

# Privacy-preserving Approximate GWAS computation based on Homomorphic Encryption

Duhyeong Kim\*, Yongha Son, Dongwoo Kim, Andrey Kim, Seungwan Hong,  
and Jung Hee Cheon

Department of Mathematical Sciences, Seoul National University  
doodoo1204@snu.ac.kr

**Abstract.** One of three tasks in a secure genome analysis competition called IDASH 2018 was to develop a solution for privacy-preserving GWAS computation based on homomorphic encryption. The scenario is that a data holder encrypts a number of individual records, each of which consists of several phenotype and genotype data, and provide the encrypted data to an untrusted server. Then, the server performs a GWAS algorithm based on homomorphic encryption without the decryption key and outputs the result in encrypted state so that there is no information leakage on the sensitive data to the server.

In this paper, we develop a privacy-preserving semi-parallel GWAS algorithm by applying an approximate homomorphic encryption scheme HEAAN. Fisher scoring and semi-parallel GWAS algorithms are modified to be efficiently computed over homomorphically encrypted data with several optimization methodologies; substitute matrix inversion by an adjoint matrix, avoid computing a superfluous matrix of super-large size, and transform the algorithm into an approximate version.

Our modified semi-parallel GWAS algorithm based on homomorphic encryption which achieves 128-bit security takes 30–40 minutes for 245 samples containing 10,000–15,000 SNPs. Compared to the true  $p$ -value from the original semi-parallel GWAS algorithm, the  $F_1$  score of our  $p$ -value result is over 0.99.

## 1 Introduction

After the successful completion of the Human Genome Project in the early 21st century, high throughput technology on genetic variations has been rapidly developed and widely studied. In particular, through the development of microarray chip with rather small computational cost, it became possible to determine the genotype of millions of single nucleotide polymorphism (SNP), a variation in a single nucleotide that occurs at a specific position in the genome, for each individual. With those statistical data of genotypes, many researches are proposed that investigate associations between SNPs and phenotypes like major human disease, and especially Genome-wide association study (GWAS) aims to find top significant SNPs relevant to a certain phenotype.

## 1.1 Motivation

Since genome analysis uses genomic data that are very sensitive and irreplaceable, privacy on genomic data has come up to be one of the most important issues in genome analysis including GWAS. The usual privacy-preserving methodology in data analysis is anonymization, perturbation, randomization, and condensation [20]; however, those methods leverage the quality of data with privacy resulting in an inaccurate analysis to some extent. The dilemma regarding the balance between personal privacy and analytical efficiency has been resolved by applying many cryptographic primitives, while homomorphic encryption (HE) is noticed as one of the ultimate cryptographic solutions for privacy-preserving data analysis. Conceptually, HE is an encryption scheme which allows computations over encrypted data without decryption. HE not only fundamentally prevents the leakage of input data during the analysis phase, but also provides an accurate result of analysis since it preserves the original data intactly. However, HE causes a significant blowup of computational cost for analysis, and optimization and modification of algorithm for efficient computation in HE is the main problem of applying HE in data analysis.

Since 2014, there has been an annual biomedical privacy competition hosted by Integrating Data for Analysis, Anonymization and SHaring (IDASH), a national center for biomedical computing in the United States. One of three tasks in IDASH 2018 [2] was to develop a solution for privacy-preserving GWAS computation based on HE, and we participated in this competition with our delicately constructed algorithms.

## 1.2 Summary of Results

In this study, we propose approximate HE algorithms for privacy-preserving GWAS computation. To be precise, we transform well-known Fisher scoring and semi-parallel GWAS algorithm into *HE-friendly* algorithms so that we can efficiently evaluate them in encrypted state. Note that the HE-friendly modified Fisher Scoring algorithm can be generally used for logistic regression, not only for GWAS.

The main challenges in transforming the semi-parallel GWAS algorithm (Algorithm 1) to an HE algorithm are complex matrix operations such as multiplication and inversion. Since matrix inversion in HE is complicated and costly, we substitute it by computation of the adjoint matrices and determinant. With this approach, the original Fisher scoring algorithm can be successfully modified to compute encrypted data efficiently. For the efficient computation of semi-parallel GWAS computation based on HE, moreover, we reduced the number of matrix multiplications as many as possible, and modified the original algorithm into an approximate version which requires much less computational cost. The details of our optimization methodologies are well described in Section 3.

We exploited an approximate HE scheme HEAAN [9,8] with a publicly available library [15] for the implementation of our modified semi-parallel GWAS

algorithm based on HE. The HE algorithm takes about 40 minutes for 245 samples each containing a binary phenotype, 3 covariates, and 14,841 SNPs on Linux with a 2.10GHz processor.

### 1.3 Related Works

Recently Cho, Wu and Berger [12] proposed a secure GWAS computation method in Nature Biotechnology, where the privacy of genomic data was preserved by a cryptographic tool called multiparty computation rather than HE. To the best of our knowledge, there have been no results on privacy-preserving GWAS computation based on HE before the IDASH 2018 competition. In perspective of HE-based genome analysis, the solutions [17,6,13,5] submitted to task 3 of IDASH 2017 [1] dealt with training the logistic regression model of genomic data based on HE.

## 2 Backgrounds

### 2.1 An Approximate Homomorphic Encryption Scheme HEAAN

For privacy-preserving GWAS computation, we applied an HE scheme called HEAAN proposed by Cheon et al. [9,8], which supports *approximate computation* of real numbers in encrypted state. Efficiency of HEAAN in the real world has been proved by showing its application in various fields including machine learning [17,18,10] and cyber physical system [7]. The winning solution of IDASH competition in 2017 also applied HEAAN as an HE scheme for privacy-preserving logistic regression on genomic data.

In detail, let  $\text{ct}$  be a HEAAN ciphertext of a message polynomial  $\mathbf{m}$ . Then, the decryption process with a secret key  $\text{sk}$  is done as

$$\text{Dec}_{\text{sk}}(\text{ct}) = \mathbf{m} + e \approx \mathbf{m}$$

where  $e$  is a small error attached to the message polynomial  $\mathbf{m}$ . For formal definitions, let  $L$  be a level parameter, and  $q_\ell := 2^\ell$  for  $1 \leq \ell \leq L$ . Let  $R := \mathbb{Z}[X]/(X^N + 1)$  for a power-of-two  $N$  and  $R_q$  be a modulo- $q$  quotient ring of  $R$ , i.e.,  $R_q := R/qR$ . The distribution  $\chi_{\text{key}} := \text{HW}(\mathbf{h})$  over  $R_q$  outputs a polynomial of  $\{-1, 0, 1\}$ -coefficient having  $h$  number of non-zero coefficients, and  $\chi_{\text{enc}}$  and  $\chi_{\text{err}}$  denote the discrete Gaussian distribution with some prefixed standard deviation. Finally,  $[\cdot]_q$  denotes a component-wise modulo  $q$  operation on each element of  $R_q$ . Note that those parameters  $N$ ,  $L$  and  $h$  satisfying a certain security level can be determined by Albrecht's security estimator [4,3]. The scheme description of HEAAN is as following:

- **KeyGen(params)**.
  - Sample  $s \leftarrow \chi_{\text{key}}$ . Set the secret key as  $\text{sk} \leftarrow (1, s)$ .
  - Sample  $a \leftarrow U(R_{q_L})$  and  $e \leftarrow \chi_{\text{err}}$ . Set the public key as  $\text{pk} \leftarrow (b, a) \in R_{q_L}^2$  where  $b \leftarrow [-a \cdot s + e]_{q_L}$ .

