

Secure genotype imputation using homomorphic encryption

Zhou, Junwei; Lei, Botian; Lang, Huile; Panaousis, Emmanouil; Liang, Kaitai; Xiang, Jianwen

DOI

[10.1016/j.jisa.2022.103386](https://doi.org/10.1016/j.jisa.2022.103386)

Publication date

2023

Document Version

Final published version

Published in

Journal of Information Security and Applications

Citation (APA)

Zhou, J., Lei, B., Lang, H., Panaousis, E., Liang, K., & Xiang, J. (2023). Secure genotype imputation using homomorphic encryption. *Journal of Information Security and Applications*, 72, Article 103386. <https://doi.org/10.1016/j.jisa.2022.103386>

Important note

To cite this publication, please use the final published version (if applicable). Please check the document version above.

Copyright

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

Takedown policy

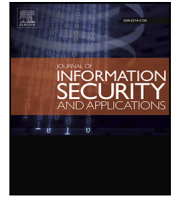
Please contact us and provide details if you believe this document breaches copyrights. We will remove access to the work immediately and investigate your claim.

Green Open Access added to TU Delft Institutional Repository

'You share, we take care!' - Taverne project

<https://www.openaccess.nl/en/you-share-we-take-care>

Otherwise as indicated in the copyright section: the publisher is the copyright holder of this work and the author uses the Dutch legislation to make this work public.



Secure genotype imputation using homomorphic encryption

Junwei Zhou^{a,*}, Botian Lei^a, Huile Lang^a, Emmanouil Panaousis^b, Kaitai Liang^c, Jianwen Xiang^a

^a Hubei Key Laboratory of Transportation of Internet of Things, School of Computer Science and Artificial Intelligence, Wuhan University of Technology, China

^b Department of Computing & Information Systems, University of Greenwich, London

^c Cybersecurity Group, EMCS, Delft University of Technology, Netherlands

ARTICLE INFO

Keywords:

Privacy-preserving
Homomorphic encryption
Genotype imputation
Privacy computing
Genetic security

ABSTRACT

Genotype imputation estimates missing genotypes from the haplotype or genotype reference panel in individual genetic sequences, which boosts the potential of genome-wide association and is essential in genetic data analysis. However, the genetic sequences involve people's privacy, confirming an individual's identification and even disease information. This work proposes a secure genotype imputation model, which uses a linear regression model and the homomorphic encryption scheme over ciphertext to impute missing genotypes. The inference model is trained with float plaintext parameters, which are round into integers to avoid high complexity homomorphic evaluation on float number operations without bootstrapping operations. Even though the rounding parameters in the inference model are not the same as those in the trained model, We find that it will no effect on the outcome of the homomorphic prediction. Thus, a high-efficiency genotype imputation inference model over the ciphertext is obtained while keeping the high-security level. The simulation results indicate that the accuracy of the secure inference model is almost the same as the original model trained on float parameters. The secure inference model's accuracy is 98.6% for a single genotype.

1. Introduction

DNA sequencing is an indispensable part of medical diagnosis and treatment, biotechnology, and genetic data analysis [1]. It can help researchers understand and study the relationship between diseases [2], distant ancestors [3], and genes. For example, the genome-wide association study (GWAS) [4] aims to study the relationship between human diseases and complex traits. However, due to individual genetic variation or genetic testing technology issues, some genotypes in genetic sequencing are missing or low-quality. The missing genotypes can affect the genetic information integrity, thereby affecting downstream analysis such as GWAS [5]. Genotype imputation uses the association between genetic variations to predict those missing or low-quality genotypes [6]. It is a foundational step of gene analysis, with a wide range of practical applications. It is now widely used in GWAS to find new risk alleles, obtain high-resolution views in fine positioning to increase the possibility of identifying causal variants and integrate Meta-analysis of research on different platforms. Currently, the state-of-the-art plaintext genotype imputation methods include IMPUTE2 [7], Minimac3 [8], and Beagle [9] and fastPHASE [10]. IMPUTE2 uses sophisticated recombination maps and dense genotype reference panels to impute missing genotypes in the research dataset. Minimac3 divides the

genome into contiguous blocks and iterates only the unique haplotypes in each genome block. It uses a reversible mapping function to reconstruct the state space used by IMPUTE2 accurately. Beagle uses the Li and Stephens haplotype frequency models with highly reduced model state space to interpolate the phased haplotypes. FastPHASE is a flexible method that allows the "blocky" pattern of linkage disequilibrium (LD), and the LD gradually decreases with distance. These methods are applied based on plaintext gene datasets to impute missing genotypes and cannot secure the gene data. The genetic data is sensitive, where a sequence larger than 75 single-nucleotide polymorphisms (SNPs) array can confirm an individual's identification [11]. A genetic sequence can reveal human ancestry, relatives, and disease type, which involves privacy concerns.

Nowadays, the pressure brought by the increasing genetic data and the huge amount of imputation computations has made people turn their attention to convenient cloud service providers [12] for imputation. However, the genetic data is stored in plaintext in the cloud, and anyone who has access to the cloud platform may obtain these plaintext data. Even the semi-trusted cloud platforms may steal the users' genetic data motivated by benefits. The privacy concerns

* Corresponding author.

E-mail addresses: junweizhou@msn.com (J. Zhou), botianlei@whut.edu.cn (B. Lei), hllang@whut.edu.cn (H. Lang), e.panaousis@greenwich.ac.uk (E. Panaousis), kaitai.liang@tudelft.nl (K. Liang), jwxiang@whut.edu.cn (J. Xiang).

<https://doi.org/10.1016/j.jisa.2022.103386>

prohibit people from trusting the third-party platform to process their genetic data [13,14].

