

Efficient FHE-Based Privacy-Enhanced Neural Network for Trustworthy AI-as-a-Service

Kwok-Yan Lam , Senior Member, IEEE, Xianhui Lu , Linru Zhang , Xiangning Wang , Huaxiong Wang , and Si Qi Goh 

Abstract—AI-as-a-Service has emerged as an important trend for supporting the growth of the digital economy. Digital service providers make use of their vast amount of customer data to train AI models (such as image recognition, financial modelling and pandemic modelling etc) and offer them as a service on the cloud. While there are convincing advantages for using such third-party models, the fact that model users are required to upload their data to the cloud is bound to raise serious privacy concerns, especially in the face of increasingly stringent privacy regulations and legislation. To promote the adoption of AI-as-a-Service while addressing privacy issues, we propose a practical approach for constructing privacy-enhanced neural networks by designing an efficient implementation of fully homomorphic encryption. With this approach, an existing neural network can be converted to process FHE-encrypted data and produce encrypted output which are only accessible by the model users, and more importantly, within an operationally acceptable time (e.g., within 1 s for facial recognition in typical border control systems). Experimental results show that in many practical tasks such as facial recognition, text classification and so on, we obtained the state-of-the-art inference accuracy in less than one second on a 16 cores CPU.

Index Terms—Fully homomorphic encryption, privacy-enhanced neural networks, look-up table algorithm, deep neural network, digital trust, secure cloud computing, data privacy, cryptographic protocol, applied cryptography.

I. INTRODUCTION

AI-AS-A-SERVICE (AIaaS) has been experiencing rapid development in recent years due to the strong demand in sharing of sophisticated and powerful AI models, which require not only technological advancement but also availability of vast amount of data resources. For digital service providers, it allows them to deploy their trained AI models on the cloud and offer the inference as a service, instead of disclosing their models to their users. For such users, with the trained AI model as a cloud

service, they simply upload their inputs to the cloud and obtain the inference results in return, instead of being required to train their own computationally expensive AI models in-house.

Besides, AIaaS plays an important role in the development of smart cities by supporting efficient implementation of intelligent distributed solutions, which typically also involve the adoption of Internet of Things (IoT) technology to gather data from the physical environment for automated control and decision-making at the cloud. However, the data handling process may lead to privacy and security concerns [1], [2], [3]. When using AIaaS, there are always non-trivial concerns that the cloud servers can access users' raw data. Either users' input or the inference result may contain privacy-sensitive information, such as biometrics, healthcare records and financial data. A cloud service that offers privacy-enhanced inference of AI models is highly desirable in order to promote the adoption of AIaaS, so that users can enjoy the service without disclosing their plaintext data to the cloud.

Typically, AI models, such as deep neural networks (DNNs), apply a sequence of evaluations on the input data and model parameters to obtain an inference output. Many techniques are studied for privacy-preserving machine learning, such as differential privacy [4], [5] and federated learning [6], [7], [8], [9]. These techniques address the concern of privacy issues of the training data and of the data sharing during model training. In this work, we focus on privacy-enhanced neural networks (PE-NN) inference and propose a Fully Homomorphic encryption (FHE) [10] based solution. FHE provides a way to encrypt data while supporting computations through the encryption envelope [11]. There are recent works [11], [12], [13], [14], [15], [16] that reported feasible implementation of privacy-enhanced neural network inferences over encrypted data by applying FHE. The flow is shown in Fig. 1.

In FHE-based PE-NN inferences, the user encrypts its sensitive data before sending it to the server. The server homomorphically evaluates the neural network over encrypted data and produces an encrypted inference output, then it returns the ciphertext to the user. The user who possesses the corresponding private key can decrypt the encrypted result. The server does not have the private key, hence it is unable to decrypt neither the input nor the output. In this model, users can enjoy AI services on the cloud without disclosing their plaintext data, and AI service providers also do not disclose the trained AI models to users.

An additional instance of Privacy-Enhancing Neural Networks (PE-NN) manifests in the realm of international

Manuscript received 8 August 2023; revised 28 December 2023; accepted 10 January 2024. Date of publication 12 January 2024; date of current version 4 September 2024. This research was supported in part by the National Research Foundation, Singapore and Infocomm Media Development Authority under its Trust Tech Funding Initiative and Strategic Capability Research Centres Funding Initiative. (Corresponding author: Kwok-Yan Lam.)

Kwok-Yan Lam, Linru Zhang, Xiangning Wang, Huaxiong Wang, and Si Qi Goh are with Nanyang Technological University, Singapore 639798 (e-mail: kwokyan.lam@ntu.edu.sg; linru.zhang@ntu.edu.sg; xiangning.wang@ntu.edu.sg; hxwang@ntu.edu.sg; siqi005@e.ntu.edu.sg).

Xianhui Lu is with the Chinese Academy of Sciences, Beijing 100045, China (e-mail: luxianhui@iie.ac.cn).

This article has supplementary downloadable material available at <https://doi.org/10.1109/TDSC.2024.3353536>, provided by the authors.

Digital Object Identifier 10.1109/TDSC.2024.3353536

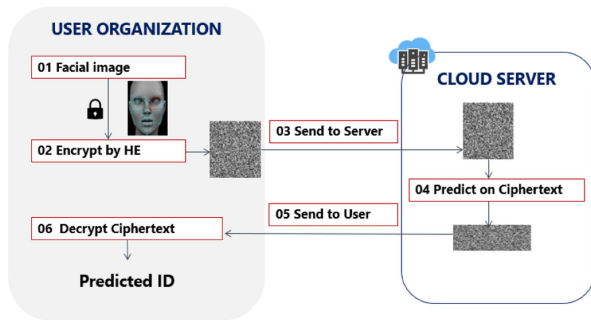


Fig. 1. Basic FHE-based privacy-enhanced neural network model (basic FHE-PE-NN).

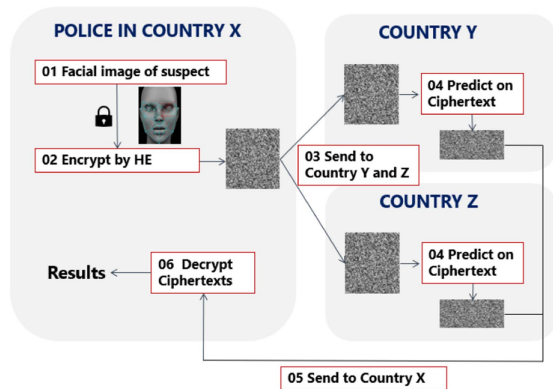


Fig. 2. Example: International collaboration to fight against transnational crime (basic model).

cooperation aimed at combating trans-border crimes. For example, when a suspect is identified within country X, law enforcement agencies (LEA) seek to ascertain whether the individual possesses a criminal record in other nations. However, the law enforcement agencies in country X are faced with limitations in accessing the datasets of other countries while transmitting unsecured data for cross-checking purposes is prohibited. Through the utilization of privacy-enhancing techniques, it is possible to employ Fully Homomorphic Encryption (FHE) to encrypt the suspect's data and transmit the resulting ciphertext to the collaborating authorities. The receiving end can subsequently perform homomorphic evaluations to determine the existence of a criminal record associated with the suspect. As the evaluation results are likewise encrypted, they can be securely transmitted back to country X. Finally, the law enforcement agency decrypts the received data using their own private key, thereby obtaining the pertinent results. Homomorphic encryption could also be utilized to perform secure data analysis and operations on other sensitive information. This could involve running encrypted machine learning algorithms on encrypted data to extract insights, identify patterns, or optimize criminal operations without exposing the underlying data. The flow is shown in Fig. 2.

However, high inference latency and low accuracy restrict existing works to be applied in real-world AI services.

a) Inference latency. High inference latency is an obstacle to applying FHE-based PE-NN in DNNs. Existing works usually

take very long time per inference in simple and shallow networks. CryptoNets [13] takes over 200 seconds per inference on encrypted image from MNIST [17]. Faster CryptoNets improved the result. It takes around 40 seconds and 20,000 seconds per inference on encrypted image from MNIST and encrypted image from CIFAR-10 [18] respectively. Recently, [14] further improved the latency and achieves around 10 seconds and 2000 seconds per inference on MNIST and CIFAR-10 respectively.

MNIST and CIFAR-10 are most commonly used datasets in existing works for benchmarks, which are simple and well studied. But in real-world applications, we usually anticipate neural network can solve more complex problems such as facial recognition and text classification. For such tasks, only [14] mentioned that it uses architectures of AlexNet, ResNet-18 and Shuffle Net for inferences on ImageNet and the result is 5 hours per inference with an accuracy of 69.4%.

To solve this problem, we analyze the architectures of various FHE schemes and concluded that there is no general-purpose FHE scheme which performs well in all homomorphic evaluations in neural networks. Evaluation of neural networks on encrypted data involves both arithmetic operations such as weighted sum and convolutions, evaluation of non-linear activations and bootstrapping after each neuron to enable further computations.

In practice, there are two main kinds of FHE schemes which are required in constructing FHE-based PE-NN. The *word-wise CKKS/BGV/BFV schemes* [19], [20], [21] (adopted in many works such as CryptoNets [13]) are good at linear evaluations but the non-linear activation evaluations are implemented by polynomial approximation, while the activations are usually chosen from non-polynomial functions such as ReLU and Sigmoid. In the other aspect, these schemes accumulate noise and it is inevitable to run the significantly slow bootstrapping in order to reduce noise. The *bit-wise FHEW/TFHE schemes* [22], [23] provides very efficient non-linear evaluations and bootstrapping operation, but the linear evaluations have to be computed gate by gate. Neither word-wise schemes nor bit-wise schemes is good at all evaluations of neural networks.

In this research, we propose that an optimized FHE scheme which combines the advantages of CKKS/BGV/BFV and TFHE/FHEW schemes to reduce the inference latency.

b) Accuracy. The poor performance of homomorphic evaluations of activations causes the low accuracy. DiNN [11] and CryptoNets [13] only support sign function and square function as the activation function separately, where both of them are not commonly used in machine learning area. n-GraphHE [24], Lola [15] and Faster CryptoNets [12] have to use the low-degree polynomial approximation activations and thus fail to obtain the state-of-the-art inference accuracy. For example, Faster CryptoNet achieves 76.72% inference accuracy on CIFAR-10, while in an unencrypted network with ReLU activations, it is 93.72%.

The bottleneck to improve the accuracy is homomorphical evaluation of non-linear activations. ReLU and Sigmoid are widely used in modern neural networks. However, existing works in FHE-based PE-NN usually use sign function, square function, or low degree polynomials to approximate activations. Although it is possible to improve the inference accuracy of FHE-based PE-NN by enlarging the degree of polynomial

approximation activations, the computing overhead increases exponentially with the degree and thus the time taken becomes unacceptable. Efficient algorithms which can evaluate the activations in modern networks accurately over encrypted data are urgently needed.

In addition to the optimizations in FHE schemes, we review the structure of DNNs. DNNs, which may consists of hundreds of layers, are considered to be computationally expensive. If the whole DNN is evaluated on encrypted data in the AI model owner’s server, then it is obvious that it will be too slow to be acceptable. Recall that in the example of trans-border crime investigation, it may take more than one day to recognize the person and check the record homomorphically. It can be accelerated significantly if all parties adopt an open human face feature extractor such as FaceNet [25]. The LEA of country X can run face feature extractor locally without encryption. Then it encrypts the feature vector before sending it out. Based on our experiment, recognizing a person from an encrypted feature vector takes less than one second, which is much faster than facial recognition from encrypted photos and yet with no accuracy loss.

A. Our Results

Towards practical FHE-based PE-NN constructions, we propose a practical approach for constructing privacy-enhanced neural networks by designing an efficient implementation of fully homomorphic encryption.

We first propose a optimized fully homomorphic encryption scheme together with an efficient design for non-linear activation evaluation. We show that our FHE scheme achieves better results in both inference accuracy and time in the benchmark MNIST dataset.

Then we propose a new model which splits a deep neural network into a plaintext evaluation part and a ciphertext evaluation part. We deploy some pre-trained feature extraction layers to user’s side and evaluate it in clear. The fine-tuned layers are retained at model owner’s server and the input to the server is encrypted. In this way, we reduce the number of expensive homomorphic evaluations without disclosing either users’ data and model providers’ trained parameters. We call this model *Hybrid FHE-based privacy-enhanced neural network* (Hybrid FHE-based PE-NN).

The proposed hybrid model prevent itself from potential information leakage risks because the plaintext evaluation part, which is evaluated within the confines of the user’s trusted computing environment. Following the completion of the plaintext evaluation part, the user encrypts the feature vector before sending it to the server. The security guarantees defined by our encryption scheme substantiates the incapacity of the server to obtain any information of the feature vector, let alone the original data. Consequently, this hybrid model satisfies the requirements of strong privacy-preserving, and meanwhile is much more efficient.

As a result, we can perform facial recognition under one second, where it requires more than a day for basic FHE-based PE-NN model. We also show that our hybrid model can be used to solve practical tasks such as speaker verification, text classification and object classification.

c) Optimized FHE scheme. Our optimized FHE scheme enables us the ability to perform weighted sums and convolutions on the approximate LWE-based additive homomorphic encryption schemes, and to evaluate non-linear activations on FHEW ciphertexts. Therefore, the evaluations of both linear and non-linear functions are very fast. Then the noise will be reduced during the evaluation of non-linear activations by applying an optimized homomorphic look-up table algorithm (LUT) [26].

We also observed the fact that the neural network is good at noise tolerating, i.e., the neural networks are usually not sensitive to the noise in less significant bits of input. This noise tolerating property also help us to reduce the size of encryption parameters and thus improves the efficiency while keeping high inference accuracy in the following ways: 1) On receiving of an input, we discard the less significant bits before encrypting; 2) In the beginning of LUT algorithm, we also discard the less significant bits of the LWE ciphertext.

