
Privacy-Preserving CNN Training with Transfer Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Privacy-preserving neural network inference has been well studied while homo-
2 morphic CNN training still remains an open challenging task. In this paper, we
3 present a practical solution to implement privacy-preserving CNN training based
4 on mere Homomorphic Encryption (HE) technique. To our best knowledge, this
5 is the first attempt successfully to crack this nut and no work ever before has
6 achieved this goal. Several techniques combine to accomplish the task: (1) with
7 transfer learning, privacy-preserving CNN training can be reduced to homomor-
8 phic neural network training, or even multiclass logistic regression (MLR) train-
9 ing; (2) via a faster gradient variant called Quadratic Gradient, an enhanced
10 gradient method for MLR with a state-of-the-art performance in convergence
11 speed is applied in this work to achieve high performance; (3) we employ the
12 thought of transformation in mathematics to transform approximating Softmax
13 function in the encryption domain to the approximation of the Sigmoid function.
14 A new type of loss function termed Squared Likelihood Error has been de-
15 veloped alongside to align with this change.; and (4) we use a simple but flexible
16 matrix-encoding method named Volley Revolver to manage the data flow in
17 the ciphertexts, which is the key factor to complete the whole homomorphic CNN
18 training. The complete, runnable C++ code to implement our work can be found
19 at: <https://anonymous.4open.science/r/HE-CNNtraining-B355/>.

20 We select REGNET_X_400MF as our pre-trained model for transfer learning. We
21 use the first 128 MNIST training images as training data and the whole MNIST
22 testing dataset as the testing data. The client only needs to upload 6 ciphertexts to
23 the cloud and it takes ~ 21 mins to perform 2 iterations on a cloud with 64 vCPUs,
24 resulting in a precision of 21.49%.

25 1 Introduction

26 1.1 Background

27 Applying machine learning to problems involving sensitive data requires not only accurate predictions
28 but also careful attention to model training. Legal and ethical requirements might limit the use of
29 machine learning solutions based on a cloud service for such tasks. As a particular encryption scheme,
30 homomorphic encryption provides the ultimate security for these machine learning applications and
31 ensures that the data remains confidential since the cloud does not need private keys to decrypt it.
32 However, it is a big challenge to train the machine learning model, such as neural networks or even
33 convolution neural networks, in such encrypted domains. Nonetheless, we will demonstrate that
34 cloud services are capable of applying neural networks over the encrypted data to make encrypted
35 training, and also return them in encrypted form.

36 1.2 Related work

37 Several studies on machine learning solutions are based on homomorphic encryption in the cloud
38 environment. Since Gilad-Bachrach et al. [1] firstly considered privacy-preserving deep learning
39 prediction models and proposed the private evaluation protocol `CryptoNets` for CNN, many other
40 approaches [2, 3, 4, 5] for privacy-preserving deep learning prediction based on HE or its combination
41 with other techniques have been developed. Also, there are several studies [6, 7, 8, 9] working on
42 logistic regression models based on homomorphic encryption.

43 However, to our best knowledge, no work ever before based on mere HE technique has presented an
44 solution to successfully perform homomorphic CNN training.

45 1.3 Contributions

46 Our specific contributions in this paper are as follows:

- 47 1. with various techniques, we initiate to propose a practical solution for privacy-preserving
48 CNN training, demonstrating the feasibility of homomorphic CNN training.
- 49 2. We suggest a new type of loss function, Squared Likelihood Error (SLE), which is
50 friendly to privacy-preserving manner. As a result, we can use the Sigmoid function to
51 replace the Softmax function which is too difficult to calculate in the encryption domain due
52 to its uncertainty.
- 53 3. We develop a new algorithm with SLE loss function for MLR using quadratic gradient.
54 Experiments show that this HE-friendly algorithm has a state-of-the-art performance in
55 convergence speed.

56 2 Preliminaries

57 We adopt “ \otimes ” to denote the kronecker product and “ \odot ” to denote the component-wise multiplication
58 between matrices.

59 2.1 Fully Homomorphic Encryption

60 Homomorphic Encryption (HE) is one type of encryption scheme with a special characteristic called
61 *Homomorphic*, which allows to compute on encrypted data without having access to the secret key.
62 Fully HE means that the scheme is fully homomorphic, namely, homomorphic with regards to both
63 addition and multiplication, and that it allows arbitrary computation on encrypted data. Since Gentry
64 proposed the first fully HE scheme [10] in 2009, some technological progress on HE has been made.
65 For example, Brakerski, Gentry and Vaikuntanathan [11] present a novel way of constructing leveled
66 fully homomorphic encryption schemes (BGV) and Smart and Vercauteren [12] introduced one of the
67 most important features of HE systems, a packing technique based on polynomial-CRT called Single
68 Instruction Multiple Data (aka SIMD) to encrypt multiple values into a single ciphertext. Another
69 great progress in terms of machine learning applications is the *rescaling* procedure [13], which can
70 manage the magnitude of plaintext effectively.

