

Secure Searching of Biomarkers Using Hybrid Homomorphic Encryption Scheme

Jung Hee Cheon¹, Miran Kim², and Yongsoo Song¹

¹ Seoul National University, Republic of Korea
{jhcheon, lucius05}@snu.ac.kr

² University of California, San Diego,
kmiran@ucsd.edu

Abstract. As genome sequencing technology develops rapidly, there has lately been an increasing need to keep genomic data secure even when stored in the cloud and still used for research. In this paper, we are interested in designing a protocol for the secure outsourcing matching problem on encrypted data. We propose an efficient method to securely search a matching position with the query data and extract some information at the position. After decryption, we only perform a small amount of comparison with the query information in plaintext. We apply this method to find a set of biomarkers in encrypted genomes.

The important feature of our method is to encode a genomic database as a single element of polynomial ring. It also requires only a single homomorphic multiplication for query computation. Thus this method has the advantage over the previous methods in parameter size, computational complexity, and communication cost.

We evaluate the performance of our method and verify that computation on large-scale personal data can be securely and practically outsourced to a cloud environment during data analysis. It takes about 3.9 seconds to *search-and-extract* the reference and alternate sequences of the queried position in a database of size 4M.

Keywords. Homomorphic encryption, Biomarkers

1 Introduction

The rapid development of genome sequencing technology enables us to access large genome dataset and it looks poised to make a significant breakthrough in medical research. While genomic data can be used for a wide range of applications including healthcare, biomedical research, and direct-to-consumer services, it has numerous special distinguishing features and it can violate personal privacy via genetic disclosure or genetic discrimination [HAHT13, EN14, NAC⁺15]. Due to these potential privacy issues, it should be managed with care.

There have been various privacy-enhancing techniques using cryptographic methods as outsourced analysis tools of genomic data. Recently, it has been suggested that we can preserve privacy through homomorphic encryption (HE), which allows computations to be carried out on ciphertexts. Yasuda et al. [YSK⁺13] gave a practical solution to find the location of a pattern in a text by computing multiple Hamming distance values on encrypted data. Lauter et al. [LLAN14] gave a solution to privately compute the basic genomic algorithms used in genome-wide association studies.

Homomorphic encryption can be applied to privacy-preserving sequence comparison, but it is still impractical for the analysis of entire human genome information. For example, Cheon et al. [CKL15] presented a protocol to compute the edit distance on homomorphically encrypted data but it took about 27 seconds even on length 8 DNA sequence. It is not easy to efficiently approximate the edit distance over encryption even though the distance to a public human DNA sequence is given [KL15]. This inefficiency comes from the difficulty of homomorphic evaluation of equality test: Encrypting the inputs bit-wise and computing over the encrypted bits yield expensive computation cost (at least linear in the data bit-length).

In this paper, we suggest an efficient method to securely search a set of biomarkers using hybrid Ring-GSW homomorphic encryption scheme. The important feature of our method is to encode a genomic database as a single element of polynomial ring. The searching operation in a database is done by a single multiplication with the query gene, thus the information of DNA sequences at the matching position with the query data is moved to the constant term of output polynomial. Then we perform the extraction procedure to obtain the DNA sequence, and after decryption we compare it with the query DNA information in plaintext. In particular, this conversion procedure not only prevents leakage of database information that has not been queried by user but also reduces the communication cost by half.

Since all the matching computation is performed on encrypted data in the cloud, the security against a semi-honest adversary follows from the semantic security of the underlying HE scheme. Furthermore, our method has the advantage over previous methods in parameter size and computational complexity since it requires only a single multiplication. Our implementation takes about 3.9 seconds to extract the reference and alternative (non-reference) sequences of the queried position when each of them consists of two single nucleotide polymorphisms (SNPs) in a database of size 4M.

2 Background

The iDASH (Integrating Data for Analysis, ‘anonymization’ and SHaring) National Center organizes the iDASH Privacy & Security challenge for secure genome analysis. This paper is based on a submission to the task 3 in 2016 iDASH challenge: secure outsourcing of testing for genetic diseases on encrypted genomes.

2.1 Problem Setting

The goal of task 3 is to privately calculate the probability of genetic diseases through matching a set of biomarkers to encrypted genomes stored in a public cloud service. The requirement is that the entire matching process needs to be carried out using homomorphic encryption so that any information about database and query should not be revealed to the server during computation.

Suppose that the client has a Variation Call Format (VCF) file which contains genotype information such as chromosome number and position in the genome. It also contains some information for each position such as reference and alternate sequences, where each base must be one of SNPs: A, T, G , and C . The client encrypts the information using homomorphic encryption and the server calculates the exact match over the encrypted data. The outcome is the absence/presence of the specified biomarkers, that is, an encryption of 1 if matched; otherwise, an encryption of 0. Finally the client decrypts the result by the secret key of homomorphic encryption.