In summary, our optimized FHE scheme has the following properties: 1) It can be applied to neural networks of arbitrary depth. 2) It supports many kinds of widely used activations, such as ReLU and Sigmoid. 3) When applied to inference of privacy-enhanced neural networks, it is fast and accurate.

d) Efficient design for non-linear activation evaluation. We further improve the evaluations of non-linear activations to speed up our system. Typically, a general method used to evaluate non-linear activations is homomorphic look-up-table algorithm (LUT) [11], [26], [27]. The core idea of LUT is to encode the “Table” containing the value of the non-linear activation into a polynomial, so that we can use ciphertext to locate the position of the desired output.

In this work, we focus on improving the LUT algorithm. Our LUT algorithm takes the characteristics of neural network which tends to be able to tolerate noise into consideration and propose an approximate LUT algorithm over integer. Compared to DiNN [11], the output of our LUT is in an integer ring, instead of binary. Compared to PEGASUS [27], we choose a small degree of LUT function to reduce time cost because of the characteristics of activation function which tends to be able to tolerate noise. We also propose two approaches to improve the efficiency of LUT algorithms:

1) *RNS polynomial multiplication.* We use number theoretic transform (NTT) multiplication to speed up the multiplication. As a brief introduction, one NTT multiplication includes: 1) (NTT) turn two polynomials into two vectors (We call it residue numeral system (RNS) form); 2) (Position-wise multiplication) Perform position-wise multiplication between two vectors; 3) (Inverse-NTT (INTT)) turn the resulted vector into a polynomial. Lastly, we reduce the number of NTT and INTT to a third of before.

2) *Modulo function.* We design different modulo functions for different operations. We also reduce the number of modulo function by half compared to the previous method.

As a result, our improved system only takes 0.14s per inference on MNIST dataset, while 0.42s is needed without improvements.

e) Hybrid FHE-based PE-NN model. The applications of basic FHE-based PE-NN is restricted by the fact that evaluation an entire DNN on encrypted data is too inefficient. In order to reduce

the inference latency, we develop the idea that divide the network into two parts: an open network and a private network. We call it Hybrid FHE-based privacy-enhanced neural network model.

In this model, we make use of edge computing in a different purpose. The open network is distributed in the users' side and the private network is remained in the cloud server. The user first runs the open network in plaintext locally, then the results are encrypted and being sent to the server. The server evaluates the private network on encrypted input and returns an inference output which is also encrypted. Only the user who has the secret key can decrypt the output and see the result. By replacing some computationally expensive ciphertext evaluations in the server side with low-cost plaintext computations in the user side, our model significantly reduce the inference latency.

The open network is made up of general feature extracting layers such as convolutional layers. Usually, it is chosen from well-known open-source networks such as FaceNet [25] (for facial recognition), TextCNN [28] (for text classification) and InceptionV3 [29] (for object classification). The private model consists of some lower layers which is highly related to the dataset and the task. The private model is trained by the AI model owner. The training method and trained parameters are often considered as critical intellectual property by AI model owner, who are typically not willing to share them. Notice that when we set the whole network as private network, our model is exactly same as basic FHE-PE-NN model, so our model can be viewed as a generalization of basic FHE-PE-NN model.

B. Overview of the paper

In Section II, we introduce backgrounds of our work. In Section III, we propose our optimized FHE scheme and propose a protocol for its application on PE-NN. Followed by Section IV, where we further optimize our FHE scheme, and focus on improving the evaluation of non-linear activations. Next up in Section V, we propose a hybrid PE-NN model so that homomorphically evaluating deep complex neural networks becomes practical. Finally in Section VI we report the performance of our system and show our evaluation results of different PE-NN such as hand-writing digits recognition, facial recognition, speaker verification, text classification and object classification.

II. PRELIMINARIES

A. Notations

We denote all real numbers by \mathbb{R} . We denote all integers by \mathbb{Z} . For a positive integer q , we use $\mathbb{Z}_q := \mathbb{Z}/q\mathbb{Z}$ to denote the multiplicative group of integers modulo q . We use upper-case letter for matrix and use lower-case bold letters for vectors. Given a vector \mathbf{x} , we write $\mathbf{x}[i-1]$ for its i th entry, i.e., $\mathbf{x}[0]$ is its first entry. We denote $\langle \mathbf{a}, \mathbf{b} \rangle$ the inner product of vectors \mathbf{a} and \mathbf{b} . We use lower-case letters for polynomials. For a positive integer n , we write $[n]$ to denote the set $\{0, 1, 2, \dots, n-1\}$. We write $a \stackrel{\$}{\leftarrow} S$ to denote that, a is sampled uniformly random from set S . Let \mathcal{E} be a distribution, we use $e \stackrel{\$}{\leftarrow} \mathcal{E}$ to denote that e is randomly sampled according to \mathcal{E} .

Finally, we use $\mathbf{1}_S$ to denote the indicator function

$$\mathbf{1}_S := \begin{cases} 1, & S \text{ is TRUE} \\ 0, & S \text{ is FALSE} \end{cases}.$$

B. (Ring) Learning With Errors

The learning with errors (LWE) problem [30] was proposed as a generalization of learning parity with noise and it is widely used in construction of many cryptosystems.

For positive integers n and $q \geq 2$, a vector $\mathbf{s} \in \mathbb{Z}_q^n$, and a probability distribution $\mathcal{E} = \mathcal{E}(n)$ over \mathbb{Z}_q , let $\mathcal{A}_{\mathbf{s}, \mathcal{E}}$ be the distribution obtained by choosing a vector $\mathbf{a} \stackrel{\$}{\leftarrow} \mathbb{Z}_q^n$ uniformly at random and a noise term $e \stackrel{\$}{\leftarrow} \mathcal{E}$, and outputting $(\mathbf{a}, \langle \mathbf{a}, \mathbf{s} \rangle + e) \in \mathbb{Z}_q^n \times \mathbb{Z}_q$. The LWE problem is defined as follows.

Definition II.1. (LWE) For an integer $q = q(n)$ and an error distribution $\mathcal{E} = \mathcal{E}(n)$ over \mathbb{Z}_q , the LWE problem $\text{LWE}_{n,m,q,\mathcal{E}}$ is defined as: Given m independent samples from $\mathcal{A}_{\mathbf{s}, \mathcal{E}}$, output \mathbf{s} with non-negligible probability.

The decisional version is to distinguish between m samples chosen according to $\mathcal{A}_{\mathbf{s}, \mathcal{E}}$ for some uniformly random \mathbf{s} and m samples from the uniform distribution over $\mathbb{Z}_q^n \times \mathbb{Z}_q$.

The Ring learning with errors (RLWE) problem [31] is a variant of LWE which is widely used to design homomorphic encryption schemes [19], [20], [23], [32]. The secret s is chosen from ring R . An RLWE sample $(a, as + e) \in R_q^2$ is generated by choosing $a \stackrel{\$}{\leftarrow} R_q$ uniformly at random and noise e from the error distribution \mathcal{E} over R . Here $q \geq 2$ is an integer modulus. The decisional version is to distinguish between the RLWE sample for some secret s and a sample from uniform distribution over R_q^2 .

C. Homomorphic Encryption (HE)

A cryptosystem that supports computation on ciphertext without decryption is known as homomorphic encryption (HE) [22]. A symmetric key HE scheme consists of following PPT algorithms.

- Key generation. The algorithm takes security parameter λ as input and outputs a secret key sk , and a evaluation key evk .
- Encryption. The algorithm takes a message m and the secret key sk as input and outputs a ciphertext $ct = \text{Enc}(sk, m)$.
- Decryption. The algorithm takes a ciphertext ct and the secret key sk as input and output a message $m' = \text{Dec}(sk, ct)$.
- Homomorphic evaluation. The algorithm takes the evaluation key evk , a function f and a set of ciphertexts as input and output a ciphertext $ct = \text{Eval}(f, evk, ct_1, \dots, ct_l)$.

Roughly speaking, an encryption scheme is *homomorphic* in an operation \circ , if there is another operation \diamond such that $\text{Enc}(m_1) \diamond \text{Enc}(m_2) = \text{Enc}(m_1 \circ m_2)$.

The security notion we consider is indistinguishability under chosen plaintext attack (IND-CPA) security. [33]

Definition II.2 (IND-CPA security of symmetric key HE scheme). A symmetric key HE scheme with message space \mathcal{M} is IND-CPA secure if for any polynomial time adversary \mathcal{A} , it

Algorithm 1: NTT Multiplication.

-
- 1: **Input:** Polynomials $a, b \in \mathbb{Z}_q[x]/(x^n + 1)$.
 - 2: Find $\bar{a} = NTT(a), \bar{b} = NTT(b)$. (One $NTT(\cdot)$ costs $O(n \log n)$.)
 - 3: Calculate $\bar{c} = \bar{a} \times \bar{b}$. (Position-wise multiplication (PMUL), $O(n)$.)
 - 4: Find $c = INTT(\bar{c})$. (One $INTT(\cdot)$ costs $O(n \log n)$.)
 - 5: **Output:** Polynomial c .
-

holds that

$$\left| \Pr[\mathcal{A}_{evk}(\text{Enc}(sk, m_0)) = 1] - \Pr[\mathcal{A}_{evk}(\text{Enc}(sk, m_1)) = 1] \right| = \text{negl}(\lambda),$$

where $(sk, evk) \leftarrow \text{KeyGen}(1^\lambda)$, and $m_0, m_1 \in \mathcal{M}$.

Briefly speaking, an encryption scheme is *homomorphic* in an operation \circ , if there is another operation \diamond such that $\text{Enc}(m_1) \diamond \text{Enc}(m_2) = \text{Enc}(m_1 \circ m_2)$. There are many sub-classes of homomorphic encryption scheme, depending on the types and number of operations it supports.

- *Partially homomorphic encryption*. It only supports the evaluation of circuits consisting of one type of gate. It is usually simple and fast. For example, *ElGamal cryptosystem* [34] can evaluate unbounded number of modular multiplications, and *Paillier cryptosystem* [35] can evaluate unbounded number of modular additions.
- *Somewhat homomorphic encryption*. It supports the evaluation of arbitrary circuits with multiple types of gates of pre-determined bounded depth. Typically, ciphertexts are “noisy”, and the noise grows along with the increment of homomorphic computations, where ultimately the noise will make the resulting ciphertext indecipherable.
- *Fully homomorphic encryption (FHE)*. It supports the evaluation of arbitrary circuits with multiple types of gates of unbounded depth. FHE solves the noise problem by using bootstrapping technique. In bootstrapping, the ciphertext is “refreshed” to reduce its noise level [10].

D. Number-Theoretic Transform (NTT)

In our encryption scheme, almost all computations involve high-degree polynomial multiplications. Naive polynomial multiplication costs $O(n^2)$ time, where n is the degree of polynomials. In field \mathbb{R} or \mathbb{C} , using Faster Fourier Transforms (FFT) [36] for polynomial multiplication is a common technique. It changes the time cost from $O(n^2)$ to $O(n \log n)$, and this will have significant improvement when n is large, especially in the area of cryptography. Note that we only consider integers in cryptography, so a variant of FFT algorithm is widely used on finite field, which is called number-theoretic transform (NTT) [37].

The basic idea of NTT is that for some appropriately chosen prime q , $\mathbb{Z}_q[x]/(x^n + 1)$ ¹ and \mathbb{Z}_q^n are isomorphic, as

¹We briefly introduce the choice of q and n here. We choose n is a power of 2 in $x^n + 1$, i.e., $x^n + 1$ is the $2n$ th cyclotomic polynomial. In this case $x^n + 1$

stated in [31]. Therefore we can define a mapping $NTT(\cdot) : \mathbb{Z}_q[x]/(x^n + 1) \rightarrow \mathbb{Z}_q^n$, which turns a polynomial into a integer vector. By definition of isomorphism, an inverse mapping is $INTT(\cdot) : \mathbb{Z}_q^n \rightarrow \mathbb{Z}_q[x]/(x^n + 1)$, which turns a vector back to the polynomial. Based on the definition of isomorphism, we have:

- $NTT(a + b) = NTT(a) + NTT(b)$.
- $NTT(a \times b) = NTT(a) \times NTT(b)$.

Here the multiplication of vectors is position-wise multiplication. The vector $NTT(a)$ is called RNS form of polynomial a , and we use \bar{a} to represent it in the following. Now we are ready to propose NTT multiplication. The time cost is $2O(n \log n) + O(n) + O(n \log n) = O(n \log n)$, which is better than the naive $O(n^2)$.

III. OPTIMIZED FHE SCHEME

In FHE-based privacy-enhanced neural network inferences, the user encrypts the data before sends it to the server. The server homomorphically evaluates the neural network over encrypted data, and produces an encrypted inference output. It then returns the result’s ciphertext to the user. Evaluation of neural networks on encrypted data involves both arithmetic operations such as weighted sum and convolutions, evaluation of non-linear activations and bootstrapping after each neuron to enable further computations. However, existing FHE schemes are not designed for evaluation neural networks on encrypted data. Most of them focus on high-precision evaluations and are very slow in inferring of neural network. Some of them only support specific non-linear activations (e.g., [11], [13]). Our optimized FHE scheme which combines the advantages of CKKS and TFHE schemes will help to reduce the inference latency, and achieve the state-of-the-art inference accuracy.

A. Background: Neural Network

Neural network is a artificial model inspired by biological brains. A neural network can be considered as population of artificial neurons arranged in layers. The raw data is encoded properly and fed into the input layer and the output layer will output the inference result. Each internal layer (i.e., hidden layer) receives the data generated by its previous layer and outputs the processed data for the next layer.