71 Modern fully HE schemes, such as HEAAN, usually support several common homomorphic operations:
72 the encryption algorithm `Enc` encrypting a vector, the decryption algorithm `Dec` decrypting
73 a ciphertext, the homomorphic addition `Add` and multiplication `Mult` between two ciphertexts, the
74 multiplication `cMult` of a constant vector with a ciphertext, the rescaling operation `ReScale` to reduce
75 the magnitude of a plaintext to an appropriate level, the rotation operation `Rot` generating a new
76 ciphertext encrypting the shifted plaintext vector, and the bootstrapping operation `bootstrap` to
77 refresh a ciphertext usually with a small ciphertext modulus.

78 2.2 Database Encoding Method

79 For a given database Z , Kim et al. [6] first developed an efficient database encoding method, in order
80 to make full use of the HE computation and storage resources. They first expand the matrix database
81 to a vector form V in a row-by-row manner and then encrypt this vector V to obtain a ciphertext
82 $Z = Enc(V)$. Also, based on this database encoding, they mentioned two simple operations via

83 shifting the encrypted vector by two different positions, respectively: the complete row shifting
 84 and the *incomplete* column shifting. These two operations performing on the matrix Z output the
 85 matrices Z' and Z'' , as follows:

$$Z = \begin{bmatrix} x_{10} & x_{11} & \dots & x_{1d} \\ x_{20} & x_{21} & \dots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n0} & x_{n1} & \dots & x_{nd} \end{bmatrix}, \quad Z' = Enc \begin{bmatrix} x_{20} & x_{21} & \dots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n0} & x_{n1} & \dots & x_{nd} \\ x_{10} & x_{11} & \dots & x_{1d} \end{bmatrix},$$

$$Z'' = Enc \begin{bmatrix} x_{11} & \dots & x_{1d} & x_{20} \\ x_{21} & \dots & x_{2d} & x_{30} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{nd} & x_{10} \end{bmatrix}, \quad Z''' = Enc \begin{bmatrix} x_{11} & \dots & x_{1d} & x_{10} \\ x_{21} & \dots & x_{2d} & x_{20} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{nd} & x_{n0} \end{bmatrix}.$$

86 The complete column shifting to obtain the matrix Z''' can also be achieved by two Rot, two cMult,
 87 and an Add.

88 Other works [14, 4] using the same encoding method also developed some other procedures, such
 89 as SumRowVec and SumColVec to calculate the summation of each row and column, respectively.
 90 Such basic common and simple operations consisting of a series of HE operations are significantly
 91 important for more complex calculations such as the homomorphic evaluation of gradient.

92 2.3 Convolutional Neural Network

93 Inspired by biological processes, Convolutional Neural Networks (CNN) are a type of artificial neural
 94 network most commonly used to analyze visual images. CNNs play a significant role in image
 95 recognition due to their powerful performance. It is also worth mentioning that the CNN model is
 96 one of a few deep learning models built with reference to the visual organization of the human brain.

97 2.3.1 Transfer Learning

98 Transfer learning in machine learning is a class of methods in which a pretrained model can be used
 99 as an optimization for a new model on a related task, allowing rapid progress in modeling the new
 100 task. In real-world applications, very few researchers train entire convolutional neural networks
 101 from scratch for image processing-related tasks. Instead, it is common to use a well-trained CNN
 102 as a fixed feature extractor for the task of interest. In our case, we freeze all the weights of the
 103 selected pre-trained CNN except that of the final fully-connected layer. We then replace the last
 104 fully-connected layer with a new layer with random weights (such as zeros) and only train this layer.

105 **REGNET_X_400MF** To use transfer learning in our privacy-preserving CNN training, we adopt
 106 a new network design paradigm called RegNet, recently introduced by Facebook AI researchers,
 107 as our pre-trained model. RegNet is a low-dimensional design space consisting of simple, regular
 108 networks. In particular, we apply REGNET_X_400MF as a fixed feature extractor and replaced the final
 109 fully connected layer with a new one of zero weights. CNN training in this case can be simplified
 110 to multiclass logistic regression training. Since REGNET_X_400MF only receive color images of size
 111 224×224 , the grayscale images will be stacked threefold and images of different sizes will be resized
 112 to the same size in advance. These two transformations can be done by using PyTorch.

113 2.3.2 Datasets

114 We adopt three common datasets in our experiments: MNIST, USPS, and CIFAR10. Table 1 describes
 115 the three datasets.