2.2 Practical Homomorphic Encryption

Fully Homomorphic cryptosystems allow us to homomorphically evaluate any arithmetic circuit without decryption. However, the noise of the resulting ciphertext grows during homomorphic evaluations, slightly with addition but substantially with multiplication. For efficiency reasons, for tasks which are known in advance, we use a more practical *Somewhat Homomorphic Encryption* (SHE) scheme, which evaluates functions up to a certain complexity. In particular, two techniques are used for noise management of SHE: one is the *modulus-switching* technique introduced by Brakerski, Gentry and Vaikuntanathan [BGV12],

which scales down a ciphertext during every multiplication operation and reduces the noise by its scaling factor. The other is a *scale-invariant* technique proposed by Brakerski such that the same modulus is used throughout the evaluation process [Bra12].

Let us denote by $[\cdot]_Q$ the reduction modulo Q into the interval $(-Q/2, Q/2] \cap \mathbb{Z}$ of the integer or integer polynomial (coefficient-wise). For a security parameter λ , we choose an integer $M = M(\lambda)$ that defines the M -th cyclotomic polynomial $\Phi_M(X)$. For a polynomial ring $\mathcal{R} = \mathbb{Z}[X]/(\Phi_M(X))$, set the plaintext space to $\mathcal{R}_t := \mathcal{R}/t\mathcal{R}$ for some fixed $t \geq 2$ and the ciphertext space to $\mathcal{R}_Q := \mathcal{R}/Q\mathcal{R}$ for an integer $Q = Q(\lambda)$. Let $\chi = \chi(\lambda)$ denote a noise distribution over the ring \mathcal{R} . We use the standard notation $a \leftarrow \mathcal{D}$ to denote that a is chosen from the distribution \mathcal{D} .

The Basic Scheme. The following is a description of basic homomorphic encryption scheme based on the hardness of (decisional) *Ring Learning with Errors* (RLWE) assumption, which was first introduced by Lyubashevsky, Peikert and Regev [LPR10]. The assumption is that it is infeasible to distinguish the following two distributions. The first distribution consists of pairs (a_i, u_i) , where a_i and u_i are drawn uniformly at random from \mathcal{R}_Q . The second distribution consists of pairs of the form $(a_i, b_i) = (a_i, a_i \mathbf{s} + e_i)$ where a_i is uniformly random in \mathcal{R}_Q and \mathbf{s}, e_i are drawn from the error distribution χ . To improve efficiency for HE, we use sparse secret keys \mathbf{s} with coefficients sampled from $\{0, \pm 1\}$ as in [GHS12].

- RLWE.ParamsGen(λ): Given the security parameter λ , choose an integer M , a modulus Q , a plaintext modulus t with $t|Q$, and discrete Gaussian distribution χ_{err} . Output $params \leftarrow (M, Q, t, \chi_{err})$.
- RLWE.KeyGen($params$): On the input parameters, let $N = \phi(M)$ and choose a sparse random \mathbf{s} from $\{0, \pm 1\}^N$. Generate an RLWE instance $(a, b) = (a, [-a\mathbf{s} + e]_Q)$ for $e \leftarrow \chi_{err}$. We set the secret key $\mathbf{sk} \leftarrow \mathbf{s}$ and the public key $\mathbf{pk} \leftarrow (a, b)$.
- RLWE.Enc(m, \mathbf{pk}): To encrypt $m \in \mathcal{R}_t$, choose a small polynomial v and two Gaussian polynomials e_0, e_1 over \mathcal{R} and output the ciphertext

$$\begin{aligned} \mathbf{ct} &\leftarrow (c_0, c_1) \\ &= ((Q/t)m, 0) + (bv + e_0, av + e_1) \in \mathcal{R}_Q^2. \end{aligned}$$

- RLWE.Dec(\mathbf{ct}, \mathbf{sk}): Given a ciphertext $\mathbf{ct} = (c_0, c_1)$, output $m \leftarrow \lfloor (t/Q) \cdot [c_0 + \mathbf{s} \cdot c_1]_Q \rfloor$.
- RLWE.Add($\mathbf{ct}, \mathbf{ct}'$): Given two ciphertexts $\mathbf{ct} = (c_0, c_1)$ and $\mathbf{ct}' = (c'_0, c'_1)$, the homomorphic addition is computed by $\mathbf{ct}_{\text{add}} \leftarrow ([c_0 + c'_0]_Q, [c_1 + c'_1]_Q)$.

Throughout this paper, we assume that the integer M is a power of two so that $N = M/2$ and $\phi_M(X) = X^N + 1$. We adapt the conversion and modulus-switching techniques of [DM15]. The conversion algorithm changes an RLWE encryption of $m = \sum_i m_i X^i$ into an LWE encryption of its constant term m_0 , and the modulus switching reduces the ciphertext modulus Q down to q while preserving the message. We note that an LWE ciphertext is represented as a vector in \mathbb{Z}_q for some modulus q , and the decryption procedure is done by an inner product of the ciphertext and the secret key vector.

- RLWE.Conv(\mathbf{ct}): Given a ciphertext $\mathbf{ct} = (c_0, c_1)$ with $c_0 = \sum_i c_{0,i} X^i$ and $c_1 = \sum_i c_{1,i} X^i$, output the vector $\mathbf{ct}' = (c_{0,0}, c_{1,0}, -c_{1,N-1}, \dots, -c_{1,1})$.
- LWE.ModSwitch(\mathbf{ct}): Given a ciphertext $\mathbf{ct} \in \mathbb{Z}_Q^{N+1}$, output the vector $\mathbf{ct}' \leftarrow \lfloor (q/Q) \cdot \mathbf{ct} \rfloor \in \mathbb{Z}_q^{N+1}$.