Neural networks are usually composed of layer of following types: Fully-connected layer (FC layer, every neuron of the layer takes all incoming data as inputs), convolutional layer (convolution evaluation is applied to its input), and so on.

At each neuron in layer l , the processing happens in two steps:

is irreducible over \mathbb{Z} , and thus $\mathbb{Z}[x]/(x^n + 1)$ is an integral domain. For q , we choose q such that any polynomial in $\mathbb{Z}_q[x]/(x^n + 1)$ can be embedded to a vector according to the Chinese remainder theorem (CRT). After that, polynomial multiplication can be accelerated from $O(n^2)$ to $O(n \log(n))$. This requires that $x^n + 1$ splits modulo q , i.e.,

$$x^n + 1 = \prod_{i=1}^{2n-1} (x - z^i) \pmod{q},$$

which is implied by $q \equiv 1 \pmod{2n}$.

- Linear function: It is a weighted sum of inputs, which can be described as $\gamma_l = W_l x_{l-1} + \beta_l$, where W_l is the weight matrix of layer l and β_l is the bias of layer l .
- Activation function: The calculated weighted sum is passed to the activation function. An activation function is a mathematical function which adds non-linearity to the network. There are some commonly used activation functions: Sigmoid, ReLu and Softmax. It can be described as $x_l = f_l(\gamma_l)$.

The neural network with one hidden layer can approximate any continuous function [38], [39]. Furthermore, a deep neural network with several layers of non-linearities has better ability in more complex tasks. In this work, our optimized FHE scheme is able to evaluate neural networks of arbitrary depth with possibly many hidden layers.

B. Overview

Our system is built on both a LWE-based secret key encryption scheme and a RLWE-based secret key encryption scheme. The entire part of our system run on integer, which provides us a possibility to apply faster algorithms in implementation (such as faster polynomial multiplication algorithm based on NTT).

Compared with LWE, which operates in \mathbb{Z}_q , RLWE operates in polynomial rings. (e.g., $R_q := \mathbb{Z}_q[x]/(x^N + 1)$). Therefore, for simpler homomorphic evaluations, especially linear functions, LWE can have short and efficient implementations, while RLWE has increased computational complexity. However, for more complex functions, RLWE is powerful and is widely used in homomorphic encryption schemes. In general both LWE and RLWE can be used simultaneously in one scheme. For example, the famous homomorphic look-up table (LUT) algorithm which we will elaborate later.

The LWE-based secret key encryption scheme is used to encrypt the input test. Assume the length of input vector is l_{in} , we will generate l_{in} LWE ciphertexts. In each ciphertext, according to the noise-tolerance of neural network, we can add small noise to the main message. As a result, we can choose relatively small ciphertext modulus and dimension. In fact, it is a widely-used way to include small noise and approximate processions when applying homomorphic encryption scheme to neural network. For example, CKKS [20], which is one of the most competitive and efficient homomorphic encryption scheme, has added small noise in the main message. It is also noteworthy that we are not “assuming the sufficient noise tolerance of neural network”. Instead, we *observe* that the existence of small noise won't deviate the neural network result too much. Our experiments in Section VI show that inference accuracy on encrypted data is only 0.8% less than the inference accuracy on plain data, on MNIST dataset.

At each neuron, our scheme performs:

Homomorphic evaluation for linear functions on LWE ciphertext. Since our LWE secret key encryption scheme is a type of additive homomorphic, we can simply evaluate the linear functions by homomorphic scale multiplication and addition.

Homomorphic evaluation for non-linear activations. We use the homomorphic look-up table algorithm (LUT) to evaluate

non-linear activations and perform the bootstrapping at the same time. The core of LUT is to encode all possible output of the non-linear function $g(\cdot)$ into a polynomial $f(X)$ (denote by LUT function), so that we can “blindly” rotate the polynomial to locate the position of the desired output. The rotates are driven by a set of evaluation keys, which is an encryption of LWE secret key under an RLWE secret key and stored as an RGSW ciphertext. The evaluation keys are generated in the key generation phase in users' side, then they are stored in the server and used to drive the activation evaluations.² At the same time, the noise is reduced since we always “refresh” the ciphertext.

In order to achieve better efficiency, we hope the degree of polynomials in the LUT algorithm is small. So we make use of the noise tolerance of neural network again. Before entering LUT process, we squeeze the output range of the inner-product first. We only store a few number of most significant bits and discard some inaccurate least significant bits. Notice that the coefficients of the LUT function $f(X)$ store all possible value of the algorithm's output, so the number of coefficients (i.e., the degree of $f(X)$) is decided by the input range of the LUT algorithm (i.e., the output range of the inner-product). We will be able to achieve small degree of LUT function by squeezing the inner-product into a small range.

Key switching and homomorphic rounding algorithms. The LUT algorithm “re-encrypts” the input ciphertext, so that the underlying secret key and ciphertext modulus are changed. To ensure the network is extendable and can be applied to deep neural network, the output ciphertext should be in same size as the input ciphertext. Hence, we need key switching and homomorphic rounding algorithms to resize the output ciphertext.

In this way, the output of each neuron is in the same form as the input to the neuron, and the noise is reduced during bootstrapping. It ensures that the output of one neuron can be used for further computations in the next layer without an initially fixed limit on the number of layers one network has. Hence, our scheme can evaluate arbitrarily deep neural networks.

C. LWE Encryption and Linear Evaluation

We show our LWE-based secret key encryption scheme and related computations we used in our system. We write $ct \in \text{LWE}_s^{n,q}(m)$ to state that ct is a LWE ciphertext of m . Here q is the modulus, n is the dimension, and s is the secret key. And we use the same way for other encryption schemes. The detailed algorithms are Algorithm 2.

We introduce two basic linear evaluation algorithms under the LWE encryption scheme. Homomorphic addition 3 is to evaluate an addition between two ciphertexts. It takes two LWE ciphertexts (encryption of m_1 and m_2) as input and outputs a new LWE ciphertext whose decryption result is $m_1 + m_2$. Homomorphic scale multiplication 4 is to evaluate a multiplication between one plaintext number and one ciphertext, which takes an encryption of m and a plaintext number c as input and outputs a new LWE ciphertext whose decryption result is $m \cdot c$.

²It is secure for users to send evaluation keys to the server, since the evaluation keys are RGSW ciphertexts. Due to the security of cryptographic protocols, the server can learn nothing about the LWE secret key from the evaluation keys.

Algorithm 2: LWE Encryption Scheme.

Input: plaintext m , modulus q , dimension n , secret key $s \xleftarrow{\$} \{-1, 0, 1\}^n$, error distribution \mathcal{E}_{LWE} on \mathbb{Z}^q .
Output: LWE ciphertext $(a, b) \in \text{LWE}_s^{n,q}(m)$.
 Sample an error value $e \xleftarrow{\$} \mathcal{E}_{\text{LWE}}$.
 Sample vector $a \xleftarrow{\$} \mathbb{Z}_n^q$.
 Define $b := m + e - \langle a, s \rangle \pmod q$.
 Return (a, b) .

Algorithm 3: Homomorphic Addition.

Input: LWE ciphertexts $(a_1, b_1) \in \text{LWE}_s^{n,q}(m_1)$ and $(a_2, b_2) \in \text{LWE}_s^{n,q}(m_2)$.
Output: LWE ciphertext $(a', b') \in \text{LWE}_s^{n,q}(m_1 + m_2)$.
 Compute $a' := a_1 + a_2 \pmod q$ and $b' := b_1 + b_2 \pmod q$.
 Return (a', b') .

Algorithm 4: Homomorphic Scale Multiplication.

Input: LWE ciphertext $(a, b) \in \text{LWE}_s^{n,q}(m)$, plaintext $c \in \mathbb{Z}_n^q$.
Output: LWE ciphertext $(a', b') \in \text{LWE}_s^{n,q}(c \times m)$.
 Compute $a' := c \cdot a \pmod q$ and $b' := c \cdot b \pmod q$.
 Return (a', b') .

Algorithm 5: Homomorphic Inner Product Evaluation.

Input: LWE ciphertext ct of vector m : $\text{ct} \in \text{LWE}_s^{n,q}(m)$, and plaintext vector c .
Output: LWE ciphertext $(a', b') \in \text{LWE}_s^{n,q}(\langle c, m \rangle)$.
 Let $L = \dim(c) = \dim(m)$, $a' = \mathbf{0}$, $b' = 0$.
for $i = 1, \dots, L$ **do**
 Parse $(a_i, b_i) \in \text{LWE}_s^{n,q}(m_i)$ from ct.
 Compute $a' = a' + c_i a_i$ and $b' = b' + c_i b_i$.
end for
 Compute $a' = a' \pmod q$ and $b' = b' \pmod q$.
 Return (a', b') .

We extend 2 to encrypt a L dimension vector m

$$\text{LWE}_s^{n,q}(m) := (\text{LWE}_s^{n,q}(m_1), \text{LWE}_s^{n,q}(m_2), \dots, \text{LWE}_s^{n,q}(m_L)).$$

Then homomorphic addition and scale multiplication can be extended in a similar way. Using these algorithms, we are able to homomorphically evaluate the inner product of a plaintext vector and ciphertext vector (Algorithm 5). The homomorphic inner product evaluation algorithm takes an encryption of vector m and a plaintext vector c as input and outputs an encryption of $\langle m, c \rangle$.

D. RLWE Encryption and Non-Linear Evaluations

We show the RLWE-based secret key encryption scheme and related operations we used in our system. Our scheme is defined on ring $R_q := \mathbb{Z}_q[X]/(X^n + 1)$.

Algorithm 6: RLWE Encryption Scheme.

Input: plaintext m , modulus q , dimension n , secret key s whose coefficients are sampled uniformly random from $\{-1, 0, 1\}$, error distribution $\mathcal{E}_{\text{RLWE}}$.
Output: RLWE ciphertext $(a, b) \in \text{RLWE}_s^{n,q}(m)$.
 Sample error element $e \xleftarrow{\$} \mathcal{E}_{\text{RLWE}}$.
 Sample element $a \xleftarrow{\$} R_q$.
 Define $b := m + e - a \times s \in R_q$.
 Return (a, b) .

Algorithm 7: Extended RLWE Encryption.

Input: plaintext m , modulus q , dimension n , secret key $s \in \{-1, 0, 1\}^n$, error distribution $\mathcal{E}_{\text{RLWE}}$, decomposition base B .
Output: Extended RLWE ciphertext $\{(a_i, b_i)\}_{i \in [\log(q)/\log(B)]} \in \widetilde{\text{RLWE}}_s^{n,q}(m)$.
for $i = 0, 1, \dots, \log(q)/\log(B) - 1$ **do**
 Find $(a_i, b_i) \in \text{RLWE}_s^{n,q}(m \times B^i)$.
end for
 Return $\{(a_i, b_i)\}_{i \in [\log(q)/\log(B)]}$.

Algorithm 8: RGSW Encryption Based on Extended RLWE Encryption.

Input: plaintext m , modulus q , dimension n , secret key $s \in \{-1, 0, 1\}^n$.
Output: RGSW ciphertext $(\beta, \alpha) \in \left(\widetilde{\text{RLWE}}_s^{n,q}(m), \widetilde{\text{RLWE}}_s^{n,q}(s \times m) \right)$.

1) *RLWE Encryption and Related Homomorphic Operations:* The encryption scheme is Algorithm 6. And we can define scale multiplication and addition for RLWE ciphertext, in the same way as Algorithms 4 and 3.

Building blocks of evaluation key generation. The Algorithms 7 and 8 are sub-algorithms for generating evaluation keys. Since the algorithm of generating evaluation keys is related to the LUT, for convenience, we present the algorithms of generating evaluation keys together with the LUT algorithm (See Algorithms 11 and 12).

Building blocks of LUT. Besides the generation of evaluation key, we also introduce the external multiplication \odot (Algorithm 10) and its sub-operation \diamond (Algorithm 9). The external multiplication is an important component of LUT, which shows how to do multiplication between RLWE ciphertext and RGSW ciphertext.

Extract0 (Detailed algorithm is shown in Appendix, available online) is to extract the constant term of polynomial m from its RLWE ciphertext $\text{RLWE}_s^{n,q}(m)$. It will output a LWE ciphertext which encrypts the 0th coefficient of m . This extraction algorithm allows us to be able to convert a RLWE ciphertext to into LWE ciphertexts and will be in LUT algorithm.

2) *Non-Linear Evaluations - Homomorphic Look-Up Table Algorithm (LUT):* The look-up table algorithm takes a LWE

Algorithm 9: Extended RLWE Ciphertext Multiplication (\diamond).

Input: plaintext operand r , modulus q , dimension n , decomposition base B , extended RLWE ciphertext $\{(a_i, b_i)\}_{i \in [\log(q)/\log(B)]} \in \widetilde{\text{RLWE}}_s^{n,q}(m)$.
Output: RLWE ciphertext $(a, b) \in \text{RLWE}_s^{n,q}(r \times m)$.
Decompose r s.t. $r = \sum_{i=0}^{\log(q)/\log(B)-1} r_i \times B^i$.
Compute $(a, b) = \sum_{i=0}^{\log(q)/\log(B)-1} r_i \times (a_i, b_i) = \sum_{i=0}^{\log(q)/\log(B)-1} (r_i a_i, r_i b_i)$.
Ensure (a, b) in R_q^2 .
Return (a, b) .

Algorithm 10: RLWE and RGSW Multiplication (\odot).

Input: $(a, b) \in \text{RLWE}_s^{n,q}(m_1)$, $(\beta, \alpha) \in \text{RGSW}_s^{n,q}(m_2)$.
Output: $(a', b') \in \text{RLWE}_s^{n,q}(m_1 \times m_2)$.
Return $(a', b') = a \odot \alpha + b \odot \beta$ and ensure $(a', b') \in R_q^2$.

ciphertext $ct \in \text{LWE}_s^{n,q}(m)$ and an evaluation function $F(\cdot)$ as input and it outputs a LWE ciphertext which encrypts $\Delta F(m)$, where Δ is a scale factor to limit noise.

