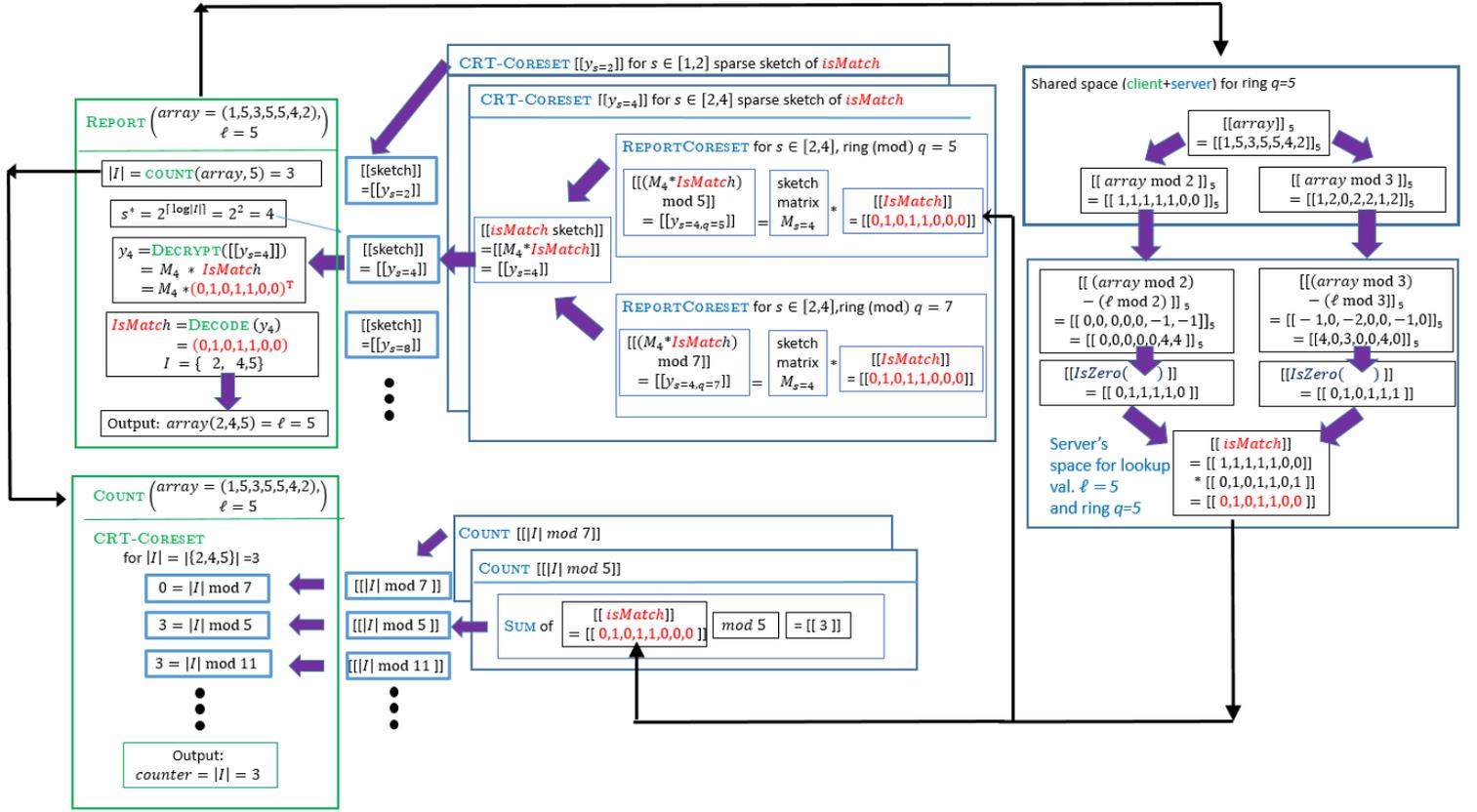


Figure 2: Overview of the suggested system for reporting the indices $I = \{2, 4, 5\}$ of the lookup value $\ell = 5$ in the array $(1, 5, 3, 5, 5, 4, 2)$. Green and blue titles correspond to algorithms that run on the client's and server's side, respectively. The rightmost boxes represent encrypted saved versions of the array. After the client calls REPORT on the right, a version is computed from the shared memory (top right) to each of the few ring values (only $q = 5$ is shown). ISMATCH (bottom right) contains binary results for each ring. The client computes the number $|I|$ of occurrences (bottom left) of ℓ in the array, by applying the CRT-CORESET on the COUNT results that are computed by the server. It then computes the closest power of 2 for $|I|$ (i.e., $s^* = 4$), decrypts the corresponding sketched vector ($y_s = 4$), and ignores the other sketches. The client then computes I efficiently from y_s . Each sketch vector is computed by the server using CRT-Coreset which computes, for every ring's size q , its REPORTCORESET (in the middle) that uses $\text{ISMATCH} = \text{array} - (\ell, \dots, \ell)$ for q (bottom right).

query $\ell \in \mathcal{Q}$, the output is the set \mathcal{I} of all entries in array matching ℓ :

$$\mathcal{I} = \{i \in \{1, \dots, m\} \mid \text{ISMATCH}(x_i, \ell) = 1\}.$$

In the context of our coreset for FHE paradigm, we relax the report problem to allow outputting a short sketch for \mathcal{I} , named, *report-coreset*, on which the client applies an efficient decoding algorithm to obtain \mathcal{I} . We assume an upper-bound $s \in [0, m]$ on the number of matches is given; See Section 4.2.3 for treatment of the case s is unknown. Looking ahead, for efficiency we'd like the report-coreset to be short and the decoding time to be efficient, specifically, we'd like them to be polynomial in s and $\log m$.

Definition 2 (Report coreset). Let m , $array$, ℓ and \mathcal{I} be as in Definition 1. Let $s \in [0, m]$ be an upper bound on the number of matches $|\mathcal{I}|$. A vector $y \in \mathbb{Z}^k$ is an s -report coreset for $(array, \ell)$ if, given only y , we can decode (compute) the set \mathcal{I} .

The secure-report problem is the report-problem where the input $(array, \ell)$ and report-coreset output y are replaced by their encryption.

Definition 3 (Secure-report coreset problem). Let m , $array$, ℓ , \mathcal{I} , s and y be as in Definition 2. In the *secure-report problem*, on input ciphertexts $(\llbracket array \rrbracket, \llbracket \ell \rrbracket)$, the output is a ciphertext $\llbracket y \rrbracket$ for an s -report coreset y for $(array, \ell)$.

Here $\llbracket array \rrbracket = (\llbracket x_1 \rrbracket, \dots, \llbracket x_m \rrbracket)$ is the entry-by-entry encryption of the data $array$, and $\llbracket \ell \rrbracket$ is the encryption of the query ℓ , where encryption is with an FHE scheme specified as parameters to the problem.

A secure-report protocol via coresets. The usage scenario is that the client first generates keys and uploads the encrypted data $\llbracket array \rrbracket$ to the server; the client can then issue repeated encrypted report queries $\llbracket \ell \rrbracket$, for which the server homomorphically computes and sends to the client the s -report coreset $\llbracket y \rrbracket$; the client can then decrypt and decode to obtain the desired outcome \mathcal{I} .

In details, for m , $array$, ℓ , \mathcal{I} , s and y as in Definition 2, the protocol for Secure report via the secure-report coreset involves two parties, a client and a server with shared input parameters (see below); a client's input $(array, \ell)$ for $array$ a length m array of data-elements and ℓ a query (or lookup value); the client's output the list of matching elements \mathcal{I} ; and no input or output for the server. The steps of the protocol are as follows First, in an initialization and data upload phase, the client generates keys for the FHE specified in the shared parameter, encrypts $array$ entry-by-entry, and uploads the corresponding ciphertexts $\llbracket array \rrbracket$ to the server. Next, in a search-phase, the client can repeatedly encrypt report queries ℓ and send the encrypted queries $\llbracket \ell \rrbracket$ to the server. The server then homomorphically computes s -report coreset $\llbracket y \rrbracket$, and sends it to the client. The client then decrypts and decodes to obtain the desired outcome \mathcal{I} .