An RLWE ciphertext $\text{ct} = (c_0, c_1)$ has the decryption structure of the form $c_0 + c_1 \cdot \mathbf{s} = (Q/t) \cdot m + e$ and its constant term is

$$c_{0,0} + c_{1,0}s_0 - \sum_{i=1}^{N-1} c_{1,N-i}s_i = (Q/t) \cdot m_0 + e_0.$$

It can be represented as an inner product of a vector $(c_{0,0}, c_{1,0}, -c_{1,N-1}, -c_{1,N-2}, \dots, -c_{1,1})$ and the desired LWE secret key $\mathbf{s} = (1, s_0, \dots, s_{N-1})$. Hence the output of the conversion algorithm can be seen as an LWE encryption of m_0 . It is also easy to check that if $\text{ct} \in \mathbb{Z}_Q^{N+1}$ satisfies $\langle \text{ct}, \mathbf{s} \rangle = (Q/t) \cdot m + e \pmod{Q}$, then the output of LWE.ModSwitch algorithm satisfies $\langle \text{ct}', \mathbf{s} \rangle = (q/t) \cdot m + e' \pmod{q}$ for some $e' \approx (q/Q) \cdot e$. These techniques have been proposed for an efficient bootstrapping [DM15], but they will play totally different roles in our application. Finally an LWE ciphertext of modulus q can be decrypted by \mathbf{s} as follows.

– LWE.Dec(ct, sk): Given a ciphertext $\text{ct} \in \mathbb{Z}_q^{N+1}$, output the value $m \leftarrow \lfloor (t/q) \cdot \langle \text{ct}, \mathbf{s} \rangle \rfloor_q$.

If $\langle \text{ct}, \mathbf{s} \rangle = (q/t) \cdot m + e \pmod{q}$ for some small enough e , it returns the correct message m modulo t . More precisely, the decryption procedure works if $|te/q| < 1/2$.

The Ring-GSW Scheme. Gentry, Sahai, and Waters [GSW13] suggested a fully homomorphic encryption based on the LWE problem, where the message is encrypted as an approximate eigenvalue of a ciphertext. Ducas and Micciancio [DM15] described its RLWE variant. The RGSW symmetric encryption scheme consists of the following algorithms.

- RGSW.ParamsGen(\cdot), RGSW.KeyGen(\cdot): Use the same parameter $params$ and secret key \mathbf{s} with the basic RLWE scheme. Additionally set the decomposition base $B_{\mathbf{g}}$ and exponent $d_{\mathbf{g}}$ satisfying $B_{\mathbf{g}}^{d_{\mathbf{g}}} \geq Q$.
- RGSW.Enc(m, sk): To encrypt $m \in \mathcal{R}_t$, pick a matrix $\mathbf{a} \in \mathcal{R}_Q^{2d_{\mathbf{g}}}$ uniformly at random, and $\mathbf{e} \in \mathcal{R}^{2d_{\mathbf{g}}} \simeq \mathbb{Z}^{2d_{\mathbf{g}} \cdot n}$ with discrete Gaussian distribution χ of parameter ς , and output the ciphertext

$$\text{CT} \leftarrow [\mathbf{b}, \mathbf{a}] + m\mathbf{G} \in \mathcal{R}_Q^{2d_{\mathbf{g}} \times 2}$$

where $\mathbf{b} = -\mathbf{a} \cdot \mathbf{s} + \mathbf{e}$ and the gadget matrix $\mathbf{G} = (\mathbf{I} \parallel B_{\mathbf{g}}\mathbf{I} \parallel \dots \parallel B_{\mathbf{g}}^{d_{\mathbf{g}}-1}\mathbf{I})^T \in \mathcal{R}_Q^{2d_{\mathbf{g}} \times 2}$ for 2×2 identity matrix \mathbf{I} .

Let $\text{WD}_{B_{\mathbf{g}}}(\cdot)$ be the decomposition with the base $B_{\mathbf{g}}$, where the dimension of input vector is multiplied by $d_{\mathbf{g}}$ through this algorithm. The RGSW encryption of m satisfies $\text{CT} \cdot (1, \mathbf{s}) = m \cdot (1, \mathbf{s}, \dots, B_{\mathbf{g}}^{d_{\mathbf{g}}-1}, B_{\mathbf{g}}^{d_{\mathbf{g}}-1}\mathbf{s}) + \mathbf{e}$. Roughly, m is an *approximate* eigenvalue of $\text{WD}_{B_{\mathbf{g}}}(\text{CT})$ with respect to the eigenvector $(1, \mathbf{s}, \dots, B_{\mathbf{g}}^{d_{\mathbf{g}}-1}, B_{\mathbf{g}}^{d_{\mathbf{g}}-1}\mathbf{s})$. In [CGGI16], the hybrid multiplication between RGSW ciphertexts and RLWE ciphertexts has been defined as follows.