Bit-by-bit Look-up table algorithm for general cases. We first show a bit-by-bit look-up table algorithm (Algorithm 11), which can be used in any cases, however not optimized.

2-bit Look-up table algorithm for single hidden layer. When the input LWE ciphertext has a secret key $s \in \{0, 1\}^n$ instead of $s \in \{-1, 0, 1\}^n$, we can use the optimized Look-up table algorithm to reduce the number of external operations from n to $n/2$ compared with Algorithm 11 while retaining the security [40], [41]. The details are in Algorithm 12. In the case only one activation layer is evaluated homomorphically, we can skip the key switch and rounding operations after the look-up table evaluation, and proceed to other linear evaluations directly. However, this method cannot be applied if more than one activation layer. This is because the output of any Look-up table algorithm must be an LWE ciphertext with secret key $s \in \{-1, 0, 1\}^n$.

2-bit Look-up table algorithm for multiple hidden layers. For multiple hidden layers, the key switching algorithm is required between each hidden layer, which consists of computations on RLWE ciphertexts. However, the binary secret for the Ring variant of LWE is still an open problem [22]. So, we cannot convert the LWE secret key from $s \in \{-1, 0, 1\}^n$ to $s \in \{0, 1\}^n$ here, but we can still construct a 2-bit look-up table algorithm. We will include the details in Appendix, available online.

E. Other Building Blocks of Our Scheme

Our framework also contains the following functions: homomorphic rounding and key switching algorithms. We will briefly introduce them in this section, and the detailed algorithms are in Appendix, available online.

Rounding. The homomorphic rounding algorithm will change a LWE ciphertext with ciphertext modulus q' to a LWE ciphertext with ciphertext modulus q .

Key switching. The key switching algorithm will change a LWE ciphertext with secret key s' and dimension n' to a LWE ciphertext with new secret key s and dimension n .

F. Security and Privacy Analysis

In our FHE-based solution, data privacy protection is to ensure that the cloud server can know neither the raw input data nor the inference results in clear. Since the cloud server only plays with encrypted data in the whole inference process, the data privacy protection relies on the security of underlying FHE schemes. In Theorem 3.1, we show that the FHE scheme achieves CPA security as in Definition 2.2.

Theorem III.1 (CPA security). The FHE scheme proposed in Section III is CPA secure under the LWE and RLWE assumptions.

Proof (sketch of proof). The view of a CPA adversary for our scheme is very similar to Regev's scheme [30], which is proved secure under LWE assumption, with the exception that our adversary also gets to see the evaluation key. However, the evaluation key contains a sequence of RLWE instances which are in the same form as CKKS ciphertexts [20]. The evaluation keys are indistinguishable from uniform according to the security of CKKS scheme. Therefore, the security of our scheme can be derived from the security of Regev's scheme. Please refer [30] and [20] for the detailed proof.

G. Our Framework

We show the full protocol for privacy-enhanced neural network using our optimized FHE scheme in Algorithm 13.³

H. Noise Analysis of Optimized FHE Scheme

In this section, we analyze the growing of noise throughout our system. Let σ_{LWE}^2 be the variance of noise used in the LWE encryption. Define $\sigma_{\text{RGSW}}^2, \sigma_{\text{KS}}^2$ in the same way.

By the widely used assumptions, we assume that in each polynomial all the coefficients behave like independent zero-mean random variables of the same variance [32] (weaker than i.i.d.), and central limit heuristic [23]. The procedures are similar as in [27]. Further, note that it suffices to find the change of error variance within one neuron of each layer. Fix Layer l , assume that $\text{LWE}_{l-1} := \{(a_i, b_i)\}_{i \in [H_{l-1}]}$ has an error whose variance is σ^2 in each LWE ciphertext (a_i, b_i) .

a) *Linear function.* $ip_h = \sum_{i \in [H_{l-1}]} W_l[h, i] \times (a_i, b_i) + \beta_l[h] \in \text{LWE}_s^{n,q}(\cdot)$ and thus the noise becomes e_{ip} with variance $\sigma_{ip}^2 := \sum_{i \in [H_{l-1}]} W_l^2[h, i] \sigma^2 \leq \|W_l\|^2 \sigma^2$. Here we define

$$\|W_l\|^2 := \max_{h \in [H_l]} \left\{ \sum_{i \in [H_{l-1}]} W_l^2[h, i] \right\}.$$

b) *LUT.* As in [27], the LUT evaluations outputs a LWE ciphertext that decrypts to $\Delta F(m + e_1) + e_2$ for input m . And

³The LWE secret key s and RLWE secret key s' are stored in the user's side and kept secret from the server. The encrypted key (i.e., evaluation keys and key switching keys) are generated by the user and stored in the server. Since the keys are encrypted, it is secure to pass them to the server.

Algorithm 11: Bit-by-Bit Look-Up Table Evaluation.

1: **Input:** LWE ciphertext $(\mathbf{a}, b) \in \text{LWE}_{\mathbf{s}}^{n,q}(m)$ s.t. $|m| < q/4$, scale factor Δ , evaluation function $F(\cdot) : \mathbb{Z}_q \rightarrow \mathbb{Z}_q$, RLWE parameter set (n', q') , RLWE secret key s' , a set of evaluation keys w.r.t. the LWE secret key $\mathbf{s} \in \{-1, 0, 1\}^n$

$$\text{EK}_{j,+} = \text{RGSW}_{s',q'}^{n',q'}(\mathbf{1}_{s[j]>0}), \text{EK}_{j,-} = \text{RGSW}_{s',q'}^{n',q'}(\mathbf{1}_{s[j]<0}), j = 0, 1, \dots, n-1.$$

2: **Output:** LWE ciphertext $ct' \in \text{LWE}_{s'}^{n',q'}(\Delta F(m))$ where s' is the trivial vector form of polynomial s' (from high degree to low degree).

3: Let $\eta_k = kq/(2n')$ for $1 \leq k \leq n'/2$. Define a polynomial $f \in R_{n',q'}$ whose coefficients are

$$f_j = \begin{cases} \lceil \Delta F(0) \rceil & \text{if } j = 0 \\ \lceil \Delta F(-\eta_j) \rceil & \text{if } 1 \leq j \leq n'/2 \\ \lceil -\Delta F(\eta_{n'-j}) \rceil & \text{if } n'/2 < j < n' \end{cases}.$$

4: Let $b' = \lceil 2n'b/q \rceil$, let $\mathbf{a}' = \lceil 2n'\mathbf{a}/q \rceil$.

5: Initialize $\text{AC} = (0, f \times X^{b'}) \in R_{n',q'}^2$.

6: **for** $j=0,1,\dots,n-1$ **do**

7: $\text{AC} += \left((X^{\mathbf{a}'[j]} - 1)\text{EK}_{j,+} + (X^{-\mathbf{a}'[j]} - 1)\text{EK}_{j,-} \right) \odot \text{AC}$, note that all calculations are in $R_{n',q'}^2$.

8: **end for**

9: Return $\text{Extract0}(\text{AC})$.

Algorithm 12: 2-Bit Look-Up Table Evaluation for Single Hidden Layer Neural Network.

1: **Input:** LWE ciphertext $(\mathbf{a}, b) \in \text{LWE}_{\mathbf{s}}^{n,q}(m)$ s.t. $|m| < q/4$, scale factor Δ , evaluation function $F(\cdot) : \mathbb{Z}_q \rightarrow \mathbb{Z}_q$, RLWE parameter set (n', q') , RLWE secret key s' , a set of evaluation keys w.r.t. the LWE secret key $\mathbf{s} \in \{0, 1\}^n$, and for $j = 0, 1, \dots, n/2 - 1$:

$$\text{EK}_{j,0} = \text{RGSW}_{s',q'}^{n',q'}(\mathbf{s}[2j]\mathbf{s}[2j+1]), \text{EK}_{j,1} = \text{RGSW}_{s',q'}^{n',q'}(\mathbf{s}[2j](1-\mathbf{s}[2j+1])), \text{EK}_{j,2} = \text{RGSW}_{s',q'}^{n',q'}(\mathbf{s}[2j+1](1-\mathbf{s}[2j])).$$

2: **Output:** LWE ciphertext $ct' \in \text{LWE}_{s'}^{n',q'}(\Delta F(m))$ where s' is the vector form of polynomial s' (from high degree coefficient to low degree coefficient).

3: Let $\eta_k = kq/(2n')$ for $1 \leq k \leq n'/2$. Define a polynomial $f \in R_{n',q'}$ whose coefficients are

$$f_j = \begin{cases} \lceil \Delta F(0) \rceil & \text{if } j = 0 \\ \lceil \Delta F(-\eta_j) \rceil & \text{if } 1 \leq j \leq n'/2 \\ \lceil -\Delta F(\eta_{n'-j}) \rceil & \text{if } n'/2 < j < n' \end{cases}.$$

4: Let $b' = \lceil 2n'b/q \rceil$, let $\mathbf{a}' = \lceil 2n'\mathbf{a}/q \rceil$.

5: Initialize $\text{AC} = (0, f \times X^{b'}) \in R_{n',q'}^2$.

6: **for** $j=0,1,\dots,n/2-1$ **do**

7: $\text{AC} += \left((X^{\mathbf{a}'[2j]+\mathbf{a}'[2j+1]} - 1)\text{EK}_{j,0} + (X^{\mathbf{a}'[2j]} - 1)\text{EK}_{j,1} + (X^{\mathbf{a}'[2j+1]} - 1)\text{EK}_{j,2} \right) \odot \text{AC}$, all calculations are in $R_{n',q'}^2$.

8: **end for**

9: Return $\text{Extract0}(\text{AC})$.

the LUT input is the above ciphertext ip_h of linear function. We first find e_1 , which is the error of look-up index. Our analysis is based on the Algorithm 11, and for other proposed algorithms the analysis is almost identical.

- By our LUT algorithm (Algorithm 11), (\mathbf{a}', b') has decryption result $b' + \langle \mathbf{a}', \mathbf{s} \rangle = \lceil 2n'm/q \rceil + e$. By the central limit heuristic, e has variance $4n'^2\sigma_{ip}^2/q^2 + (\|\mathbf{s}\|_2^2 + 1)/12$, where the second term is from the rounding $\lceil \cdot \rceil$. Note that ciphertext is generated from uniformly random distribution, thus the loss of $\lceil \cdot \rceil$ is from $U[-1/2, 1/2]$, i.e., the uniform distribution on $[-1/2, 1/2]$.
- Next, note the definition of polynomial f 's coefficients, the above error e is scaled up by a factor

of $q/(2n')$, which results in an error with variance $\sigma_{ip}^2 + q^2(\|\mathbf{s}\|_2^2 + 1)/(48n'^2)$. Also note that the error of $\lceil 2n'm/q \rceil$ is scaled up by a factor of $q/(2n')$, this error has variance $(1/12) \times q^2/(4n'^2)$ if assuming the loss of $\lceil \cdot \rceil$ is from $U[-1/2, 1/2]$. Summing up both parts we have $\text{var}(e_1) = \sigma_{ip}^2 + q^2(\|\mathbf{s}\|_2^2 + 2)/(48n'^2) = \|\mathbf{W}_l\|^2\sigma^2 + q^2(\|\mathbf{s}\|_2^2 + 2)/(48n'^2)$.

Next, we proceed to find e_2 . In each step j we calculate the external product

$$\text{AC} += \left((X^{\mathbf{a}'[j]} - 1)\text{EK}_{j,+} + (X^{-\mathbf{a}'[j]} - 1)\text{EK}_{j,-} \right) \odot \text{AC}.$$

Algorithm 13: Full Protocol for Privacy-Enhanced Neural Network.

- 1: **Input Layer parameters:** LWE parameter set (n, q, \mathbf{s}) .
- 2: **Hidden Output Layer parameters:** RLWE and RGSW related parameter set (n', q', s') where $n|n'$, scale factor Δ .
- 3: **Public parameters:** Decomposition base B , Key-Switching decomposition base B_{KS} .
- 4: Given input vector \mathbf{x} of length N , define a set of N LWE ciphertexts $\{(\mathbf{a}_i, b_i)\}$ each of which is from $\text{LWE}_{\mathbf{s}}^{n,q}(x[i])$.
- 5: Generate a set of evaluation keys w.r.t. the LWE secret key $\mathbf{s} \in \{-1, 0, 1\}^n$:
- 6: $\text{EK}_{j,+} = \text{RGSW}_{s',q'}^{n',q'}(\mathbf{1}_{\mathbf{s}[j]>0}), \text{EK}_{j,-} = \text{RGSW}_{s',q'}^{n',q'}(\mathbf{1}_{\mathbf{s}[j]<0}), j = 0, 1, \dots, n-1$.
- 7: Generate a set of LWE switching keys w.r.t. decomposition base B :
- 8: $\text{SK}_j \in \widetilde{\text{RLWE}}_{\mathbf{s}}^{n,q} \left(\sum_{l=0}^{n-1} \mathbf{s}[jn+l] X^l \right), j = 0, 1, \dots, n'/n-1$.
- 9: Let input layer be Layer 0 and let output layer be Layer L .
- 10: **for** $l = 1, 2, \dots, L$ **do**
- 11: Let H_l be the number of neurons in layer l .
- 12: Let W_l, β_l be the weight matrix and the bias vector from layer $l-1$ to layer l , respectively.
- 13: **for** $h = 0, 1, 2, \dots, H_l-1$ **do**
- 14: Let $\text{LWE}_{l-1} := \{(\mathbf{a}_{l-1,i}, b_{l-1,i})\}_{i \in [H_{l-1}]}$, each of which is a LWE ciphertext $\text{LWE}_{\mathbf{s}}^{n,q}(\cdot)$ from Layer $l-1$. For the sake of simplicity we omit the subscript of layer $l-1$ and write (\mathbf{a}_i, b_i) .
- 15: Homomorphically evaluate the linear function in the h th neuron in Layer l :
 $ip_h = \sum_{i \in [H_{l-1}]} W_l[h, i] \times (\mathbf{a}_i, b_i) + \beta_l[h]$, i.e., $ip_h \in \text{LWE}_{\mathbf{s}}^{n,q}(\gamma_l[h])$.
- 16: Evaluate the look-up-table in the h th neuron: $ct'_h = \text{LUT}(ip_h) \in \text{LWE}_{\mathbf{s}}^{n',q'}(\Delta \times \cdot)$.
- 17: Switch from n' to smaller n : $\tilde{ct}_h = \text{LWE_KS}(ct'_h, B_{\text{KS}}) \in \text{LWE}_{\mathbf{s}}^{n,q'}(\Delta \times \cdot)$.
- 18: Round from q' to smaller q and rescale: $ct_h = \text{LWE_Rounding}(\tilde{ct}_h) \in \text{LWE}_{\mathbf{s}}^{n,q}(\cdot)$.
- 19: Include ct_h in LWE_l .
- 20: **end for**
- 21: **end for**
- 22: Return LWE_L to User. User can decrypt all ciphertexts in it with secret key \mathbf{s} , and then find the inference result.

