

# Privacy-Preserving Deep Learning via Additively Homomorphic Encryption<sup>\*</sup>

Le Trieu Phong<sup>1</sup>, Yoshinori Aono<sup>1</sup>, Takuya Hayashi<sup>1,2</sup>, Lihua Wang<sup>1</sup>, and Shiho Moriai<sup>1</sup>

<sup>1</sup> National Institute of Information and Communications Technology (NICT), Japan

<sup>2</sup> Kobe University, Japan

{phong, aono, wlh, shiho.moriai}@nict.go.jp  
t-hayashi@eedept.kobe-u.ac.jp

**Abstract.** We build a privacy-preserving deep learning system in which many learning participants perform neural network-based deep learning over a combined dataset of all, without actually revealing the participants' local data. To that end, we revisit the previous work by Shokri and Shmatikov (ACM CCS 2015) and point out that local data information may be actually leaked to an honest-but-curious server. We then move on to fix that problem via building an enhanced system with following properties: (1) no information is leaked to the server; and (2) accuracy is kept intact, compared to that of the ordinary deep learning system also over the combined dataset.

Our system is a bridge between deep learning and cryptography: we utilise stochastic gradient descent (SGD) applied to neural networks, in combination with additively homomorphic encryption. We show that our usage of encryption adds tolerable overhead to the ordinary deep learning system.

**keywords:** Privacy, deep learning, neural network, additively homomorphic encryption, LWE-based encryption, Paillier encryption.

## 1 Introduction

### 1.1 Background

In recent years, *deep learning* (aka, *deep machine learning*) has produced exciting results in both academia and industry, in which deep learning systems are approaching or even surpassing human-level accuracy. This is thanks to algorithmic breakthroughs and physical parallel hardware applied to *neural networks* when processing massive amount of data.

Massive collection of data, while vital for deep learning, raises the issue of privacy. Individually, a collected photo can be permanently kept on a server of a company, out of the control of the photo's owner. At law, privacy and confidentiality worries may prevent hospitals and research centers from sharing their medical datasets, barring them from enjoying the advantage of large-scale deep learning over the joint datasets.

As a directly related work, Shokri and Shmatikov (ACM CCS 2015) [19] presented a system for privacy-preserving deep learning, allowing local datasets of several participants staying home while the learned model for the neural network over the joint dataset can be obtained by the participants. To achieve the result, the system in [19] needs the following: each learning participant, using local data, first computes gradients of a neural network; then a part (e.g. 1% ~ 10%) of those gradients must be sent to a parameter cloud server. The server is *honest-but-curious*: it is assumed to be *curious* in extracting the data of individuals; and yet, it is assumed to be *honest* in operations.

To protect privacy, the system of Shokri and Shmatikov admits an accuracy/privacy tradeoff (see Table 1): sharing no local gradients leads to perfect privacy but not desirable accuracy; on the other hand, sharing all local gradients violates privacy but leads to good accuracy. To compromise, sharing a part of local gradients is the main solution in [19] to keep as less accuracy decline as possible.

---

<sup>\*</sup> Full version of a paper at the 8-th International Conference on Applications and Techniques in Information Security (ATIS 2017) [17].

**Table 1.** Comparison of techniques.

System	Method to protect gradients (against the curious server)	Potential information leakage to server?
Shokri-Shmatikov [19]	Partial sharing [19][Sect. 5] (and Laplace noises [19][Sect. 7])	yes
Ours (Section 4)	Additively homomorphic encryption	no

## 1.2 Our contributions

We demonstrate that, in the system of Shokri and Shmatikov [19], even a small portion of the gradients stored over the cloud server can be exploited: namely, local data can be unwillingly extracted from those gradients. Illustratively, we show in Section 3 a few examples on how a small fraction of gradients leaks useful information on data.

We then propose a novel system for deep learning to protect the gradients over the *honest-but-curious* cloud server, using additively homomorphic encryption. All gradients are encrypted and stored on the cloud server. The additive homomorphic property enables the computation over the gradients. Our system is described in Section 4, and depicted in Figure 4, enjoying following properties on security and accuracy:

**Security.** *Our system leaks no information of participants to the honest-but-curious parameter (cloud) server.*

**Accuracy.** *Our system achieves identical accuracy to a corresponding deep learning system (Asynchronous SGD, see below) trained over the joint dataset of all participants.*

In short, our system enjoys the best of both worlds: security as in cryptography, and accuracy as in deep learning. See Theorem 1 and Theorem 2 in Section 4.