- Hybrid.Mult(CT, ct): Given an RGSW ciphertext $\text{CT} \in \mathcal{R}_Q^{2d_{\mathbf{g}} \times 2}$ and an RLWE ciphertext $\text{ct} \in \mathcal{R}_Q^2$ output the vector $\text{ct}' \leftarrow \text{CT}^T \cdot \text{WD}_{B_{\mathbf{g}}}(\text{ct})$.

If CT and ct are RGSW and RLWE encryptions of m and m' , respectively, their multiplication ct' is a valid RLWE encryption of mm' . For convenience, we will denote Hybrid.Mult(CT, ct) algorithm by \boxtimes , i.e., $(\text{CT}, \text{ct}) \in \mathcal{R}_Q^{2d_{\mathbf{g}} \times 2} \times \mathcal{R}_Q^2 \mapsto \text{CT} \boxtimes \text{ct} \in \mathcal{R}_Q^2$.

3 Privacy-preserving Database Searching and Extraction

Let us consider a database of a set of n tuples. Each tuple consists of pairs (d_i, α_i) for $i = 1, \dots, n$, where d_i denotes a *data-tag* in the domain $\{0, 1, \dots, \mathcal{T} - 1\}$ and α_i represents the corresponding *value attribute* in a plaintext space $\mathbb{Z}_t \setminus \{0\}$. Note that all the tags should be distinct from each other. For instance, in the case of personal information database, α_i may be the age of user whose identity number is d_i .

Given a *query tag* d from a tag domain and a *query value* α from a plaintext space, the matching problem is to determine the existence of an index i such that $(d, \alpha) = (d_i, \alpha_i)$. Now consider the following simplified search query: select α_i if there exists an index i such that $d_i = d$; otherwise zero (\perp). The purpose of this section is to store the database and carry out this search query on the public cloud. The server should learn nothing from encrypted query and any information other than the final result should not be leaked to user. Throughout this work, we will use semi-honest (honest but curious) adversary model, which is a standard assumption for evaluation of homomorphic encryption.

3.1 Our Method of Secure Search-and-extract

Our main idea is the following encoding method of database suitable for the efficient computation of equality test and extraction:

$$\text{DB}(X) = \sum_i \alpha_i X^{d_i} \in \mathbb{Z}_t[X].$$

The user encrypts this polynomial with the RLWE public-key encryption scheme and stores the ciphertext ct_{DB} in the server. At the query phase, given a query tag d , the user encrypts the monomial X^{-d} with the RGSW symmetric encryption scheme and sends the ciphertext CT_{Q} to the server. We assume that the RGSW encryption scheme has the same secret key sk as the one of RLWE encryption scheme.

Given two ciphertexts $\text{CT}_{\text{Q}} \leftarrow \text{RGSW.Enc}(X^{-d})$ and $\text{ct}_{\text{DB}} \leftarrow \text{RLWE.Enc}(\text{DB}(X))$, the server first performs their multiplication to obtain an ciphertext, denoted by $\text{ct}_{\text{mult}} = \text{CT}_{\text{Q}} \boxtimes \text{ct}_{\text{DB}}$. It follows from the previous section that ct_{mult} is a valid RLWE encryption of the polynomial

$$\text{DB}(X) \cdot X^{-d} = \sum_i \alpha_i X^{d_i - d} \in \mathcal{R}_t.$$

Since we use the cyclotomic polynomial $\phi_M(X) = X^N + 1$ of power-of-two degree, the polynomial ring \mathcal{R} has the property $X^N = -1$. Thus, for any tag d , the constant term of the polynomial $\text{DB}(X) \cdot X^{-d}$ is α_i if there is some index i satisfying $d = d_i$, otherwise zero.

Now the server applies the RLWE.Conv algorithm on ct_{mult} to compute an LWE encryption ct_{conv} of this constant term. This conversion procedure not only prevents the leakage of information that has not been queried but also reduces the size of output ciphertext by half. In addition, the (optional) modulus-switching procedure can be considered to get a ciphertext ct_{res} with a smaller modulus size and reduce the communication cost. Finally the user decrypts this LWE ciphertext and gets the desired value α_i or zero (\perp). Algorithm 1 summarizes the procedure of secure *search-and-extraction*.

Our method can be modified to support a secure comparison of data values using a hash (one-way) function. If hashed values of α_i are used as polynomial coefficients, our method will return a hashed value of α_i to the user instead of α_i . The user may check whether the resulting value and the hashed query value are the same or not without knowing information about database.

Algorithm 1 Procedure of secure search-and-extraction

- 1: Database encryption: The data owner encodes the genomic information as $DB(X)$ and submits its encryption to the server:

$$ct_{DB} \leftarrow RLWE.Enc(DB(X)).$$

- 2: Query encryption: The user encodes the query tag d and sends its encryption to the server:

$$CT_Q \leftarrow RGSW.Enc(X^{-d}).$$

- 3: Evaluation phase: The server computes their multiplication, and carries out the conversion and modulus-switching operations:

$$ct_{mult} \leftarrow Hybrid.Mult(CT_Q, ct_{DB}).$$

$$ct_{conv} \leftarrow RLWE.Convert(ct_{mult}).$$

$$ct_{res} \leftarrow LWE.ModSwitch(ct_{conv}).$$

Return the resulting ciphertext ct_{res} to the user.