The shared input parameters for our protocol (Protocol 3) are λ , \mathcal{E} , p , \mathcal{M} , \mathcal{Q} , ISMATCH, m and s defined as follows: λ is the security parameter, p the smallest prime larger than m , $\mathcal{E} = (\text{Gen}, \text{Enc}, \text{Dec}, \text{Eval})$ an FHE scheme that can compute homomorphic evaluations in $GF(p)$, \mathcal{M} the data elements space, \mathcal{Q} the queries-space, ISMATCH a polynomial realizing the matching criterion $\text{ISMATCH}: \mathcal{M} \times \mathcal{Q} \rightarrow \{0, 1\}$, m the number of data elements (for an upper-bound pad with zeros), and s an upper-bound on the number of matching data elements.

In our extension for the case that an upper bound s on the number of matches is unknown (Section 4.2.3), we omit s from the shared parameters.

In our optimization utilizing our CRT-coreset to reduce the ring size for the homomorphic-evaluations from linear to logarithmic in m (Section 4.2.2), we replace the large prime p by a small set $\mathcal{P} = \{p_1, \dots, p_\kappa\}$ of small primes. Specifically, we set \mathcal{P} to consist of the first κ primes larger than $\log_2 m$ for $\kappa = \lceil \log m / \log \log m \rceil$.

Remarks and extensions. For simplicity of the presentation we assume the entire data $array$ is uploaded in a single execution of the protocol, together with the keys generation. This can be easily modified. The keys can be generated at a separate time, possibly by a separate entity distributing the public key to the clients who will upload the data (data-sources) and to the server, and distributing the secret key to the clients who will issue search queries (search-clients). The upload can occur over time and from multiple data-sources, gradually uploading new encrypted data records for the server to append to its array of ciphertexts.

The question motivating this work, which we affirmatively answer is as follows:

Can we solve the secure-report problem via a server homomorphically evaluating of a polynomial of degree independent of the number of data elements m , and where the communication and client's is complexity sub-linear in m ?

3.2 Main Theoretical Result

The main result of this paper is the first solution for the Secure report-coreset problem on FHE encrypted data (see Definition 3) via polynomial of degree independent of the number m of elements in the data *array*. We employ this polynomial to give a secure-report protocol via coreset, in a single-server, single-round protocol with no pre-processing of cleartext data (no sorting, indexing, or search oriented data-structures are used). The protocol attains our complexity goals: The server homomorphically evaluates a polynomial of degree independent of the number m of elements in the data *array*, and the communication and client's complexity is sub-linear (in fact, poly-logarithmic) in m and quadratic in the upper-bound s on the number of matches.

Theorem 1 (basic settings, see extension below). *There exists a protocol for secure-report via coresets satisfying the following, when executed on client's input array $= (x_1, \dots, x_m) \in \mathcal{M}^m$ and $\ell \in \mathcal{Q}$, and shared parameters $(\lambda, p, \mathcal{E}, \mathcal{M}, \mathcal{Q}, \text{ISMATCH}, m, s)$, as defined in Section 3.1:*

- **Structural requirements:** *The protocol is a single-server, single-round protocol, with no preprocessing on cleartext data.*
- **Correctness:** *If $|\mathcal{I}| \leq s$, then the client's output is $\mathcal{I} = \{i \mid \text{ISMATCH}(x_i, \ell) = 1\}$.*
- **Complexity:** *The complexity of the protocol is as follows, when denoting by d, μ the degree and overall multiplication in the ISMATCH polynomial specified in the shared parameters:*
 - *The server homomorphically-evaluates over $GF(p)$ a polynomial of degree d and overall multiplications $m \cdot \mu$ (for p the smallest prime larger than m).*
 - *The communication consists of $|\ell|$ ciphertexts sent from the client and $O(s^2 \log m)$ ciphertexts sent from the server (for $|\ell|$ denoting the representation length of ℓ).*
 - *The client's running-time is polynomial in $|\ell|, \lambda, s$ and $\log m$.*
- **Security:** *The protocol is secure against semi-honest adversary controlling the server.*

Our subsequent protocols offer some extension and optimizations and tradeoffs on the above theorem, as detailed next.

3.2.1 Extensions and Optimizations

The above theorem states our basic result that can be optimized and extended in several forms.

Outputting \mathcal{L} . The output can be extended to consist of (index, value) pairs for all matches, i.e., outputting

$$\mathcal{L} = \{(i, x_i) \mid \text{ISMATCH}(x_i, \ell) = 1\}.$$

This extension increased the client's run-time, communication complexity, and number of polynomials evaluated by the server by a factor w , for w the binary representation length of records x_i . There is no degradation however in the structural requirements or the complexity polynomials evaluated by the server.

Reducing the prime to logarithmic in m . The prime p can be reduced to $p = \tilde{O}(\log m)$ using our CRT-coreset. In this case we essentially repeat our protocol in parallel over distinct fields $GF(p_1), \dots, GF(p_\kappa)$ for $\kappa = \lceil \log_{\log_2 m} m \rceil$ and p_1, \dots, p_κ the first κ primes larger than $\log_2 m$. This incurs a $\kappa = o(\log_2 m)$ increase on the client's run-time, communication complexity, and number of polynomials evaluated by the server. There is no degradation however in the structural requirements, or the complexity of the polynomials evaluated by the server.

Handling unknown number of matches. If a sufficiently small upper-bound s on the number of matches is unknown we offer two incomparable alternatives.

The first alternative is to add a communication round in which the client and server engage in a secure computation protocol to compute the number of matches $|\mathcal{I}|$, and then executing our Secure report protocol setting the upper-bound s to be at least $|\mathcal{I}|$ (possibly $s > |\mathcal{I}|$ to not reveal the exact count). Computing the number of matches can be done in a single-round protocol, where the server evaluates a polynomial of degree d and overall multiplications $\tilde{O}(m \cdot \mu)$. Elaborating on the above, if allowing homomorphic evaluations over a large field $GF(p)$ for $p > m$, then counting the number of matches is straightforward (Section 4.2.3). To reduce the field size to logarithmic $p = \tilde{O}(\log m)$, we propose employing our CRT-coreset (Section 4.2.2).

The second alternative is request that the server homomorphically evaluates the report coreset for $s \in \{2^0, 2^1, 2^2, \dots, m\}$ (sending each resulting encrypted coresets as soon as its computation has completed). This does not increase the degree of the polynomials evaluated by the server's, but does have a high overall number of multiplications and servers communication complexity. Nonetheless, this approach might be relevant if the client can abort (preferably, covertly without the server knowing when the client aborts); this results in client's received communication and run-time that are still polynomial in λ , $|\ell|$, $|\mathcal{I}|$ and $\log m$, as desired.

3.3 Details: Compatibility Requirements

We next set some notation (below) and give the technical details on the compatibility requirements from ISMATCH (Section 3.3.1) and the requirements from the used FHE scheme (Section 3.3.2).

Notations. For a data record z or a lookup value, we denote by $|z|$ the binary representation length of z . For an integer m , we denote $[m] = \{1, \dots, m\}$. We follow the convention that array indexes start from 1, with entry i denoted $array(i)$. Vectors are column vectors unless stated otherwise. We assume the array size m is a power of two (otherwise pad with zero). Logarithms are in base 2 unless explicitly stated otherwise.

3.3.1 Compatibility Requirements for isMatch

We next elaborate on the compatibility requirements for ensuring that (a) the input representation is compatible to the message space of the FHE scheme, and (b) the provided polynomial for ISMATCH is compatible with both the input representation and the ring for the FHE evaluations.

The polynomial ISMATCH specifies the matching criteria by mapping pairs $(x_i, \ell) \in \mathcal{M} \times \mathcal{Q}$ of data element and query to a Boolean value $ISMATCH(x_i, \ell) \in \{0, 1\}$ accepting value 1 if-and-only-if the element x_i is a match to the query ℓ . To plug-into our protocol the polynomial ISMATCH, we make few requirements for compatible with the input representation, compatibility with the underlying FHE, and for security.