116 3 Technical details

117 3.1 Multiclass Logistic Regression

118 Multiclass Logistic Regression, or Multinomial Logistic Regression, can be seen as an extension of
 119 logistic regression for multi-class classification problems. Supposing that the matrix $X \in \mathbb{R}^{n \times (1+d)}$,

Table 1: Characteristics of the several datasets used in our experiments

Dataset	No. Samples (training)	No. Samples (testing)	No. Features	No. Classes
USPS	7,291	2,007	16×16	10
MNIST	60,000	10,000	28×28	10
CIFAR-10	50,000	10,000	3×32×32	10

120 the column vector $Y \in \mathbb{N}^{n \times 1}$, the matrix $\bar{Y} \in \mathbb{R}^{n \times c}$, and the matrix $W \in \mathbb{R}^{c \times (1+d)}$ represent
 121 the dataset, class labels, the one-hot encoding of the class labels, and the MLR model parameter,
 122 respectively:

$$\begin{aligned}
 X &= \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} x_{[1][0]} & x_{[1][1]} & \cdots & x_{[1][d]} \\ x_{[2][0]} & x_{[2][1]} & \cdots & x_{[2][d]} \\ \vdots & \vdots & \ddots & \vdots \\ x_{[n][0]} & x_{[n][1]} & \cdots & x_{[n][d]} \end{bmatrix}, \\
 Y &= \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \xrightarrow{\text{one-hot encoding}} \bar{Y} = \begin{bmatrix} \bar{y}_1 \\ \bar{y}_2 \\ \vdots \\ \bar{y}_n \end{bmatrix} = \begin{bmatrix} y_{[1][1]} & y_{[1][2]} & \cdots & y_{[1][c-1]} \\ y_{[2][1]} & y_{[2][2]} & \cdots & y_{[2][c-1]} \\ \vdots & \vdots & \ddots & \vdots \\ y_{[n][1]} & y_{[n][2]} & \cdots & y_{[n][c-1]} \end{bmatrix}, \\
 W &= \begin{bmatrix} w_{[0]} \\ w_{[1]} \\ \vdots \\ w_{[c-1]} \end{bmatrix} = \begin{bmatrix} w_{[0][0]} & w_{[0][1]} & \cdots & w_{[0][d]} \\ w_{[1][0]} & w_{[1][1]} & \cdots & w_{[1][d]} \\ \vdots & \vdots & \ddots & \vdots \\ w_{[c-1][0]} & w_{[c-1][1]} & \cdots & w_{[c-1][d]} \end{bmatrix}.
 \end{aligned}$$

MLR aims to maximize L or $\ln L$:

$$L = \prod_{i=1}^n \frac{\exp(x_i \cdot w_{[y_i]}^T)}{\sum_{k=0}^{c-1} \exp(x_i \cdot w_{[k]}^T)} \mapsto \ln L = \sum_{i=1}^n [x_i \cdot w_{[y_i]}^T - \ln \sum_{k=0}^{c-1} \exp(x_i \cdot w_{[k]}^T)].$$

123 The loss function $\ln L$ is a multivariate function of $[(1+c)(1+d)]$ variables, which has its column-
 124 vector gradient ∇ of size $[(1+c)(1+d)]$ and Hessian square matrix ∇^2 of order $[(1+c)(1+d)]$ as
 125 follows:

$$\begin{aligned}
 \nabla &= \frac{\partial \ln L}{\partial \pi} = \left[\frac{\partial \ln L}{\partial w_{[0]}}, \frac{\partial \ln L}{\partial w_{[1]}}, \dots, \frac{\partial \ln L}{\partial w_{[c-1]}} \right]^\top, \\
 \nabla^2 &= \begin{bmatrix} \frac{\partial^2 \ln L}{\partial w_{[0]} \partial w_{[0]}} & \frac{\partial^2 \ln L}{\partial w_{[0]} \partial w_{[1]}} & \cdots & \frac{\partial^2 \ln L}{\partial w_{[0]} \partial w_{[c-1]}} \\ \frac{\partial^2 \ln L}{\partial w_{[1]} \partial w_{[0]}} & \frac{\partial^2 \ln L}{\partial w_{[1]} \partial w_{[1]}} & \cdots & \frac{\partial^2 \ln L}{\partial w_{[1]} \partial w_{[c-1]}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 \ln L}{\partial w_{[c-1]} \partial w_{[0]}} & \frac{\partial^2 \ln L}{\partial w_{[c-1]} \partial w_{[1]}} & \cdots & \frac{\partial^2 \ln L}{\partial w_{[c-1]} \partial w_{[c-1]}} \end{bmatrix}.
 \end{aligned}$$

126 **Nesterov's Accelerated Gradient** With ∇ or ∇^2 , first-order gradient algorithms or second-order
 127 Newton–Raphson method are commonly applied in MLE to maximise $\ln L$. In particular, Nesterov's
 128 Accelerated Gradient (NAG) is a practical solution for homomorphic MLR without frequent inversion
 129 operations. It seems plausible that the NAG method is probably the best choice for privacy-preserving
 130 model training.