- 4: Decryption phase: The user decrypts the ciphertext with the secret key and gets the desired value:

$$\alpha \leftarrow LWE.Dec(ct_{res}).$$

3.2 Comparison with Related Work

Equality test has been traditionally considered difficult to perform on homomorphic encryption, because of its large circuit depth [CKK15, KL15, CKK16]. They evaluate the equality test on each encrypted tuple of database, so at least $\Omega(n)$ homomorphic operations are required for searching on database of size n . In addition, Boneh et al. [BGH⁺13] does not protect the database information to the users, that is, the whole database can be recovered by the resulting ciphertext of a query. However, our method is very efficient in parameter size and complexity since it requires only a single hybrid multiplication.

One limitation of this method is that the tags d_i should be bounded by ciphertext dimension N to construct the encoding polynomial $DB(X)$. Since the dimension N has a significant influence on the performance of HE scheme, too large value of N has an impractical impact on the performance. In the next section, we will describe how to overcome this problem in terms of the application to genomic data.

4 Secure Searching of Biomarkers

We return to our main goal of task3: secure outsourcing matching of a set of biomarkers to encrypted genomes. We describe how to encode and encrypt the genotype information of VCF file in order to apply the privacy-preserving database searching and extraction.

4.1 Encoding of Genomic Data

VCF file contains multiple genotype information lines, where each of them consists of a triple $(ch_i, pos_i, SNPs_i)$ of chromosome number, position, and a sequence of SNP alleles. A chromosome identifier ch ranges from 1 to 22, X, and Y. A non-negative integer pos represents the reference position with the first base having position 1, and $SNPs$ is a reference or

alternate sequence in $\{A, T, G, C\}^*$. A query from user is also a triple of the same form and we aim to decide absence/presence of this biomarker in the database file.

We represent the sex chromosomes X and Y as 0 and 23, respectively. Then we define an encoding function $\mathcal{E} : \mathbb{Z} \times \mathbb{Z} \rightarrow \mathbb{Z}$ by

$$(\text{ch}, \text{pos}) \mapsto d = \text{ch} + 24 \cdot \text{pos}.$$

In the following, we describe how to encode the SNPs. For convenience we set the upper bound for the length of SNPs, so let n_{SNP} be the maximal number of reference (or alternate) alleles to be compared between the query genome and user genome in the target database. Each of SNP is represented by two bits as

$$A \mapsto 00, T \mapsto 01, G \mapsto 10, C \mapsto 11,$$

and then concatenated with each other. Next we pad with 1 to the left of the bit string in order to express the starting position of SNPs. Finally it is zero-padded into a binary string of length $\ell_{\text{SNP}} = 2 \cdot n_{\text{SNP}} + 1$, and we convert it into an integer value, denoted by α_i . If a single nucleotide variant at the given locus is not known, then it is encoded as 0-string. For example, ‘GC’ is encoded as a bit string 1|10|11, which will be represented as an integer $1|10|11_{(2)} = 27$.

Now consider the case that we wish to encode the reference and alternate alleles together. Let α_i^{ref} and α_i^{alt} denote the integer encodings of n_{SNP} reference alleles and n_{SNP} alternate alleles, respectively. Then we define an encoding α_i by the concatenation of two encodings, *i.e.*, $\alpha_i = 2^{\ell_{\text{SNP}}} \cdot \alpha_i^{\text{ref}} + \alpha_i^{\text{alt}}$ as an integer. Table 1 shows the format of database file and illustrates some examples of encoded genomic data.

Table 1. The format of genome data and its encoding with $n_{\text{SNP}} = 10$

CHROM	POS	d	REF	ALT	α
1	161235340	3869648161	<i>G</i>	<i>A</i>	12582916
1	161235596	3869654305	<i>C</i>	<i>T</i>	14680069
1	161235657	3869655769	<i>G</i>	<i>T</i>	12582917
1	161235981	3869663545	<i>G</i>	<i>A</i>	12582916
1	161237503	3869700073	.	<i>TTTTTGT</i>	21849
1	161237891	3869709385	<i>G</i>	<i>A</i>	12582916
1	161238009	3869712217	<i>G</i>	.	12582912
1	161238488	3869723713	<i>A</i>	<i>G</i>	8388614
1	161238683	3869728393	<i>G</i>	<i>A</i>	12582916
1	161238856	3869732545	<i>T</i>	.	10485760
1	161239028	3869736673	<i>AG</i>	.	37748736
1	161239142	3869739409	<i>A</i>	<i>G</i>	8388614
1	161239346	3869744305	<i>G</i>	<i>T</i>	12582917
1	161239470	3869747281	<i>C</i>	<i>T</i>	14680069
1	161239788	3869754913	.	<i>AA</i>	16
1	161239978	3869759473	<i>C</i>	<i>T</i>	12582917
1	161240641	3869775385	<i>TGAT</i>	.	740294656

4.2 Encryption of Genomic Data

A database file is encoded as a set of pair (d_i, α_i) for $i = 1, \dots, n$ such that $d_i = \mathcal{E}(\text{ch}_i, \text{pos}_i)$ and α_i is the encoded integer of the i -th SNP allele string. Then the encodings d_i and α_i are regarded as data-tag and value attribute, respectively. The data user constructs a polynomial $\text{DB}(X) = \sum_k c_k X^k$ such that

$$c_k = \begin{cases} \alpha_i & \text{if } k = d_i \text{ for some } i, \\ \alpha \leftarrow \mathbb{Z}_t & \text{otherwise.} \end{cases}$$

The user encrypts the polynomial with the RLWE public-key encryption scheme as described above.