- Sample  $a' \leftarrow U(R_{q_L^2})$  and  $e' \leftarrow \chi_{\text{err}}$ . Set the evaluation key as  $\text{evk} \leftarrow (b', a') \in R_{q_L^2}^2$  where  $b' \leftarrow [-a's + e' + q_L \cdot s^2]_{q_L^2}$ .
- $\text{Enc}_{\text{pk}}(\mathbf{m})$ . For a message  $\mathbf{m} \in R$ , sample  $v \leftarrow \chi_{\text{enc}}$  and  $e_0, e_1 \leftarrow \chi_{\text{err}}$ . Output the ciphertext  $\text{ct} = [v \cdot \text{pk} + (\mathbf{m} + e_0, e_1)]_{q_L}$ .
- $\text{Dec}_{\text{sk}}(\text{ct})$ . For a ciphertext  $\text{ct} = (c_0, c_1) \in R_{q_\ell}^2$ , output a message  $\mathbf{m}' = \lfloor c_0 + c_1 \cdot s \rfloor_{q_\ell}$ .
- $\text{C.Add}(\text{ct}, \text{ct}')$ . For  $\text{ct}, \text{ct}' \in R_{q_\ell}^2$ , output  $\text{ct}_{\text{add}} \leftarrow [\text{ct} + \text{ct}']_{q_\ell}$ .
- $\text{C.Sub}(\text{ct}, \text{ct}')$ . For  $\text{ct}, \text{ct}' \in R_{q_\ell}^2$ , output  $\text{ct}_{\text{sub}} \leftarrow [\text{ct} - \text{ct}']_{q_\ell}$ .
- $\text{C.Mult}_{\text{evk}}(\text{ct}, \text{ct}')$ . For  $\text{ct} = (c_0, c_1), \text{ct}' = (c'_0, c'_1) \in R_{q_\ell}^2$ , let  $(d_0, d_1, d_2) = (c_0 c'_0, c_0 c'_1 + c_1 c'_0, c_1 c'_1)$ . Output  $\text{ct}_{\text{mult}} \leftarrow [(d_0, d_1) + \lfloor q_L^{-1} \cdot d_2 \cdot \text{evk} \rfloor]_{q_\ell}$ .

Let  $\text{ct}_1$  and  $\text{ct}_2$  be ciphertexts of message polynomials  $\mathbf{m}_1$  and  $\mathbf{m}_2$ . Then, the homomorphic evaluation algorithms  $\text{C.Add}$  and  $\text{C.Mult}$  satisfy

$$\begin{aligned} \text{Dec}_{\text{sk}}(\text{C.Add}(\text{ct}_1, \text{ct}_2)) &\approx \mathbf{m}_1 + \mathbf{m}_2, \\ \text{Dec}_{\text{sk}}(\text{C.Mult}_{\text{evk}}(\text{ct}_1, \text{ct}_2)) &\approx \mathbf{m}_1 \cdot \mathbf{m}_2, \end{aligned}$$

i.e., addition and multiplication can be *internally* done even in encrypted state. For more details of the scheme including the correctness and security analysis, we refer the readers to [9].

Since each message  $\mathbf{m} \in R$  is a  $\mathbb{Z}$ -coefficient polynomial, not a real number, there should be a conversion between polynomials and real numbers to encrypt real numbers. In this regard, we use a (field) isomorphism  $\tau$  from  $\mathbb{R}[X]/(X^N + 1)$  to  $\mathbb{C}^{N/2}$  called canonical embedding. A plaintext vector  $\mathbf{m} = (m_0, \dots, m_{N/2-1})$  is first transformed into  $\tau^{-1}(\mathbf{m}) \in \mathbb{R}[X]/(X^N + 1)$ , and then rounded off to an integer-coefficient polynomial. However, the naive rounding-off  $\lfloor \tau^{-1}(\mathbf{m}) \rfloor$  can derive quite large relative error on the plaintext. To control the error, we round it off after scaling up by  $p$  bits for some integer  $p$ , i.e.,  $\lfloor 2^p \cdot \tau^{-1}(\mathbf{m}) \rfloor$ , so that the relative error is reduced. Clearly, a decoding algorithm for  $\mathbf{m}$  would be  $2^{-p} \cdot \tau(\mathbf{m})$ :

- $\text{Ecd}(\mathbf{m}; p)$ . For  $\mathbf{m} = (m_0, \dots, m_{N/2-1})$  in  $\mathbb{C}^{N/2}$  and a precision bit  $p > 0$ , output a polynomial  $\mathbf{m} \leftarrow \lfloor 2^p \cdot \tau^{-1}(\mathbf{m}) \rfloor \in R$  where the rounding  $\lfloor \cdot \rfloor$  is done coefficient-wisely.
- $\text{Dcd}(\mathbf{m}; p)$ . For  $\mathbf{m} \in R$ , output a plaintext vector  $\mathbf{m}' = 2^{-p} \cdot \tau(\mathbf{m}) \in \mathbb{C}^{N/2}$ .

To sum up, to encrypt a plaintext vector of real (complex) numbers  $\mathbf{m}$ , we first encode  $\mathbf{m}$  into  $\mathbf{m} \leftarrow \text{Ecd}(\mathbf{m}; p)$  with a certain precision bit  $p$ , and then generate a ciphertext  $\text{ct} \leftarrow \text{Enc}_{\text{pk}}(\mathbf{m})$  with the public key  $\text{pk}$ .

Now consider  $\text{ct}_1$  and  $\text{ct}_2$  be ciphertexts of  $\mathbf{m}_1$  and  $\mathbf{m}_2$  in  $\mathbb{C}^{N/2}$ . Since our encoding method scales each plaintext vector up by  $2^p$ , the plaintext vector of a ciphertext  $\text{ct}' \leftarrow \text{C.Mult}_{\text{evk}}(\text{ct}_1, \text{ct}_2)$  is (approximately)  $2^p \cdot \mathbf{m}_1 \odot \mathbf{m}_2$ , not  $\mathbf{m}_1 \odot \mathbf{m}_2$ , which will result in exponential growth of plaintexts. To deal with this problem, we adjust the scaling factor by the following procedure so-called rescaling:

- $\text{RS}_{\ell \rightarrow \ell'}(\text{ct})$ . For a ciphertext  $\text{ct} \in R_{q_\ell}^2$ , output  $\text{ct}' \leftarrow \lfloor [(q_{\ell'}/q_\ell) \cdot \text{ct}] \rfloor_{q_{\ell'}}$ .

After the rescaling procedure  $\text{ct}_{\text{mult}} \leftarrow \text{RS}(\text{ct}')$ , the plaintext vector of the output  $\text{ct}_{\text{mult}}$  is (approximately)  $\mathbf{m}_1 \odot \mathbf{m}_2$ , and the ciphertext modulus  $q_L$  is reduced by  $2^p$ . As a result, the level parameter  $L$  should be carefully chosen according to the multiplicative depth of a target computation. In order to present our algorithm in a simple form, we will not describe these rescaling procedures for the rest of the paper, but we remark that in the actual use of HEAAN, there should be delicate consideration on scaling of message.

To deal with a plaintext vector of the form  $\mathbf{m} \in \mathbb{C}^k$  having length  $K \leq N/2$  for some power-of-two divisor  $K$  of  $N/2$ , HEAAN encrypts  $\mathbf{m}$  into a ciphertext of an  $N/2$ -dimensional vector  $(\mathbf{m} \parallel \dots \parallel \mathbf{m}) \in \mathbb{C}^{N/2}$ . This implies that a ciphertext of  $\mathbf{m} \in \mathbb{C}^K$  can be understood as a ciphertext of  $(\mathbf{m} \parallel \dots \parallel \mathbf{m}) \in \mathbb{C}^{K'}$  for powers-of-two  $K$  and  $K'$  satisfying  $K \leq K' \leq N/2$ .

Finally, the HEAAN scheme provides the rotation operation on plaintext slots, i.e., it enables us to securely obtain an encryption of the shifted plaintext vector  $(m_r, \dots, m_{N/2-1}, m_0, \dots, m_{r-1})$  from an encryption of  $(m_0, \dots, m_{N/2-1})$ . It is necessary to generate an additional public information  $\text{rk}$ , called the rotation key. We denote the rotation operation as follows.

- $\text{Rot}_{\text{rk}}(\text{ct}; r)$ . For the rotation key  $\text{rk}$ , output a ciphertext  $\text{ct}'$  encrypting the (left) rotated plaintext vector of  $\text{ct}$  by  $r (> 0)$  positions as above example. If  $r < 0$ , it denotes the right rotation by  $(-r)$  positions.

We omit a subscript of each algorithm of HEAAN for convenience if it is obvious.

## 2.2 Matrix packing method and Rotate function

In this subsection, we describe an encoding method to encrypt a matrix structure in a ciphertext which was also introduced in [17]. Consider an  $n \times m$  matrix  $Z$

$$Z = \begin{bmatrix} z_{0,0} & \cdots & z_{0,m-1} \\ \vdots & \ddots & \vdots \\ z_{n-1,0} & \cdots & z_{n-1,m-1} \end{bmatrix}.$$

We first pad zeros to set the number of rows and columns to be powers-of-two, say  $\bar{n}$  and  $\bar{m}$ , and assume that  $\log \bar{n} + \log \bar{m} \leq \log(N/2)$ . Then we pack the whole matrix in a single ciphertext  $\text{ct}_Z$  in a column-by-column manner. As described above, the algorithm  $\text{Rot}(\text{ct}_Z; r)$  can shift the encrypted vector by  $r$  positions. In particular, we can perform row and column rotations of an encrypted matrix with this operation. When  $r = \bar{n} \cdot j$ , and the result will be the (left) column rotation of the encrypted matrix  $Z$  by  $j$  columns.