As in [32] and [23], under their assumptions we have: Given polynomials a, b , whose variances of the coefficients are σ_a^2 and σ_b^2 respectively, then:

- the variance of coefficients of polynomial $a + b$ is $\sigma_a^2 + \sigma_b^2$;
- the variance of coefficients of polynomial ab is $n'\sigma_a^2\sigma_b^2$.

Based on the above we are able to find the variance of polynomial calculations. For the sake of simplicity, when we say the variance of a polynomial we mean the variance of the coefficients of this polynomial.

- Let $\text{RGSW}_j := (X^{a[j]} - 1)\text{EK}_{j,+} + (X^{-a[j]} - 1)\text{EK}_{j,-}$. First we find the variances of the coefficients of the polynomial in each slot of RGSW_j . By our definition of algorithm, $\text{EK}_{j,+}$ and $\text{EK}_{j,-}$ have the same error variance: σ_{RGSW}^2 in each slot. Then note that $(X^{a[j]} - 1)\text{EK}_{j,+}$ is a sum of $X^{a[j]}\text{EK}_{j,+}$ and $-\text{EK}_{j,+}$. They both have variances σ_{RGSW}^2 . The other part $(X^{-a[j]} - 1)\text{EK}_{j,-}$ is the same, so the variance in each slot of $(X^{a[j]} - 1)\text{EK}_{j,+} + (X^{-a[j]} - 1)\text{EK}_{j,-}$ is $4\sigma_{\text{RGSW}}^2$.
- Suppose in step j , $\text{AC} = (a_j, b_j)$ where both a_j and b_j are polynomials from the same ring. Given decomposition base B , they have decomposition

$$a_j \rightarrow (\hat{a}_0, \hat{a}_1, \dots, \hat{a}_{d-1}), b_j \rightarrow (\hat{b}_0, \hat{b}_1, \dots, \hat{b}_{d-1}),$$

where $a_j = \sum_{i=0}^{d-1} B^i \hat{a}_i, b_j = \sum_{i=0}^{d-1} B^i \hat{b}_i$. We also write

$$\begin{aligned} \text{RGSW}_j &= (\beta, \alpha) \\ &= \left((\beta[0], \dots, \beta[d-1]), (\alpha[0], \dots, \alpha[d-1]) \right). \end{aligned}$$

By definition of our RGSW each $\beta[j]$ or $\alpha[j]$ is a pair of polynomials. Then

$$\text{RGSW}_j \odot \text{AC} = a_j \diamond \alpha + b_j \diamond \beta = \sum_{i=0}^{d-1} \left(\hat{a}_i \alpha[i] + \hat{b}_i \beta[i] \right).$$

Since all \hat{a}_j and \hat{b}_j have B -bounded coefficients, the variance of each slot in $\text{RGSW} \odot \text{AC}$ is bounded by

$$d \times 2 \times n' B^2 \times (4\sigma_{\text{RGSW}}^2) = 8n' dB^2 \sigma_{\text{RGSW}}^2.$$

- Finally, $\text{RGSW}_j = (X^{a[j]} - 1)\text{EK}_{j,+} + (X^{-a[j]} - 1)\text{EK}_{j,-}$ in each step j has the same distribution, and thus has the same error variance. The total variance in n steps is $\text{var}(e_2) = 8nn' dB^2 \sigma_{\text{RGSW}}^2$.

Key Switch. By our definitions of Key-Switching key $\{\text{SK}_j\}$ and $\widetilde{\text{RLWE}}$, each SK_j has d RLWE ciphertexts, each of which has error of variance σ_{KS}^2 . The error e_{KS} is from $\sum_{j=0}^{n'/n-1} \tilde{a}_j \diamond \text{SK}_j$. Similarly as the \diamond operation in the LUT part, note that $\text{deg}(\tilde{a}_j) = n$, the variance of e_{KS} is bounded by

$$n'/n \times n B_{\text{KS}}^2 \sigma_{\text{KS}}^2 = n' B_{\text{KS}}^2 \sigma_{\text{KS}}^2.$$

c) Rounding. Suppose the rounding factor is z , simply round down the variance by a factor of z^2 . To round the ciphertext to integer, $\lceil \cdot \rceil$ results in an additional variance of $\text{var}(e_{\text{RD}}) = (\|\mathbf{s}\|_2^2 + 1)/12$.

Like in [27] and many other works, we assume $F(\cdot)$ is L -lipschitz and then we will have $|F(m + e_1) - F(m)| \leq L|e_1|$. Next, $|e_1|, |e_2|$ and $|e_{\text{KS}}|$ can be bounded w.h.p. by

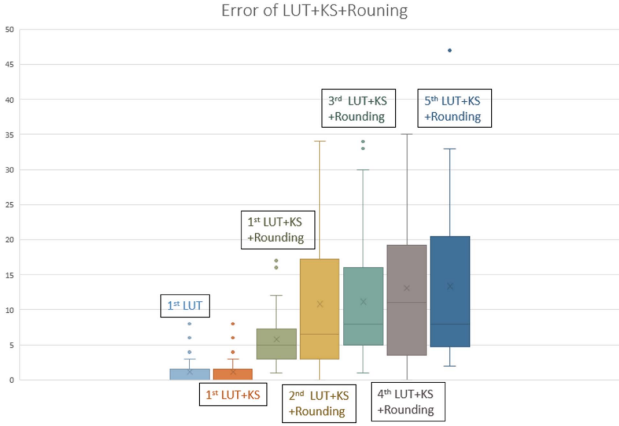


Fig. 3. Test result on noise growing in LUT+KS+Rounding.

$O(\sqrt{\text{var}(e_i)}), i = 1, 2, \text{KS}$ under central limit heuristic [23]. To sum up, after scaling down Δ , the error between $F(m)$ and $F(m + e_1) + e_2/\Delta + e_{\text{KS}}/\Delta + e_{\text{RD}}$ is bounded by

$$\begin{aligned} & O\left(L\sqrt{\text{var}(e_1)} + \frac{\sqrt{\text{var}(e_2)} + \sqrt{\text{var}(e_{\text{KS}})}}{\Delta} + \sqrt{\text{var}(e_{\text{RD}})}\right) \\ &= O\left(L\sqrt{\|W_l\|^2\sigma^2 + q^2(\|s\|_2^2 + 2)/(48n^2)}\right. \\ & \quad \left. + \frac{\sqrt{8nn'dB^2\sigma_{\text{RGSW}}^2} + \sqrt{n'B_{\text{KS}}^2\sigma_{\text{KS}}^2}}{\Delta} + \sqrt{(\|s\|_2^2 + 1)/12}\right). \end{aligned}$$

Also, we give the error variance of the LWE ciphertext that output to the next layer

$$\begin{aligned} \sigma_{i+1}^2 &\leq \left(L^2\text{var}(e_1) + \frac{\text{var}(e_2) + \text{var}(e_{\text{KS}})}{\Delta^2} + \text{var}(e_{\text{RD}})\right) \\ &\leq \left(L^2(\|W_l\|^2\sigma_l^2 + q^2(\|s\|_2^2 + 2)/(48n^2))\right. \\ & \quad \left.+ \frac{8nn'dB^2\sigma_{\text{RGSW}}^2 + n'B_{\text{KS}}^2\sigma_{\text{KS}}^2}{\Delta^2} + (\|s\|_2^2 + 1)/12\right). \end{aligned}$$

1) *Experiments on Noise Growing in Multiple Layers:* After showing that our system performed well in single hidden layer neural network, we move our focus to multiple layers. The key to achieve good results in multiple layers is to ensure that the noise is always in a suitable range. In Section III-H, we showed a theoretical analysis on the growing of noise. Now, we show some experimental results about it. The noise is growing in two steps. The first is in the inner-product computation. In this part, the growing of noise is linear and easy to control by choosing suitable parameters. So we focus on the second step: Look-up table (LUT), key switching (KS) and rounding. We ask the program to run one inner-product computation and then to perform LUT+KS+Rounding 5 times continuously and get the following Fig. 3 (the evaluation function in LUT is ReLu). The picture shows the difference between the result of “*i*th LUT+KS+Rounding” and the real value. 0 means that the result in “*i*th LUT+KS+Rounding” is the same as the real value, which means no error exists. From the picture we can see that, most of

points fall in (0,20), which is relatively small compared to the inner-product range (−1000, 1000).

IV. EFFICIENT DESIGN FOR NON-LINEAR ACTIVATION EVALUATION

A. Overview

In this section, we further optimize our FHE scheme and focus on improving the evaluation of non-linear activations. Many previous work on FHE-based PE-NN use polynomial approximation activations to evaluate non-linear activations. However, it does not perform well. n-GraphHE [24], Lola [15] and Faster CryptoNet [12] have to use the low-degree polynomial approximation activations and thus fail to obtain the state-of-the-art inference accuracy. High-degree polynomial approximation activations can improve the accuracy, but the computing overhead increases exponentially with the degree, so the inference becomes inefficient.

We adopt homomorphic look-up-table (LUT) algorithm to evaluate non-linear activations. Although LUT algorithm is faster than polynomial approximation, it is the slowest part among all homomorphic evaluations in FHE-based PE-NN. To achieve better performance, we propose an efficient design for LUT-based non-linear activation evaluation.

We find that the NTT/INTT and the modulo calculations take around 70% time in one LUT. Therefore, we reduce the number of NTT and INTT in one LUT evaluation and optimize modulo calculations. As a result, our system becomes 3 times faster than before. Our improved system only takes 0.14s per input on MNIST dataset, while 0.42s is needed without improvements. We list our optimizations in below sections and an efficiency analysis to compare the difference in terms of the amount of time where we managed to bring down significantly.

B. Reduce the Number of NTT and INTT

In this section, we show our optimizations on polynomial multiplications, which takes around 45% time (computed from Table I) in one LUT. We introduce 2 methods to reduce the number of NTT/INTTs in LUT algorithm. Section IV-B1 is a general method and can be applied in any cases. Section IV-B2 is designed for the neural network where very large modulus is necessary, such as hundreds of bits. In practical scenario, we should choose suitable method according to the concrete problems and neural networks.

1) *General Method to Reduce Number of NTT/INTTs in LUT:* Taking our 2-bit look-up table evaluation for single hidden layer algorithm (Algorithm 12) as an example, we will show how to reduce the number of times where NTT/INTTs are called.

First, we count the number of NTT/INTTs in one LUT. We use d to denote that in the beginning of each external product \odot , each polynomial in AC is decomposed to d polynomials. Regarding the computations of $(X^{\alpha[2j] + \alpha'[2j+1]} - 1)EK_{j,0}$, $(X^{\alpha[2j]} - 1)EK_{j,1}$ and $(X^{\alpha'[2j+1]} - 1)EK_{j,2}$, we discover that instead of calling NTT multiplication $3d$ times, we can simply implement it by rotating the polynomial in $EK_{j,\cdot}$, which only

TABLE I
ANALYSIS ON LUT ALGORITHM

Operations	LUT in our previous work			Improved LUT			Time saved
	Number of calculation	Time per calculation	Time in total	Number of calculation	Time per calculation	Time in total	
NTT	4096	15.29us	62628us	1792	15.29us	27399 us	35229us
INTT	2048	16.05us	32870us	512	16.05us	8192us	24678us
PMUL	2048	2.84us	5816us	8193	2.84us	23268us	-17452us
Modulo Calculation (After additions)	10625	5.9us	62687us	4994	4.5us	22473us	40214us
Polynomial addition	8192	1.48us	12124us	6144	1.48us	9093us	3031us
Bit-decomposition	512	9.5us	4864us	512	9.5us	4864us	0
Quick multiplication	6145	8us	49160us	0	-	0	49160us
Total time	230149us=230ms			95289us=95ms			135ms (60%)

involves polynomial addition and subtraction and is faster than NTT multiplication. We call it ‘quick multiplication’ and it was used in the implementation of the non-optimized LUT algorithm. Therefore, the NTT multiplications only appear in \odot . One \odot includes $4d$ NTT polynomial multiplications. So one LUT calculation includes $4d \cdot (n/2) = 2dn$ NTT polynomial multiplications. Each NTT polynomial multiplication includes 3 NTT/INTT transformations, so one LUT has $6dn$ NTT/INTT transformations.