The query genes are also encoded as a pair of integers (d, α) , however, we consider only the information of d is encrypted using the RGSW symmetric encryption scheme, that is, the user encrypts the monomial X^{-d} .

5 Implementation and Experimental Results

In this section, we explain how to set the parameters and describe our optimization techniques for the implementation. We also present our results using the techniques. The dataset was randomly selected from Personal Genome Project. Our implementation is publicly available on github [KS16].

5.1 How to Set Parameters

The security of the underlying homomorphic encryption scheme relies on the hardness of the RLWE assumption. We derive a lower-bound on the ring dimension as $N \geq \frac{\lambda+110}{7.2} \cdot \log_2 Q$ to get λ -bit security level from the security analysis of [GHS12].

Given the ciphertext modulus Q , it follows from the estimation of noise growth during evaluations [DM15] and decryption condition that we get the upper bound on the plaintext modulus t to ensure the correctness of decryption after computation. So we set t as the largest power-of-two integer less than the upper bound. If the encodings of the allele strings are too large, we divide them into smaller integers so that each of them is smaller than t . Then we repeat the algorithm to construct the corresponding polynomials of each integer.

5.2 Optimization Techniques

As we mentioned before, the ring dimension N needs to be larger than the encoded integers d_i 's. However, the encoded integers d_i from VCF files have bits size about 32, while a dimension N with about $11 \leq \log_2 N \leq 16$ is considered appropriate for implementation of HE schemes to achieve both security and efficiency. Hence direct application of our method to the VCF file would yield an impractical result.

For compression of tag data and its re-randomization, we make the use of a pseudo random number generator $H(\cdot)$ which transforms a tag d_i into a pair of two non-negative integers d_i^* and d_i^\dagger less than N . Our implementation adopts SHA-3 and extracts $\log_2 N = 11$ bits of the hashed value for each of d_i^* and d_i^\dagger .

We construct two polynomials

$$\text{DB}^*(X) = \sum_k c_k^* X^k, \quad \text{DB}^\dagger(X) = \sum_k c_k^\dagger X^k$$

by the Algorithm 2. Note that for any $1 \leq i \leq n$ and $H(d_i) = (d_i^*, d_i^\dagger) \in \{0, \dots, N-1\}^2$, the pair of constructed polynomials DB^* and DB^\dagger satisfy $\alpha_i = c_{d_i^*} + c_{d_i^\dagger}$. The procedure of database encoding for secure search of biomarkers is described in Algorithm 2.

Algorithm 2 Encoding genomic data

```

1:  $c_{d_1^*}^* \leftarrow \alpha_1 \in \mathbb{Z}_q, c_{d_1^\dagger}^\dagger \leftarrow \alpha_1 - c_{d_1^*}^*$ 
2:  $d_1^* \in \mathcal{D}^*, d_1^\dagger \in \mathcal{D}^\dagger$ 
3: for  $i \in \{2, \dots, n\}$  do
4:   if  $d_i^* \notin \mathcal{D}^*$  and  $d_i^\dagger \notin \mathcal{D}^\dagger$  then
5:      $c_{d_i^*}^* \leftarrow \alpha_i \in \mathbb{Z}_q, c_{d_i^\dagger}^\dagger \leftarrow \alpha_i - c_{d_i^*}^*$ 
6:   else if  $d_i^* \in \mathcal{D}^*$  and  $d_i^\dagger \notin \mathcal{D}^\dagger$  then
7:      $c_{d_i^\dagger}^\dagger \leftarrow \alpha_i - c_{d_i^*}^*$ 
8:   else if  $d_i^* \notin \mathcal{D}^*$  and  $d_i^\dagger \in \mathcal{D}^\dagger$  then
9:      $c_{d_i^*}^* \leftarrow \alpha_i - c_{d_i^\dagger}^\dagger$ 
10:  end if
11:   $d_i^* \in \mathcal{D}^*, d_i^\dagger \in \mathcal{D}^\dagger$ 
12: end for
13: return  $\text{DB}^*(X) = \sum_k c_k^* X^k, \text{DB}^\dagger(X) = \sum_k c_k^\dagger X^k$ 

```

Let ct_{DB}^* and $\text{ct}_{\text{DB}}^\dagger$ denote the ciphertexts of the polynomials DB^* and DB^\dagger , respectively. Similarly, given the query encoding d , the user computes its randomized value $H(d) = (d^*, d^\dagger)$ and encrypts the two polynomials X^{-d^*} and X^{-d^\dagger} . We denote the ciphertexts by CT_Q^* and CT_Q^\dagger . The server computes the hybrid multiplication to obtain the ciphertexts

$$\text{ct}_{\text{mult}}^* = \text{CT}_Q^* \boxtimes \text{ct}_{\text{DB}}^*, \text{ct}_{\text{mult}}^\dagger = \text{CT}_Q^\dagger \boxtimes \text{ct}_{\text{DB}}^\dagger.$$

Now let ct denote the ciphertext computed by the homomorphic addition between $\text{ct}_{\text{mult}}^*$ and $\text{ct}_{\text{mult}}^\dagger$. Finally the server converts it into an LWE ciphertext and performs the modulus-switching procedure as described above. The Algorithm 3 describes the procedure of secure *search-and-extraction* using our proposed optimization techniques.