131 3.2 Chiang's Quadratic Gradient

132 Chiang's Quadratic Gradient (CQG) [15, 16, 9] is a faster, promising gradient variant that can
 133 combine the first-order gradient descent/ascent algorithms and the second-order Newton–Raphson
 134 method, accelerating the raw Newton–Raphson method with various gradient algorithms and probably

135 helpful to build super-quadratic algorithms. For a function $F(x)$ with its gradient g and Hessian
 136 matrix H , to build CQG, we first construct a diagonal matrix \bar{B} from the Hessian H itself:

$$\bar{B} = \begin{bmatrix} \frac{1}{\varepsilon + \sum_{i=0}^d |\bar{h}_{0i}|} & 0 & \dots & 0 \\ 0 & \frac{1}{\varepsilon + \sum_{i=0}^d |\bar{h}_{1i}|} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{1}{\varepsilon + \sum_{i=0}^d |\bar{h}_{di}|} \end{bmatrix},$$

137 where \bar{h}_{ji} is the elements of the matrix H and ε is a small constant positive number.

138 CQG for the function $F(x)$, defined as $G = \bar{B} \cdot g$, has the same dimension as the raw gradient g . To
 139 apply CQG in practice, we can use it in the same way as the first-order gradient algorithms, except
 140 that we need to replace the naive gradient with the quadratic gradient and adopt a new learning rate
 141 (usually by increasing 1 to the original learning rate).

142 For efficiency in applying CQG, a good bound matrix should be attempted to obtain in order to
 143 replace the Hessian itself. Chiang has proposed the enhanced NAG method via CQG for MLR with a
 144 fixed Hessian [17, 7, 18] substitute built from $\frac{1}{2}X^T X$.

145 3.3 Approximating Softmax Function

146 It might be impractical to perfectly approximate Softmax function in the privacy-preserving domain
 147 due to its uncertainty. To address this issue, we employ the thought of transformation from mathemat-
 148 ics: transforming one tough problem into another easier one. That is, instead of trying to approximate
 149 the Softmax function, we attempt to approximate the Sigmoid function in the encryption domain,
 150 which has been well-studied by several works using the least-square method.

In line with standard practice of the log-likelihood loss function involving the Softmax function, we
 should try to maximize the new loss function

$$L_1 = \prod_{i=1}^n \frac{1}{1 + \exp(-x_i \cdot w_{[y_i]}^T)}.$$

151 We can prove that $\ln L_1$ is concave and deduce that $\frac{1}{4}E \otimes X^T X$ can be used to build the CQG for
 152 $\ln L_1$. However, the performance of this loss function $\ln L_1$ is not ideal, probably because for the
 153 individual example its gradient and Hessian contain no information about any other class weights not
 154 related to this example.

Squared Likelihood Error After many attempts to finding a proper loss function, we develop
 a novel loss function that can have a competitive performance to the log-likelihood loss function,
 which we term Squared Likelihood Error (SLE):

$$L_2 = \prod_{i=1}^n \prod_{j=0}^{c-1} (\bar{y}_i - \text{Sigmoid}(x_i \cdot w_{[y_i]}^T))^2 \mapsto \ln L_2 = \sum_{i=1}^n \sum_{j=0}^{c-1} \ln |\bar{y}_i - \text{Sigmoid}(x_i \cdot w_{[y_i]}^T)|.$$

155 We can also prove that $\ln L_2$ is concave and that $\frac{1}{4}E \otimes X^T X$ can be used to build the CQG for $\ln L_2$.
 156 The loss function SLE might be related to Mean Squared Error (MSE): the MSE loss function sums
 157 all the squared errors while SLE calculates the cumulative product of all the squared likelihood errors.

158 Combining together all the techniques above, we now have the enhanced NAG method with the SLE
 159 loss function for MLR training, described in detail in Algorithm 1.

160 Performance Evaluation We test the convergence speed of the raw NAG method with log-
 161 likelihood loss function (denoted as RawNAG), the NAG method with SLE loss function (denoted
 162 as SigmoidNAG), and the enhanced NAG method via CQG with SLE loss function (denoted as
 163 SigmoidNAGQG) on the three datasets described above: USPS, MNIST, and CIFAR10. Since two
 164 different types of loss functions are used in these three methods, the loss function directly measuring
 165 the performance of various methods will not be selected as the indicator. Instead, we select precision
 166 as the only indicator in the following Python experiments. Note that we use REGNET_X_400MF to in

Algorithm 1 The Enhanced NAG method with the SLE loss function for MLR Training

Input: training dataset $X \in \mathbb{R}^{n \times (1+d)}$; one-hot encoding training label $Y \in \mathbb{R}^{n \times c}$; and the number κ of iterations;