For the row rotation of an encrypted matrix, we use so-called *masking approach*. Consider  $n \times m$  matrices  $M_i$  and  $\bar{M}_i$ , where the first  $n-i$  rows of  $M_i$  (resp.  $\bar{M}_i$ ) are filled with 1 (resp. 0) and the last  $i$  rows of  $M_i$  (resp.  $\bar{M}_i$ ) are filled with 0 (resp. 1). Let  $\text{msk}_i$  and  $\bar{\text{msk}}_i$  be the ciphertext of  $M_i$  and  $\bar{M}_i$ , respectively. For row rotation of an encrypted matrix  $Z$  by  $i$  rows, we first compute  $\text{ct}_1 \leftarrow \text{Rot}(\text{ct}_Z, i)$  and  $\text{ct}_2 \leftarrow \text{Rot}(\text{ct}_Z, i - n)$ . Then, we mask them as  $\text{ct}_1 \leftarrow \text{C.Mult}(\text{ct}_1, \text{msk}_i)$  and

$\text{ct}_2 \leftarrow \mathbf{C}.\text{Mult}(\text{ct}_2, \overline{\text{msk}_i})$ . As a result, the output of  $\mathbf{C}.\text{Add}(\text{ct}_1, \text{ct}_2)$  is a ciphertext of upper row rotation of  $Z$  by  $i$  rows.

Those row and column rotations of an encrypted matrix are denoted as follows:

- $\mathbf{C}.\text{ColumnRot}(\text{ct}_Z, j)$ . For a ciphertext  $\text{ct}_Z$  of a matrix  $Z$  and an integer  $j$ , output a ciphertext  $\text{ct}$  of left column rotation of  $Z$  by  $j$  columns.
- $\mathbf{C}.\text{RowRot}(\text{ct}_Z, i)$ . For a ciphertext  $\text{ct}_Z$  of a matrix  $Z$  and an integer  $j$ , output a ciphertext  $\text{ct}$  of upper row rotation of  $Z$  by  $i$  rows.

### 2.3 Semi-parallel GWAS Algorithm

A naive application of GWAS analysis can be done by running a logistic regression for each SNP, which resulting in high computational cost since the number of SNPs can be usually hundred thousands or more. To overcome this problem, Sikorska et al. [23] proposed a semi-parallel GWAS algorithm which reduces the required computation time from 6 hours to 10-15 minutes using projections.

Let  $n$  be the number of samples each of which consists of  $m$  (binary) SNP data and  $k'$  covariate data. Then the whole SNP data and covariate data can be organized as an  $n \times m$  matrix  $S$  and an  $n \times k'$  matrix  $X_0$ , respectively. For  $k := k' + 1$ , we define a matrix  $X$  as the concatenation of a vector whose components are 1 and  $X_0$ , denoted by  $X = (\mathbf{1} || X_0)$ . Let  $\mathbf{y}$  be a target binary phenotype vector of length  $n$ . With these  $\overrightarrow{\text{inputs}}$ , the semi-parallel GWAS algorithm outputs the  $m$ -dimensional vector  $\overrightarrow{\text{pval}}$  which indicates the  $p$ -value of each SNP with respect to the target phenotype. The detail of the algorithm is described in Algorithm 1.

---

#### Algorithm 1 The Original Semi-Parallel GWAS

---

**Input:** SNP matrix  $S \in \{0, 1\}^{n \times m}$ , covariate matrix  $X \in \mathbb{Q}^{n \times k}$ , phenotype vector  $\mathbf{y} \in \mathbb{Q}^n$ , and # iteration iter.

**Output:**  $p$ -value vector  $\overrightarrow{\text{pval}} = (\text{pval}_1, \dots, \text{pval}_m) \in \mathbb{Q}^m$ .

- 1:  $(\beta, \mathbf{p}, W) \leftarrow \text{FisherScoring}(X, \mathbf{y}; \text{iter})$  ▷ FisherScoring( $\cdot$ ) is described in Algorithm 2
  - 2:  $\mathbf{v} \leftarrow \log(\mathbf{p}/(\mathbf{1} - \mathbf{p})) + (\mathbf{y} - \mathbf{p})/\text{diag}(W)$
  - 3:  $S^* \leftarrow S - X(X^T W X)^{-1} X^T W S$
  - 4:  $\mathbf{v}^* \leftarrow \mathbf{v} - X(X^T W X)^{-1} X^T W \mathbf{z}$
  - 5:  $\mathbf{c} \leftarrow S^{*T} W \mathbf{v}^* \in \mathbb{Q}^m$
  - 6:  $\mathbf{d} \leftarrow \text{diag}(S^{*T} W S^*) \in \mathbb{Q}^m$
  - 7: **for**  $i = 1$  to  $m$  **do**
  - 8:      $a_i \leftarrow -|c_i/\sqrt{d_i}|$
  - 9:      $\text{pval}_i = 2 \cdot \int_{-\infty}^{a_i} \rho(x) dx$  ▷  $\rho(x) := \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}x^2)$
  - 10: **end for**
  - 11: **return**  $\overrightarrow{\text{pval}}$
- 

The semi-parallel GWAS algorithm involves logistic regression on  $(X, \mathbf{y})$  in the first step, and Fisher Scoring [19] described in Algorithm 2 is one of the most highly efficient algorithm for logistic regression.

---

**Algorithm 2** FisherScoring

---

**Input:** Covariate matrix  $X \in \mathbb{Q}^{n \times k}$ , phenotype vector  $\mathbf{y} \in \mathbb{Q}^n$ , and # iteration  $\text{iter}$ .

**Output:** Coefficient vector  $\boldsymbol{\beta} \in \mathbb{Q}^k$ , fitted vector  $\mathbf{p} \in \mathbb{Q}^n$ , weight matrix  $W \in \mathbb{Q}^{n \times n}$

- 1:  $\boldsymbol{\beta}^{(0)} = \mathbf{0} \in \mathbb{Q}^k$ ,  $\mathbf{p}^{(0)} = \overline{0.5} \in \mathbb{Q}^n$ ,  $W^{(0)} = 0.25 \cdot I_n$
  - 2: **for**  $t = 0$  to  $\text{iter} - 1$  **do**
  - 3:    $\mathbf{v}^{(t)} \leftarrow \log(\mathbf{p}^{(t)}/(\mathbf{1} - \mathbf{p}^{(t)})) + (\mathbf{y} - \mathbf{p}^{(t)})/\text{diag}(W^{(t)})$
  - 4:    $\boldsymbol{\beta}^{(t+1)} \leftarrow (X^T W^{(t)} X)^{-1} X^T W^{(t)} \mathbf{v}^{(t)}$
  - 5:    $\mathbf{p}^{(t+1)} \leftarrow \sigma(X \boldsymbol{\beta}^{(t+1)})$
  - 6:    $W^{(t+1)} \leftarrow \text{diagonal matrix of } \mathbf{p}^{(t+1)} \odot (\mathbf{1} - \mathbf{p}^{(t+1)})$
  - 7: **end for**
  - 8: **return**  $(\boldsymbol{\beta}^{(\text{iter})}, \mathbf{p}^{(\text{iter})}, W^{(\text{iter})})$
- 

In both algorithms,  $\sigma(x) := 1/(1 + \exp(-x))$  is called sigmoid function. For both algorithms, if the input of functions such as logarithm ( $\log$ ), division ( $/$ ) and sigmoid ( $\sigma$ ) is a vector, then it means to apply the function component-wisely resulting in the output vector of the same length. The notation  $T$  in the superscript denotes the matrix transpose. The operation  $\odot$  denotes the Hadamard (component-wise) multiplication of two vectors, and the notation  $\text{diag}(\cdot)$  with an input of a square matrix means the diagonal vector of the input.

### 3 Our Optimization Methodology

The aim of this study is to construct an HE algorithm for privacy-preserving semi-parallel GWAS computation. Since non-polynomial operations such as matrix inverse or real number inversion is a challenging stuff in HE, we need to modify the original semi-parallel GWAS algorithm into HE-friendly form for efficiency. Moreover, the super-large data size of GWAS requires too much computational cost in encrypted state, and this issue also should be resolved. In this regard, we introduce our optimization methodology to the algorithm.