The homomorphic encryption (HE) [15,16] allows computing over encrypted data directly without revealing data. It provides mathematically provable security guarantees for protecting genotype data while performing imputations in an untrusted cloud platform. We note that other cryptologies like multiparty computation (MPC) [17,18] require interaction between multiple parties that hold the data, and the cloud platform performs imputation operations. Although the functional performance of MPC-based methods is impressive, they may cause problems such as network latency and high bandwidth usage. Thus, we rely on the HE scheme to secure the gene data and perform the imputation evaluation over encrypted genetic data. However, the HE scheme limits the computational circuit depth or requires time-consuming bootstrapping operations to refresh the ciphertext, making the HE-based applications computationally inefficient. The main challenge of the HE-based genotype imputation is to achieve efficient and accurate genotype imputation under the premise of ensuring security.

To explore the practical feasibility of the cryptographic methods for genotype imputation, IDASH organized the genotype imputation track in IDASH2019 Genomic Privacy Challenges. The participating teams focused on the three most advanced HE cryptosystems, which are namely Brakerski/Fan-Vercauteren (BFV) [19], Cheon-Kim-Kim-Song (CKKS) [20], and fast fully homomorphic encryption over the torus (TFHE) [21]. The highest accuracy of the participating teams was 95.5% with a 128-bit security level from the HE standardization workshop paper [22]. Since the imputation accuracy of these secure genotype imputation methods is lower than those of plaintext methods. We need higher accuracy to improve the downstream analysis, such as GWAS.

In this work, we propose a fast and secure genotype imputation [23] inference model based on the TFHE [21]. The main contributions are as follows:

- Based on the iDASH Secure Genome Analysis Challenge 2019 dataset, we provided several secure genotype imputation models of a single variant, exhibiting the connection between the tag and target variants. Simulation results indicate that the accuracy for a single genotype of a variant of 1000 individuals reached 98.6%.
- We used an LWE-based linear regression model without bootstrapping or key-switching operation, improving genotype imputation speed. Our experimental results show that the imputation time for 1000 individuals' single genotype is approximately 0.269 s.
- We found that the rounding error in the imputation vanished the noise brought by HE encryption. The obtained secure genotype imputation model has a similar performance to the original plaintext model.

We trained the model with genetic data in plaintext. The genotypes of tag variants used for the inference model are encrypted into ciphertext. We set the message space of ciphertext in advance by calculating the maximum multi-sum between the trained parameters and genotype value. The homomorphic evaluation is based on the security of the LWE problem. The encryption phase uses the LWE security concept and there did not use bootstrapping operation and key-switching key. The key distribution will not cause any time costs. There is no additional overhead (except for linearly increasing input size) to extend the interpolation calculation. The rounding parameters will not affect homomorphic computation on the ciphertext and the security level.

The remaining parts of the paper are as follows. Section 2 gives the related work of secure genotype imputation. Section 3 describes the model's training process on the plaintext and the implementation of secure inference models. Section 4 introduces our experiments. The conclusion and expectations are provided in the last section.

2. Related work

In 1978, Rivest [24] first proposed the HE scheme's assumption, which allowed various calculations over encrypted data. Nowadays, HE schemes are generally classified into three types: Partially Homomorphic Encryption (PHE) [25], Somewhat Homomorphic Encryption (SWHE) [26–28], and Fully Homomorphic Encryption (FHE) [21,29] according to calculation depth and capacity. PHE can allow homomorphic multiplications or homomorphic additions with limited calculation depth. SWHE supports homomorphic multiplications and additions with limited calculation depth. FHE supports the arbitrary depth of any calculations on the ciphertext by using bootstrapping. The bootstrapping operation [29–31] can refresh ciphertext and reduce noise in ciphertext, which allows FHE to achieve the arbitrary depths of circuits and maintain the decryption's correctness.

The current popular HE schemes include BFV [19,26], CKKS [20], BGV [28] and TFHE. These schemes are implemented based on the ring learning with error (R-LWE) problem [32], while TFHE is based on LWE and GSW [33] problems. Both BFV and BGV allow homomorphic calculations on vectors of finite field elements, and the CKKS scheme allows approximate homomorphic calculations on real or complex numbers.

In the era of cloud computing and machine learning, HE provides a solution to protect users' outsourced data [34–37]. The user uploads encrypted data to the cloud service without decryption, and the cloud service directly performs homomorphic addition and multiplication on the ciphertext. The other computing on ciphertext can be constructed using homomorphic addition and multiplication.

HE&Genotype imputation. Advances in information technology and bioinformation have made people and institutions use third-party cloud platforms to store and process data, such as online health status monitoring [38], disease diagnosis [39,40], and genotype imputation [41,42]. However, once users upload their data in plaintext to the third-party platform, they will lose control of their sensitive data. Anyone who can access the third-party platform can steal users' genetic data.

Kocabas [38] proposed a secure health monitoring system (real-time monitoring of heartbeat frequency) with the FHE scheme. When the monitor system obtained the user's heartbeat frequency, it encrypted the frequency data locally using HELib library [43] and uploaded encrypted data to the cloud platform for analysis. The cloud platform analyzed the encrypted data and transmitted the encrypted result to the user. Then the user decrypted it to check the result. In this process, since the secret key was in the users' hands, the cloud platform could not obtain any information about the user's data. Thus the privacy of the user's health status was guaranteed. Meehan [39] proposed a secure model using the TFHE library to diagnose whether users had breast cancer. It guaranteed the privacy and safety of users' health status and avoided disease discrimination. Kim etc. [41] proposed to combine genotype imputation and HE scheme to protect people's genetic data without data leakage. The UTMSR team presented a fast and secure linear regression model based on BFV and CKKS schemes, and the accuracy achieved 95.4%. EPFL team applied a multinomial logistic model to impute missing genotypes based on the CKKS scheme homomorphically and got 95.5% accuracy. The Chimera team presented a TFHE-based logistic regression model, and the accuracy was 95.1%. The SNU team applied a one-hidden layer neural network with the CKKS scheme, and the accuracy was 95.0%. The models only took 380 microseconds to predict the genotype of a variant of 1000 individuals. They output each type of genotype's probability to determine the imputed genotype. Compared with the plaintext models (higher than 97.1% [41]), the accuracy of the secure models can still be improved.