Output: the parameter matrix $V \in \mathbb{R}^{c \times (1+d)}$ of the MLR

- 1: Set $\bar{H} \leftarrow -\frac{1}{4}X^\top X$ $\triangleright \bar{H} \in \mathbb{R}^{(1+d) \times (1+d)}$
- 2: Set $V \leftarrow \mathbf{0}$, $W \leftarrow \mathbf{0}$, $\bar{B} \leftarrow \mathbf{0}$ $\triangleright V \in \mathbb{R}^{c \times (1+d)}$, $W \in \mathbb{R}^{c \times (1+d)}$, $\bar{B} \in \mathbb{R}^{c \times (1+d)}$
- 3: **for** $j := 0$ to d **do**
- 4: $\bar{B}[0][j] \leftarrow \varepsilon$ $\triangleright \varepsilon$ is a small positive constant such as $1e - 10$
- 5: **for** $i := 0$ to d **do**
- 6: $\bar{B}[0][j] \leftarrow \bar{B}[0][j] + |\bar{H}[i][j]|$
- 7: **end for**
- 8: **for** $i := 1$ to $c - 1$ **do**
- 9: $\bar{B}[i][j] \leftarrow \bar{B}[0][j]$
- 10: **end for**
- 11: **for** $i := 0$ to $c - 1$ **do**
- 12: $\bar{B}[i][j] \leftarrow 1.0/\bar{B}[i][j]$
- 13: **end for**
- 14: **end for**
- 15: Set $\alpha_0 \leftarrow 0.01$, $\alpha_1 \leftarrow 0.5 \times (1 + \sqrt{1 + 4 \times \alpha_0^2})$
- 16: **for** $count := 1$ to κ **do**
- 17: Set $Z \leftarrow X \times V^\top$ $\triangleright Z \in \mathbb{R}^{n \times c}$ and V^\top means the transpose of matrix V
- 18: **for** $i := 1$ to n **do** $\triangleright Z$ is going to store the inputs to the Sigmoid function
- 19: **for** $j := 0$ to d **do**
- 20: $Z[i][j] \leftarrow 1/(1 + e^{-Z[i][j]})$
- 21: **end for**
- 22: **end for**
- 23: Set $g \leftarrow (Y - Z)^\top \times X$ $\triangleright g \in \mathbb{R}^{c \times (1+d)}$
- 24: Set $G \leftarrow \mathbf{0}$
- 25: **for** $i := 0$ to $c - 1$ **do**
- 26: **for** $j := 0$ to d **do**
- 27: $G[i][j] \leftarrow \bar{B}[i][j] \times g[i][j]$
- 28: **end for**
- 29: **end for**
- 30: Set $\eta \leftarrow (1 - \alpha_0)/\alpha_1$, $\gamma \leftarrow 1/(n \times count)$ $\triangleright n$ is the size of training data
- 31: $w_{temp} \leftarrow W + (1 + \gamma) \times G$
- 32: $W \leftarrow (1 - \eta) \times w_{temp} + \eta \times V$
- 33: $V \leftarrow w_{temp}$
- 34: $\alpha_0 \leftarrow \alpha_1$, $\alpha_1 \leftarrow 0.5 \times (1 + \sqrt{1 + 4 \times \alpha_0^2})$
- 35: **end for**
- 36: **return** W

167 advance extract the features of USPS, MNIST, and CIFAR10, resulting in a new same-size dataset
168 with 401 features of each example. Figure 1 shows that our enhanced methods all converge faster
169 than other algorithms on the three datasets.

170 3.4 Double Volley Revolver

171 Unlike those efficient, complex encoding methods [3], Volley Revolver is a simple, flexible
172 matrix-encoding method specialized for privacy-preserving machine-learning applications, whose
173 basic idea in a simple version is to encrypt the transpose of the second matrix for two matrices to
174 perform multiplication. Figure 2 describes a simple case for the algorithm adopted in this encoding
175 method.

176 The encoding method actually plays a significant role in implementing privacy-preserving CNN
177 training. Just as Chiang mentioned in [4], we show that Volley Revolver can indeed be used to
178 implement homomorphic CNN training. This simple encoding method can help to control and
179 manage the data flow through ciphertexts.