#### 3.1 Modification of Fisher Scoring

The main obstacle of Fisher Scoring (Algorithm 2) is a matrix inversion for  $U = X^T W X \in \mathbb{Q}^{k \times k}$ . Our main idea to overcome this is the fact  $U^{-1} = \frac{1}{\det(U)} \cdot \text{adj}(U)$ . To be precise, observe that

$$\mathbf{v}^{(t)} = \log\left(\frac{\mathbf{p}^{(t)}}{\mathbf{1} - \mathbf{p}^{(t)}}\right) + \frac{\mathbf{y} - \mathbf{p}^{(t)}}{\text{diag}(W^{(t)})} = X \boldsymbol{\beta}^{(t)} + \frac{\mathbf{y} - \mathbf{p}^{(t)}}{\text{diag}(W^{(t)})},$$

from which we obtain an iterative updating equation on  $\boldsymbol{\beta}^{(t)}$  as follows:

$$\begin{aligned} \boldsymbol{\beta}^{(t+1)} &= U^{-1} X^T W^{(t)} \left( X \boldsymbol{\beta}^{(t)} + \frac{\mathbf{y} - \mathbf{p}^{(t)}}{\text{diag}(W^{(t)})} \right) \\ &= \boldsymbol{\beta}^{(t)} + U^{-1} X^T (\mathbf{y} - \mathbf{p}^{(t)}) \\ &= \boldsymbol{\beta}^{(t)} + \frac{1}{\det(U)} \cdot \text{adj}(U) X^T (\mathbf{y} - \mathbf{p}^{(t)}). \end{aligned}$$

Here one needs to compute the inverse of  $\det(U)$ , but this non-polynomial operation is rather expensive in HE. The key observation on this equation is that the second term  $U^{-1}X^T(\mathbf{y} - \mathbf{p}^{(t)})$  essentially converges to 0 as  $t \rightarrow \infty$  since  $\beta^{(t)}$  converges to some point. From this, we may expect that the convergence would be still valid even when we neglect the term  $\det(U)^{-1}$  and substitute it by some appropriate constant. Namely, we can modify the equation as

$$\beta^{(t+1)} = \beta^{(t)} + \alpha \cdot \text{adj}(U)X^T(\mathbf{y} - \mathbf{p}^{(t)}).$$

for some constant  $\alpha > 0$ . In practice, this approximate version of the Fisher scoring algorithm works quite well with slightly slower convergence rate.

### 3.2 Computation of $\text{diag}(S^{*T}WS^*)$ without computing $S^*$

The main observation of this subsection is that computation of  $n \times m$  matrix  $S^*$  in Algorithm 1 is superfluous for obtaining  $p$ -values. Indeed, we directly compute  $\det(U) \cdot \text{diag}(S^{*T}WS^*)$ , which can be obtained without computing (matrix) inversion and super-large matrix  $S^*$ . To be precise, using the fact that  $S^* = S - XU^{-1}V$  for  $V := X^TWS$ , we get

$$\begin{aligned} S^{*T}WS^* &= (S - XU^{-1}V)^T W (S - XU^{-1}V) \\ &= S^TWS - V^T U^{-1}V. \end{aligned}$$

Based on this observation, we compute  $\det(U) \cdot \text{diag}(S^{*T}WS^*)$  by following:

1. Compute  $U = X^TWX$  and  $V = X^TWS$ .
2. Compute  $\text{adj}(U)$  and  $\det(U)$ .
3. Compute  $\det(U) \cdot \text{diag}(S^TWS) - \text{diag}(V^T \text{adj}(U)V)$ .

### 3.3 Approximate Computation of $S^{*T}W\mathbf{v}^*$

We also take the main observation of the previous subsection so that we do not compute  $S^*$ . From the definition of  $S^*$  and  $\mathbf{v}^*$ , it holds that  $S^{*T}W\mathbf{v}^* = S^T W\mathbf{v}^*$ . Then we have the following equations:

$$\begin{aligned} S^T W\mathbf{v}^* &= S^T W(I - XU^{-1}X^T W)\mathbf{v} \\ &= S^T W(I - XU^{-1}X^T W) \frac{\mathbf{y} - \mathbf{p}}{\text{diag}(W)} \\ &= S^T(\mathbf{y} - \mathbf{p}) - S^T W X U^{-1} X^T (\mathbf{y} - \mathbf{p}) \\ &\simeq S^T(\mathbf{y} - \mathbf{p}) \end{aligned}$$

where the last approximation is valid since the term  $X^T(\mathbf{y} - \mathbf{p})$  is sufficiently close to the zero vector, which is resulted from the Fisher Scoring. Therefore, we compute  $\det(U) \cdot S^T(\mathbf{y} - \mathbf{p})$  which is a reliable approximation of  $\det(U) \cdot S^{*T}W\mathbf{v}^*$ , in much less computational costs.

### 3.4 Our Modified semi-parallel GWAS Algorithm

To sum up all our algorithmic optimization techniques described in above, we have Algorithm 3 and 4, which are modified Fisher Scoring and semi-parallel GWAS algorithms, respectively.

---

#### Algorithm 3 ModifiedFisherScoring

---

**Input:** Covariate matrix  $X \in \mathbb{Q}^{n \times k}$ , phenotype vector  $\mathbf{y} \in \mathbb{Q}^n$ , # iteration  $\text{iter}$ , and constant  $\alpha > 0$ .  
**Output:** Coefficient vector  $\beta \in \mathbb{Q}^k$ , fitted vector  $\mathbf{p} \in \mathbb{Q}^n$ , weight matrix  $W \in \mathbb{Q}^{n \times n}$   
1:  $\beta^{(0)} = \mathbf{0} \in \mathbb{Q}^k$ ,  $\mathbf{p}^{(0)} = \mathbf{0.5} \in \mathbb{Q}^n$ ,  $W^{(0)} = 0.25 \cdot I_n$   
2: **for**  $t = 0$  to  $\text{iter} - 1$  **do**  
3:    $U^{(t)} \leftarrow X^T W^{(t)} X$   
4:    $\beta^{(t+1)} \leftarrow \beta^{(t)} + \alpha \cdot \text{adj}(U^{(t)}) X^T (\mathbf{y} - \mathbf{p}^{(t)})$   
5:    $\mathbf{p}^{(t+1)} \leftarrow \sigma(X \beta^{(t+1)})$   
6:    $W^{(t+1)} \leftarrow \text{diagonal matrix of } \mathbf{p}^{(t+1)} \odot (\mathbf{1} - \mathbf{p}^{(t+1)})$   
7: **end for**  
8: **return**  $\beta^{(\text{iter})}$ ,  $\mathbf{p}^{(\text{iter})}$ ,  $W^{(\text{iter})}$

---

In Algorithm 3, the constant  $\alpha$  takes a similar role to the learning rate in gradient descent algorithm [22], which can be adjusted if necessary.

---

#### Algorithm 4 Modified Semi-Parallel GWAS

---

**Input:** SNP matrix  $S \in \{0, 1\}^{n \times m}$ , covariate matrix  $X \in \mathbb{Q}^{n \times k}$ , phenotype vector  $\mathbf{y} \in \mathbb{Q}^n$ , # iteration  $\text{iter}$ , and constant  $\alpha > 0$ .  
**Output:**  $\overrightarrow{p\text{val}}$  =  $(\text{pval}_1, \dots, \text{pval}_m) \in \mathbb{Q}^m$ .  
1:  $\beta, \mathbf{p}, W \leftarrow \text{ModifiedFisherScoring}(X, \mathbf{y}; \text{iter}, \alpha)$   
2:  $U \leftarrow X^T W X$ , and compute  $\text{adj}(U)$ ,  $\det(U)$   
3:  $V \leftarrow X^T W S$   
4:  $\mathbf{c} \leftarrow S^T (\mathbf{y} - \mathbf{p})$   
5:  $\mathbf{d} \leftarrow \det(U) \cdot \text{diag}(S^T W S) - \text{diag}(V^T \text{adj}(U) V)$   
6: **for**  $i = 1$  to  $m$  **do**  
7:    $z_i \leftarrow \det(U) \cdot c_i^2 / d_i$   
8:    $\text{pval}_i = 2 \cdot \int_{-\infty}^{-\sqrt{z_i}} \rho(x) dx$  ▷ Done in unencrypted state  
9: **end for**  $\overrightarrow{p\text{val}}$   
10: **return**  $\overrightarrow{p\text{val}}$

---

We remark that step 8 of Algorithm 4, a conversion procedure from the squared  $z$ -score  $z_i$  to the  $p$ -value  $\text{pval}_i$ , is done in unencrypted state. Namely, we decrypt the ciphertext of  $z_i$  for  $1 \leq i \leq m$  after step 7 so that  $z_i$ 's are publicized. We stress that the squared  $z$ -score has exactly *the same information* as the  $p$ -value, so publishing squared  $z$ -scores does not leak any additional information more than publishing  $p$ -values.