We propose a secure linear regression inference model with the TFHE library based on the LWE problem to secure genetic data and efficiently impute the missing SNPs. The inference model is trained with plaintext float parameters, which are round into integers to avoid high complexity HE evaluation on float number operations. We performed

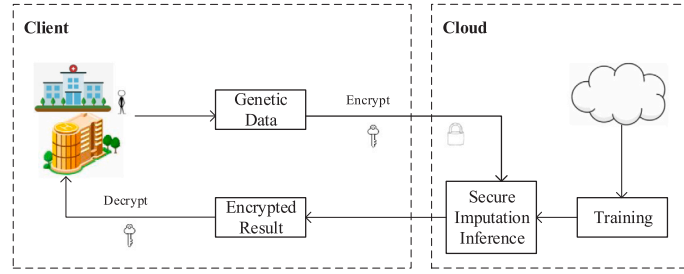


Fig. 1. The overview of our privacy-preserving genotype imputation inference.

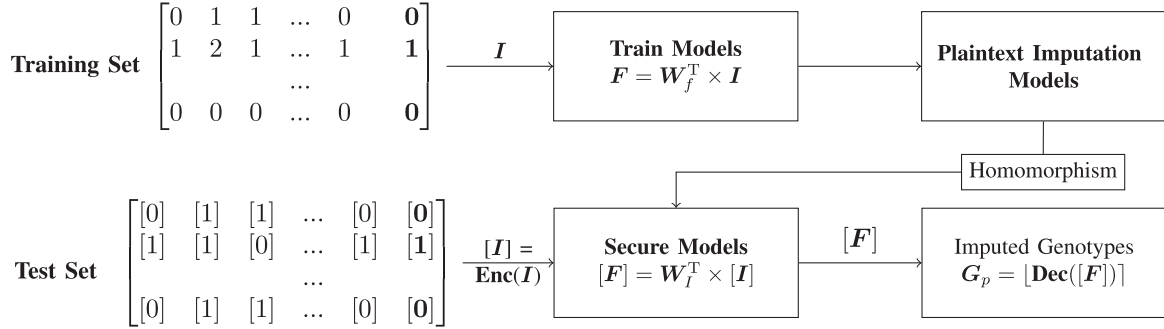


Fig. 2. Training and inference on the cloud.

detailed experiments on the time and memory requirements of the HE-based imputation model and demonstrated the feasibility of large-scale secure imputation. We found comparable performance (without decrease) in imputation accuracy with total genomic data security benefit. Our experimental results provide evidence that HE-based methods can perform efficient calculations to analyze massive genetic data.

3. The proposed model

The secure inference model over encrypted data is shown in Fig. 1. The client encrypts genetic data with a secret key and then uploads the encrypted data to the cloud service provider (CSP) for genotype imputation. CSP returns the encrypted result to the client after imputation computations. Even if an untrusted third party steals the encrypted genetic data in the entire system imputation process, he cannot obtain any information because he does not have the key. This section will introduce the models in detail, including data encryption and decryption process in client and model training and secure imputation inference model in the cloud.

3.1. Preliminaries

This subsection gives the notations, definitions of the LWE problem and the LWE encryption scheme, and homomorphic computations used in the paper.

Notation. A vector is denoted by a bold letter and $\langle a, b \rangle$ denotes the inner product between two vector a and b . $\| \cdot \|_1$ denotes the L_1 norm of a vector, $\| \cdot \|_2$ denotes the L_2 norm of a vector, and $\| \cdot \|_\infty$ denotes the infinite norm of a vector. \mathbb{R} denotes the real numbers, \mathbb{Z} denotes the integers, and \mathbb{T} is torus \mathbb{R}/\mathbb{Z} . Furthermore, given a set \mathcal{Q} , $a \xleftarrow{\$} \mathcal{Q}$ represents that a is chosen uniformly and randomly from \mathcal{Q} .

Learning With Errors. Regev [44] introduced the Learning With Errors (LWE) problem in 2005. Let n be a positive integer ($n \geq 1$), any vector $s \xleftarrow{\$} \mathbb{Z}^n$, ϕ be a distribution over \mathbb{R} , and the noise e is sampled from distribution ϕ . For any vector s , we define the LWE distribution $LWE_{s,\phi}$ as (a, b) , where the vector $a \xleftarrow{\$} \mathbb{T}^n$ and $b = \langle s, a \rangle + e$.

Regev defined the LWE problem is: for a fixed $s \xleftarrow{\$} \mathbb{Z}^n$, it is hard to distinguish between $LWE_{s,\phi}$ and the uniform distribution over \mathbb{T}^{n+1} .

Regev stated that the LWE problem is as asymptotically difficult as the worst-case lattice problem.

LWE-based encryption. Define a positive integer B related to the message space. Let $m \in [-B, B]$ be an integer message. The torus is split into $2B + 1$ slices, and each slice denotes one possible integer value.

Let n denotes the security parameter, $s \xleftarrow{\$} \mathbb{Z}^n$. ϕ is a Gaussian distribution.

Enc(m): Return (a, b) , with $a \xleftarrow{\$} \mathbb{T}^n$, and $b = \langle s, a \rangle + \frac{m}{2B+1} + e$, where $e \leftarrow \phi$.

Dec(s, (a, b)): Return $m = \lfloor (b - \langle s, a \rangle) \times (2B + 1) \rfloor$.

Homomorphic addition. Suppose two messages $m_1, m_2 \in [-B, B]$ with a randomly generated secret key s_1, a_1 and a_2 are randomly chosen from \mathbb{T}^n . $c_1 = \text{Enc}(m_1)$, $b_1 = \langle s_1, a_1 \rangle + \frac{m_1}{2B+1} + e_1$. $c_2 = \text{Enc}(m_2)$, $b_2 = \langle s_1, a_2 \rangle + \frac{m_2}{2B+1} + e_2$, and we can get $c_1 + c_2 = (a_1 + a_2, b_1 + b_2)$.