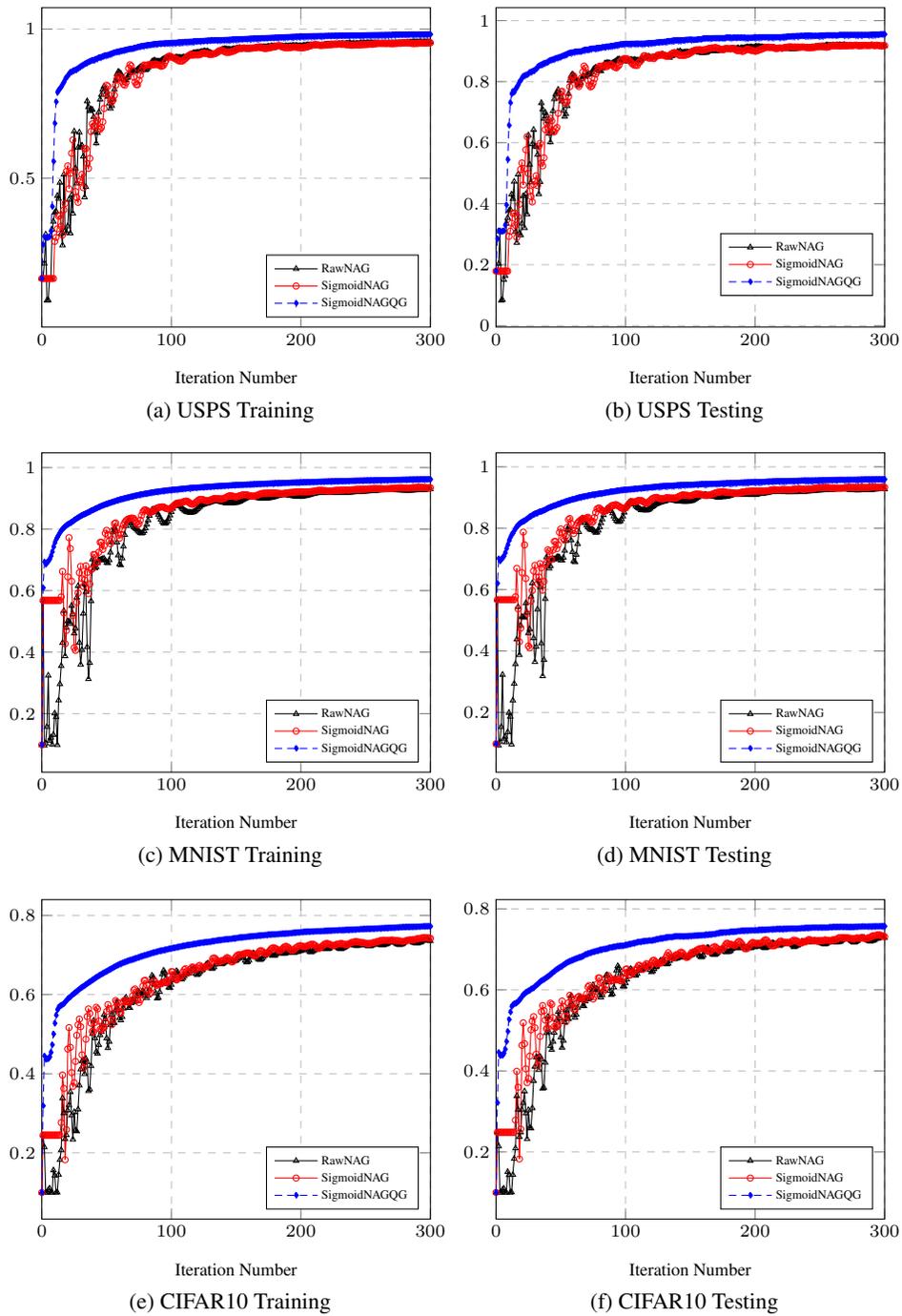


Figure 1: Training and Testing precision results for raw NAG vs. NAG with SLE vs. The enhanced NAG with SLE

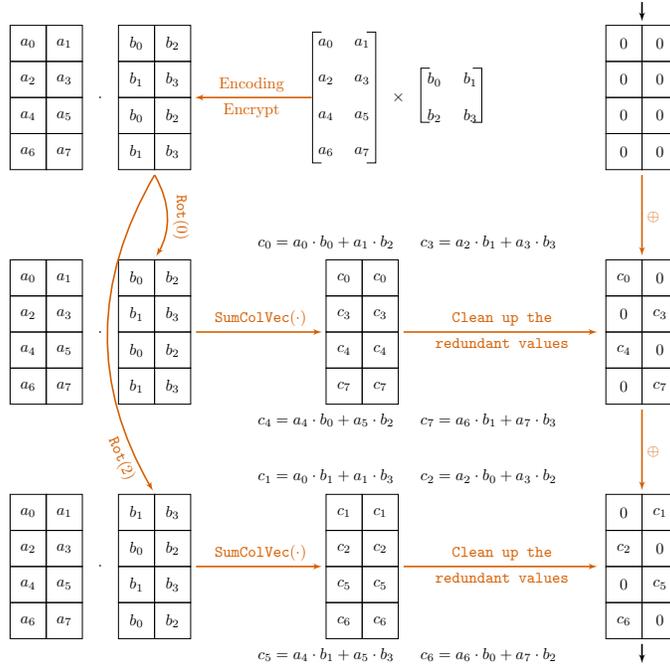


Figure 2: The matrix multiplication algorithm of Volley Revolver for the 4×2 matrix A and the matrix B of size 2×2

180 However, we don't need to stick to encrypting the transpose of the second matrix. Instead, either of
 181 the two matrices is transposed would do the trick: we could also encrypt the transpose of the first
 182 matrix, and the corresponding multiplication algorithm due to this change is similar to the Algorithm
 183 2 from [4].