## 4 Homomorphic Evaluation of the Modified semi-parallel GWAS algorithm

Upon HE-friendly algorithms discussed in the previous section, there still remain computational issues regarding more fundamental operations. Recall that HEAAN basically supports component-wise addition and multiplication along with data slot rotations. However, we encrypt the data matrix by column-by-column manner, and our algorithms include complex operations such as matrix multiplication, evaluation of the adjoint matrix, a sigmoid function, and so on. In this regard, we specify how we can deal with such operations efficiently, reducing the number of multiplications or the total depth of multiplications required which are the main bottleneck of HE.

Note that this section consists of rather technical contents related to HE, since it includes HE algorithms of all building blocks for Algorithm 3 and Algorithm 4. One can simply embrace the fact that every operation required in Algorithm 3 and Algorithm 4 can be efficiently done based on HE, if not really interested in the details.

Hereafter,  $[a]_k$  with an integer  $a$  denotes a residue number in  $[0, k-1]$  modulo  $k$ . An  $n$ -dimensional vector  $\mathbf{a} = (a_1, \dots, a_n)$  is simply denoted by  $(a_i)_{1 \leq i \leq n}$ , and an  $n \times m$  matrix  $A$  having  $(i, j)$ -entry  $a_{i,j}$  is denoted by  $[a_{i,j}]_{1 \leq i \leq n, 1 \leq j \leq m}$ . For both cases, if the size is obvious from context, we simply write a vector by  $(a_i)$ , and a matrix by  $[a_{i,j}]$ . Every vector and matrix in this section is assumed to be of size power-of-two, which is in line with our packing method introduced in Section 2.2.

### 4.1 Adjoint matrix and Determinant

In step 4 of Algorithm 3 and step 2 of Algorithm 4, we need to compute the adjoint matrix and the determinant of the matrix  $U = [u_{i,j}]_{i,j} := X^T W X$ . Basically, we exploit the following facts for  $(i, j)$ -minor  $M_{i,j} \in \mathbb{Q}$  of  $U$  from basic linear algebra:

$$\begin{aligned} \text{adj}(U) &= [(-1)^{i+j} \cdot M_{i,j}] \in \mathbb{Q}^{k \times k}, \\ \det(U) &= \sum_{i=0}^{k-1} u_{i,0} \cdot (-1)^i \cdot M_{i,0}. \end{aligned}$$

Given an encryption of  $U$ , denoted by  $C_U$ , we generate  $(k-1)^2$  ciphertexts  $C_{i,j}$  for  $1 \leq i, j \leq k-1$  from  $C_U$ , whose plaintext is an  $i$ -row (upper) rotation and  $j$ -column (left) rotation of  $U$ . We first consider the 0-th plaintext slot, i.e., the  $(0, 0)$ -position of the plaintext matrix, of the ciphertexts. Since every  $u_{a,b}$  for  $1 \leq a, b \leq k-1$  is a  $(0, 0)$ -entry of plaintext matrix for one and only one ciphertext  $C_{a,b}$ , we can compute a ciphertext whose  $(0, 0)$ -entry of the plaintext matrix is  $M_{0,0}$  from  $C_{a,b}$ 's, by homomorphically evaluating the polynomial  $f$  which outputs  $M_{0,0}$  with input  $u_{a,b}$ 's.

Now observe that  $M_{a,b}$  and  $M_{a',b'}$  have a formula of the same form where the subscript indices of  $u_{i,j}$  are shifted by  $(a' - a, b' - b)$  modulo  $k$ . Thanks to this index-shifting property, the homomorphic evaluation of the polynomial  $f$  with input  $C_{a,b}$ 's essentially outputs a ciphertext of which the plaintext matrix is  $[M_{i,j}]$ .

After computing the ciphertext of  $[M_{i,j}]$  as above, we can obtain the ciphertext  $C_{\text{adj}}$  of  $\text{adj}(U)$  by multiplying a ciphertext  $C_{\text{sgn}}$  of  $[(-1)^{i+j}]$ . Finally, the ciphertext  $C_{\text{det}}$  of determinant  $\det(U)$  is easily obtained from the homomorphic multiplication of  $C_U$  and  $C_{\text{adj}}$ , followed by  $\log k$  rotations and summations.

In case of  $k = 4$ , for example, the polynomial  $f$  is defined as  $f([u_{i,j}]_{1 \leq i,j \leq 3}) = u_{1,1}u_{2,2}u_{3,3} - u_{1,1}u_{2,3}u_{3,2} - u_{2,1}u_{1,2}u_{3,3} + u_{2,1}u_{1,3}u_{3,2} + u_{3,1}u_{1,2}u_{2,3} - u_{3,1}u_{1,3}u_{2,2}$ . Then, the homomorphic evaluation to obtain  $C_{\text{adj}}$  is done as Algorithm 5.

---

**Algorithm 5** C.Adj( $C_U; C_{\text{sgn}}$ ) for  $k = 4$

---

**Input:** Ciphertexts  $C_U$  of  $U$  and  $C_{\text{sgn}}$  of  $[(-1)^{i+j}]_{0 \leq i,j \leq 3}$

**Output:** Ciphertext  $C_{\text{adj}}$  of  $\text{adj}(U)$

```

1: for  $i = 1$  to  $3$  do
2:    $C_{i,0} \leftarrow \text{C.RowRot}(C_U, i)$ 
3:   for  $j = 1$  to  $3$  do
4:      $C_{i,j} \leftarrow \text{C.ColumnRot}(C_{i,0}, j)$ 
5:   end for
6: end for
7:  $C_1 \leftarrow \text{C.Mult}(\text{C.Mult}(C_{1,1}, C_{2,2}), C_{3,3})$ 
8:  $C_2 \leftarrow \text{C.Mult}(\text{C.Mult}(C_{1,1}, C_{2,3}), C_{3,2})$ 
9:  $C_3 \leftarrow \text{C.Mult}(\text{C.Mult}(C_{2,1}, C_{1,2}), C_{3,3})$ 
10:  $C_4 \leftarrow \text{C.Mult}(\text{C.Mult}(C_{2,1}, C_{1,3}), C_{3,2})$ 
11:  $C_5 \leftarrow \text{C.Mult}(\text{C.Mult}(C_{3,1}, C_{1,2}), C_{2,3})$ 
12:  $C_6 \leftarrow \text{C.Mult}(\text{C.Mult}(C_{3,1}, C_{1,3}), C_{2,2})$ 
13: for  $i = 2$  to  $6$  do
14:   if  $i$  is even then
15:      $C_1 \leftarrow \text{C.Sub}(C_1, C_i)$ 
16:   else
17:      $C_1 \leftarrow \text{C.Add}(C_1, C_i)$ 
18:   end if
19:  $C_{\text{adj}} \leftarrow \text{C.Mult}(C_1, C_{\text{sgn}})$ 
20: end for
21: return  $C_{\text{adj}}$ 

```

---

## 4.2 Matrix Multiplications

Let  $A = [a_{i,j}]$  be an  $n \times k$  matrix with  $n \geq k$  and  $B = [b_{i,j}]$  be an  $n \times m$  matrix. We use an Algorithm 6 computing a ciphertext  $C_{A^T B}$  of  $A^T B$  from ciphertexts  $C_A$  and  $C_B$  of  $A$  and  $B$ , which is inspired from the hybrid method by Juvekar et al. [16].