Compatibility of ISMATCH to the input representation mean that data element and query pair (x_i, ℓ) must be in the domain of ISMATCH. A few examples follow. When the

matching criterion is an equality-test polynomial $\text{ISMATCH}: \{0, 1\}^w \times \{0, 1\}^w \rightarrow \{0, 1\}$ (Section 4.2.5), we require that data elements x_i and queries (or lookup value) ℓ are specified in binary representation of length w . When the matching criterion is an equality-test polynomial $\text{ISMATCH}: \{0, \dots, p-1\}^w \times \{0, \dots, p-1\}^w \rightarrow \{0, 1\}$ over $GF(p)$ for a prime $p > 2$ (Section 4.2.5), we require that data elements x_i and query (or lookup value) ℓ are specified using w digits in base p . When the matching criterion is a range query, the elements and queries come from different domains, where for example elements may be integer values x_i and queries are a pair of integer endpoints (a, b) so that $\text{ISMATCH}(x_i, (a, b)) = 1$ if-and-only-if x_i is in the ranger (a, b) . In general, for \mathcal{M} and \mathcal{Q} the spaces of data elements and queries respectively, the matching criterion is a function $\text{ISMATCH}: \mathcal{M} \times \mathcal{Q} \rightarrow \{0, 1\}$.

Compatibility of ISMATCH with the underlying FHE means the following. First, the data elements space \mathcal{M} and queries space \mathcal{Q} must be contained in the message space for the FHE so that they can be encrypted. Second, the polynomial provided for ISMATCH must be over the same ring where the FHE computations are executed so that it can be homomorphically evaluated.

For security, we require that \mathcal{M} consists of equal length data elements and \mathcal{Q} consists of equal length queries, because the encryption does not hide length information; use padding as needed.

3.3.2 Compatibility Requirements for our Black-Box usage of FHE

Our protocols employ a fully (or, leveled) homomorphic encryption (FHE) in a black-box fashion: we require only a black-box usage of the standard algorithms for FHE (key generation, encryption, decryption, and evaluation). The only requirement we make on the scheme is that we can choose as a parameter the *plaintext modulus* to be a prime number p of our choice, so that the homomorphic operations are additions and multiplications modulo p . This is the case in many of the FHE candidates, for example, [BGV12]. For security of our scheme we require that the FHE scheme is semantically secure.

To emphasize the plaintext modulus p we use the following notations for the standard algorithms specifying an FHE scheme $\mathcal{E} = (\text{Gen}, \text{Enc}, \text{Dec}, \text{Eval})$:

- Gen is a randomized algorithm that takes a security parameter λ as input and a prime p , and outputs a public key $pk_p = (p, pk)$ and a secret key $sk_p = (p, sk)$ for plaintext modulus p , denoted:

$$(pk_p = (p, pk), sk_p = (p, sk)) \leftarrow \text{Gen}(1^\lambda; p).$$

- Enc is a randomized algorithm that takes pk_p and a plaintext message msg , and outputs a ciphertext $\llbracket msg \rrbracket_p$ for plaintext modulus p , denoted:

$$\llbracket msg \rrbracket_p \leftarrow \text{Enc}_{pk_p}(msg).$$

- Dec is an algorithm that takes sk_p and a ciphertext $\llbracket msg \rrbracket_p$ as input, and outputs a plaintext msg' , denoted:

$$msg' \leftarrow \text{Dec}_{sk_p}(\llbracket msg \rrbracket_p).$$

Correctness is the standard requirement that $msg' = msg$.

- Eval is a (possibly randomized) algorithm takes pk_p , a polynomial $f(x_1, \dots, x_t)$, and a tuple of ciphertexts $(\llbracket m_1 \rrbracket_p, \dots, \llbracket m_t \rrbracket_p)$, and outputs a ciphertext c , denoted:

$$c \leftarrow \text{Eval}_{pk_p}(f, \llbracket m_1 \rrbracket_p, \dots, \llbracket m_t \rrbracket_p).$$

Correctness is the requirement that decryption would return the message resulting from evaluating (modulo p) the polynomial $f()$ on inputs m_1, \dots, m_t ,

$\text{Dec}_{sk_p}(\text{Eval}_{pk_p}(f, \llbracket m_1 \rrbracket_p, \dots, \llbracket m_t \rrbracket_p)) = f(m_1, \dots, m_t) \bmod p$. *Semantic security* implies that the resulting ciphertext c is computationally indistinguishable from a fresh ciphertext $\llbracket f(m_1, \dots, m_t) \rrbracket_p$.

Shorthand notations we use in this writeup are detailed next. (a) We encrypt messages $m = (m_1, \dots, m_t)$ digit-by-digit, and encrypt arrays entry-by-entry. We abuse notation by denoting the resulting tuple of ciphertexts by $\llbracket msg \rrbracket_p = (\llbracket m_1 \rrbracket_p, \dots, \llbracket m_t \rrbracket_p)$ for $\llbracket msg(i) \rrbracket_p \leftarrow \text{Enc}_{pk_p}(m_i)$ and $\llbracket array \rrbracket_p = (\llbracket array(1) \rrbracket_p, \dots, \llbracket array(m) \rrbracket_p)$. (b) When the modulus p is for a ciphertext $\llbracket m \rrbracket_p$ clear from context we omit it and write $\llbracket m \rrbracket$. (c) When the public key pk_p is clear from the context we use a more natural presentation of homomorphic evaluation:

4 The Secure Report via Coreset Protocol

In this section we give the details of our upload and report protocol, starting with the basic version of the protocols and continuing with the extensions, optimizations and trade-offs.

4.1 The Basic Protocol

In this section we give the details of our upload and report protocol (basic version), and analyze the protocol to prove Theorem 1; See a protocol summary in Fig. 3.

The Secure report protocol (basic version). The protocol begins with an initialization and data upload phase. In this phase the client generates keys for the FHE, encrypts the data elements x_1, \dots, x_m , and sends the encrypted data and the public key to the server. The FHE is initialized here with plaintext modulus p so that $p > m$ for m the provided upper-bound on the number of data elements.

Next the protocol enters the report queries phase. In this phase the client can repeatedly issue encrypted report queries ℓ . For each such encrypted query $\llbracket \ell \rrbracket$, the server does the following. First the server homomorphically evaluates the specified matching polynomial ISMATCH to obtain the encrypted indicator vector $\chi = (\chi(1), \dots, \chi(m))$, where

$$\llbracket \chi(i) \rrbracket_p := \text{ISMATCH}(\llbracket x_i \rrbracket_p, \llbracket \ell \rrbracket_p)$$

for all $i \in [m]$. Next, the server homomorphically multiplies χ by an (s, m) -sketch matrix $A_{s,m}$, to obtain and send to the client the encryption of $y = A_{s,m} \cdot \chi$:

$$\llbracket y \rrbracket_p := A_{s,m} \cdot \llbracket \chi \rrbracket_p.$$

The client decrypts using her secret key to obtain y , and then decodes y using the decoding algorithms for the sketch to obtain χ .

Proof of Theorem 1. The protocol is a single-server, single-round protocol, that performs no pre-processing on the cleartext data.

Correctness is argued as follows. First observe that χ satisfies that $\chi(i) = 1$ if-and-only-if $\text{ISMATCH}(x_i, \ell) = 1$. This implies that χ is a binary length m vector with at most s positive entries (by the premise that the number of matches is upper-bounded by s). By the correctness of the sketch $A_{s,m}$, we know that when computing $y = A_{s,m} \chi$ over the real number, the sketch decoding algorithm returns the set of entries where χ is positive, that is, the correct output \mathcal{I} . In our case however computation is over $GF(p)$. Nonetheless, our choice of $p > m$ ensures that no overflow occurs, and the obtained result is identical to the result when computing over the reals. Namely, the output is the desired set $\mathcal{I} = \{i \in [m] \mid \text{ISMATCH}(x_i, \ell) = 1\}$.