**Our tradeoff.** Protecting the gradients against the cloud server comes with the cost of increased communication between the learning participants and the cloud server. We show in Table 2 that the increased factors are not big: less than 3 for concrete datasets MNIST [3] and SVHN [15]. For example, in the case of MNIST, if each learning participant needs to communicate 0.56 MB<sup>3</sup> of plain gradients to the server at each upload or download; then in our system with LWE-based encryption, the corresponding communication cost at each upload or download becomes

$$2.47 \text{ (Table 2's factor)} \times 0.437 \text{ (original MB)} \approx 1 \text{ MB}$$

which needs around 8 milliseconds to be transmitted over a 1 Gbps channel. Technical details are in Sections 5 and 6.

On the computational side, we estimate that our system finishes in around 2.25 hours to obtain around 97.14% accuracy when training and testing over the MNIST dataset.

**Discussion on the tradeoffs.** With our system using additively homomorphic encryption against the curious server, we show that the trade-off between *accuracy/privacy* in [19] can be shifted to *efficiency/privacy* in ours. The accuracy/privacy tradeoff of [19] may make privacy-preserving deep learning less attractive compared to ordinary deep learning, as accuracy is the main appeal in the field. Our efficiency/privacy tradeoff, keeping ordinary deep learning accuracy intact, can be solved if more processing units and more dedicated programming codes are employed.

<sup>3</sup> Size of 109386 gradients each of 32 bits; the size is computed via the formula  $\frac{109386 \times 32}{8 \times 10^6} \approx 0.437$ .

### 1.3 Technical overviews

A succinct comparison is in Table 1. Below we present the underlying technicalities.

**Asynchronous SGD (ASGD) [8,18], no privacy protection.** Both our system and that of [19] rely on the fact that neural networks can be trained via a variant of SGD called asynchronous SGD [8,18] with *data parallelism* and *model parallelism*. Specifically, first a global weight vector  $W_{\text{global}}$  for the neural network is initialised randomly. Then, at each iteration, replicas of the neural network are run over local datasets (data parallelism), and the corresponding local gradient vector  $G_{\text{local}}$  is sent to a cloud server. For each  $G_{\text{local}}$ , the cloud server then updates the global parameters as follows:

$$W_{\text{global}} := W_{\text{global}} - \alpha \cdot G_{\text{local}} \quad (1)$$

where  $\alpha$  is a learning rate. The updated global parameters  $W_{\text{global}}$  are broadcasted to all the replicas, who then use them to replace their old weight parameters. The process of updating and broadcasting  $W_{\text{global}}$  is repeated until a desired minimum for a pre-defined cost function (based on cross-entropy or squared-error) is reached. For model parallelism, the update at (1) is computed in parallel via components of the vectors  $W_{\text{global}}$  and  $G_{\text{local}}$ .

**Shokri-Shmatikov systems.** The system in [19][Sect. 5] can be called as *gradients-selective* ASGD, by following reasons. In [19][Sect. 5], the update rule at (1) is modified as follows:

$$W_{\text{global}} := W_{\text{global}} - \alpha \cdot G_{\text{local}}^{\text{selective}} \quad (2)$$

in which vector  $G_{\text{local}}^{\text{selective}}$  contains selective (say 1%  $\sim$  10%) gradients of  $G_{\text{local}}$ . The update using (2) allows each participant to choose which gradients to share globally, with the hope of reducing the risk of leaking sensitive information on the participant’s local dataset to the cloud server. However, showed in Section 3, even a small part of gradients leaks information to the server.

In [19][Sect. 7], Shokri-Shmatikov showed an additional technique on using differential privacy to counter-measure indirect leakage from gradients. Their strategy is to add Laplace noises into  $G_{\text{local}}^{\text{selective}}$  at (2). Due to noises, this method declines the learning accuracy which is the main appeal in deep learning.