As the first step, we compute  $k$  ciphertexts of

$$\text{diag}_t(A) = (a_{i,[i-t]_k})_{0 \leq i \leq n-1}$$

for  $0 \leq t \leq k-1$ . For this, we use ciphertexts  $\text{dmsk}_t$  of  $n \times k$  masking matrix, of which the  $(i, j)$ -entry is  $\delta_{[i+t]_k, j}$ . Here  $\delta_{i, j}$  denotes the Kronecker Delta. By summing column rotations of  $A \odot \text{dmsk}_t$ , we get ciphertexts of  $n \times m$  matrices  $\text{Expdiag}_t(A)$  having  $m$  identical columns  $\text{diag}_t(A)$ . Then, we compute the following matrix  $M$ :

$$\begin{aligned} M &= \sum_{t=0}^{k-1} \rho_t(\text{Expdiag}_t(A) \odot B) \\ &= \sum_{t=0}^{k-1} \rho_t([a_{i,[i-t]_k} \cdot b_{i,j}]) \\ &= \sum_{t=0}^{k-1} [a_{[i+t]_n, [i]_k} \cdot b_{[i+t]_n, j}] \in \mathbb{Q}^{n \times m}, \end{aligned}$$

where  $\rho_t$  is an (upward)  $t$ -rotation of matrix by row. Thus, by properly summing rows of  $M$ , we obtain  $A^T B = \sum_{t=0}^{k-1} [a_{t,i} \cdot b_{t,j}] \in \mathbb{Q}^{k \times m}$ . The detail of this algorithm is described in Algorithm 6.

---

**Algorithm 6**  $\text{C.MatMul}(C_A, C_B)$

---

**Input:** Ciphertexts  $C_A$  of  $A \in \mathbb{Q}^{n \times k}$  and  $C_B$  of  $B \in \mathbb{Q}^{n \times m}$  with  $n \geq k$ , and masking ciphertexts  $\{\text{dmsk}_t\}_t$

**Output:** A ciphertext of  $A^T B$ .

```

1:  $C_M \leftarrow 0$ 
2: for  $t = 0$  to  $k - 1$  do
3:    $C_t \leftarrow \text{C.Mult}(\text{dmsk}_t, C_A)$ 
4:   for  $i = 0$  to  $\log k - 1$  do
5:      $C_t \leftarrow \text{C.Add}(C_t, \text{C.ColumnRot}(C_t, 2^i))$ 
6:   end for
7:    $C_M \leftarrow \text{C.Add}(C_M, \text{C.RowRot}(\text{C.Mult}(C_t, C_B), t))$ 
8: end for
9: for  $i = 0$  to  $\log(n/k) - 1$  do
10:   $C_M \leftarrow \text{C.Add}(C_M, \text{C.RowRot}(C_M, 2^i \cdot k))$ 
11: end for
12: return  $C_M$ 

```

---

Indeed, our GWAS algorithm contains several matrix multiplications of the form  $A^T DB$  for the diagonal matrix  $D$ . For this, we first compute  $B' = DB$  by  $\text{Expdiag}_0(D) \odot B$ , and then obtain  $A^T DB$  by computing  $A^T B'$  with the above method. Note that this requires only one additional Hadamard multiplication.

### 4.3 Matrix-vector multiplications

By understanding a vector by an one column matrix, we can perform all matrix-vector multiplications in our algorithms, except  $X\boldsymbol{\beta}$  that appears in step 5 of Algorithm 3.

In fact, we can also compute  $X\boldsymbol{\beta}$  by changing Algorithm 6 a little bit. Recall that  $X = [x_{i,j}]$  is an  $n \cdot k$  matrix and  $\boldsymbol{\beta} = (\beta_i)$  is a  $k$ -length vector, and it holds that  $n \geq k$ . We now again compute  $k$  ciphertexts of  $\text{diag}_i(X)$ , and then compute

$$\begin{aligned} & \sum_{t=0}^{k-1} \text{diag}_t(X) \odot \rho_{-t}((\boldsymbol{\beta} || \cdots || \boldsymbol{\beta})) \\ &= \sum_{t=0}^{k-1} (x_{i,[i-t]_k}) \odot (\beta_{[i-t]_k}) \\ &= \sum_{t=0}^{k-1} (x_{i,[i-t]_k} \cdot \beta_{[i-t]_k}) = X\boldsymbol{\beta}, \end{aligned}$$

whose algorithm is described by Algorithm 7.

---

#### Algorithm 7 $\text{C.MatVecMul}(C_X, C_\beta)$

---

**Input:** Ciphertexts  $C_X$  of  $X \in \mathbb{Q}^{n \times k}$  and  $C_\beta$  of  $\boldsymbol{\beta} \in \mathbb{Q}^k$  and masking ciphertexts  $\{\text{dmsk}_t\}_t$

**Output:** A ciphertext of  $X\boldsymbol{\beta}$ .

- 1:  $C_M \leftarrow 0$
  - 2: **for**  $t = 0$  to  $k - 1$  **do**
  - 3:      $C_t \leftarrow \text{C.Mult}(\text{dmsk}_t, C_X)$
  - 4:     **for**  $i = 0$  to  $\log k - 1$  **do**
  - 5:          $C_t \leftarrow \text{C.Add}(C_t, \text{C.ColumnRot}(C_t, 2^i))$
  - 6:     **end for**
  - 7:      $C_M \leftarrow \text{C.Add}(C_M, \text{C.Mult}(C_t, \text{C.RowRot}(C_\beta, -t)))$
  - 8: **end for**
  - 9: **return**  $C_M$
- 

We remark that there is another simple method for matrix-vector multiplication that we use for  $\mathbf{c} = S^T(\mathbf{y} - \mathbf{p})$ . For simplicity, let  $\mathbf{x} = (x_i) := \mathbf{y} - \mathbf{p}$ . Then by rotating and summing all rows of  $S \odot [\mathbf{x} || \cdots || \mathbf{x}] = [s_{i,j} \cdot x_i]$ , we obtain a matrix having the same size with  $S$  and consisting of identical rows  $\mathbf{c} = S^T \mathbf{x}$ . This requires only one Hadamard multiplication and  $\log n$  rotations. However, strictly speaking, this resulting ciphertext is not a ciphertext of  $S^T \mathbf{x}$ , since it encrypts a matrix having  $\mathbf{c}$  row-wisely, not column-wisely. Thus we can only use this simple method only for  $S^T \mathbf{x}$ , where this row-wise packing does not matter after then. The detail of this algorithm is described in Algorithm 8.

---

**Algorithm 8**  $\mathbf{C.MatVecMul}_2(C_S, C_x)$ 

---

**Input:** Ciphertexts  $C_S$  of a matrix  $S \in \mathbb{Q}^{n \times m}$  and  $C_x$  of a vector  $x \in \mathbb{Q}^n$

**Output:** A ciphertext of a matrix having identical rows  $S^T x$ .

- 1:  $C_M \leftarrow \mathbf{C.Mult}(C_S, C_x)$
  - 2: **for**  $i = 0$  to  $\log n - 1$  **do**
  - 3:      $C_M \leftarrow \mathbf{C.Add}(C_M, \mathbf{C.RowRot}(C_M, 2^i))$
  - 4: **end for**
  - 5: **return**  $C_M$
- 

#### 4.4 Fast $\mathbf{diag}(A^T BA)$ computations

To obtain  $\mathbf{diag}(A^T BA)$ , one can perform matrix multiplication followed by diagonal extraction, but it is obviously not optimal since this computes unnecessary entries of  $A^T BC$  other than diagonal entries. Thus we use another method that only compute the diagonal entries.

Let  $A = [a_{i,j}]$  be an  $n \times m$  matrix. As an incremental step, we first consider  $\mathbf{diag}(A^T DC)$  where  $D = [d_{i,j}]$  is an  $n \times n$  diagonal matrix and  $C = [c_{i,j}]$  is an  $n \times m$  matrix. Then it holds that  $\mathbf{diag}(A^T DC)_j = \sum_{i=0}^{n-1} d_{i,i} \cdot a_{i,j} \cdot c_{i,j}$ . Now, from an encryption of  $D$ , we compute an encryption of  $n \times m$  matrix  $\mathbf{Expdiag}_0(D)$  and then by rotating and summing

$$A \odot \mathbf{Expdiag}_0(D) \odot C = [d_{i,i} \cdot a_{i,j} \cdot c_{i,j}]_{i,j}$$

through all rows, we obtain a matrix consisting of identical rows  $\mathbf{diag}(A^T DC)$ . One can easily check that Algorithm 8 with input  $C_A$  and  $\mathbf{C.Mult}(C_{\mathbf{diag}(D)}, C_C)$  exactly performs this computation, and then we omit the explicit algorithm. Note that this can be directly applied for  $\mathbf{diag}(S^T WS)$  computation of step 5 of Algorithm 4.