5.3 Implementation Results

The use of variable type ‘*int32_t*’ accelerates the speed of implementations and basic C++ std libraries, so we set $Q = 2^{32}$ as the ciphertext modulus. We also set $t = 2^{11}$ as the modulus parameter of the plaintext space to ensure the correctness for the output ciphertext. We take the following parameters for Gadget matrix \mathbf{G} : $B_g = 128$ and $d_g = 5$, so that they satisfy the condition $B_g^{d_g} \geq Q$.

Each coefficient of the secret key sk is chosen at random from $\{0, \pm 1\}$ and we set 64 as the number of nonzero coefficients in the secret key. As in the work of [DM15], we considered the Gaussian distribution of standard deviation $\sigma = 1.4$ to sample random error polynomials.

For the efficiency of homomorphic multiplication, we also used the optimized library for complex FFT, *i.e.*, the *Fast Fourier Transform in the West* [FJ05]. That is, we use the complex primitive $2N$ -th root of unity rather than a primitive root in a prime field of order

Algorithm 3 Procedure of optimized secure search of biomarkers

- 1: Database encryption: The data owner encodes the genomic information as $DB^*(X)$ and $DB^\dagger(X)$. Then the user submits the ciphertexts to the server:

$$\begin{aligned} ct_{DB}^* &\leftarrow \text{RLWE.Enc}(DB^*(X)), \\ ct_{DB}^\dagger &\leftarrow \text{RLWE.Enc}(DB^\dagger(X)). \end{aligned}$$

- 2: Query encryption: The user encodes the query as X^{-d^*} and X^{-d^\dagger} . Then the user sends the ciphertexts to the server:

$$\begin{aligned} CT_Q^* &\leftarrow \text{RGSW.Enc}(X^{-d^*}), \\ CT_Q^\dagger &\leftarrow \text{RGSW.Enc}(X^{-d^\dagger}). \end{aligned}$$

- 3: Evaluation phase: The server computes their multiplications:

$$ct_{\text{mult}}^* \leftarrow CT_Q^* \boxtimes ct_{DB}^*, \quad ct_{\text{mult}}^\dagger \leftarrow CT_Q^\dagger \boxtimes ct_{DB}^\dagger.$$

Let $ct \leftarrow ct_{\text{mult}}^* + ct_{\text{mult}}^\dagger$. The server converts it into an LWE ciphertext and performs modulus-switching operations:

$$\begin{aligned} ct_{\text{conv}} &\leftarrow \text{RLWE.Convert}(ct), \\ ct_{\text{res}} &\leftarrow \text{LWE.ModSwitch}(ct_{\text{conv}}). \end{aligned}$$

Return the resulting ciphertext ct_{res} to the user.

- 4: Decryption phase: The user decrypts the ciphertext with the secret key and gets the desired value:

$$\alpha \leftarrow \text{LWE.Dec}(ct_{\text{res}}).$$

Q. We measure a running time of 0.804 seconds to set up the FFT environment at dimension $2N = 2^{12}$. The key generation of two schemes takes about 0.247 ms in total.

Table 2 presents the time complexity and storage for the evaluation of secure searching of biomarkers. All the experiments were performed on a single Intel Core i5 running at 2.9 GHz processor. The chosen parameters provide $\lambda = 128$ bits of security level.

6 Conclusions

In this work, we suggested an efficient method to securely search the query tag and extract the corresponding value from a database over hybrid GSW homomorphic encryption scheme. We came up with a solution to the secure outsourcing matching problem by using polynomial encoding and extraction of desired value based on the multiplication of an RGSW ciphertext and an ordinary RLWE ciphertext. And then we applied this method to find a set of biomarkers in DNA sequences.

Our solution shows the progress of cryptographic techniques in terms of their capability can support real-world genome data analysis in a cloud environment. We list a few fascinating open problems to remain. First, we only considered the semi-honest adversary model in this work. Other tools such as homomorphic authenticated scheme may lead to more efficient protocols in the malicious settings. Another issue is to support k multiple queries while

Table 2. Implementation results of secure searching of biomarkers

DB size	n_{SNP}	Complexity				Storage		
		Query-enc	DB-enc	Eval	Dec	Query	DB	Result
10K	2	3.247ms	3.563ms	0.018s	0.004ms	160KB	3MB	0.75MB
	5		7.212ms	0.039s	0.011ms		6MB	1.5MB
	10		14.813ms	0.079s	0.027ms		12MB	3MB
100K	2		21.424ms	0.111s	0.034ms		17MB	4.25MB
	5		42.415ms	0.227s	0.064ms		34MB	8.5MB
	10		99.921ms	0.454s	0.139ms		68MB	17MB
4M	2		0.745s	3.954s	1.171ms		593MB	148MB
	5		1.506s	7.911s	1.949ms		1185MB	296MB
	10		3.001s	15.442s	3.795ms		2370MB	593MB

maintaining the performance and communication cost less than k times of a single query case. We expect to have much faster performance by applying more efficient method.