**Our system.** The system we designed can be called as *gradients-encrypted* ASGD, by following reasons. In our system in Section 4, we make use of the following update fomular

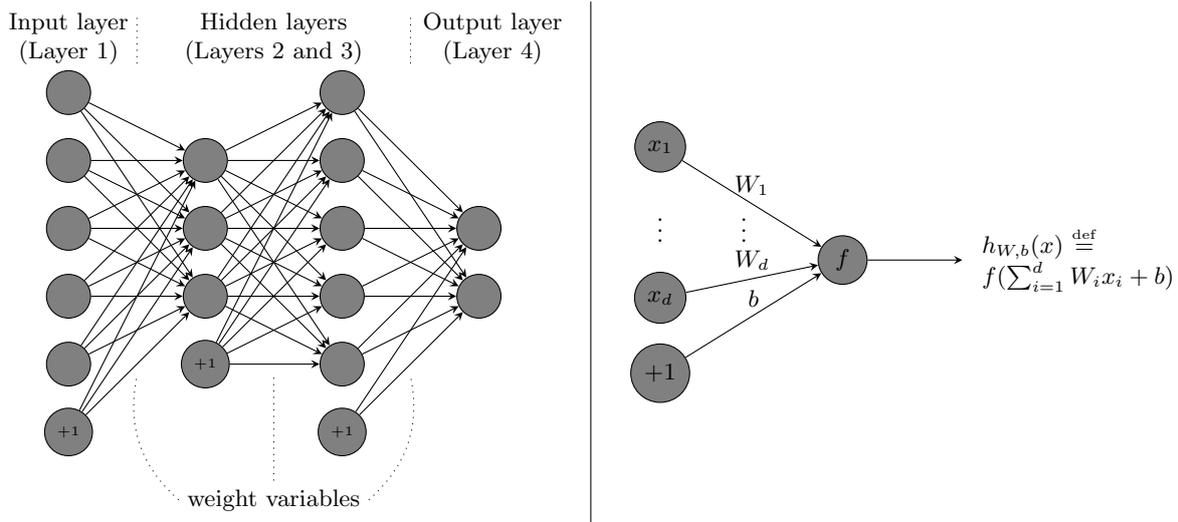
$$\mathbf{E}(W_{\text{global}}) := \mathbf{E}(W_{\text{global}}) + \mathbf{E}(-\alpha \cdot G_{\text{local}}) \quad (3)$$

in which  $\mathbf{E}$  is homomorphic encryption supporting addition over ciphertexts. The decryption key is only known to the participants and not to the cloud server. Therefore, the honest-but-curious cloud server knows nothing about each  $G_{\text{local}}$ , and hence obtains no information on each local dataset of participants. Nonetheless, as

$$\mathbf{E}(W_{\text{global}}) + \mathbf{E}(-\alpha \cdot G_{\text{local}}) = \mathbf{E}(W_{\text{global}} - \alpha \cdot G_{\text{local}})$$

by the additively homomorphic property of  $\mathbf{E}$ , each participant will get the correctly updated  $W_{\text{global}}$  by decryption. Moreover, when the original update at (1) is parallelised in components of the vectors  $W_{\text{global}}$  and  $G_{\text{local}}$ , our system applies homomorphic encryption to each of the components accordingly.

In addition, to ensure the integrity of the homomorphic ciphertexts, each client will use a secure channel such as TLS/SSL (distinct from each other) to communicate the homomorphic ciphertexts with the server.



**Fig. 1.** (left) a neural network with 5 inputs, 2 hidden layers, 2 outputs; (right) a network with one neuron.

## 1.4 More related works

Gilad-Bachrach et al. [10] present a system called *CryptoNets*, which allows homomorphically encrypted data feedforwarding an already-trained neural network. Because *CryptoNets* assumes that the weights in the neural network have been trained beforehand, the system aims at making prediction for individual data item.

The goal of our paper and [19] differ from that of [10], as our system and Shokri-Shmatikov’s exactly aim at training the weights utilising multiple data sources, while *CryptoNets* [10] does not.

We consider the cloud server as the adversary in this paper while learning participants are seen as honest entities. This is because in our scenario learning participants are considered as organisations such as financial institutions or hospitals acting with responsibilities by laws. Our scenario and adversary model is different from Hitaj et al. [12] which examines dishonest learning participants.