Toward our goal  $\mathbf{diag}(A^T BA)$  with a full matrix  $B$ , we exploit the above diagonal-case method after decomposing  $B$  into diagonal matrices. Let  $B_t$  be a diagonal matrix with the diagonal  $\mathbf{diag}_t(B)$  for  $0 \leq t \leq n - 1$ , then it holds that

$$\mathbf{diag}(A^T BA) = \sum_{t=0}^{n-1} \mathbf{diag}(A^T \cdot B_t \cdot \rho_t(A)).$$

Therefore, after obtaining encryptions of  $\mathbf{Expdiag}_t(B)$  and  $\rho_t(A)$  from encryptions of  $B$  and  $A$ , we can directly apply the diagonal-case method on each  $\mathbf{diag}(A^T \cdot B_t \cdot \rho_t(A))$  computation for  $1 \leq t \leq n$  and finally obtain the encryption of  $\mathbf{diag}(A^T BA)$ .

Here we again remark that, since these methods use Algorithm 8, they also ruin the column-wise packing as we already pointed out. Hence after applying these methods, it would be hard to perform another matrix operation. Indeed, one can check that the diagonal extractions are required for step 5 of Algorithm 4, which is the last part of algorithm that uses matrix structure.

## 4.5 Approximate Computation of Sigmoid

Since the sigmoid function  $\sigma(x) = 1/(1+\exp(-x))$  is not a polynomial, we exploit an approximate polynomial of the function to evaluate based on HE. Following the methodology of [17,18], we used least square approximation method over the interval  $[-8, 8]$ . The approximate polynomials  $g(x)$  of degree 7 is computed as

$$0.5 + 1.735 \cdot \frac{x}{8} - 4.194 \cdot \left(\frac{x}{8}\right)^3 + 5.434 \cdot \left(\frac{x}{8}\right)^5 - 2.507 \cdot \left(\frac{x}{8}\right)^7.$$

The maximal error between  $\sigma(x)$  and  $g(x)$  is approximately 0.032.

## 4.6 Inverse of Real Numbers

In step 7 of Algorithm 4, we need to compute the inverse of  $d_i$  for  $1 \leq i \leq m$ . To compute the inverse of real numbers, we exploit the Goldschmidt's division algorithm [14], which outputs an approximate value of the inverse through iterative polynomial evaluations. Refer to [21] for more details of the algorithm.

# 5 Results

In this section, we present the experimental results of our modified semi-parallel GWAS algorithm based on HEAAN with a publicly available library [15]. All experiments were implemented in C++ 11 standard, and performed on Linux with Intel Xeon CPU E5-2620 v4 at 2.10GHz processor with multi-threading (8 threads) turned on.

## 5.1 Dataset Description

We used a dataset of 245 samples where each sample contains a binary phenotype, 3 covariates (height, weight, age), and 25,484 SNP data provided by IDASH 2018 competition. The dataset is divided into two sets named by `iDash_Test` and `iDash_Eval` each composed of 245 samples containing common phenotype and 3 covariates but different number of SNPs; 10,643 and 14,841 SNPs, respectively. We used `iDash_Test` to set optimal parameters, and `iDash_Eval` was used to evaluate our algorithm in the competition. Note that the first column of the covariate matrix  $X \in \mathbb{Q}^{n \times k}$  is a vector of which all the components are 1. Therefore, the parameters are  $(n, m, k) = (245, 10643, 4)$  for `iDash_Test` and  $(n, m, k) = (245, 14841, 4)$  for `iDash_Eval`.

## 5.2 Experimental Setting and Parameter Selection

We propose two HEAAN parameter sets achieving 128-bit or higher security for two experiments denoted by `Exp I` and `Exp II` in Table 1. The security levels of HEAAN parameter sets were estimated with Albrecht's security estimator [4,3] of which inputs are the ring dimension  $N$ , the modulus  $Q$ , the Hamming weight  $h$  of a secret polynomial, and the error distribution  $\chi_{\text{err}}$ . Note that since the

Table 1: Parameters for HEAAN, and Running time of KeyGen, Enc and Dec

Exp	HE parameters				Time (sec)		
	$\log N$	$L$	$p$	$h$	KeyGen	Enc	Dec
I	17	1300	50	56	157	68	0.28
II	17	1700	50	78	197	90	0.15

modulus of the evaluation key  $evk$  is  $2^{2L}$ , the security of HEAAN is estimated with input  $(N, Q = 2^{2L}, h, \chi_{\text{err}})$ .

Exp I is a streamlined version operating Algorithm 4 only until step 5; that is, it does not perform the last division process. Exp II includes the division process (step 7 of Algorithm 4), so it naturally requires larger  $L$  than Exp I.

We first set the scaling parameter  $p$  to be sufficiently large so that errors derived from HEAAN do not effect on significant bits of plaintexts. The level parameter  $L$  is chosen by considering the fact that  $p$  levels are consumed for each homomorphic multiplication. See Section 2.1 for specific definitions of HEAAN parameters. Besides those HEAAN parameters, one also needs to select an appropriate constant  $\alpha > 0$  in Algorithm 3. After experimenting on several values, we set  $\alpha = 8$ . Note that the choice of  $\alpha$  merely depends on the size of  $X$ , but not on  $S$ .

### 5.3 Experimental Results and Evaluation

We demonstrated our modified semi-parallel GWAS algorithm in encrypted state, and evaluated the accuracy of our algorithm comparing it to that of the original algorithm which is performed in unencrypted state. The comparison result of  $p$ -value is described as a (log-scale) graph in Figure 1.

We plotted each SNPs according to the  $p$ -values computed by original algorithm and ours denoted by `True` and `Enc`, respectively. The diagonal line represents the line  $y = x$ , and closer distribution of points to this line implies higher accuracy. The Figure 1-(a) shows that the accuracy of our algorithm increases with the number of iterations for Fisher scoring, where the data set `iDash_Test` is used. The Figure 1-(b) shows that the accuracy of Exp II, which includes a division procedure, is comparable to that of Exp I without the division, where the data set `iDash_Eval` is used. Comparing Exp I and Exp II, there exists a trade-off between computational time and information leakage. The output of Exp I is the vector of squared statistics  $(z_i)_{1 \leq i \leq m}$  which has exactly same information with the  $p$ -value vector `pval`, but it takes 20 minutes longer than Exp I. On the other hand, since Exp I outputs the numerator  $\det(U) \cdot c_i^2$  and the denominator  $d_i$  (in Algorithm 4) separately, it leaks some information more than  $p$ -values. However,

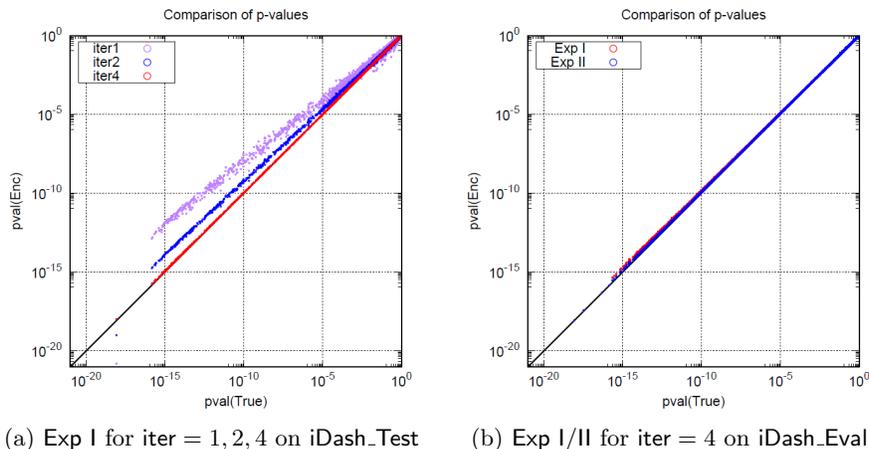


Fig. 1: Comparison of  $p$ -values on IDASH datasets

it still seems to be very hard to extract any important information of input data from the numerator and denominator.

For more concrete evaluation, we classified each SNP as positive or negative depending on whether the corresponding  $p$ -value is larger or smaller than the given threshold (e.g.  $10^{-2}$ ,  $10^{-5}$ , or  $10^{-12}$ ). Then, the accuracy of our algorithm compared to the original algorithm can be checked by a well-known statistical measure called  $F_1$  score. The  $F_1$  score of our algorithm is calculated by regarding the positive SNPs classified by the original GWAS algorithm to be the correct positive samples. For the formal definition of  $F_1$  score, we refer readers to [11].

The performance of our algorithm including the computation time and the  $F_1$  score on each parameter set is described in Table 2, where *iter* denotes the number of iterations in Fisher Scoring, *Comp. time* denotes the running time of our algorithm in encrypted state, and *TH* denotes the threshold of  $p$ -values for classification.