Complexity analysis follows. The server homomorphically-evaluates over $GF(p)$ the composition of two polynomials: the ISMATCH polynomial and the linear polynomial that computes the product of a plaintext matrix and an encrypted vector. The former has degree d , whereas the latter add nothing to the degree because the linear combinations

specifying the entries $y(i)$ are computed using only homomorphic addition (by summing the encrypted entries $ind(j)$ for all j where the $A_{s,m}(i,j) = 1$). Thus, the server computes a polynomial of degree d (the degree of ISMATCH). The overall number multiplications is μ multiplications to compute ISMATCH on each of the m data elements, resulting in a total of $m \cdot \mu$ multiplications. The communication consists $\llbracket \ell \rrbracket$ and $\llbracket y \rrbracket$. The latter consists of $|\ell|$, the latter consists of a ciphertext for each entry of y , that is, $O(s^2 \log m)$ ciphertexts as the number of rows in $A_{m,s}$. The client's running-time is the time to encrypt and decrypt the sent and received ciphertexts respectively, plus the time to decode y . The time to encrypt and decrypt is polynomial in the number of ciphertexts and the security parameter λ , the decoding time is polynomial in s and $\log m$. The overall running time is therefore polynomial in $|\ell|$, λ , s and $\log m$.

Security against a semi-honest adversary controlling the server follows from the semantic-security on the underlying FHE scheme, because the entire protocol can be simulated by an efficient simulator oblivious of the client's input and output. \square

Shared Parameters: The security parameter λ , the FHE scheme $\mathcal{E} = (\text{Gen}, \text{Enc}, \text{Dec}, \text{Eval})$, a prime $p > m$, the records space \mathcal{M} and lookup values space \mathcal{Q} , a polynomial realizing the matching criterion $\text{ISMATCH}: \mathcal{M} \times \mathcal{Q} \rightarrow \{0, 1\}$, an upper-bound on the number of data records m , and an upper-bound on the number of matches $s \in [0, m]$.

Inputs: The client's input is a data *array* $= (x_1, \dots, x_m) \in \mathcal{M}^m$ and queries $\ell \in \mathcal{Q}$; The server has no input.

Outputs: The client's output is the set $\mathcal{I} = \{i \in [m] \mid \text{ISMATCH}(x_i, \ell) = 1\}$; The server has no output.

Initialization and upload phase The client does the following:

- Generate keys $(pk_p, sk_p) \leftarrow \text{Gen}(1^\lambda; p)$.
- Encrypt the data elements $\llbracket x_i \rrbracket_p \leftarrow \text{Enc}_{pk_p}(x_i)$ for all $i \in [m]$.
- Send to server pk_p and $(\llbracket x_1 \rrbracket_p, \dots, \llbracket x_m \rrbracket_p)$.

Report queries phase:

1. To issue each report query ℓ , the client encrypts the query $\llbracket \ell \rrbracket_p \leftarrow \text{Enc}_{pk_p}(\ell)$ and sends $\llbracket \ell \rrbracket_p$ to the server.
2. The server then does the following:
 - (a) Compute $\llbracket \chi(i) \rrbracket_p := \text{ISMATCH}(\llbracket x_i \rrbracket_p, \llbracket \ell \rrbracket_p)$ for all $i \in [m]$.
 - (b) Compute $\llbracket y \rrbracket_p := A_{s,m} \cdot \llbracket \chi \rrbracket_p$ for $A_{s,m}$ an (s, m) -sketch matrix (see Section 2) and $\llbracket \chi \rrbracket_p = (\llbracket \chi(1) \rrbracket_p, \dots, \llbracket \chi(m) \rrbracket_p)$; send $\llbracket y \rrbracket_p$ to the client.
3. The client decrypts $y \leftarrow \text{Dec}_{sk_p}(\llbracket y \rrbracket_p)$ and decodes y (see Section 2) to obtain the output \mathcal{I} .

Figure 3: Secure report protocol via report-coreset.

4.2 Extensions and Optimizations

We discuss extensions and optimizations of our protocol for reporting both value and index; handling the case that s is unknown; reducing the ring size for the homomorphic evaluation from linear to logarithmic in m ; and handling dynamic data management.

4.2.1 Reporting values on top of indexes (\mathcal{L})

To report the matching elements values $x_i \in \{0, 1\}^w$ (\mathcal{L}) on top of their indexes (\mathcal{I}) we do the following. The server homomorphically evaluates $w + 1$ sketches, instead of a single sketch: One sketch for χ (this is the same as in the basic protocol), plus w sketches – a sketch for each digit $j \in [w]$ in the binary representation of the matching data elements.

The sketch for the j -th digit is computed as follows. Denote the bits in the binary representation of x_i by $x_i(1), \dots, x_i(w)$ (for $i \in [m]$). Consider the vector consisting of all the j -th bits of the data elements:

$$allData_j = (x_1(j), \dots, x_m(j))$$

Observe that the server implicitly holds the encrypted version of $allData_j$ (encrypted entry-by-entry). To compute the sketch for the j -bit of the matching data elements the server first homomorphically maps $allData_j$ to the s -sparse vector $sparseData_j$ whose i -th entries ($i \in [m]$) are the product of $\chi(i)$ and $allData_j(i)$:

$$\llbracket sparseData_j(i) \rrbracket_p := \llbracket \chi(i) \rrbracket_p \cdot \llbracket allData_j(i) \rrbracket_p.$$

Next, the server homomorphically evaluates the sketch for $sparseData_j(i)$ by left-multiplying its encrypted version by the cleartext binary matrix $A_{s,m}$

$$\llbracket sketch_j \rrbracket_p := A_{s,m} \cdot \llbracket sparseData_j \rrbracket$$

The server then sends to the client all the encrypted sketches $\llbracket y \rrbracket_p$ and $\llbracket sketch_j \rrbracket_p$ for all $j \in [w]$. The client decrypts and decodes $\llbracket y \rrbracket_p$ to obtain \mathcal{I} , and decrypts and decodes each $\llbracket sketch_j \rrbracket_p$ to obtain the j -th entry of the matching elements. Putting it all together, the protocol outputs

$$\mathcal{L} = \{(i, x_i) \mid i \in [m], \text{ISMATCH}(x_i, \ell) = 1\}.$$

The complexity is affected as follows. There is a growth by a factor w in the server's overall multiplications, the communication complexity and the client's time. The degree of the homomorphically evaluated polynomial however does not suffer from any degradation.

4.2.2 Reducing the Modulus p to Quasi-Logarithmic in m

Our basic protocol requires that the homomorphic evaluations are done over $GF(p)$ for large prime modulus $p > m$. In this section we employ the CRT-coreset to reduce the modulus to nearly logarithmic in m , while repeating the small modulus protocol for $\kappa = \lceil \log m / \log \log m \rceil$ times.

The small modulus protocol is similar to our basic protocol, except for the following. First, the initialization, encryption and homomorphic evaluations are with respect to a small modulus $p = O(\log m)$. Second, the initialization upload and homomorphic evaluations are done in parallel for κ primes p_1, \dots, p_κ (specifically, the κ first primes larger than $\log m$). Third, the client receiving κ sketches:

$$\llbracket y \rrbracket_{p_1}, \dots, \llbracket y \rrbracket_{p_\kappa}$$

does the following: decrypt all sketches to obtain the residues $y_j = A_{s,m} \chi \pmod{p_j}$ for $j \in [\kappa]$, apply the CRT-decoding algorithm on these residues to obtain the value $y = A_{s,m} \chi$ over the reals, now apply the sketch-decoding algorithm (as in the basic protocol) to obtain \mathcal{I} .

In this modified protocol, the degree of the polynomial evaluated by the server is unchanged while the modulus is reduced to logarithmic $p = \tilde{O}(\log m)$ (by the prime numbers theorem). The overall number of multiplications by the server and the communication and client's complexity grow by a factor of κ .

4.2.3 Handling Unknown Number of Matches

We discussed two approaches to address settings when an upper-bound s on the number of matches is unknown.

The first approach, yielding a 2-rounds protocol revealing an upper-bound s on $|\mathcal{I}|$. In this approach the client and server engage in a preliminary round of communication for the client to learn $|\mathcal{I}|$ (the server learns nothing). The client then sends to the server an upper-bound s on $|\mathcal{I}|$ as she pleases, and engages in the report protocol for known upper-bound s .