## 2 Preliminaries

### 2.1 On additively homomorphic encryption

**Definition 1. (Homomorphic encryption)** *Public key additively homomorphic encryption (PHE) schemes consist of the following (possibly probabilistic) poly-time algorithms.*

- $\text{ParamGen}(1^\lambda) \rightarrow pp$ :  $\lambda$  is the security parameter and the public parameter  $pp$  is implicitly fed in following algorithms.
- $\text{KeyGen}(1^\lambda) \rightarrow (pk, sk)$ :  $pk$  is the public key, while  $sk$  is the secret key.
- $\text{Enc}(pk, m) \rightarrow c$ : probabilistic encryption algorithm produces  $c$ , the ciphertext of message  $m$ .
- $\text{Dec}(sk, c) \rightarrow m$ : decryption algorithm returns message  $m$  encrypted in  $c$ .
- $\text{Add}(c, c')$ : for ciphertexts  $c$  and  $c'$ , the output is the encryption of plaintext addition  $c_{\text{add}}$ .
- $\text{DecA}(sk, c_{\text{add}})$ : decrypting  $c_{\text{add}}$  to obtain an addition of plaintexts.

Ciphertext indistinguishability against chosen plaintext attacks [11] (or CPA security for short below) ensures that no bit of information is leaked from ciphertexts.

## 2.2 On deep machine learning

**Some concepts and notations.** Deep machine learning can be seen as a set of techniques applied to neural networks. In Figure 1 is a neural network with 5 inputs, 2 hidden layers, and 2 outputs. The node with +1 represents the bias term. The neuron nodes are connected via weight variables. In a deep learning structure of neural network, there can be multiple layers each with thousands of neurons.

Each neuron node (except the bias node) is associated with an *activation function*  $f$ . Examples of  $f$  in deep learning are  $f(z) = \max\{0, z\}$  (rectified linear),  $f(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$  (hyperbolic tangent), and  $f(z) = (1 + e^{-z})^{-1}$  (sigmoid). The output at layer  $l + 1$ , denoted as  $a^{(l+1)}$ , is computed as  $a^{(l+1)} = f(W^{(l)}a^{(l)} + b^{(l)})$  in which  $(W^{(l)}, b^{(l)})$  is the weights connecting layers  $l$  and  $l + 1$ , and  $a^{(l)}$  is the output at layer  $l$ .

The learning task is, given a training dataset, to determine these weight variables to minimise a pre-defined cost function such as the cross-entropy or the squared-error cost function [2]. The cost function can be computed over all data items in the training dataset; or over a subset (called mini-batch) of  $t$  elements from the training dataset. Denote the cost function for the latter case as  $J_{|\text{batch}|=t}$ . In the extreme case of  $t = 1$ , corresponding to maximum stochasticity,  $J_{|\text{batch}|=1}$  is the cost function defined over 1 single data item.

**Stochastic gradient descent (SGD).** Let  $W$  be the flattened vector consisting all weight variables, namely we take all weights in the neural network and arrange them consecutively to form the vector  $W$ . Denote  $W = (W_1, \dots, W_{n_{gd}}) \in \mathbb{R}^{n_{gd}}$ . Let

$$G = \left( \frac{\delta J_{|\text{batch}|=t}}{\delta W_1}, \dots, \frac{\delta J_{|\text{batch}|=t}}{\delta W_{n_{gd}}} \right) \quad (4)$$

be the gradients of the cost function  $J_{|\text{batch}|=t}$  corresponding to variables  $W_1, \dots, W_{n_{gd}}$ . The variable update rule in SGD is as follows, for a learning rate  $\alpha \in \mathbb{R}$ :

$$W := W - \alpha \cdot G \quad (5)$$

in which  $\alpha \cdot G$  is component-wise multiplication, namely  $\alpha \cdot G = (\alpha G_1, \dots, \alpha G_{n_{gd}}) \in \mathbb{R}^{n_{gd}}$ . The learning rate  $\alpha$  can also be changed adaptively as described in [2, 9].

**Asynchronous (aka. Downpour) SGD [8, 18].** By (4) and (5), as long as the gradients  $G$  can be computed, the weights  $W$  can be updated. Therefore, the data used in computing  $G$  can be distributed (namely, data parallelism). Moreover, the update process can be parallelised by considering separate components of the vectors (namely, model parallelism).