As we have seen in Figure 1, more iterations of Fisher scoring provides higher accuracy measured by higher  $F_1$  score. Note that 4 iterations suffice to provide high  $F_1$  score even in a very small threshold such as  $10^{-12}$ . Also, Exp II calculating approximate inverse in encrypted state provides almost similar  $F_1$  score to Exp I without such approximation. It implies that the error of inverse approximation does not seriously impact the whole approximation. Furthermore, Exp II shows even higher  $F_1$  score than Exp I due to the cancellation of errors from our algorithmic approximation and that from the inverse approximation.

For about 15,000 SNP data, our algorithm works in less than 40 minutes when we exclude step 7 of Algorithm 4 in encrypted state, or in about 60 minutes otherwise. We emphasize that each iteration of Fisher scoring takes about 3

Table 2: Experimental Results for each parameter set.

Data	Params		Comp. time	$F_1$ Score		
	Exp	iter		TH: $10^{-2}$	TH: $10^{-5}$	TH: $10^{-12}$
iDash_Test	I	1	13 min*	0.960	0.969	0.243
	I	2	27 min	0.985	0.985	0.955
	<b>I</b>	<b>4</b>	<b>32 min</b>	<b>1.000</b>	<b>0.999</b>	<b>0.997</b>
	<b>II</b>	<b>4</b>	<b>52 min</b>	<b>1.000</b>	<b>0.999</b>	<b>0.998</b>
iDash_Eval	<b>I</b>	<b>4</b>	<b>38 min</b>	<b>0.998</b>	<b>0.995</b>	<b>0.992</b>
	<b>II</b>	<b>4</b>	<b>62 min</b>	<b>0.998</b>	<b>0.996</b>	<b>0.994</b>

\*: We used more streamlined parameter;  $\log N = 16, L = 950, p = 50, h = 91$ .

minutes while the Goldschmidt’s division algorithm takes less than 30 seconds. Exp II takes much longer time than Exp I due to the larger level parameter  $L$ .

## 6 Discussion

*Scalability.* Our algorithm is executed and evaluated with about hundreds of samples each containing ten thousand SNPs, and 3 covariates which can be seen as a small-size data in usual GWAS analysis. We emphasize that our algorithm is highly scalable in the number of samples or SNPs, since we circumvent the naive execution of large-sized matrix operations through the proper algorithmic modification. To test the scalability of our algorithm in practice, we randomly generated 500 samples each of which consist of 3 covariates and 30,000 SNPs.<sup>1</sup> The experiment Exp 1 on this random dataset encrypted with properly chosen hyperparameters  $\text{iter} = 4$  and  $\alpha = 2^{-7}$  still showed quite accurate  $p$ -value result compared to the result obtained by running Algorithm 2 in unencrypted state within 2 hours.

*Fisher Scoring.* Our HE-friendly modified Fisher scoring (Algorithm 3) works quite well in practice, but there still remains to obtain some theoretical results on the convergence of the algorithm with respect to the new parameter  $\alpha$ . Furthermore, we should consider an error in every operation derived from HEAAN when homomorphically evaluate the algorithm. As a result, research on the convergence of the *erroneous* version of our modified Fisher scoring algorithm should

<sup>1</sup> Each column of the covariate matrix was uniform randomly generated in the interval  $[150, 200]$ ,  $[40, 100]$  and  $[20, 80]$  considering height, weight and age, respectively. Each element of the SNP matrix was uniform randomly chosen as a binary matrix.

be very interesting topic as a further work. In addition, we note that our modified Fisher scoring algorithm can be generally used for logistic regression, not restricted to GWAS algorithm.

## 7 Conclusions

Interest on privacy-preserving genome data analysis based on HE has grown up very rapidly since the annual IDASH competition was launched, and GWAS is one of the most important technologies in this area which was also selected as one of three tasks in IDASH 2018 competition. Our HE-friendly modified semi-parallel GWAS algorithm was successfully implemented based on an approximate HE scheme HEAAN, and we could obtain the  $p$ -value result in about 30–40 minutes for 10,000–15,000 SNP data with sufficiently high accuracy compared to the result obtained in unencrypted state.

## References

1. IDASH 2017 Competition. <http://www.humangenomeprivacy.org/2017/>.
2. IDASH 2018 Competition. <http://www.humangenomeprivacy.org/2018/>.
3. M. R. Albrecht. A Sage Module for estimating the concrete security of Learning with Errors instances., 2017. <https://bitbucket.org/malb/lwe-estimator>.
4. M. R. Albrecht, R. Player, and S. Scott. On the concrete hardness of learning with errors. *Journal of Mathematical Cryptology*, 9(3):169–203, 2015.
5. C. Bonte and F. Vercauteren. Privacy-preserving logistic regression training. Cryptology ePrint Archive, Report 2018/233, 2018. <https://eprint.iacr.org/2018/233>.
6. H. Chen, R. Gilad-Bachrach, K. Han, Z. Huang, A. Jalali, K. Laine, and K. Lauter. Logistic regression over encrypted data from fully homomorphic encryption. *BMC medical genomics*, 11(4):81, 2018.
7. J. H. Cheon, K. Han, S. M. Hong, H. J. Kim, J. Kim, S. Kim, H. Seo, H. Shim, and Y. Song. Toward a secure drone system: Flying with real-time homomorphic authenticated encryption. *IEEE Access*, 6:24325–24339, 2018.
8. J. H. Cheon, K. Han, A. Kim, M. Kim, and Y. Song. Bootstrapping for approximate homomorphic encryption. In *Annual International Conference on the Theory and Applications of Cryptographic Techniques*, pages 360–384. Springer, 2018.
9. J. H. Cheon, A. Kim, M. Kim, and Y. Song. Homomorphic encryption for arithmetic of approximate numbers. In *Advances in Cryptology–ASIACRYPT 2017: 23rd International Conference on the Theory and Application of Cryptology and Information Security*, pages 409–437. Springer, 2017.
10. J. H. Cheon, D. Kim, Y. Kim, and Y. Song. Ensemble method for privacy-preserving logistic regression based on homomorphic encryption. *IEEE Access*, 6:46938–46948, 2018.
11. N. Chinchor. Muc-4 evaluation metrics. In *Proceedings of the 4th conference on Message understanding*, pages 22–29. Association for Computational Linguistics, 1992.
12. H. Cho, D. J. Wu, and B. Berger. Secure genome-wide association analysis using multiparty computation. *Nature biotechnology*, 36(6):547, 2018.

13. J. L. Crawford, C. Gentry, S. Halevi, D. Platt, and V. Shoup. Doing real work with fhe: The case of logistic regression. Cryptology ePrint Archive, Report 2018/202, 2018. <https://eprint.iacr.org/2018/202>.
14. R. E. Goldschmidt. *Applications of division by convergence*. PhD thesis, Massachusetts Institute of Technology, 1964.
15. K. Han, A. Kim, M. Kim, and Y. Song. Implementation of HEAAN, 2016. <https://github.com/snucrypto/HEAAN>.
16. C. Juvekar, V. Vaikuntanathan, and A. Chandrakasan. Gazelle: A low latency framework for secure neural network inference. In *27th USENIX Security Symposium (USENIX Security 18)*. USENIX Association, 2018.
17. A. Kim, Y. Song, M. Kim, K. Lee, and J. H. Cheon. Logistic regression model training based on the approximate homomorphic encryption. *BMC Medical Genomics*, 11(4):83, Oct 2018.
18. M. Kim, Y. Song, S. Wang, Y. Xia, and X. Jiang. Secure logistic regression based on homomorphic encryption: Design and evaluation. *JMIR medical informatics*, 6(2), 2018.
19. N. T. Longford. A fast scoring algorithm for maximum likelihood estimation in unbalanced mixed models with nested random effects. *Biometrika*, 74(4):817–827, 1987.
20. M. B. Malik, M. A. Ghazi, and R. Ali. Privacy preserving data mining techniques: current scenario and future prospects. In *Third International Conference on Computer and Communication Technology (IC3CT)*, pages 26–32. IEEE, 2012.
21. P. Markstein. Software division and square root using goldschmidt’s algorithms. In *Proceedings of the 6th Conference on Real Numbers and Computers (RNC’6)*, volume 123, pages 146–157, 2004.
22. S. Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*, 2016.
23. K. Sikorska, E. Lesaffre, P. F. Groenen, and P. H. Eilers. Gwas on your notebook: fast semi-parallel linear and logistic regression for genome-wide association studies. *BMC bioinformatics*, 14(1):166, 2013.