Computing $|\mathcal{I}|$ is as follows. For the simple case when that modulus p is larger than m , computing $|\mathcal{I}|$ is very simple: The server homomorphically evaluates

$$\llbracket cnt \rrbracket_p := \sum_{i \in [m]} \llbracket \chi(i) \rrbracket_p$$

and sends to the client who decrypts and learns $|\mathcal{I}|$. For the case where we require the use of small moduli $p = \tilde{O}(\log m)$ the protocol employs the CRT-coreset: The server homomorphically evaluates the above sum modulo p_j the $\kappa = O(\log m)$ small primes p_j :

$$\llbracket cnt \rrbracket_{p_j} := \sum_{i \in [m]} \llbracket \chi(i) \rrbracket_{p_j}$$

and sends the tuple of resulting ciphertexts to the client. The client decrypts $\llbracket cnt \rrbracket_{p_j}$ for all $j \in [\kappa]$ to obtain the residues modulo p_j of $|\mathcal{I}|$. From these residues the client computes $|\mathcal{I}|$ using the standard CRT-decode algorithm.

The second approach yields a 1-round protocol, albeit with high overall multiplications and outgoing communication on the server's side (communication that the client mostly can ignore). In this case, the server sends to the client $\llbracket |\mathcal{I}| \rrbracket_p$ as well as the messages sent in the report protocol when executed with number of matches upper-bounds $s = 2^0, 2^1, \dots, m$ (assuming here w.l.o.g that m is a power of two). The client decrypts $\llbracket |\mathcal{I}| \rrbracket_p$ to learn $|\mathcal{I}|$ and then decrypts only the sketch corresponding to the smallest s larger than $|\mathcal{I}|$. Although the Cloud has communication complexity of $\Omega(m)$, the user can abort the protocol after receiving communication that is polynomial in $|\mathcal{I}|$ and $\log m$. We remark that we assumed here that $p > m$; extension to small moduli p is via using the CRT-coreset analogously to the above.

4.2.4 Dynamic Data Management

We next elaborate on how to extending the Secure report functionality to dynamic data management: Insert, Update and Delete.

Insert. Insertion of additional data elements is simply be the client encrypting and sending the element for the server to append at the end of the encrypted array (and update the array size parameter m accordingly).

Update. To update the i -th data element x_i , the client issues an update request (UPDATE, $\llbracket i \rrbracket$, $\llbracket diff \rrbracket$) with encrypted index i and encrypted difference value $diff = new - old$ for new the new value for x_i and old the current value. (If the client doesn't have old she can use PIR to retrieve it.) The server then homomorphically adds to each element $\llbracket x_j \rrbracket$ of the stored array the value $\text{ISEQUAL}(\llbracket i \rrbracket, j) \cdot \llbracket diff \rrbracket$, where i and j are in binary representation and $\text{ISEQUAL}: \{0, 1\}^{\log m} \times \{0, 1\}^{\log m} \rightarrow \{0, 1\}$ is the equality-test polynomial (see Section 4.2.5). This results in a new encrypted array $\llbracket x' \rrbracket$ satisfying $x'_i = new$ and $\forall j \neq i, x'_j$ is unchanged.

Delete. Deletion of elements can be implemented by updating them to a reserved “Deleted” symbol. Another option is switching the value of the element we wish to delete to that of the last element in the array and reducing the number of elements m by 1 (for cases when the dynamic size of the data is either maintained by the client, or is not a secret and can be maintained by the server).

4.2.5 isMatch Instantiation in our Experimental Results

For our experimental results we used an equality-test instantiation of ISMATCH for exact-match Secure report. The data and query are specified in binary representation.

Specifically we instantiate ISMATCH to using the following equality-test polynomial over $GF(p)$ for data and query in binary representation. The polynomial $\text{ISMATCH}: \{0, 1\}^w \times \{0, 1\}^w \rightarrow \{0, 1\}$ is defined by:

$$\text{ISMATCH}((a_1, \dots, a_w), (b_1, \dots, b_w)) = \prod_{i=1}^w (1 - (a_i - b_i)^2) \pmod{p}.$$

of degree $2w$.

We remark that for $p = 2$ a degree w polynomial suffice, as there is no need to square the difference $a_i - b_i$. For large $p > 2$ the standard polynomial raises the difference by $p - 1$ using Fermat’s Little Theorem to guarantee that all non-zero values map to one. Nonetheless for data in binary representation, as considered here, taking the square suffice as we are guaranteed that the difference is in $\{-1, 0, 1\}$.

5 System & Experimental Results

In this section we describe experiments on the Secure report that we implemented based on our algorithms (Protocol 3). For example, it can report the locations of all 1’s in a 10-sparse vector of size $3 \cdot 10^9 = 3,000,000,000$ entries in less than one minute by using a single machine on Amazon EC2 cloud. The system is fully open sourced, and all our experiments are reproducible. We hope to extend and improve the system in future papers together with both the theoretical and practical community.

5.1 The System

System Overview. The system maintains an encrypted database that is stored on Amazon’s AWS cloud. The system gets from the user an encrypted lookup value ℓ to search for, and a column name $array$ in a database table of length m . The encryption is computed on the user’s side using a secret key that is unknown to the cloud. The user can send the request through a web-browser, that can be run e.g. from a smart-phone or a laptop. The system then runs our Secure report coreset algorithm on the cloud, and returns a report coreset for $(array, \ell)$. The web browser then decrypts this coreset on the user’s machine and uses it to compute the solution to the report query, which is the indices i_1, \dots, i_{cnt} in $array$ that matches ℓ . Database updates can be maintained between search calls, and support multiple users that share the same security key.

Hardware. Our system is generic but in this section we evaluate it on Amazon’s AWS cloud. We use one of the standard suggested EC2 x1.32xlarge server, each with 64 2.4 GHz Intel Xeon E5-2676 v3 (Haswell) cores and 1,952 GigaByte of RAM.

Open Software and Security. The algorithms were implemented in $C++$. HELib library [Hal13] was used for the FHE commands. The source of our system is open under the GNU v3 license and can be found in [AFS18a]. For our experiments below we use a security key of 80 bits of security.

5.2 Experimental Results

In this sub-section we describe our preliminary experiments with our system and explain the results. Due to lack of space we omit more results and description that can be found in the fuller version [AFS18b].

Data. We ran the system on a lookup value $\ell = 1$ in an array of m integers $array \in \{0, \dots, m-1\}^m$. The vector was all zeroes except for s random indices, and different values for m and s were used. As expected, the actual values behind the encrypted records had no effect on the running times.

The Experiment. We ran Algorithm 3 for database table columns ranging from $m = 10$ to $m = 3,000,000,000 = 3 \cdot 10^9$ records, and $cnt \in \{10, 20, 40\}$.

Results. Our experimental results for a single machine on the cloud are shown as the circle points in Fig 1. The table of exact values appears in the full version [AFS18b].

The user’s decoding time was negligible in all the experiments, so the cloud’s time equals to the overall running time. For example, the graph shows that a single machine can report in about 2 minutes all 1’s in a 20-sparse column of 3,000,000,000 binary entries.

Comparison to the direct approach. Our theoretical results proves that the running time of our new algorithm is only poly-logarithmic in the number m of entries compared to the direct approach which is polynomial in m . However, it assumes that both our and the direct algorithm may use m machines in parallel, i.e., a machine for each record. The goal of our experiment was to show a significant time reduction even using a single machine on the cloud where a running time that is linear in m is expected.

The graph in Fig. 1 is log scaled so a linear curve shows a polynomial relation, and its slope is the degree of the polynomial. Our experiments indeed show that our algorithm is linear in m as expected. The direct algorithm perform worse, both on small number of entries (the ratio is approximately 1000), but also asymptotically: The slope of the direct approach is about 1.5, indicating a running time of $O(m^{1.5})$ (as oppose to our $O(m)$). The main reason for the improvement in running time of our algorithm is that unlike the direct algorithm, our algorithm does not use multiplications to compute the sketch.

6 Conclusion

In this paper we give a secure protocol for the report problem on encrypted data elements and query. Our protocol retrieves all matching elements in a data array, where both data elements and query are encrypted using Fully Homomorphic Encryption (FHE).