$\text{Dec}(s_1, c_1 + c_2) = \lfloor (b_1 + b_2 - \langle s_1, a_1 + a_2 \rangle) \times (2B + 1) \rfloor$
 $= \lfloor (m_1 + m_2) + (e_1 + e_2) \times (2B + 1) \rfloor$

Suppose $m_1 + m_2 \in [-B, B]$, and the noise $(e_1 + e_2) \times (2B + 1)$ is within the controllable range, and the ciphertext is expanded in the ciphertext space after the homomorphic addition, then the decryption of homomorphic addition result will be correct with overwhelming probability: $\text{Dec}(s_1, (c_1 + c_2)) = m_1 + m_2$. Therefore, addition over ciphertext is homomorphic.

Homomorphic multiplication. Homomorphic multiplication supports the calculation between ciphertext and an integer plaintext. Let k be an integer constant in plaintext, the message $m \in [-B, B]$, a is randomly chosen from \mathbb{T}^n , $c = \text{Enc}(m)$, $b = \langle s, a \rangle + \frac{m}{2B+1} + e$. The multiplication over ciphertext and plaintext and decryption processes are as follows.

$k \times c = k \times (a, b) = (k \times a, k \times b)$

$\text{Dec}(s, k \times c) = \lfloor (k \times b - \langle s, k \times a \rangle) \times (2B + 1) \rfloor$
 $= \lfloor k \times m + k \times e \times (2B + 1) \rfloor$

If $k \times m \in [-B, B]$ and the noise $k \times e \times (2B + 1)$ is within the controllable range, and the expanded ciphertext size is within the ciphertext space after the homomorphic multiplications, the decryption of multiplication will succeed with overwhelming probability: $\text{Dec}(s, k \times$

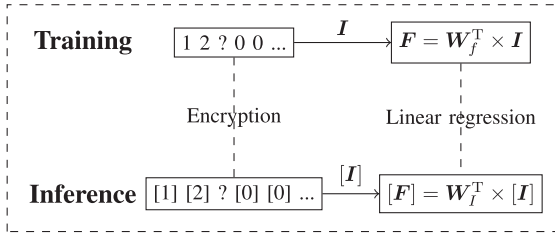


Fig. 3. Linear regression model. The CSP trains the model on the plaintext, rounds the model's weight, then constructs the ciphertext model to accept input ciphertexts from the client.

$c) = k \times m$. Therefore, we believe that the multiplication over ciphertext and plaintext is homomorphic.

Here is a toy example where insecure parameters are used for straightforward explanation. Let us choose $n = 4, B = 12, m = 5$ and $a = (0.1, 0, 0.1, 0), s = (0, 1, 0, 1), e = 0.001$. To encrypt, we need to compute $b = \langle s, a \rangle + \frac{m}{2B+1} + e = 0.3001$, thus the ciphertext is $(a, b) = (0.1, 0, 0, ., 0, 0.3001)$. We can use s to decrypt and then $m = \lfloor (b - \langle s, a \rangle) \times (2B + 1) \rfloor = \lfloor 5.0025 \rfloor = 5$ can be obtained.

3.2. Client

Clients are limited by high computation capacity and genotype imputation technology; thus, they are willing to upload genetic data to a third-party platform for convenient imputation. Before clients send their data to CSP, they first encrypt the genetic data with a secret key, and CSP will impute missing genotypes over these encrypted data.

We use the subsets of genetic sequences [45] to reduce the cost of large-scale genotype imputation and enhance the power of genetic data analysis. The genotypes in the subset are called tag variants, and we study the relationship between these tag variants to impute missing or low-quality variant genotypes (called target variants).

Suppose user's genetic data in plaintext is $I (n \times 1)$, each genotype value $I_i \in I$ is a discrete value which is defined as 0, 1, or 2. 0 denotes homozygous reference genotype, 1 denotes heterozygous genotype, and 2 indicates homozygous alternate genotype. The genetic data I is encrypted by $\text{Enc}(\cdot)$ operation into encrypted data $[I]$, and then client inputs $[I]$ into the secure imputation inference model.

$$I : (1, 2, 0, 0, \dots, 2) \xrightarrow{\text{Enc}(\cdot)} [I]([1], [2], [0], [0], \dots, [2]) \quad (1)$$

3.3. Cloud service provider

The cloud service provider (CSP) uses the linear correlation between genetic data to train and get the model's parameters by the public genomic dataset such as 1000 genomic project [46]. CSP takes the client's encrypted genetic data as the input of the secure inference model and returns the encrypted imputed result to the client.

CSP uses the forward propagation and backpropagation [47] to train the model on a public genetic dataset. The forward propagation calculates each neuron's output, and the backpropagation updates and optimizes the weight parameters. The training process usually involves thousands of iterations. We use stochastic gradient descent, and the parameters will be optimized after a new sample is trained. CSP trains the model on the public genetic dataset and imputes missing SNPs over the encrypted data. Fig. 2 illustrated the process of training and inference.

3.3.1. Model training

As shown in Fig. 3, the linear regression model's input is $I (n \times 1)$, and $I_i \in \{0, 1, 2\}$, while its output is a float vector $F (m \times 1)$. Before the model starts training, weight parameters $W_f (n \times m)$ are randomly initialized as float numbers between (0, 1). In the forward propagation, $F = W_f^T \times I$. During the backpropagation process, we use the mean-square error function as the cost function J , where F' denotes the target genotype, F is the model's output.

$$J = \frac{1}{2}(F - F')^2 \quad (2)$$

The stochastic gradient descent algorithm is performed according to Eq. (3), where α is the learning rate and W_f represents the weight parameters.

$$W_f = W_f - \alpha \frac{\partial J}{\partial W_f} \quad (3)$$

The parameters of the model are iterative optimized until J convergence to a minimal value or the iteration number reaches the limitation. After training the model, CSP constructs a homomorphic linear regression inference model based on the trained parameters.