A general method is to store $EK_{j,0}$, $EK_{j,1}$ and $EK_{j,2}$ in RNS form when generating them. Such evaluation keys are generated in the initialization phase and can be used repeatedly. When generating the evaluation keys, we can do NTT transformations after RGSW encryption to convert them into RNS form and store them.

In LUT calculation, when calculating $(X^{\alpha'[2j]+\alpha'[2j+1]} - 1)EK_{j,0} + (X^{\alpha'[2j]} - 1)EK_{j,1} + (X^{\alpha'[2j+1]} - 1)EK_{j,2}$, we first perform NTT transformation on $X^{\alpha'[2j]+\alpha'[2j+1]} - 1$, $X^{\alpha'[2j]} - 1$ and $X^{\alpha'[2j+1]} - 1$. Then we proceed to compute position-wise multiplications and additions. The output of this part is $4d$ RNS form polynomials.

Next, we move to the calculation of \odot . The left side of \odot is already in RNS form. So we first decompose the right side from 2 polynomials to $2d$ polynomials and do NTT transformations $2d$ times. The other calculations in \odot can be done by position-wise multiplications and additions. Now the output of \odot is 2 RNS form polynomials. We apply 2 INTT transformations on them and get 2 regular polynomials.

Therefore, in each loop, we only need $3 + 2d$ NTT transformations and 2 INTT transformations. So one improved LUT only has $(3 + 2d + 2) \cdot (n/2) = dn + 2.5n$ NTT/INTT transformations.

For example, if we take $n = 512$ and $d = 2$ (same as our experiments), then the number of NTT/INTT transformations reduced from 6144 to 2304.

2) *Further Improvements on Neural Network With Large Modulus*: When it is necessary to use hundreds bits modulus Q , a well known technique is to set Q as a product of L distinct and machine-word-sized primes: $Q = \prod_{i=0}^{L-1} q_i$. Each q_i is also appropriate chosen so that NTT multiplication can be applied. The propose of this decomposition is, when modulus is very large, the multiplications and mod algorithm on computer/server become very slow. Using machine-word-sized is many times faster.

Algorithm 14: CRT Based LUT for Large Modulus.

- 1: First turn AC_0 , all RGSW ciphertexts into RNS form.
 - 2: **for** $i = 0, 1, \dots, n/2 - 1$ **do**
 - 3: Calculate \overline{EK}_i
 - 4: $\overline{AC}_{i+1} = \overline{EK}_i \odot \overline{AC}_i + \overline{AC}_i$
 - 5: **end for**
 - 6: **Output**: $\text{Extract0}(\text{INTT}(\overline{AC}_{n/2}))$, i.e., the LWE ciphertext of the constant term in the plaintext of $\text{INTT}(\overline{AC}_{n/2})$.
-

Algorithm 15: (\diamond) Operator in RNS Form, $\bar{\alpha} \diamond \bar{f}$.

- 1: **Input**: Extended RLWE ciphertext $\bar{\alpha}$ which includes L RLWE ciphertexts, i.e., $2L$ polynomials (in RNS form). According to CRT, find gadget vector $\vec{g} = (g_0, g_1, \dots, g_{L-1}) \in \mathbb{Z}^L$ based on $\{q_i\}$. Let $\bar{\alpha} = \{(\bar{a}_i, \bar{b}_i)\}_{i=0}^{L-1}$, where each $(a_i, b_i) \in \text{RLWE}_{s', Q}^{n', Q}(g_i \times \cdot)$. A RNS form polynomial $\bar{f} \in \prod_{i=0}^{L-1} \mathbb{Z}_{q_i}^{n'}$.
 - 2: Decompose \bar{f} into $\{\bar{f}_i = \bar{f} \bmod q_i, i = 0, 1, \dots, L - 1\}$.
 - 3: Calculate and output $\sum_{i=0}^{L-1} (\bar{f}_i \bar{a}_i, \bar{f}_i \bar{b}_i)$. (Position-wise multiplication)
-

By Chinese remainder theorem (CRT), $R_{n,Q}$ and $\prod_{i=0}^{L-1} R_{n,q_i}$ are isomorphic, which means that for each polynomial f in $R_{n,Q}$, it is equivalently L polynomials in $R_{n,q_i}, i = 0, \dots, L - 1$, so its RNS form \bar{f} is L vectors in $\mathbb{Z}_{q_i}^n, i = 0, \dots, L - 1$.

For simplicity we write $EK_i := (X^{\alpha'[2i]+\alpha'[2i+1]} - 1) \cdot EK_{i,0} + (X^{\alpha'[2i]} - 1)EK_{i,1} + (X^{\alpha'[2i+1]} - 1)EK_{i,2}$.

Recall in Algorithm 12, in each loop we calculate $AC_{i+1} = EK_i \odot AC_i$. We use overline to represent RNS form

$$f(x) \in R_{n',Q} \rightarrow \bar{f} \in \prod_{i=0}^{L-1} \mathbb{Z}_{q_i}^{n'},$$

which includes L vectors. The detailed algorithm is in Algorithm 14.

If we write $\overline{EK}_i = (\bar{\alpha}_i, \bar{\beta}_i)$, and $\overline{AC}_i = (\bar{a}_i, \bar{b}_i)$, then we have

$$\overline{AC}_{i+1} = \overline{EK}_i \odot \overline{AC}_i = \bar{\alpha}_i \diamond \bar{a}_i + \bar{\beta}_i \diamond \bar{b}_i.$$

Define (\diamond) operator in RNS form Algorithm 15.

C. Optimize Modulo Calculations

In this section, we present our optimizations on modulo calculations, which takes around 25% time (computed from Table I) in one LUT.

Different modulo calculations for different operations. Notice that the modulo calculation after addition is much easier than the modulo calculation after multiplication. The reason is that the size of coefficients grows slowly in addition. For example, we have two 59-bit coefficients c_1 and c_2 , then $c_1 + c_2$ is at most 60-bit, while $c_1 \times c_2$ could be 118-bit. So we write two modulo calculation functions: one is for addition, the other is for multiplication.

Reduce the number of modulo calculation. We separate the modulo calculation from the polynomial addition. In most works, modulo calculation is followed by every polynomial addition, which means that at any time, when we do a polynomial addition, then we will do a modulo calculation. In fact, many modulus calculations are not necessary. For example, in our previous face recognition experiment, the modulus is a 59-bit prime. However in our implementation, we use 64-bit integer data type to store the coefficients.

From $AC+ = ((X^{a'[2j]+a'[2j+1]} - 1)EK_{j,0} + (X^{a'[2j]} - 1)EK_{j,1} + (X^{a'[2j+1]} - 1)EK_{j,2}) \odot AC$, we can see that we only need compute one modulo calculation after two additions in the left side of \odot . Similarly, by the definition of \odot (Algorithm 10) and \diamond (Algorithm 9), we only need one modulo calculation after 2-3 additions. As a result, the number of modulo calculations is half of before.

D. Efficiency Analysis

In this section, we test and record the time taken of each part in one LUT, and to determine how much time our algorithm could save.

The parameters are same as those we used in our experiments. Length of LWE secret key is 512. Degree of LUT function and RLWE ciphertext is 2048. The modulus in our experiment is a 59-bit prime, which is not very large, so the method in Section IV-B2 does not fit for our case. Therefore, we apply the general method (Section IV-B1) to speed up our system. We summarize our results in the Table I. To observe the comparison clearly, we only ran it in single thread. From the table we can see that compared to our non-optimized LUT algorithm, our improved LUT saves 60% time.

V. HYBRID FHE-BASED PE-NN MODEL

A. Overview

The applications of basic FHE-based PE-NN is restricted by computationally expensive evaluations of DNNs on encrypted data. However, DNNs are widely used to solve many practical applications, such as facial recognition. Aiming for real-world scenarios, we develop Hybrid FHE-based PE-NN model, which could solve many practical problems within 1 second, such as facial recognition, text classification and object classification.

In many multi-party cases, such as the example of international collaboration to fight against transnational crime, if all

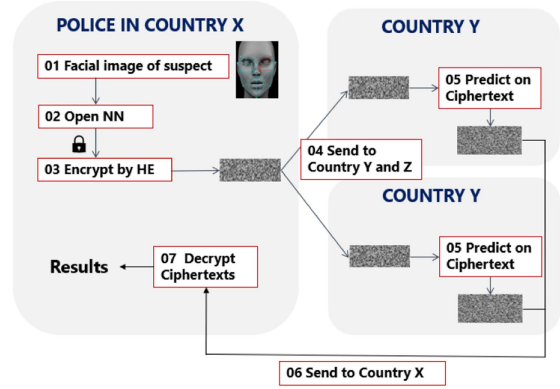


Fig. 4. Example: International collaboration to fight against transnational crime (hybrid model).

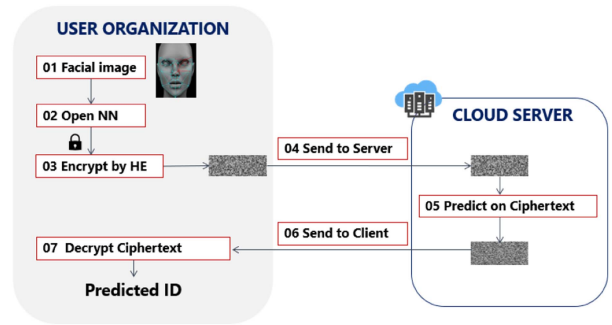


Fig. 5. Hybrid FHE-based PE-NN model.

parties commit on a feature extractor, then only the extracted vectors need to be encrypted and sent out, instead of the original photos. The flow is shown in Fig. 4 and the details can be found in Section I It can reduce the inference and communication latency significantly.

In AIaaS, there are two privacy problems in this model: 1. For AI solution providers, private networks are often trained by private datasets, and owners do not want to share the parameters with others. 2. On users' side as well, the users of AI models are not willing to disclose both the input data and the inference results to the server. To solve problem 1, model owners only publish open networks, and keep private networks in providers' cloud servers which allow users to make queries. To solve problem 2, users encrypt their data by homomorphic encryption, then the cloud servers evaluate the private networks on encrypted data and return encrypted inference results.

In summary, we use edge computing to drive the model in two steps (see Fig. 5): 1. User first runs the open network in plaintext locally; 2. The user encrypts the output from open network and sends the ciphertext to the server; 3. The server evaluates the private network on encrypted data and returns an encrypted inference output. 4. Only the user who has the secret key can decrypt and see the result.

In practical applications, such as facial recognition (FaceNet [25]) and object classification (InceptionV3 [29]), there are many open-source pre-trained neural network projects. Since both the parameters and models of these projects are public, so we can set them as the open network (usually discard

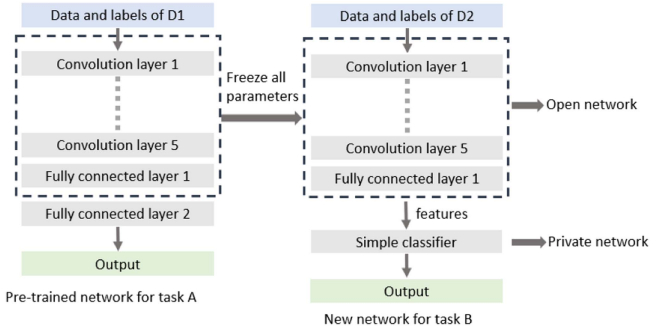


Fig. 6. Freezing all parameters in pre-trained network.

the fully connected layers). Meanwhile, the open project is usually not sufficient enough to solve their problems perfectly. There are two scenarios that might happen, either users' tasks differ too much from the open projects' tasks, or the open projects are not specific enough to solve the users' problem while running the inference on users' dataset. In the other words, the performance of open projects might tie to their own dataset. Inspired by transfer learning, AI solution provider can train a shallow network follows open network to solve user's task. Since the AI solution provider who trains the private network is usually not willing to share the parameters with users, so we set it as the private network and remain unseen in provider's cloud server. By transferring some computationally expensive ciphertext evaluations in the server side to lower-cost plaintext computations in the user side, our model significantly reduce the inference latency.

B. Training Method and Network Structure

Transfer learning [42] is a good technique to train the network in our hybrid model. Transfer learning focuses on storing knowledge gained while solving problem and applying it to a different but related problem. From practical standpoints, transferring information from previously learned tasks for the learning of new tasks could significantly improve the sample efficiency.

We now talk about how to use transfer learning to train the network in our hybrid model. Assuming that we want to build a network to solve task B , and we have a dataset \mathcal{D}_2 , we first search for a open-source pre-trained network on a related task A (trained by dataset \mathcal{D}_1). If we can find such a network, then we can apply transfer learning to obtain our network. There are two training methods to obtain the base of our network on task B from the pre-trained network on task A .

Freeze all parameters in pre-trained network. The idea of this method is that we use the pre-trained network as a feature extractor. Then we add a simple and shallow fully connected network (usually consists of 1 or 2 fully connected layer) and train the shallow network on these features. In this case, since the pre-trained network is open-source and we do not modify its parameters, so it can be set as the open network in our hybrid model. The shallow fully connected network is trained by ourselves and use our own dataset, so we can set it as the private network. The training method and network structure can be found in Fig. 6.

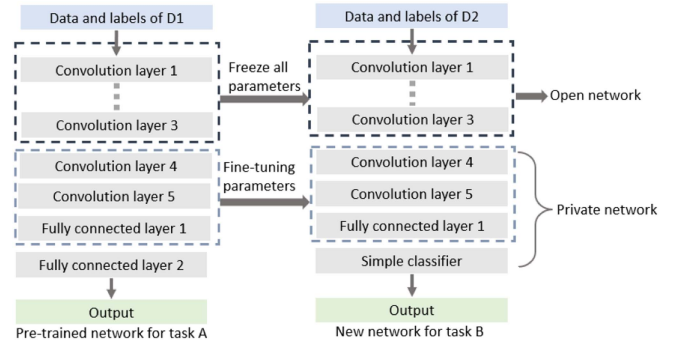


Fig. 7. Fine-tuning some layers in pre-trained network.