184 Also, if each of the two matrices are too large to be encrypted into a single ciphertext, we could also
 185 encrypt the two matrices into two teams A and B of multiple ciphertexts. In this case, we can see this
 186 encoding method as Double Volley Revolver, which has two loops: the outside loop deals with
 187 the calculations between ciphertexts from two teams while the inside loop literally calculates two
 188 sub-matrices encrypted by two ciphertexts $A_{[i]}$ and $B_{[j]}$ using the raw algorithm of Volley Revolver.

189 4 Privacy-preserving CNN Training

190 4.1 Polynomial Approximation

191 Although Algorithm 1 enables us to avoid computing the Softmax function in the encryption domain,
 192 we still need to calculate the Sigmoid function using HE technique. This problem has been well
 193 studied by several works and we adopt a simple one [19], that is (1) we first use the least-square method
 194 to perfectly approximate the sigmoid function over the range $[-8, +8]$, obtaining a polynomial Z_{11}
 195 of degree 11; and (2) we use a polynomial Z_3 of degree 3 to approximate the Sigmoid by minimizing
 196 the cost function F including the squared gradient difference:

$$F = \lambda_0 \cdot \int_{-8}^{+8} (Z_{11} - Z_3)^2 dx + \lambda_1 \cdot \int_{-8}^{+8} (Z'_{11} - Z'_3)^2 dx,$$

197 where λ_0 and λ_1 are two positive float numbers to control the shape of the polynomial to approximate.

198 Setting $\lambda_0 = 128$ and $\lambda_1 = 1$ would result in the polynomial we used in our privacy-preserving CNN
 199 training: $Z_3 = 0.5 + 0.106795345032 \cdot x - 0.000385032598 \cdot x^3$.

200 4.2 Homomorphic Evaluation

201 Before the homomorphic CNN training starts, the client needs to encrypt the dataset X , the data
202 labels \bar{Y} , the matrix \bar{B} and the weight W into ciphertexts $Enc(X)$, $Enc(\bar{Y})$, $Enc(\bar{B})$ and $Enc(W)$,
203 respectively, and upload them to the cloud. For simplicity in presentation, we can just regard
204 the whole pipeline of homomorphic evaluation of Algorithm 1 as updating the weight ciphertext:
205 $W = W + B \odot (\bar{Y} - Z_3(X \times W^\top))^\top \times X$, regardless of the subtle control of the enhanced NAG
206 method with the SLE loss function.

207 Since `Volley Revolver` only needs one of the two matrices to be transposed ahead before en-
208 cryption and $(\bar{Y} - Z_3(X \times W^\top))^\top \times X$ happened to suffice this situation between any matrix
209 multiplication, we can complete the homomorphic evaluation of CQG for MLR.

210 5 Experiments

211 The C++ source code to implement the experiments in this section is openly available at:
212 <https://anonymous.4open.science/r/HE-CNNtraining-B355/>.

213 **Implementation** We implement the enhanced NAG with the SLE loss function based on HE with
214 the library HEAAN. All the experiments on the ciphertexts were conducted on a public cloud with 64
215 vCPUs and 192 GB RAM.

216 We adopt the first 128 MNIST training images as the training data and the whole test dataset as the
217 testing data. Both the training images and testing images have been processed in advance with the
218 pre-trained model `REGNET_X_400MF`, resulting in a new dataset with each example of size 401.

219 5.1 Parameters

220 The parameters of HEAAN we selected are: $\log N = 16$, $\log Q = 990$, $\log p = 45$, $slots = 32768$,
221 which ensure the security level $\lambda = 128$. Refer [6] for the details of these parameters. We didn't
222 use bootstrapping to refresh the weight ciphertexts and thus it can only perform 2 iterations of our
223 algorithm. Each iteration takes ~ 11 mins. The maximum runtime memory in this case is ~ 18 GB.
224 The 128 MNIST training images are encrypted into 2 ciphertexts. The client who own the private data
225 has to upload these two ciphertexts, two ciphertexts encrypting the one-hot labels \bar{Y} , one ciphertext
226 encrypting the \bar{B} and one ciphertext encrypting the weight W to the cloud. The initial weight matrix
227 W_0 we adopted is the zero matrix. The resulting MLR model after 2-iteration training has reached a
228 precision of 21.49% and obtain the loss of -147206 , which are consistent with the Python simulation
229 experiment.

230 6 Conclusion

231 In this work, we initiated to implement privacy-persevering CNN training based on mere HE tech-
232 niques by presenting a faster HE-friendly algorithm.