3.3.2. Secure inference model

As shown in Fig. 2, the inference model takes the encrypted genetic data $[I] (n \times 1)$ as the input. Since the trained parameters on plaintext are float and homomorphic multiplication requires the plaintext to be an integer, we convert W_f into integer W_I . We keep two numbers after the decimal point for the W_f and scale them by 100 times. We use integers to approximate the floating-point weight, without affecting the encrypted genetic data and homomorphic computation on the ciphertext, as seen from the homomorphic multiplication formula in Section 3. Thus, the conversion will not reduce the complexity of the LWE-based homomorphic computation and security level.

$$W_I = \lfloor W_f \times 100 \rfloor \quad (4)$$

$$[F] = W_I^T \times [I] \quad (5)$$

Where $i \in [1, n]$ is an integer, and $[F]$ is encrypted imputation result. The inference of missing SNPs is shown as Eq. (5). After CSP infers the missing SNPs, it returns back the encrypted imputation result $[F]$ to the client. The client decrypts $[F]$ with the secret key, then decodes by scaling it by $1/100$. The decoded results are the predicted genotypes of missing SNPs.

To correctly evaluate the model's multi-sum, we need to include all possible values of $W_I^T \times I$ in the message space B [48], $W_I^T \times I \in [-B, B]$. Otherwise, the decryption of the multi-sum will fail. CSP first chooses models to train and obtains optimized model parameters W_I after training. The maximum of I_i is 2, and CSP can get the possible maximum multi-sum of the model: $\|W_I\|_1 \times 2$. As long as B satisfies the following formula, the multi-sum of the model will be in the range of $[-B, B]$, and homomorphic decryption will succeed correctly with an overwhelming possibility.

$$B \geq \|W_I\|_1 \times 2 \quad (6)$$

3.4. Analysis of noise

With homomorphic addition and multiplication calculation, the ciphertext's noise will expand. The expanded noise will result in decryption failure if calculation times are not limited. Thus, we should limit calculation times by reducing model's input units and the size of $\|W_I\|_2$. The standard deviation is used to evaluate whether ciphertext's noise is out of bounds, which is σ^2 in a fresh ciphertext. With every multiplication, the standard deviation gets larger by the square of the multiplier. When a ciphertext is multiplied by an integer p , the noise's standard deviation of obtained ciphertext will be expanded by p^2 times. To decrypt correctly (noise will not overflow), the following inequation needs to be satisfied:

$$\|W_I\|_2^2 \times \sigma^2 < \frac{1}{4B} \quad (7)$$

Table 1

Distance	Dataset	Tag SNPs	Target SNPs
1k	sorted_tag_SNPs_1k_genotypes	9764	500
10k	sorted_tag_SNPs_10k_genotypes	1045	

3.5. Analysis of parameter rounding

Theorem 1. Let $W_\xi = W_f - W_I/100$ be difference of between the float and integer weight parameters, $F_f = W_f^T \times I$ is the original model's output, $F_I = W_I^T \times I$ is the output of model with integer weights, $\text{Dec}([F])$ denotes the decrypted result of secure inference model's output. Since the prediction result on the plaintext is close to an integer, the noise generated by rounding the weight parameter is very small and will not affect the prediction accuracy.

Proof. W_I takes the accuracy of W_f 's two decimal places, the values in W_ξ : $W_\xi^i < 5 \times 10^{-3}$.

$$\begin{aligned} & W_f^T \times I - (W_I^T \times I)/100 \\ &= (W_f^T/100 + W_\xi^T) \times I - (W_I^T \times I)/100 \\ &= W_\xi^T \times I \end{aligned} \quad (8)$$

$$\begin{aligned} & W_f^T \times I - \text{Dec}([F])/100 \\ &= (W_f^T/100 + W_\xi^T) \times I - \\ & \lfloor W_f^T \times I + W_I^T \times e \times (2B + 1) \rfloor / 100 \\ &= W_\xi^T \times I \end{aligned} \quad (9)$$

Eq. (8) represents the output difference between models on the plaintext. From the Eq. (9), We can find that the noise in the encryption is taken from a Gaussian sample centered on the input message, with the standard deviation sd . It will not affect decryption. At the same time, since the output value on the plaintext is close to an integer and the $W_\xi < 5 \times 10^{-3}$ is small, which will not affect the imputation accuracy. Thus, even though the secure model's decrypted result contains noise, it does not affect the imputation result. \square

4. Experiments

Our models are implemented in C++. The models are run on a PC with i7-6700 CPU and 8G RAM. This section describes the experimental dataset, parameter settings, imputation accuracy, resource usage, and time consumption. The code address: https://github.com/tfhe-genotype-imputation/HE_genotype.

4.1. Dataset

The simulation datasets include two datasets that come from the iDASH Secure Genome Analysis Challenge 2019, containing 2504 individuals' genetic data.

As shown in Table 1, in the "sorted_tag_SNPs_1k_genotypes" dataset, it includes each individual's 9764 SNPs, and the distance between two nearby genotypes is 1k. "sorted_tag_SNPs_1k_genotypes" dataset includes each individual's 1045 SNPs, and the distance is 10k. The 500 target SNPs are the missing SNPs to be imputed. In experiments, 1500 individuals are used as the training set and 1004 as the test set (3:2). The dataset can be found in the following URL.¹

Table 2

Imputation accuracy of homomorphic models.