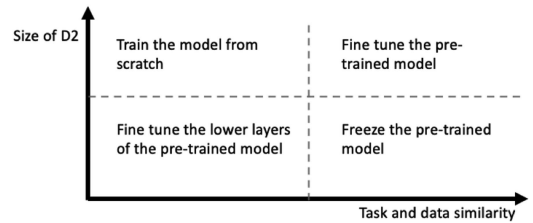


Fig. 8. Training method instruction.

Fine-Tune some layers in pre-trained network. The method will be used when the difference gap of A and task B are huge, or our own dataset \mathcal{D}_2 is large, then we can consider this method. The idea of this method is that we only freeze the parameters of the top layers in the pre-trained network, and we add a simple and shallow fully connected network (usually consists of 1 or 2 fully connected layer). Then we train the lower layers in pre-trained network together with the shallow fully connected network on our own dataset. In this case, the frozen part of the open-source pre-trained network can be set as the open network in our hybrid model, since we do not modify it. But the parameters in the lower layers are re-trained by our own dataset, so both the lower layers and the shallow fully connected network should be set as the private network. The training method and network structure can be found in Fig. 7.

Summary. When making a choice on training method, there are two criteria that we need to consider: (1) Task and data similarity; (2) Size of our dataset \mathcal{D}_2 , and follow the instruction in Fig. 8. The training method depends on the problem and sometimes should be decided by experiments in plaintext.

The worst case is that we cannot find a suitable open-source pre-trained network. Then we need to train the whole network by ourselves and determine which part can be open to public.

C. Time Complexity and Space Cost Analysis

The time complexity of one inference includes: 1) time of evaluating the open network in plaintext; 2) time of encrypting the extracted feature vector; 3) time of all homomorphic linear evaluations (inner-products in each layer of private network); 4) time of all homomorphic non-linear evaluations (activations in each layer of private network); 5) time of decrypting the inference result.

The time of open network evaluation depends on the complexity of open network and the computation resources in user side. For example, as shown in Section VI, in a 16-core CPU, some open network only takes 0.01s (such as text classification and object classification) while others may take 0.6s (speaker verification). It will be much faster if we use a powerful GPU to evaluate the open network. The encryption and decryption parts are very fast, which usually take less than 0.1 second in a laptop. Time of homomorphic linear evaluation grows linearly of the size of weight matrices, which typically takes 0.01s to evaluate a weight matrix with about 1,000 elements in private network. Evaluation of non-linear activation is the slowest part among all computations in our system. Time of one homomorphic non-linear evaluation takes about 0.1s in single-thread model. In our experiments, we apply multi-thread model to do homomorphic non-linear evaluation in a parallel way. Therefore, the total time of all homomorphic non-linear evaluations depends on the number of CPU cores and the number of activations.

The space cost consists of two parts: 1) the size of evaluation keys; 2) the size of input and output ciphertexts. Size of both parts depends on the security parameters of the FHE scheme. Evaluation keys are large (about 32 MB in our experiments) but the user only needs to send evaluation keys to the server once before all inferences, instead of sending them to the server in every inference. Input and output ciphertexts, whose size grows linearly of the length of input/output vectors, are small (less than 2 MB in our experiments).

VI. EVALUATION RESULTS

Besides the LUT evaluation in Section IV-D and noise growing evaluation in Section III-H1, in this section, we report the performance of our system. We evaluated all schemes on an AMD Milan 7313P 3.0 GHz CPU which owns 16 cores. The security level is at least 80 bits.⁴ Our experiments include one of the most popular benchmark dataset, MNIST and practical applications such as facial recognition, speaker verification, text classification and object classification. For MNIST dataset, we train a BP network with one hidden layer and 30 hidden nodes, which is a commonly used network structure in the area of inference on encrypted data [11]. The average time for one inference is 0.14s and the accuracy loss compared to inference in plaintext is only less than 1%.

Beside that, we also show that our system can be used to solve practical problems, such as facial recognition and so on. Tasks like facial recognition are more difficult and requires a deep neural network, so we apply our hybrid FHE-based PE-NN for inference. The average time for one facial recognition is 0.18s, while basic FHE-based PE-NN model takes around 28 hours in the same server environment.

We also achieve excellent results in speaker verification, text classification and object classification, which show that our system can be used in real-world tasks. To the best of our knowledge, this is the first work to solve real-world problems by applying FHE-based privacy-enhanced neural networks.

⁴We use BKZ simulator with core-(Q)sieving model to estimate the security for our scheme [43]. 80 bits security means that the attacker would have to perform 2^{80} operations to break it.

TABLE II
INFERENCE ON ENCRYPTED DATA IN MNIST DATASET

FHE scheme	Activation function	Plaintext inference	Ciphertext inference	
		Accuracy	Accuracy	Time
Our Scheme with LUT optimization	ReLU (support others)	94.80%	94.04%	0.14s
Our Scheme without LUT optimization	ReLU (support others)	94.80%	94.04%	0.42s
FHE-DiNN [11]	Sign function	94.76%	93.71%	0.49s

A. Experiments on Optimized FHE Scheme

MNIST dataset. There are 60,000 training samples and 10,000 test samples in MNIST dataset for handwritten digit recognition. Each input is a 28×28 gray-level image, which is represented by a vector with length 784. For each point of the image, we set the value is 1 if the original value is > 0 and set it to be 0 otherwise.

We first train a BP network with one hidden layer and 30 hidden nodes, whose input is 784-dim vector and output has 10 classes. The training phase is in clear and the accuracy is 94.80%.

We present our results in Table II. The first line is the result that we evaluate our FHE scheme which equips with our LUT optimization. The second line is the result that we evaluate our FHE scheme without the LUT optimization. The results show that LUT optimization in Section IV improves the efficiency for almost 3 times. Our inference accuracy on encrypted data is only 0.8% less than the inference accuracy on plain data. The third line shows the results in FHE-DiNN [11], which also evaluate their system in BP network with one hidden layer and 30 hidden nodes on MNIST dataset. Our system achieves better results in both inference accuracy and time on encrypted data.

B. Experiments on Hybrid PE-NN Model

We also show that our hybrid FHE-based PE-NN system can be used to solve practical problems with help of transfer learning. In this section, we report performance of our system on facial recognition, speaker ver, text classification and object classification. Compared to basic PE-NN model (i.e., the input is encrypted at the beginning and the whole DNN is evaluated at the server homomorphically), our hybrid PE-NN model achieves the state-of-the-art inference accuracy efficiently.

1) Facial Recognition: In this experiment, we build a hybrid network to do facial recognition in a group of people. The network will identify the identity of the input photo. We first establish our own training and test dataset, which contains photos of 30 people. Note that our test dataset is different from the training dataset.

We use a pre-trained FaceNet (trained by VGGFace2) as the open network and train a private fully connected network on our own dataset by freezing all parameters in pre-trained network. Then we implement this system by our optimized FHE scheme and test it. The structure of the recognition network is shown in Fig. 9.

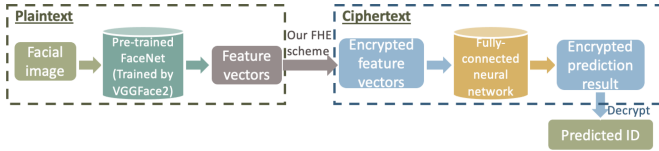


Fig. 9. Structure of face recognition network.

TABLE III
INFERENCE ON ENCRYPTED DATA IN FACIAL RECOGNITION

Model	Private network		Time	Accuracy
	Structure	Number of Activations		
Hybrid PE-NN	2 FC layers	30	Open: 0.04s Private: 0.21s	100%
Traditional NN	none	0	0.04s	100%
Basic PE-NN	FaceNet + 2 FC layers	$\approx 6M$	≈ 8 hours (estimated)	-

We present our evaluation results in Table III. In this table, we compare 3 neural network models for face recognition: 1. Our hybrid neural network, where the open network runs in plaintext and the private network runs in ciphertext; 2. Traditional neural network, which is not privacy-enhanced and the whole network runs in plaintext; 3. Basic privacy-enhanced neural network, where the whole network is privacy-enhanced and runs in ciphertext.

We can observe that traditional neural network is very fast, but it does not consider the privacy problem. Basic privacy-enhanced neural network is does it well at privacy protection, but it is too slow to be applied to real-world. Our hybrid neural network only needs less than one second per recognition, while basic privacy-enhanced neural network needs 28 hours in the same server environment. Therefore, our hybrid neural network achieved good balance between privacy protection and efficiency, and can be used in real applications.

2) *Speaker Verification*: Our scheme can be applied to process voice files as well. This can be applied when we do not want the neural network acquire the details of speech record while at the same time would require the running of neural network through the data to identify the identity of speaker. Our hybrid privacy-enhanced system can verify the owner of a speech record without knowing the speech content or the raw audio file. The dataset we used is *VoxCeleb1*, whose test set consists of 4,874 utterances from totally 40 speakers. The test set contain 37,611 trial pairs, and the task is to verify whether the trial pair of utterances are from the same speaker. We use MFA-Conformer model⁵ to do the feature extraction. After obtaining the embedding feature vector of the input audio, the vector is encrypted and send to the private network. We homomorphically evaluate the verification with pre-stored embedding vector in the private network. We present our results in Table IV. Equal error rate (EER) is used as the performance measure. We can see that the when homomorphic encryption is enabled in our hybrid PE-NN, the accuracy loss is very small, less than 0.1%. The time cost is satisfactory.

⁵Github repository: https://github.com/ductuantruong/mfa_conformer_sv.TABLE IV
INFERENCE ON ENCRYPTED DATA IN SPEAKER VERIFICATION

Model	Network structure	Time	Accuracy (1-EER)
Hybrid PE-NN	Open: MFA-Conformer Private: 1 FC layer	Open: 0.6s Private: 0.003s	99.22%
Traditional NN	Open: MFA-Conformer +1 FC layer No private network	0.6s	99.27%
Basic PE-NN	No open network Private: MFA-Conformer +1 FC layer	> 1 year (estimated)	-

TABLE V
INFERENCE ON ENCRYPTED DATA IN TEXT CLASSIFICATION

Model	Network structure	Time	Accuracy
Hybrid PE-NN	Open: TextCNN Private: 1 FC layer	Open: 0.013s Private: 0.003s	84%
Traditional NN	Open: TextCNN +1 FC layer No private network	0.013s	84%
Basic PE-NN	No open network Private: TextCNN +1 FC layer	≈ 2 minutes (estimated)	-

TABLE VI
INFERENCE ON ENCRYPTED DATA IN OBJECT CLASSIFICATION

Model	Network structure	Time	Accuracy
Hybrid PE-NN	Open: InceptionV3 Private: 1 FC layer	Open: 0.012s Private: 0.007s	99%
Traditional NN	Open: InceptionV3 +1 FC layer No private network	0.012s	99%
Traditional NN without TL	Open: InceptionV3 No private network	0.012s	85%
Basic PE-NN	No open network Private: InceptionV3 +1 FC layer	≈ 8 hours (estimated)	-

3) *Text Classification*: Text files can also be processed by our hybrid PE-NN scheme. Nowadays, we receive lots of emails and SMS. Many of them are advertisement or spam mail. It takes us lots of time to check such messages everyday. Text classification is a process of providing labels to the set of texts or words and those labels will tell us about the sentiment of the set of words. We can build a hybrid privacy-enhanced text classification system, which can add labels to text files without knowing the plaintext.

For this application, we tested our scheme on *Movie Review dataset*, which consists of positive and negative sentences/snippets. We use TextCNN [28] as the open network.

We present our results in Table V. Our scheme takes 0.016s per classification and does not lose the original accuracy compared to test in traditional neural network. Basic PE-NN takes around 6 minutes per inference, which is much slower. Notice that 84% accuracy is acceptable since the classifications are subjective and even human cannot achieve very high accuracy, e.g., [44].

4) *Object Classification*: Finally we show our system works well in object classification. For this application, we tested our scheme on *Cat and Dog dataset* (A famous competition on Kaggle.com). We use InceptionV3 [29] as the open network to extract the features. Results are in Table VI. Our scheme takes

0.019s per classification and achieved 99% accuracy. Training the 1-layer fully-connected network by transfer learning improve the inference accuracy significantly. When we use InceptionV3 directly, the inference accuracy is only 85%. After training by transfer learning, the inference accuracy becomes 99%.

VII. CONCLUSION AND FUTURE WORKS

This paper presented a practical approach for constructing privacy-enhanced neural networks by designing an efficient implementation of fully homomorphic encryption. As part of the efforts towards building a trusted digital economy, we aim to promote the adoption of AIaaS, which has emerged as an important trend for supporting the growth of the digital economy, by enabling AI models to process encrypted data.

In the global trend of digitalization, digital service providers make use of their vast amount of customer data to train AI models (such as image recognitions, financial modelling and pandemic modelling etc) and offer them as a service on the cloud. While there are convincing advantages for using such third-party models, the fact that model users are required to upload their data to the cloud is bound to raise serious privacy concerns, especially in the face of increasingly stringent privacy regulations and legislations [1].

With the proposed approach, an existing neural network can be converted to process FHE-encrypted data and produce encrypted output, which are only accessible by the model users, and more importantly, within an operationally acceptable time (e.g., within 1 s for facial recognition in typical border control systems). In our work, we apply our FHE technique to existing proven neural networks instead of building proprietary neural networks. Allowing privacy issues to be addressed separately from the accuracy of the AI models.