Specifically, asynchronous SGD uses multiple replicas of a neural network. Before each execution, each replica will download the newest weights from the parameter server; and each replica is run over a data shard, which is a subset of the training dataset. To use the power of parallel computation when the server has multiple processing units  $\text{PU}_1, \dots, \text{PU}_{n_{pu}}$ , asynchronous SGD splits the weight vector  $W$  and gradient vector  $G$  into  $n_{pu}$  parts, namely  $W = (W^{(1)}, \dots, W^{(n_{pu})})$  and  $G = (G^{(1)}, \dots, G^{(n_{pu})})$ , so that the update rule at (5) becomes as follows:

$$W^{(i)} := W^{(i)} - \alpha \cdot G^{(i)} \quad (6)$$

which is computed at processing unit  $\text{PU}_i$ . As the processing units  $\text{PU}_1, \dots, \text{PU}_{n_{pu}}$  can run in parallel, asynchronous SGD significantly increases the scale and speed of deep network training, as experimentally shown in [8].

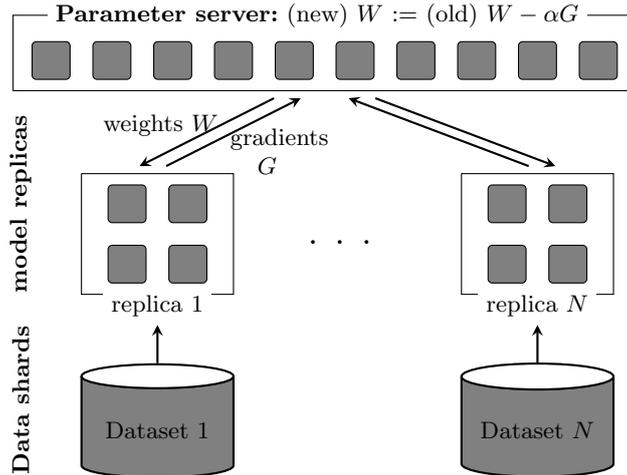


Fig. 2. Asynchronous SGD [8, 18].

### 3 Gradients leak information

This section shows that a small portion of gradients may reveal information on local data.

**Example 1 (one neuron).** For illustration of how gradients leak information on data, we first use the neural network in Figure 1, with only one neuron. In the figure, real numbers  $x_i$  ( $1 \leq i \leq d$ ) are the input data, with a corresponding truth label  $y$ ; real numbers  $W_i$  ( $1 \leq i \leq d$ ) are the weight parameters to be learned; and  $b$  is the bias. The function  $f$  is an activation function (either sigmoid, rectified linear, or hyperbolic tangent as described in Section 2.2). The cost function is defined as the distance between the predicted value  $h_{W,b}(x) \stackrel{\text{def}}{=} f(\sum_{i=1}^d W_i x_i + b)$  and the truth value  $y$ :

$$J(W, b, x, y) \stackrel{\text{def}}{=} (h_{W,b}(x) - y)^2$$

and hence the gradients are

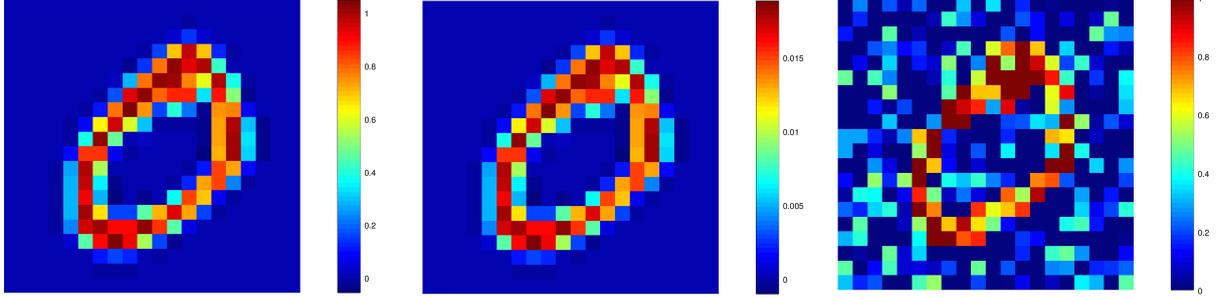
$$\begin{aligned} \eta_k &\stackrel{\text{def}}{=} \frac{\delta J(W, b, x, y)}{\delta W_k} = 2(h_{W,b}(x) - y) \frac{\delta h_{W,b}(x)}{\delta W_k} = 2(h_{W,b}(x) - y) \frac{\delta f(\sum_{i=1}^d W_i x_i + b)}{\delta W_k} \\ &= 2(h_{W,b}(x) - y) f'(\sum_{i=1}^d W_i x_i + b) \cdot x_k \end{aligned} \quad (7)$$