Models	Size	Accuracy
UTMSR [41]	32→1	95.40%
EPFL [41]	–	95.50%
CHIMERA [41]	45→1	95.10%
SNU [41]	24→1	95.00%
Ours	10→1	96.66%
	30→1	98.55%
	70→1	98.60%

4.2. Parameters

The homomorphic evaluation is based on the security of the LWE problem. The setting followed the notations of [21]. The encryption phase uses LWE security notions with no bootstrapping operations and no key-switching key. We estimated the security level from the attack models by the LWE estimator from [49] that computes the computational costs of state-of-art (R)LWE attack algorithms. We employed the LWE estimator to estimate hardness for the standard deviation $sd(2^{-25})$ and dimension $n(1024)$ and get an estimated 130-bits of security. The parameters related to the LWE estimator are the following:

- Ciphertext dimension: $n = 1024$;
- Noise standard deviation: $sd = \text{pow}(2., -25)$;
- Noise rate: $\alpha = \text{sqrt}(2 * \pi) * stdev$;
- Compatibility: $q = \text{pow}(2., 32)$;
- The attacker can use any number of samples: $m = \infty$

Inspired by the methods in [41], we conducted experiments on models of increasing input sizes. In the dataset of the different distances between variants, the models include 10, 30, and 70 tag SNPs for a single target SNP, and we represent the models as "10 → 1", "30 → 1", and "70 → 1" models.

Finally, we calculated the message space: $B = \|W_I^T\|_1 \times 2$, W_I denotes weight parameters of each model. In the experiment, we set $B = 700$, which is slightly smaller than the calculated value, and we found that it did not affect the results of the homomorphic evaluation.

4.3. Imputation accuracy

As shown in Table 2, we give the accuracy of models based on the HE scheme. Our secure linear regression reference model is similar to the UTMSR's model in [41], logistic regression models in CHIMERA and EPFL. But we select the most appropriate nearby genotypes for each missing gene. Experimental results show that our "30→1" model is about 3% higher than their model on accuracy. Unlike one-hidden layer neural network in SNU, we applied a more straightforward structure, and the results show our model is 1%–3% higher than theirs.

We construct three models on two datasets to illustrate our models' performance under different input sizes and security levels. We refer to the original models with float weights as float plaintext models, the model with integer weights as integer plaintext models, and the models on ciphertext as secure inference models for convenience.

In Tables 3 and 4, the second value refers to the accuracy of float plaintext models, and the third value refers to the integer plaintext models' accuracy, the last value refers to the accuracy gap between float plaintext models and secure inference models. The tables show that parameter rounding operation does affect the imputation accuracy. On the 1k dataset, the accuracy of the models with integer parameters is about 0.1% lower than float plaintext models. On the 10k dataset, the accuracy difference is about 0.2%. As seen from Tables 3 and 4, we can find that the accuracy of the secure inference model is close to that of the integer plaintext model. There is a small improvement because some results of the plaintext model are close to the middle of the label, which becomes correct after adding the noise of the rounding weights.

¹ <http://www.humangenomeprivacy.org/2019/competition-tasks.html>

Table 3
Test accuracy on 1k dataset.

Input→Output	Float	Integer	Security (bits)	Ciphertext	Gap
10→1	96.74%	96.65%	80	96.65%	0.09%
			130	96.65%	0.09%
30→1	98.57%	98.55%	80	98.55%	0.02%
			130	98.55%	0.02%
70→1	98.65%	98.60%	80	98.60%	0.05%
			130	98.60%	0.05%

Table 4
Test accuracy on 10k dataset.

Input→Output	Float	Integer	Security (bits)	Ciphertext	Gap
10→1	87.17%	86.87%	80	86.87%	0.30%
			130	86.87%	0.30%
30→1	88.48%	88.29%	80	88.29%	0.19%
			130	88.29%	0.19%
70→1	88.40%	88.52%	80	88.52%	0.12%
			130	88.52%	0.12%

Table 5
Memory Usage.

Models	Size	Memory (gigabytes)
UTMSR	32→1	0.03
CHIMERA	45→1	0.02
SNU	24→1	0.13
EPFL	-	0.06
Ours	10→1	0.26
	30→1	0.30
	70→1	0.39

On the 1k dataset, the secure model’s imputation accuracy for a single genotype achieved 98.6%. The “70→1” homomorphic model maintained the highest accuracy and is slightly 0.04% higher than the “30→1” homomorphic model. Furthermore, on the 10k dataset, the “70→1” homomorphic model is slightly higher than the “30→1” homomorphic model. The experimental results show that the number of nearby genotypes influences imputation accuracy. The more nearby genotypes (larger input size), the higher the accuracy.

Under the same model size, the imputation accuracy of the model with a nearby genotype distance of 1k is almost 10% higher than that of the model with 10k. The experimental results show that the accuracy of the inference model is the same under the security level of 130 bits and 80 bits.

4.4. Resource usage

With 80 bits or 130 bits security level, each genotype (0, 1, or 2) in the clear takes 2 bits, and each LWE ciphertext takes 8 bytes (64 bits). Therefore, the storage of ciphertext is 32 times that of plaintext.

As shown in Table 5, the memory usage of the proposed scheme is larger than that in the [41]. Our scheme required less than 0.39 GB.

4.5. Time consumption

We divided homomorphic evaluation into three processes: encryption, homomorphic calculation, and decryption. Table 6 refers to each process’s time consumption per 1000 individuals with a different security level. We can see from the table that the encryption operation consumes the most time, accounting for more than 95% of the total time. The whole homomorphic calculation did not use bootstrapping operations and key-switching keys, so the homomorphic calculation only spent about 5% of the total time. The decryption step took the

Table 6
Time consumption of each process.

Security (bits)	Input → Output	Enc(s)	Calculation(s)	Dec(s)
80	10 → 1	0.1941	0.00237	0.00069
		0.2438	0.00298	0.00075
130	30 → 1	0.5788	0.00619	0.00069
		0.7323	0.0078	0.00076
80	70 → 1	1.3583	0.0139	0.00072
		1.6967	0.0175	0.00076

Table 7
Total time consumption.

Input→ Output	Total Time (s)	Time (ms/variant)
10 → 1	0.269	0.268
30 → 1	0.78	0.777
70 → 1	1.71	1.708

Table 8
Average time consumption of each process for each SNP after data packaging.