Our protocol involves a single-server, single-round, and requires no pre-processing of the cleartext data. The complexity of our protocol, given an upper-bound s on the number of matches, is as follows: The server evaluates a polynomial of degree independent of the array size m ; the communication and client’s complexity is polynomial in s and $\log m$. The protocol is secure against semi-honest adversaries controlling the server. We presented extensions and optimizations of our protocol: handling the case that s is unknown; reducing the ring size for the homomorphic evaluation from linear to logarithmic in m ;

Table 2: Server’s running time of Report Coreset (Algorithm 3) as measured on a single machine on Amazon’s cloud for different database *array* size (1st column). In these experiments $cnt = 40$, i.e. at most 40 elements matched the report criteria. The 3rd column shows the time in minute to compute report with the direct algorithm and The 4th column shows the time in minute to compute report with our coreset algorithm. The 2nd column shows the the speedup of our report coreset. Since we could not run the direct approach for too large databases we give estimations In parentheses are values based on our analysis.

Records	$\frac{time(direct)}{time(coreset)}$	Direct Repot (minutes)	Report Coreset (minutes)
89,600	48	4	0.083
192,000	47	7	0.15
396,800	73	22	0.3
806,400	109	67	0.61
1,625,600	164	199	1.21
3,264,000	268	689	2.56
6,540,800	450	2,281	10.16
13,094,400	(754)	(7,661)	20.55
26,201,600	(441)	(9,068)	40.58
52,416,000	(542)	(21,990)	82.38
104,844,800	(577)	(47,568)	163.56
209,702,400	(3,107)	(508,126)	328.08
419,417,600	(3,502)	(1,148,875)	653.07

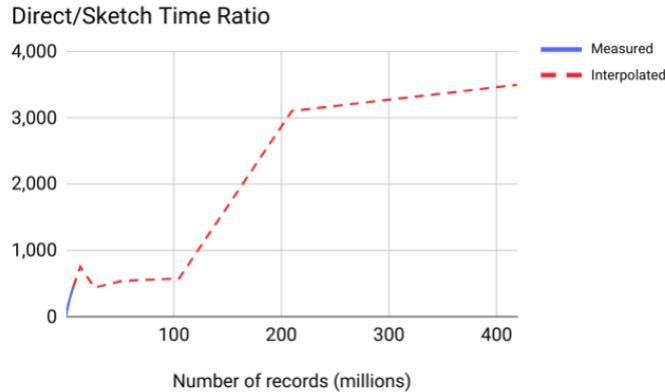


Figure 4: The speedup of using our report coreset (ratio between the direct algorithm and the report coreset) for arrays with $cnt = 40$, i.e. at most 40 records matched the query. The x -axis shows the array size (in millions) and the y -axis shows the ratio.

and handling dynamic data management. Furthermore, our protocol is compatible with versatile matching conditions ISMATCH, and can be utilized for example for exact match, similarity search, wild-card matching, range queries, Boolean queries and more.

We implemented our protocol in an open source library based on HELib, and ran experiments on Amazon’s AWS EC2 cloud. Our experiments show that we can search in *billions of encrypted data elements for an encrypted query in a minute*.

To achieve our results we put forth our coreset paradigm for homomorphic encryption, and present a low-degree polynomial for Secure report. Both contributions could be useful for future works. The coreset for homomorphic encryption paradigm has already been employed in follow-up works. The low-degree polynomial could be useful in the broader

context of secure computation, using other techniques such as secret sharing beyond FHE.

The simplicity of our algorithm makes it a good candidate to run on small embedded systems. Motivating examples are data sources such as sensors of wearable devices uploading data to a server while keeping it hidden from the Server but still allowing to query the data. In this case the sensors take the role of the data source client in the running the upload phase of our protocol. With our protocol a network of weak sensor devices can store their encrypted data in a server for the search client to issue the report queries.

References

- [AAG17] Mohamed Ahmed Abdelraheem, Tobias Andersson, and Christian Gehrman. Inference and record-injection attacks on searchable encrypted relational databases. *IACR Cryptology ePrint Archive*, 2017:24, 2017.
- [AFS18a] Adi Akavia, Dan Feldman, and Hayim Shaul. ReportLib: Open library for FHE report, 2018.
- [AFS18b] Adi Akavia, Dan Feldman, and Hayim Shaul. Secure data retrieval on the cloud homomorphic encryption meets coresets. *Cryptology ePrint Archive*, Report 2018/1003, 2018. <https://eprint.iacr.org/2018/1003>.
- [AFS18c] Adi Akavia, Dan Feldman, and Hayim Shaul. Secure search via sketching for homomorphic encryption. In *Proceedings of the 25th ACM Conference on Computer and Communications Security*. ACM, 2018.
- [AGHL18] Adi Akavia, Craig Gentry, Shai Halevi, and Max Leibovich. Setup-free secure search on encrypted data: Faster and post-processing free. *Cryptology ePrint Archive*, Report 2018/1235, 2018. <https://eprint.iacr.org/2018/1235>.
- [AHPV05] Pankaj K Agarwal, Sariel Har-Peled, and Kasturi R Varadarajan. Geometric approximation via coresets. *Combinatorial and computational geometry*, 52:1–30, 2005.
- [BFL⁺17] Vladimir Braverman, Gereon Frahling, Harry Lang, Christian Sohler, and Lin F Yang. Clustering high dimensional dynamic data streams. *arXiv preprint arXiv:1706.03887*, 2017.
- [BGH⁺13] Dan Boneh, Craig Gentry, Shai Halevi, Frank Wang, and David J Wu. Private database queries using somewhat homomorphic encryption. In *International Conference on Applied Cryptography and Network Security*, pages 102–118. Springer, 2013.
- [BGV12] Zvika Brakerski, Craig Gentry, and Vinod Vaikuntanathan. (leveled) fully homomorphic encryption without bootstrapping. In *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, ITCS '12, pages 309–325, New York, NY, USA, 2012. ACM.
- [BI05] Omer Barkol and Yuval Ishai. Secure computation of constant-depth circuits with applications to database search problems. In *Annual International Cryptology Conference*, pages 395–411. Springer, 2005.
- [BMO⁺11] Vladimir Braverman, Adam Meyerson, Rafail Ostrovsky, Alan Roytman, Michael Shindler, and Brian Tagiku. Streaming k-means on well-clusterable data. In *Proceedings of the twenty-second annual ACM-SIAM symposium on Discrete Algorithms*, pages 26–40. Society for Industrial and Applied Mathematics, 2011.

- [BV11] Zvika Brakerski and Vinod Vaikuntanathan. Efficient fully homomorphic encryption from (standard) LWE. In *IEEE 52nd Annual Symposium on Foundations of Computer Science, FOCS 2011, Palm Springs, CA, USA, October 22-25, 2011*, pages 97–106, 2011.
- [BZD10] Christos Boutsidis, Anastasios Zouzias, and Petros Drineas. Random projections for k -means clustering. In *Advances in Neural Information Processing Systems*, pages 298–306, 2010.
- [qCDD⁺16] Gizem S. qCetin, Wei Dai, Yarkin Doröz, William J. Martin, and Berk Sunar. Blind web search: How far are we from a privacy preserving search engine? *IACR Cryptology ePrint Archive*, 2016:801, 2016.
- [CG05] Graham Cormode and Minos Garofalakis. Sketching streams through the net: Distributed approximate query tracking. In *Proceedings of the 31st international conference on Very large data bases*, pages 13–24. VLDB Endowment, 2005.
- [CG07] Graham Cormode and Minos Garofalakis. Sketching probabilistic data streams. In *Proceedings of the 2007 ACM SIGMOD international conference on Management of data*, pages 281–292. ACM, 2007.
- [CGHJ12] Graham Cormode, Minos Garofalakis, Peter J Haas, and Chris Jermaine. Synopses for massive data: Samples, histograms, wavelets, sketches. *Foundations and Trends in Databases*, 4(1–3):1–294, 2012.
- [CGKS95] Benny Chor, Oded Goldreich, Eyal Kushilevitz, and Madhu Sudan. Private information retrieval. In *Foundations of Computer Science, 1995. Proceedings., 36th Annual Symposium on*, pages 41–50. IEEE, 1995.
- [CGN97] Benny Chor, Niv Gilboa, and Moni Naor. Private information retrieval by keywords, 1997.
- [CGPR15] David Cash, Paul Grubbs, Jason Perry, and Thomas Ristenpart. Leakage-abuse attacks against searchable encryption. In *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security*, pages 668–679. ACM, 2015.
- [CK08] Edith Cohen and Haim Kaplan. Tighter estimation using bottom k sketches. *Proceedings of the VLDB Endowment*, 1(1):213–224, 2008.
- [CK10] Melissa Chase and Seny Kamara. Structured encryption and controlled disclosure. In *International Conference on the Theory and Application of Cryptology and Information Security*, pages 577–594. Springer, 2010.
- [CKK16] Jung Hee Cheon, Miran Kim, and Myungsun Kim. Optimized search-and-compute circuits and their application to query evaluation on encrypted data. *IEEE Trans. Information Forensics and Security*, 11(1):188–199, 2016.
- [CKKS17] Jung Hee Cheon, Andrey Kim, Miran Kim, and Yong Soo Song. Homomorphic encryption for arithmetic of approximate numbers. In Tsuyoshi Takagi and Thomas Peyrin, editors, *Advances in Cryptology - ASIACRYPT 2017 - 23rd International Conference on the Theory and Applications of Cryptology and Information Security, Hong Kong, China, December 3-7, 2017, Proceedings, Part I*, volume 10624 of *Lecture Notes in Computer Science*, pages 409–437. Springer, 2017.