Acknowledgements. This work was supported by IT R&D program of MSIP/KEIT [No. B0717-16-0098]. The authors would like to thank the referee for helpful comments. The authors would also like to thank the iDASH Secure Genome Analysis Contest organizers, in particular Xiaoqian Jiang and Shuang Wang, for running the contest and providing the opportunity to submit competing implementations for these important tasks.

References

- [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*, pages 309–325. ACM, 2012.
- [Bra12] Zvika Brakerski. Fully homomorphic encryption without modulus swithing from classical gapsvp. In *Advances in Cryptology–CRYPTO 2012*, volume 7417, pages 868–886. Springer, 2012.
- [CGGI16] Ilaria Chillotti, Nicolas Gama, Mariya Georgieva, and Malika Izabachène. Faster fully homomorphic encryption: Bootstrapping in less than 0.1 seconds. In *Advances in Cryptology–ASIACRYPT 2016: 22nd International Conference on the Theory and Application of Cryptology and Information Security, Hanoi, Vietnam, December 4-8, 2016, Proceedings, Part I*, pages 3–33. Springer, 2016.
- [CKK15] Jung Hee Cheon, Miran Kim, and Myungsun Kim. Search-and-compute on encrypted data. In *International Conference on Financial Cryptography and Data Security*, pages 142–159. Springer, 2015.
- [CKK16] Jung Hee Cheon, Miran Kim, and Myungsun Kim. Optimized search-and-compute circuits and their application to query evaluation on encrypted data. *IEEE Transactions on Information Forensics and Security*, 11(1):188–199, 2016.
- [CKL15] Jung Hee Cheon, Miran Kim, and Kristin Lauter. Homomorphic computation of edit distance. In *International Conference on Financial Cryptography and Data Security*, pages 194–212. Springer, 2015.
- [DM15] Léo Ducas and Daniele Micciancio. Fhew: Bootstrapping homomorphic encryption in less than a second. In *Advances in Cryptology–EUROCRYPT 2015*, volume 9056, pages 617–640. Springer, 2015.

- [EN14] Yaniv Erlich and Arvind Narayanan. Routes for breaching and protecting genetic privacy. *Nature Reviews Genetics*, 15(6):409–421, 2014.
- [FJ05] Matteo Frigo and Steven G Johnson. The design and implementation of fftw3. *Proceedings of the IEEE*, 93(2):216–231, 2005.
- [GHS12] Craig Gentry, Shai Halevi, and Nigel P Smart. Homomorphic evaluation of the aes circuit. In *Advances in Cryptology–CRYPTO 2012*, volume 7417, pages 850–867. Springer, 2012.
- [GSW13] Craig Gentry, Amit Sahai, and Brent Waters. Homomorphic encryption from learning with errors: Conceptually-simpler, asymptotically-faster, attribute-based. In *Advances in Cryptology–CRYPTO 2013*, volume 8042, pages 75–92. Springer, 2013.
- [HAHT13] Mathias Humbert, Erman Ayday, Jean-Pierre Hubaux, and Amalio Telenti. Addressing the concerns of the lacks family: quantification of kin genomic privacy. In *Proceedings of the 2013 ACM SIGSAC conference on Computer & communications security*, pages 1141–1152. ACM, 2013.
- [KL15] Miran Kim and Kristin Lauter. Private genome analysis through homomorphic encryption. *BMC medical informatics and decision making*, 15(Suppl 5):S3, 2015.
- [KS16] Miran Kim and Yongsoo Song. *Implementation of secure searching of biomarkers*, 2016. <http://github.com/amedonis/HybridHE>.
- [LLAN14] Kristin Lauter, Adriana López-Alt, and Michael Naehrig. Private computation on encrypted genomic data. In *International Conference on Cryptology and Information Security in Latin America*, pages 3–27. Springer, 2014.
- [LPR10] Vadim Lyubashevsky, Chris Peikert, and Oded Regev. On ideal lattices and learning with errors over rings. In *Advances in Cryptology - EUROCRYPT 2010*, pages 1–23, 2010.
- [NAC⁺15] Muhammad Naveed, Erman Ayday, Ellen W Clayton, Jacques Fellay, Carl A Gunter, Jean-Pierre Hubaux, Bradley A Malin, and XiaoFeng Wang. Privacy in the genomic era. *ACM Computing Surveys (CSUR)*, 48(1):6, 2015.
- [YSK⁺13] M Yasuda, T Shimoyama, J Kogure, K Yokoyama, and T Koshihara. Secure pattern matching using somewhat homomorphic encryption. In *Proceedings of the 2013 ACM Cloud Computing Security Workshop*, pages 65–76, 2013.