Security (bits)	Input → Output	Enc(ms)	Calculation(ms)	Dec(ms)
80	10 → 1	0.507	0.01	0.011
		0.513	0.011	0.012
130	30 → 1	1.578	0.031	0.011
		1.6	0.04	0.012
80	70 → 1	3.43	0.075	0.012
		3.47	0.077	0.012

least time. The experimental results also show that each step’s time consumption has a linear relationship with the model’s input size.

In addition, Table 6 shows that the secure inference model with 80 bits security level consumes nearly 25% less time than the model with 130 bits security level.

Table 7 describes the total time consumption of homomorphic models with the security level of 130 bits. The test dataset has 1004 individuals. In the “10 → 1” model, the homomorphic evaluation time for each variant of 1004 individuals is approximately 0.269 s, the “30 → 1” model is about 0.78 s per variant per 1004 individuals, and the “70 → 1” model is approximately 1.71 s. The experimental results show that the secure linear model with 30 tag SNPs as the model’s input for a single genotype shows the most balanced performance in terms of timing and imputation accuracy.

Since there are 1004 individuals in the test dataset with the same calculation in homomorphic linear regression, the genes with the same SNP’s identifier of different individuals can be packaged into one ciphertext, which can calculate repeated operations in parallel. As a result, the time consumption will be significantly reduced. For the prediction of the same SNP of 1004 individuals in the “30-1” model, the input changed from 1004 × 30 LWE ciphertexts to 30 TLWE ciphertexts, and the output with 1004 LWE ciphertexts are packaged in one TLWE ciphertext. Table 8 shows that packaging multiple inputs into a single ciphertext can reduce the time consumption of our model nearly 300 times, and the accuracy remains the same.

5. Conclusion

In this work, we propose a secure and fast linear regression inference model to impute the missing genotypes. Homomorphic encryption is time-consuming and only allows simple addition or binary gates on the ciphertext. Thus we design the secure inference model carefully to maintain high imputation accuracy and efficiency. We round the trained weights to integers and reduce the time consumption of homomorphic evaluation without affecting the security level. The secure genotype imputation inference model reduces the cost of large-scale gene sequencing and guarantees the safety of genes. If new variants must be imputed, training and homomorphic imputation are performed

independently, with fast training and prediction. Since our basic model is a simple linear regression model, we consider using neural network models and activation functions to improve accuracy in the future. We also consider using ciphertext packaging technology to reduce data encryption time.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Data availability

The authors do not have permission to share data.

Acknowledgment

The work described in this paper was supported in part by the Key Research and Development Program of Hubei Province (Grant No. 2022BAA050), and the Key Research and Development Program of Hainan Province (Grant No. ZDYF2021GXJS014).