- [CKL15] Jung Hee Cheon, Miran Kim, and Kristin E Lauter. Homomorphic computation of edit distance. In *Financial Cryptography Workshops*, pages 194–212, 2015.
- [Cla10] Kenneth L Clarkson. Coresets, sparse greedy approximation, and the Frank-Wolfe algorithm. *ACM Transactions on Algorithms (TALG)*, 6(4):63, 2010.
- [CLMS13] Artur Czumaj, Christiane Lammersen, Morteza Monemizadeh, and Christian Sohler. $(1 + \epsilon)$ -approximation for facility location in data streams. In *Proceedings of the twenty-fourth annual ACM-SIAM symposium on Discrete algorithms*, pages 1710–1728. Society for Industrial and Applied Mathematics, 2013.
- [CSZ05] Richard Cole, Dennis Shasha, and Xiaojian Zhao. Fast window correlations over uncooperative time series. In *Proceedings of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining*, pages 743–749. ACM, 2005.
- [CZ17] Jiecao Chen and Qin Zhang. Bias-aware sketches. *Proceedings of the VLDB Endowment*, 10(9):961–972, 2017.
- [DHH00] Dingzhu Du, Frank K Hwang, and Frank Hwang. *Combinatorial group testing and its applications*, volume 12. World Scientific, 2000.
- [DSH14] Yarkin Doröz, Berk Sunar, and Ghaith Hammouri. Bandwidth efficient PIR from NTRU. In *Financial Cryptography and Data Security - FC 2014 Workshops, BITCOIN and WAHC 2014, Christ Church, Barbados, March 7, 2014, Revised Selected Papers*, pages 195–207, 2014.
- [GAYBL17] Matthieu Giraud, Alexandre Anzala-Yamajako, Olivier Bernard, and Pascal Lafourcade. Practical passive leakage-abuse attacks against symmetric searchable encryption. In *14th International Conference on Security and Cryptography SECRYPT 2017*. SCITEPRESS-Science and Technology Publications, 2017.
- [Gen09a] Craig Gentry. *A Fully Homomorphic Encryption Scheme*. PhD thesis, Stanford University, Stanford, CA, USA, 2009. AAI3382729.
- [Gen09b] Craig Gentry. Fully homomorphic encryption using ideal lattices. In *Proceedings of the Forty-first Annual ACM Symposium on Theory of Computing, STOC '09*, pages 169–178, New York, NY, USA, 2009. ACM.
- [GKMS01] Anna C Gilbert, Yannis Kotidis, S Muthukrishnan, and Martin Strauss. Surfing wavelets on streams: One-pass summaries for approximate aggregate queries. In *VLDB*, volume 1, pages 79–88, 2001.
- [GKS13] Minos Garofalakis, Daniel Keren, and Vasilis Samoladas. Sketch-based geometric monitoring of distributed stream queries. *Proceedings of the VLDB Endowment*, 6(10):937–948, 2013.
- [GM12] Sudipto Guha and Andrew McGregor. Graph synopses, sketches, and streams: A survey. *Proceedings of the VLDB Endowment*, 5(12):2030–2031, 2012.
- [GMN⁺16] Paul Grubbs, Richard McPherson, Muhammad Naveed, Thomas Ristenpart, and Vitaly Shmatikov. Breaking web applications built on top of encrypted data. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, pages 1353–1364. ACM, 2016.

- [GMW87a] O. Goldreich, S. Micali, and A. Wigderson. How to play any mental game. In *Proceedings of the Nineteenth Annual ACM Symposium on Theory of Computing*, STOC '87, pages 218–229, New York, NY, USA, 1987. ACM.
- [GMW87b] O. Goldreich, S. Micali, and A. Wigderson. How to play any mental game. In *Proceedings of the Nineteenth Annual ACM Symposium on Theory of Computing*, STOC '87, pages 218–229, New York, NY, USA, 1987. ACM.
- [Gol87] O. Goldreich. Towards a theory of software protection and simulation by oblivious rams. In *Proceedings of the Nineteenth Annual ACM Symposium on Theory of Computing*, STOC '87, pages 182–194, New York, NY, USA, 1987. ACM.
- [GSB⁺17] Paul Grubbs, Kevin Sekniqi, Vincent Bindschaedler, Muhammad Naveed, and Thomas Ristenpart. Leakage-abuse attacks against order-revealing encryption. In *Security and Privacy (SP), 2017 IEEE Symposium on*, pages 655–672. IEEE, 2017.
- [Hal13] Shai Halevi. Helib - an implementation of homomorphic encryption. <https://github.com/shaih/HElib/>, 2013.
- [HCB16] Jonathan Huggins, Trevor Campbell, and Tamara Broderick. Coresets for scalable bayesian logistic regression. In *Advances in Neural Information Processing Systems*, pages 4080–4088, 2016.
- [HK15] Hao Huang and Shiva Prasad Kasiviswanathan. Streaming anomaly detection using randomized matrix sketching. *Proceedings of the VLDB Endowment*, 9(3):192–203, 2015.
- [IK00] Yuval Ishai and Eyal Kushilevitz. Randomizing polynomials: A new representation with applications to round-efficient secure computation. In *focs*, page 294. IEEE, 2000.
- [IKK12] Mohammad Saiful Islam, Mehmet Kuzu, and Murat Kantarcioglu. Access pattern disclosure on searchable encryption: Ramification, attack and mitigation. In *Ndss*, volume 20, page 12, 2012.
- [IKM00] Piotr Indyk, Nick Koudas, and S Muthukrishnan. Identifying representative trends in massive time series data sets using sketches. In *VLDB*, pages 363–372, 2000.
- [INR10] Piotr Indyk, Hung Q Ngo, and Atri Rudra. Efficiently decodable non-adaptive group testing. In *Proceedings of the twenty-first annual ACM-SIAM symposium on Discrete Algorithms*, pages 1126–1142. SIAM, 2010.
- [KLL⁺16] Myungsun Kim, Hyung Tae Lee, San Ling, Shu Qin Ren, Benjamin Hong Meng Tan, and Huaxiong Wang. Better security for queries on encrypted databases. *IACR Cryptology ePrint Archive*, 2016:470, 2016.
- [KLL⁺17] Myungsun Kim, Hyung Tae Lee, San Ling, Benjamin Hong Meng Tan, and Huaxiong Wang. Private compound wildcard queries using fully homomorphic encryption. *IEEE Transactions on Dependable and Secure Computing*, 2017.
- [LLN15] Kristin E. Lauter, Adriana López-Alt, and Michael Naehrig. Private computation on encrypted genomic data. *IACR Cryptology ePrint Archive*, 2015:133, 2015.