and

$$\begin{aligned} \eta &\stackrel{\text{def}}{=} \frac{\delta J(W, b, x, y)}{\delta b} = 2(h_{W,b}(x) - y) \frac{\delta h_{W,b}(x)}{\delta b} = 2(h_{W,b}(x) - y) \frac{\delta f(\sum_{i=1}^d W_i x_i + b)}{\delta b} \\ &= 2(h_{W,b}(x) - y) f'(\sum_{i=1}^d W_i x_i + b) \cdot 1. \end{aligned} \quad (8)$$

The  $k$ -th component  $x_k$  of  $x = (x_1, \dots, x_d) \in \mathbb{R}^d$  or the truth label  $y$  can be inferred from the gradients by one of the following means:

- (O1) Observe that  $\eta_k/\eta = x_k$ . Therefore,  $x_k$  is completely leaked if  $\eta_k$  and  $\eta$  are shared to the cloud server. For example, if 1% of local gradients, chosen randomly as suggested in [19], are shared to the server, then the probability that both  $\eta_k$  and  $\eta$  are shared is  $(1/100) \times (1/100) = 1/10^4$ , which is not negligible.



(a) Original 20x20 image of handwritten number 0, seen as a vector over  $\mathbb{R}^{400}$  fed to a neural network.

(b) Recovered image using 400/10285 (3.89%) gradients (see Sect.3, Example 2). The difference with the original (a) is only at the value bar.

(c) Recovered image using 400/10285 (3.89%) gradients (see Sect.3, Example 3). There are noises but the truth label 0 can still be seen.

**Fig. 3.** Original data (a) vs. leakage information (b), (c) from a small part of gradients in a neural network.

(O2) Observe that the gradient  $\eta_k$  is proportional to the input  $x_k$  for all  $1 \leq i \leq d$ . Therefore, when  $x = (x_1, \dots, x_d)$  is an image, one can use the gradients to produce a related “proportional” image, and then obtain the truth value  $y$  by guessing.

**Example 2 (general neural networks, cf. Fig.3(b)).** The above observations (O1) and (O2) similarly hold for general neural networks, with both cross-entropy and squared-error cost functions [2]. In particular, following [2],

$$\eta_{ik} \stackrel{\text{def}}{=} \frac{\delta J(W, b, x, y)}{\delta W_{ik}^{(1)}} = \xi_i \cdot x_k \quad (9)$$

where  $W_{ik}^{(1)}$  is the weight parameter connecting layer 1’s input  $x_k$  with hidden node  $i$  of layer 2;  $\xi_i$  is a real number.

In Figure 3(b), using a neural network on [1], we demonstrate that gradients at (9) are indeed *proportional* to the original data, as Figure 3(b) only differs from Figure 3(a) at the value bar. The original data is a 20x20 image, reshaped into a vector of  $(x_1, \dots, x_{400}) \in \mathbb{R}^{400}$ . The vector is an input to a neural network of 1 hidden layer of 25 nodes; and the output layer contains 10 nodes. The total number of gradients in the neural network is  $(400 + 1) \times 25 + (25 + 1) \times 10 = 10285$ . At (9), we have  $1 \leq k \leq 400$  and  $1 \leq i \leq 25$ . We then use a small part of the gradients at (9), namely  $(\eta_{1k})_{1 \leq k \leq 400}$ , reshaped into a 20x20 image, to draw Figure 3(b). It is clear that the part (namely  $400/10285 \approx 3.89\%$ ) of the gradients reveals the truth label 0 of the original data.