Experimental results show that in many practical tasks such as facial recognition, text classification and so on, we obtained the state-of-the-art inference accuracy in less than one second on a 16 cores CPU. In conclusion, our experiments show that in various practical tasks, our scheme achieved the state-of-the-art inference accuracy efficiently. It shows the feasibility of applying FHE-based PE-NN in real-world AI services to protect users' data.

As a future work, we will further improve our proposed techniques, and apply them in the area of privacy-enhancing machine learning to train AI models using encrypted data efficiently. Some improvements in implementation aspects can help to further reduce the time cost, such as an efficient implementation of homomorphic linear evaluation in multi-thread model. Packing techniques in homomorphic encryption are widely discussed recently. We will consider packing several data in one ciphertext to further increase the efficiency in inference and especially in training.

ACKNOWLEDGMENT

Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not reflect the views of National Research Foundation, Singapore and Infocomm Media Development Authority.

REFERENCES

- [1] K.-L. Tan, C.-H. Chi, and K.-Y. Lam, "Analysis of digital sovereignty and identity: From digitization to digitalization," 2022, *arXiv:2202.10069*.
- [2] K.-Y. Lam, S. Mitra, F. Gondesen, and X. Yi, "Ant-centric IoT security reference architecture—security-by-design for satellite-enabled smart cities," *IEEE Internet Things J.*, vol. 9, no. 8, pp. 5895–5908, Apr. 2022.
- [3] J. Fan et al., "Understanding security in smart city domains from the ant-centric perspective," *IEEE Internet Things J.*, vol. 10, no. 13, pp. 11199–11223, Jul. 2023.
- [4] C. Dwork, "Differential privacy: A survey of results," in *Proc. Int. Conf. Theory Appl. Models Comput.*, Springer, 2008, pp. 1–19.
- [5] M. Yang, I. Tjuawinata, and K.-Y. Lam, "K-means clustering with local -privacy for privacy-preserving data analysis," *IEEE Trans. Inf. Forensics Secur.*, vol. 17, pp. 2524–2537, 2022.
- [6] Z. Liu, J. Guo, W. Yang, J. Fan, K.-Y. Lam, and J. Zhao, "Privacy-preserving aggregation in federated learning: A survey," *IEEE Trans. Big Data*, to be published, doi: [10.1109/TBDDATA.2022.3190835](https://doi.org/10.1109/TBDDATA.2022.3190835).
- [7] H. Yang et al., "Lead federated neuromorphic learning for wireless edge artificial intelligence," *Nature Commun.*, vol. 13, no. 1, 2022, Art. no. 4269.
- [8] Z. Liu, J. Guo, K.-Y. Lam, and J. Zhao, "Efficient dropout-resilient aggregation for privacy-preserving machine learning," *IEEE Trans. Inf. Forensics Secur.*, vol. 18, pp. 1839–1854, 2022.
- [9] H. Yang, J. Zhao, Z. Xiong, K.-Y. Lam, S. Sun, and L. Xiao, "Privacy-preserving federated learning for UAV-enabled networks: Learning-based joint scheduling and resource management," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 10, pp. 3144–3159, Oct. 2021.
- [10] C. Gentry, "Fully homomorphic encryption using ideal lattices," in *Proc. 41st Annu. ACM Symp. Theory Comput.*, 2009, pp. 169–178.
- [11] F. Bourse, M. Minelli, M. Minihold, and P. Paillier, "Fast homomorphic evaluation of deep discretized neural networks," in *Proc. Annu. Int. Cryptol. Conf.*, Springer, 2018, pp. 483–512.
- [12] E. Chou, J. Beal, D. Levy, S. Yeung, A. Haque, and L. Fei-Fei, "Faster cryptonets: Leveraging sparsity for real-world encrypted inference," 2018, *arXiv:1811.09953*.
- [13] R. Gilad-Bachrach, N. Dowlin, K. Laine, K. Lauter, M. Naehrig, and J. Wernsing, "CryptoNets: Applying neural networks to encrypted data with high throughput and accuracy," in *Proc. Int. Conf. Mach. Learn.*, PMLR, 2016, pp. 201–210.
- [14] Q. Lou and L. Jiang, "SHE: A fast and accurate deep neural network for encrypted data," in *Proc. Adv. Neural Inf. Process. Syst.*, 2019, Art. no. 900.
- [15] A. Brutzkus, R. Gilad-Bachrach, and O. Elisha, "Low latency privacy preserving inference," in *Proc. Int. Conf. Mach. Learn.*, PMLR, 2019, pp. 812–821.
- [16] P. Mishra, R. Lehmkuhl, A. Srinivasan, W. Zheng, and R. A. Popa, "Delphi: A cryptographic inference service for neural networks," in *Proc. 29th USENIX Secur. Symp.*, 2020, pp. 2505–2522.
- [17] Y. LeCun, "The MNIST database of handwritten digits," 1998. [Online]. Available: <http://yann.lecun.com/exdb/mnist/>
- [18] A. Krizhevsky and G. Hinton, "Learning multiple layers of features from tiny images," 2009.
- [19] Z. Brakerski, C. Gentry, and V. Vaikuntanathan, "(Leveled) fully homomorphic encryption without bootstrapping," *ACM Trans. Comput. Theory*, vol. 6, no. 3, pp. 1–36, 2014.
- [20] J. H. Cheon, A. Kim, M. Kim, and Y. Song, "Homomorphic encryption for arithmetic of approximate numbers," in *Proc. Int. Conf. Theory Appl. Cryptol. Inf. Secur.*, Springer, 2017, pp. 409–437.
- [21] S. Halevi, Y. Polyakov, and V. Shoup, "An improved RNS variant of the BFV homomorphic encryption scheme," in *Proc. Cryptographers' Track RSA Conf.*, San Francisco, CA, USA, Springer, 2019, pp. 83–105.
- [22] I. Chillotti, N. Gama, M. Georgieva, and M. Izabachène, "TFHE: Fast fully homomorphic encryption over the torus," *J. Cryptol.*, vol. 33, no. 1, pp. 34–91, 2020.
- [23] L. Ducas and D. Micciancio, "FHEW: Bootstrapping homomorphic encryption in less than a second," in *Proc. Annu. Int. Conf. Theory Appl. Cryptographic Techn.*, Springer, 2015, pp. 617–640.
- [24] F. Boemer, Y. Lao, R. Cammarota, and C. Wierzynski, "nGraph-HE: A graph compiler for deep learning on homomorphically encrypted data," in *Proc. 16th ACM Int. Conf. Comput. Front.*, 2019, pp. 3–13.
- [25] F. Schroff, D. Kalenichenko, and J. Philbin, "FaceNet: A unified embedding for face recognition and clustering," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2015, pp. 815–823.
- [26] D. Micciancio and Y. Polyakov, "Bootstrapping in fhe-like cryptosystems," in *Proc. 9th Workshop Encrypted Comput. Appl. Homomorphic Cryptogr.*, 2021, pp. 17–28.

- [27] W.-J. Lu, Z. Huang, C. Hong, Y. Ma, and H. Qu, "PEGASUS: Bridging polynomial and non-polynomial evaluations in homomorphic encryption," in *Proc. IEEE Symp. Secur. Privacy*, 2021, pp. 1057–1073.
- [28] B. Guo, C. Zhang, J. Liu, and X. Ma, "Improving text classification with weighted word embeddings via a multi-channel TextCNN model," *Neurocomputing*, vol. 363, pp. 366–374, 2019.
- [29] X. Xia, C. Xu, and B. Nan, "Inception-V3 for flower classification," in *Proc. 2nd Int. Conf. Image Vis. Comput.*, 2017, pp. 783–787.
- [30] O. Regev, "On lattices, learning with errors, random linear codes, and cryptography," *J. ACM*, vol. 56, no. 6, pp. 1–40, 2009.
- [31] V. Lyubashevsky, C. Peikert, and O. Regev, "On ideal lattices and learning with errors over rings," in *Proc. Annu. Int. Conf. Theory Appl. Cryptographic Techn.*, Springer, 2010, pp. 1–23.
- [32] A. Costache and N. P. Smart, "Which ring based somewhat homomorphic encryption scheme is best?," in *Proc. Cryptographers' Track RSA Conf.*, Springer, 2016, pp. 325–340.
- [33] Z. Brakerski and V. Vaikuntanathan, "Efficient fully homomorphic encryption from (standard) LWE," *SIAM J. Comput.*, vol. 43, no. 2, pp. 831–871, 2014.
- [34] T. ElGamal, "A public key cryptosystem and a signature scheme based on discrete logarithms," *IEEE Trans. Inf. Theory*, vol. 31, no. 4, pp. 469–472, Jul. 1985.
- [35] P. Paillier, "Public-key cryptosystems based on composite degree residuosity classes," in *Proc. Int. Conf. Theory Appl. Cryptographic Techn.*, Springer, 1999, pp. 223–238.
- [36] J. W. Cooley and J. W. Tukey, "An algorithm for the machine calculation of complex Fourier series," *Math. Comput.*, vol. 19, no. 90, pp. 297–301, 1965.
- [37] R. Agarwal and C. Burrus, "Fast convolution using fermat number transforms with applications to digital filtering," *IEEE Trans. Acoust. Speech Signal Process.*, vol. 22, no. 2, pp. 87–97, Apr. 1974.
- [38] K. Hornik, "Approximation capabilities of multilayer feedforward networks," *Neural Netw.*, vol. 4, no. 2, pp. 251–257, 1991.
- [39] G. Cybenko, "Approximation by superpositions of a sigmoidal function," *Math. Control Signals Syst.*, vol. 2, no. 4, pp. 303–314, 1989.
- [40] Z. Brakerski, A. Langlois, C. Peikert, O. Regev, and D. Stehlé, "Classical hardness of learning with errors," in *Proc. 45th Annu. ACM Symp. Theory Comput.*, 2013, pp. 575–584.
- [41] D. Micciancio, "On the hardness of learning with errors with binary secrets," *Theory Comput.*, vol. 14, no. 1, pp. 1–17, 2018.
- [42] S. J. Pan and Q. Yang, "A survey on transfer learning," *IEEE Trans. Knowl. Data Eng.*, vol. 22, no. 10, pp. 1345–1359, Oct. 2010.
- [43] X. Lu et al., "LAC: Practical ring-LWE based public-key encryption with byte-level modulus," *Cryptol. ePrint Arch.*, 2018.
- [44] A. K. Sharma, S. Chaurasia, and D. K. Srivastava, "Sentimental short sentences classification by using CNN deep learning model with fine tuned Word2Vec," *Procedia Comput. Sci.*, vol. 167, pp. 1139–1147, 2020.



Kwok-Yan Lam (Senior Member, IEEE) received the BSc degree (1st Class Hons.) from the University of London, in 1987, and the PhD degree from the University of Cambridge, in 1990. He is the associate vice president (Strategy and Partnerships) in the President's Office, and professor with the School of Computer Science and Engineering, Nanyang Technological University, Singapore. He is currently also the executive director of the National Centre for Research in Digital Trust, and director of the Strategic Centre for Research in Privacy-Preserving Technologies and Systems (SCRiPTS). From 2020, he is on part-time secondment to the INTERPOL as a consultant with Cyber and New Technology Innovation. Prior to joining NTU, he has been a professor of the Tsinghua University, PR China (2002–2010) and a faculty member of the National University of Singapore and the University of London since 1990. He was a visiting scientist with the Isaac Newton Institute, Cambridge University, and a visiting professor with the European Institute for Systems Security. In 1998, he received the Singapore Foundation Award from the Japanese Chamber of Commerce and Industry in recognition of his research and development achievement in information security in Singapore. He is the recipient of the Singapore Cybersecurity Hall of Fame Award, in 2022. His research interests include distributed systems, intelligent systems, IoT security, distributed protocols for blockchain, homeland security and cybersecurity.



Xianhui Lu received the PhD degree in information security from Southwest Jiaotong University, in 2009. He is a professor with the Institute of Information Engineering, Chinese Academy of Sciences. His current research focuses on homomorphic encryption, post-quantum cryptography and physical layer cryptography. He is one of the editors of the post-quantum study project of ISO/IEC SC27 WG2. He is the Principal author of the algorithm LAC, which is currently one of the 26 s round candidates of the NIST Post-Quantum Cryptography Competition.



Linru Zhang received the BSc degree from Sun Yat-Sen University, in 2016, and the PhD degree in computer science from the University of Hong Kong, in 2021. She is currently a research fellow with Nanyang Technological University, Singapore. Her research interest includes cryptography and privacy-preserving machine learning.



Xiangning Wang received the BSc degree from Peking University, in 2016, and the PhD degree in computer science from the University of Hong Kong, in 2021. He is currently a research fellow with Nanyang Technological University, Singapore. His research interest includes privacy-preserving machine learning and differential privacy.



Huaxiong Wang received the PhD degree in mathematics from the University of Haifa, Israel, in 1996, and the PhD degree in computer science from the University of Wollongong, Australia, in 2001. He has been with Nanyang Technological University in Singapore since 2006, where he is a professor with the Division of Mathematical Sciences. Currently he is also the co-director of National Centre for Research in Digital Trust and the deputy director of Strategic Centre for Research in Privacy-Preserving Technologies and Systems with NTU. Prior to NTU, he held faculty positions with Macquarie University and University of Wollongong in Australia, and visiting positions with ENS de Lyon in France, City University of Hong Kong, National University of Singapore and Kobe University in Japan. His research interest is in cryptography and cybersecurity. He was the program co-chair of Asiacrypt 2020 and 2021.



Si Qi Goh received the BE degree in computer science from Monash University, Malaysia with first class honors, in 2020. She is currently working toward the PhD degree with the School of Computer Science and Engineering, Nanyang Technological University, Singapore. Her research interest includes privacy-preserving machine learning, explainable artificial intelligence, and digital trust.