- [Phi16] Jeff M Phillips. Coresets and sketches. *arXiv preprint arXiv:1601.00617*, 2016.
- [PW16] David Pouliot and Charles V Wright. The shadow nemesis: Inference attacks on efficiently deployable, efficiently searchable encryption. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, pages 1341–1352. ACM, 2016.
- [RAD78] R L Rivest, L Adleman, and M L Dertouzos. On data banks and privacy homomorphisms. *Foundations of Secure Computation*, Academia Press, pages 169–179, 1978.
- [Raz87] Alexander A Razborov. Lower bounds on the size of bounded depth circuits over a complete basis with logical addition. *Mathematical Notes of the Academy of Sciences of the USSR*, 41(4):333–338, 1987.
- [RLNZ09] Galen Reeves, Jie Liu, Suman Nath, and Feng Zhao. Managing massive time series streams with multi-scale compressed trickles. *Proceedings of the VLDB Endowment*, 2(1):97–108, 2009.
- [RVVV17] Sujoy Sinha Roy, Frederik Vercauteren, Jo Vliegen, and Ingrid Verbauwhede. Hardware assisted fully homomorphic function evaluation and encrypted search. *IEEE Transactions on Computers*, 2017.
- [Smo87] Roman Smolensky. Algebraic methods in the theory of lower bounds for boolean circuit complexity. In *Proceedings of the nineteenth annual ACM symposium on Theory of computing*, pages 77–82. ACM, 1987.
- [SWP00] Dawn Xiaoding Song, David Wagner, and Adrian Perrig. Practical techniques for searches on encrypted data. In *Security and Privacy, 2000. S&P 2000. Proceedings. 2000 IEEE Symposium on*, pages 44–55. IEEE, 2000.
- [TJW⁺16] Haixu Tang, Xiaoqian Jiang, Xiaofeng Wang, Shuang Wang, Heidi Sofia, Dov Fox, Kristin Lauter, Bradley Malin, Amalio Telenti, Li Xiong, and Lucila Ohno-Machado. Protecting genomic data analytics in the cloud: state of the art and opportunities. *BMC Medical Genomics*, 9(1):63, Oct 2016.
- [W⁺14] David P Woodruff et al. Sketching as a tool for numerical linear algebra. *Foundations and Trends® in Theoretical Computer Science*, 10(1–2):1–157, 2014.
- [WYG⁺17] Frank Wang, Catherine Yun, Shafi Goldwasser, Vinod Vaikuntanathan, and Matei Zaharia. Splinter: Practical private queries on public data. In *NSDI*, pages 299–313. USENIX Association, 2017.
- [Yao86] Andrew Chi-Chih Yao. How to generate and exchange secrets. In *Proceedings of the 27th Annual Symposium on Foundations of Computer Science, SFCS '86*, pages 162–167, Washington, DC, USA, 1986. IEEE Computer Society.
- [YSK⁺13] Masaya Yasuda, Takeshi Shimoyama, Jun Kogure, Kazuhiro Yokoyama, and Takeshi Koshiba. Secure pattern matching using somewhat homomorphic encryption. In *Proceedings of the 2013 ACM Workshop on Cloud Computing Security Workshop, CCSW '13*, pages 65–76, New York, NY, USA, 2013. ACM.
- [ZKP16] Yupeng Zhang, Jonathan Katz, and Charalampos Papamanthou. All your queries are belong to us: The power of file-injection attacks on searchable encryption. In *USENIX Security Symposium*, pages 707–720, 2016.

A Direct Polynomial

We present the direct polynomial for returning all matches. We present here the polynomial as defined on cleartext data. To compute on encrypted data the polynomial is replaced by its homomorphic version, as standard.

Specifications for the direct polynomial. We first define the direct polynomial's parameters, input and output.

- The direct polynomial is parameterized by the number of elements m , a prime $p > m$ so that computations are in $GF(p)$, the number of matches to be returned $cnt \in \{1, \dots, m\}$, the domain of data-elements \mathcal{M} and of queries \mathcal{Q} , and a polynomial realizing the matching criterion $ISMATCH: \mathcal{M} \times \mathcal{Q} \rightarrow \{0, 1\}$.
- The input to the polynomial is the data $array = (x_1, \dots, x_m) \in \mathcal{M}^m$ and query $\ell \in \mathcal{Q}$.
- The output is the indexes of the matching elements $output \in \{0, \dots, m\}^{cnt}$, where for every $k \in [cnt]$, $output(k)$ is the index $i \in [m]$ of the k -th match for ℓ in $array$. That is, $output(k) = i$ if-and-only-if both the following holds:
 1. The i -th element is a match: $ISMATCH(x_i, \ell) = 1$, and
 2. There are $k - 1$ preceding matches: $|\{j \in [i - 1] \mid ISMATCH(x_j, \ell) = 1\}| = k - 1$.

In case $cnt > |\mathcal{I}|$, for \mathcal{I} the set of all matches as defined in Section 3, the output on entries $k \in \{|\mathcal{I}| + 1, \dots, cnt\}$ is 0.

Realization of the direct polynomial. We next specify how to realized the direct polynomial. Computation in the following is over $GF(p)$ for a prime $p > m$. For every $k \in [cnt]$, we define:

$$output(k) = \sum_{i=1}^m i \cdot ISMATCH(x_i, \ell) \cdot \left(1 - (\text{PREVMATCHES}(i) - (k - 1))^{p-1}\right)$$

where $\text{PREVMATCHES}(i)$ is the number of matches in the sub-array (x_1, \dots, x_{i-1}) , as computed by the following polynomial:

$$\text{PREVMATCHES}(i) = \sum_{j=1}^{i-1} ISMATCH(x_j, \ell).$$

Analysis of the direct polynomial. The analysis for the direct polynomial follows.

Correctness. To argue correctness observer first that $\text{PREVMATCHES}(i)$ is the number of elements in x_1, \dots, x_{i-1} that match ℓ , namely:

$$\text{PREVMATCHES}(i) = |\{\mathcal{I} \cap \{1, \dots, i - 1\}\}|.$$

Next, for every $k \in \{1, \dots, |\mathcal{I}|\}$, by Fermat's Little Theorem (FLT), we have that

$$1 - (\text{PREVMATCHES}(i) - (k - 1))^{p-1} = \begin{cases} 1 & \text{if } \text{PREVMATCHES}(i) = k - 1 \\ 0 & \text{otherwise.} \end{cases}$$

Therefore, $\text{ISMATCH}(x_i, \ell) \cdot \left(1 - (\text{PREVMATCHES}(i) - (k - 1))^{p-1}\right)$ evaluates to 1 if-and-only-if both conditions (1) and (2) from the output definition above hold, and it is 0 otherwise. So,

$$i \cdot \text{ISMATCH}(\text{array}(i), \ell) \cdot \left(1 - (\text{PREVMATCHES}(i) - (k - 1))^{p-1}\right)$$

is equal to i if-and-only-if i is the index of the k -th matching element to ℓ in array , and it is 0 otherwise. We conclude that the $\text{output}[k] = i$ for i the index of the k -th matching element to ℓ in array .

Complexity. We show that the direct polynomial has degree $\Omega(m \cdot d)$ for d the degree of ISMATCH . The direct polynomial for computing $\text{output}(k)$ is the product of several sub-polynomials: (i) The plaintext coefficient i , (ii) The ISMATCH polynomial, and (iii) The product of $(p - 1)$ polynomials, each of degree identical to the degree of PREVMATCHES . The latter is equal to the degree of ISMATCH (as it is the sum of ISMATCH evaluations). The degree of a product of polynomials is the sum of their degrees. So the degree of the direct polynomial is $d + (p - 1)d = pd$. We conclude, by assigning $p > m$, that the degree of the direct polynomial is $\Omega(m \cdot d)$.

Remarks. The requirement that $p > m$ is necessary for correctness. Otherwise, the output is incorrect due to two reasons. First, for $p \leq |\mathcal{I}|$, PREVMATCHES would have “false-zero” value whenever the number of previous-matches is a multiple of p , leading to incorrect output. Second, for $p \in (|\mathcal{I}|, m)$, the output values are the residues “ $i \bmod p$ ” of the desired indexes i , which is insufficient information for reconstructing i when $i \in [p + 1, m]$.

We note that there are other possible implementations for the direct polynomial. In particular, the output of PREVMATCHES could be specified in binary representation. In the case, computations can be done over $GF(2)$, and the *comparison* of PREVMATCHES to $k - 1$ can be done in low-degree $\log m$. However, the degree of *counting* the number of previous matches PREVMATCHES becomes $\Omega(m \cdot d)$.