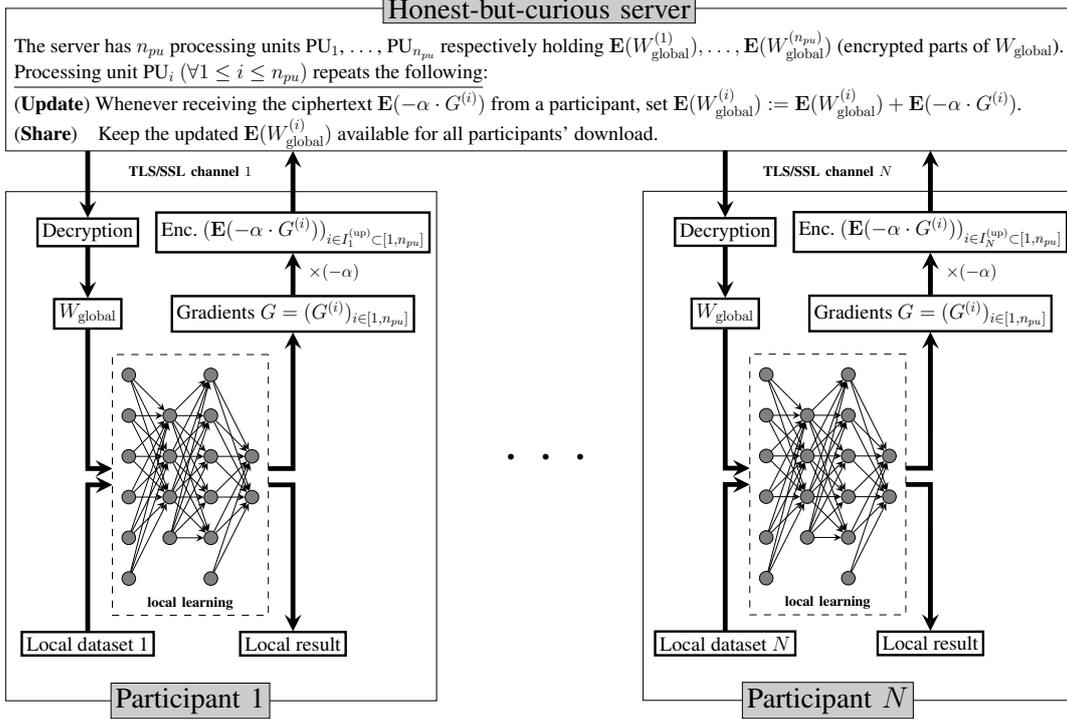
**Example 3 (general neural networks, with regularization, cf. Fig.3(c)).** In a neural network with regularization, following [2] we have

$$\begin{aligned} \eta_{ik} &\stackrel{\text{def}}{=} \frac{\delta J(W, b, x, y)}{\delta W_{ik}^{(1)}} = \xi_i \cdot x_k + \lambda W_{ik}^{(1)} \\ \eta_i &\stackrel{\text{def}}{=} \frac{\delta J(W, b, x, y)}{\delta b_i^{(1)}} = \xi_i \end{aligned}$$

where notations are as in Example 2 above,  $b_i^{(1)}$  is the bias associated with node  $i$  of layer 2; and  $\lambda \geq 0$  is a regularization term.

As in observation (O1), both  $\eta_{ik}$  and  $\eta_i$  is known to server with a non-negligible probability. Additionally, in Figure 3(c), we use the following observation:

$$\frac{\eta_{ik}}{\eta_i} = x_k + \frac{\lambda W_{ik}^{(1)}}{\xi_i}$$



**Fig. 4.** Our system (*gradients-encrypted Asynchronous SGD*) for privacy-preserving deep learning, with a cloud server and  $N$  participants.

which is an approximation of the data  $x_k$ . In Figure 3(c), we take  $\lambda = 0.1$ , and other details are identical to Example 2 above. Due to the term  $\lambda W_{ik}^{(1)}/\xi_i$ , there are noises in the figure, but the truth value (number 0) of the original data can still be seen.

#### 4 Our system: privacy-preserving deep learning without accuracy decline

Our system is depicted in Figure 4, consisting of a common cloud server and  $N$  (e.g. = 10x) learning participants.

**Learning Participants.** The participants jointly set up the public key  $pk$  and secret key  $sk$  for an additively homomorphic encryption scheme. The secret key  $sk$  is kept confidential against the cloud server, but is known to all learning participants. Each participant will establish a TLS/SSL secure channel, different from each other, to communicate and protect the integrity of the homomorphic ciphertexts.