References

- [1] Shendure J, Balasubramanian S, Church GM, Gilbert W, Rogers J, Schloss JA, et al. DNA sequencing at 40: Past, present and future. *Nature* 2017;550(7676):345–53.
- [2] Allen HL, Estrada K, Lettre G, Berndt SI, Weedon MN, Rivadeneira F, et al. Hundreds of variants clustered in genomic loci and biological pathways affect human height. *Nature* 2010;467(7317):832–8.
- [3] Locke AE, Kahali B, Berndt SI, Justice AE, Pers TH, Day FR, et al. Genetic studies of body mass index yield new insights for obesity biology. *Nature* 2015;518(7538):197–206.
- [4] Tam V, Patel N, Turcotte M, Bossé Y, Paré G, Meyre D. Benefits and limitations of genome-wide association studies. *Nature Rev Genet* 2019;20(8):467–84.
- [5] Evangelou E, Ioannidis JP. Meta-analysis methods for genome-wide association studies and beyond. *Nature Rev Genet* 2013;14(6):379–89.
- [6] Schaid DJ, Chen W, Larson NB. From genome-wide associations to candidate causal variants by statistical fine-mapping. *Nature Rev Genet* 2018;19(8):491–504.
- [7] Howie BN, Donnelly P, Marchini J. A flexible and accurate genotype imputation method for the next generation of genome-wide association studies. *PLoS Genetics* 2009;5(6):e1000529.
- [8] Das S, Forer L, Schönherr S, Sidore C, Locke AE, Kwong A, et al. Next-generation genotype imputation service and methods. *Nature Genet* 2016;48(10):1284–7.
- [9] Browning BL, Zhou Y, Browning SR. A one-penny imputed genome from next-generation reference panels. *Am J Hum Genet* 2018;103(3):338–48.
- [10] Servin B, Stephens M. Imputation-based analysis of association studies: Candidate regions and quantitative traits. *PLoS Genetics* 2007;3(7):e114.
- [11] Lin Z, Owen AB, Altman RB. Genomic research and human subject privacy. *Science* 2004;305(5681):183.
- [12] Zhang L, Cui Y, Mu Y. Improving security and privacy attribute based data sharing in cloud computing. *IEEE Syst J* 2019;14(1):387–97.
- [13] Naveed M, Ayday E, Clayton EW, Fellay J, Gunter CA, Hubaux J-P, et al. Privacy in the genomic era. *ACM Comput Surv* 2015;48(1):1–44.
- [14] Berger B, Cho H. Emerging technologies towards enhancing privacy in genomic data sharing. *Genome Biol* 2019;20(1):1–3.
- [15] Acar A, Aksu H, Uluagac AS, Conti M. A survey on homomorphic encryption schemes: Theory and implementation. *ACM Comput Surv* 2018;51(4):1–35.
- [16] Gürsoy G, Chielie E, Brannon CM, Maniatakos M, Gerstein M. Privacy-preserving genotype imputation with fully homomorphic encryption. *Cell Syst* 2022;13(2):173–82.
- [17] Bost R, Popa RA, Tu S, Goldwasser S. Machine learning classification over encrypted data. 2014, Cryptology EPrint Archive.
- [18] Mohassel P, Zhang Y. Secureml: A system for scalable privacy-preserving machine learning. In: 2017 IEEE symposium on security and privacy. IEEE; 2017, p. 19–38.
- [19] Brakerski Z. Fully homomorphic encryption without modulus switching from classical GapSVP. In: Annual cryptography conference. Springer; 2012, p. 868–86.
- [20] Cheon JH, Kim A, Kim M, Song Y. Homomorphic encryption for arithmetic of approximate numbers. In: International conference on the theory and application of cryptology and information security. Springer; 2017, p. 409–37.
- [21] Chillotti I, Gama N, Georgieva M, Izabachène M. TFHE: Fast fully homomorphic encryption over the torus. *J Cryptol* 2020;33(1):34–91.
- [22] Albrecht M, Chase M, Chen H, Ding J, Goldwasser S, Gorbunov S, et al. Homomorphic encryption standard. In: Protecting privacy through homomorphic encryption. Springer; 2021, p. 31–62.
- [23] Michie D, Spiegelhalter DJ, Taylor C, et al. Machine learning. *Neural Stat Classification* 1994;13(1994):1–298.
- [24] Rivest RL, Adleman L, Dertouzos ML, et al. On data banks and privacy homomorphisms. *Found Secure Comput* 1978;4(11):169–80.
- [25] Cominetti EL, Simplicio MA. Fast additive partially homomorphic encryption from the approximate common divisor problem. *IEEE Trans Inf Forensics Secur* 2020;15:2988–98.
- [26] Fan J, Vercauteren F. Somewhat practical fully homomorphic encryption. *IACR Cryptol ePrint Arch* 2012;2012:144.
- [27] Chen H, Iliashenko I, Laine K. When HEAAN meets FV: A new somewhat homomorphic encryption with reduced memory overhead. *IACR Cryptol EPrint Arch* 2020;2020:121.
- [28] Brakerski Z, Gentry C, Vaikuntanathan V. (Leveled) fully homomorphic encryption without bootstrapping. *ACM Trans Comput Theory (TOCT)* 2014;6(3):1–36.
- [29] Gentry C, Halevi S. Implementing gentry’s fully-homomorphic encryption scheme. In: Annual international conference on the theory and applications of cryptographic techniques. Springer; 2011, p. 129–48.
- [30] Van Dijk M, Gentry C, Halevi S, Vaikuntanathan V. Fully homomorphic encryption over the integers. In: Annual international conference on the theory and applications of cryptographic techniques. Springer; 2010, p. 24–43.
- [31] Chillotti I, Gama N, Georgieva M, Izabachene M. Faster fully homomorphic encryption: Bootstrapping in less than 0.1 seconds. In: International conference on the theory and application of cryptology and information security. Springer; 2016, p. 3–33.
- [32] Liu Z, Seo H, Roy SS, Großschädl J, Kim H, Verbauwhede I. Efficient ring-LWE encryption on 8-bit AVR processors. In: International workshop on cryptographic hardware and embedded systems. Springer; 2015, p. 663–82.
- [33] Gentry C, Sahai A, Waters B. Homomorphic encryption from learning with errors: Conceptually-simpler, asymptotically-faster, attribute-based. In: Annual cryptography conference. Springer; 2013, p. 75–92.
- [34] Chillotti I, Joye M, Paillier P. New challenges for fully homomorphic encryption. In: Privacy-preserving machine learning (PPML-PrIML 2020) NeurIPS 2020 workshop. 2020.
- [35] Hong S, Park JH, Cho W, Choe H, Cheon JH. Secure multi-label tumor classification using homomorphic encryption. *Res Square* 2021.
- [36] Pereira HVL. Bootstrapping fully homomorphic encryption over the integers in less than one second. In: IACR international conference on public-key cryptography. Springer; 2021, p. 331–59.
- [37] Chillotti I, Joye M, Paillier P. Programmable bootstrapping enables efficient homomorphic inference of deep neural networks. *IACR Cryptol EPrint Arch* 2021;2021:91.
- [38] Kocabas O, Soyata T. Utilizing homomorphic encryption to implement secure and private medical cloud computing. In: 2015 IEEE 8th international conference on cloud computing. IEEE; 2015, p. 540–7.
- [39] Meehan A, Ko RK, Holmes G. Deep learning inferences with hybrid homomorphic encryption. 2018.
- [40] Wood A, Najarian K, Kahrobaei D. Homomorphic encryption for machine learning in medicine and bioinformatics. *ACM Comput Surv* 2020;53(4):1–35.
- [41] Kim M, Harmanci AO, Bossuat J-P, Carpov S, Cheon JH, Chillotti I, et al. Ultrafast homomorphic encryption models enable secure outsourcing of genotype imputation. *Cell Syst* 2021;12(11):1108–20.
- [42] Dokmai N, Kockan C, Zhu K, Wang X, Sahinalp SC, Cho H. Privacy-preserving genotype imputation in a trusted execution environment. 2021, BioRxiv, Cold Spring Harbor Laboratory.
- [43] Halevi S, Shoup V. HELIB-An Implementation of homomorphic encryption. *Cryptology EPrint Archive, Report* 2014/039, 2014.
- [44] Regev O. On lattices, learning with errors, random linear codes, and cryptography. *J ACM* 2009;56(6):1–40.
- [45] Das S, Abecasis GR, Browning BL. Genotype imputation from large reference panels. *Annu Rev Genom Hum Genet* 2018;19:73–96.
- [46] Siva N. 1000 Genomes project. *Nature Biotechnol* 2008;26(3):256–7.
- [47] Li J, Cheng J-h, Shi J-y, Huang F. Brief introduction of back propagation (BP) neural network algorithm and its improvement. In: Advances in computer science and information engineering. Springer; 2012, p. 553–8.
- [48] Bourse F, Minelli M, Minihold M, Paillier P. Fast homomorphic evaluation of deep discretized neural networks. In: Annual international cryptology conference. Springer; 2018, p. 483–512.
- [49] Albrecht MR, Player R, Scott S. On the concrete hardness of learning with errors. *J Math Cryptol* 2015;9(3):169–203.