233 The HE operation bootstrapping could be adopted to refresh the weight ciphertexts. Python exper-
234 iments imitating the privacy-preserving CNN training using Z_3 as Sigmoid substitution showed
235 that using a large amount of data such as 8,192 images to train the MLE model for hundreds of
236 iterations would finally reach 95% precision. The real experiments over ciphertexts conducted on a
237 high-performance cloud with many vCPUs would take weeks to complete this test, if not months.

238 References

- 239 [1] Ran Gilad-Bachrach, Nathan Dowlin, Kim Laine, Kristin Lauter, Michael Naehrig, and John
240 Wernsing. Cryptonets: Applying neural networks to encrypted data with high throughput and
241 accuracy. In *International conference on machine learning*, pages 201–210. PMLR, 2016.
- 242 [2] Hervé Chabanne, Amaury De Wargny, Jonathan Milgram, Constance Morel, and Emmanuel
243 Prouff. Privacy-preserving classification on deep neural network. *Cryptology ePrint Archive*,
244 2017.

- 245 [3] Xiaoqian Jiang, Miran Kim, Kristin Lauter, and Yongsoo Song. Secure outsourced matrix
246 computation and application to neural networks. In *Proceedings of the 2018 ACM SIGSAC*
247 *Conference on Computer and Communications Security*, pages 1209–1222, 2018.
- 248 [4] John Chiang. A novel matrix-encoding method for privacy-preserving neural networks (infer-
249 ence). *arXiv preprint arXiv:2201.12577*, 2022.
- 250 [5] Florian Bourse, Michele Minelli, Matthias Minihold, and Pascal Paillier. Fast homomorphic
251 evaluation of deep discretized neural networks. In *Advances in Cryptology–CRYPTO 2018:*
252 *38th Annual International Cryptology Conference, Santa Barbara, CA, USA, August 19–23,*
253 *2018, Proceedings, Part III 38*, pages 483–512. Springer, 2018.
- 254 [6] Andrey Kim, Yongsoo Song, Miran Kim, Keewoo Lee, and Jung Hee Cheon. Logistic regression
255 model training based on the approximate homomorphic encryption. *BMC medical genomics*,
256 11(4):83, 2018.
- 257 [7] Charlotte Bonte and Frederik Vercauteren. Privacy-preserving logistic regression training. *BMC*
258 *medical genomics*, 11(4):86, 2018.
- 259 [8] Miran Kim, Yongsoo Song, Shuang Wang, Yuhou Xia, and Xiaoqian Jiang. Secure logistic
260 regression based on homomorphic encryption: Design and evaluation. *JMIR medical informatics*,
261 6(2):e19, 2018.
- 262 [9] John Chiang. Privacy-preserving logistic regression training with a faster gradient variant. *arXiv*
263 *preprint arXiv:2201.10838*, 2022.
- 264 [10] Craig Gentry. Fully homomorphic encryption using ideal lattices. In *Proceedings of the*
265 *forty-first annual ACM symposium on Theory of computing*, pages 169–178, 2009.
- 266 [11] Zvika Brakerski, Craig Gentry, and Vinod Vaikuntanathan. (leveled) fully homomorphic
267 encryption without bootstrapping. *ACM Transactions on Computation Theory (TOCT)*, 6(3):1–
268 36, 2014.
- 269 [12] N.P. Smart and F. Vercauteren. Fully homomorphic simd operations. Cryptology ePrint Archive,
270 Report 2011/133, 2011. <https://ia.cr/2011/133>.
- 271 [13] Jung Hee Cheon, Andrey Kim, Miran Kim, and Yongsoo Song. Homomorphic encryption for
272 arithmetic of approximate numbers. In *International Conference on the Theory and Application*
273 *of Cryptology and Information Security*, pages 409–437. Springer, 2017.
- 274 [14] Kyoohyung Han, Seungwan Hong, Jung Hee Cheon, and Daejun Park. Logistic regression on
275 homomorphic encrypted data at scale. In *Proceedings of the AAAI Conference on Artificial*
276 *Intelligence*, volume 33, pages 9466–9471, 2019.
- 277 [15] John Chiang. Multinomial logistic regression algorithms via quadratic gradient, 2023.
- 278 [16] John Chiang. Quadratic gradient: Uniting gradient algorithm and newton method as one. *arXiv*
279 *preprint arXiv:2209.03282*, 2022.
- 280 [17] Dankmar Böhning and Bruce G Lindsay. Monotonicity of quadratic-approximation algorithms.
281 *Annals of the Institute of Statistical Mathematics*, 40(4):641–663, 1988.
- 282 [18] Dankmar Böhning. Multinomial logistic regression algorithm. *Annals of the institute of*
283 *Statistical Mathematics*, 44(1):197–200, 1992.
- 284 [19] John Chiang. On polynomial approximation of activation function. *arXiv preprint*
285 *arXiv:2202.00004*, 2022.