Then, the participants locally hold their datasets and run replicas of a deep learning based neural network. The initial (random) weight  $W_{\text{global}}$  to run the local deep learning is initialised by Participant 1 who also sends  $\mathbf{E}(W_{\text{global}}^{(1)}), \dots, \mathbf{E}(W_{\text{global}}^{(n_{pu})})$  to the server initially, in which  $W_{\text{global}}^{(i)}$  is also a vector constituting the  $i$ -th part of  $W_{\text{global}}$ . The gradient vector  $G$  obtained after each execution of the neural network is split into  $n_{pu}$  parts, namely  $G = (G^{(1)}, \dots, G^{(n_{pu})})$ , multiplied by the learning rate  $\alpha$ , and then encrypted using the public key  $pk$ . The resulting encryption  $\mathbf{E}(-\alpha \cdot G^{(i)})$  ( $\forall 1 \leq i \leq n_{pu}$ ) from each learning participant is sent to the processing unit  $PU_i$  of the server. It is also worth noting that the learning rate  $\alpha$  can be adaptively changed locally at each learning participant as described in [8].

As seen in Figure 4, each participant  $1 \leq k \leq N$  will perform the following steps:

1. Download the ciphertexts  $\mathbf{E}(W_{\text{global}}^{(j)})$  stored at the processing units  $PU_j$  of the server for all  $j \in I_k^{(\text{down})} \subset [1, n_{pu}]$ . Usually  $I_k^{(\text{down})} = [1, n_{pu}]$  namely the participant will download all

encrypted parts of the global weight, but it is possible that  $I_k^{(\text{down})} \not\subseteq [1, n_{pu}]$  if the learning participant has restricted download bandwidth.

2. Decrypt the above ciphertexts using the secret key  $sk$  to obtain  $W_{\text{global}}^{(j)}$  for all  $j \in I_k^{(\text{down})}$ , and replacing those values into the corresponding places of  $W_{\text{global}} = (W_{\text{global}}^{(1)}, \dots, W_{\text{global}}^{(n_{pu})})$ .
3. Get a mini-batch of data from its local dataset.
4. Using the values of  $W_{\text{global}}$  and data items at steps **2** and **3**, compute the gradients  $G = (G^{(1)}, \dots, G^{(n_{pu})})$  with respect to variable  $W_{\text{global}}$ .
5. Encrypt and send back the ciphertexts  $\mathbf{E}(-\alpha \cdot G^{(i)}) \forall i \in I_k^{(\text{up})} \subset [1, n_{pu}]$  to the corresponding processing unit  $\text{PU}_i$  of the server. The subset for uploading  $I_k^{(\text{up})} \subset [1, n_{pu}] = \{1, \dots, n_{pu}\}$  depends on the choice of participant  $k$ . For a full upload,  $I_k^{(\text{up})} = [1, n_{pu}]$  so that all encrypted gradients are uploaded to the server.

The downloads and uploads of the encrypted parts of  $W_{\text{global}}$  can be *asynchronous* in two aspects: the participants are independent with each other; and the processing units are also independent with each other.

**Cloud Server.** The cloud server is a common place to recursively update the encrypted weight parameters. In particular, each processing unit  $\text{PU}_i$  at the server, after receiving any encryption  $\mathbf{E}(\alpha \cdot G^{(i)})$ , computes

$$\mathbf{E}(W_{\text{global}}^{(i)}) + \mathbf{E}(-\alpha \cdot G^{(i)}) \left( \text{which is } = \mathbf{E}(W_{\text{global}}^{(i)} - \alpha \cdot G^{(i)}) \right)$$

where the equality is thanks to the additively homomorphic property of encryption. Therefore, the part  $W_{\text{global}}^{(i)}$  at the processing unit  $\text{PU}_i$  is updated to  $W_{\text{global}}^{(i)} - \alpha \cdot G^{(i)}$ , or notationally  $W_{\text{global}}^{(i)} := W_{\text{global}}^{(i)} - \alpha \cdot G^{(i)}$ .

**Theorem 1 (Security against the cloud server).** *Our system in Figure 4 leaks no information of the learning participants to the honest-but-curious cloud server, provided that the underlying homomorphic encryption scheme is CPA-secure.*

*Proof.* The participants only send encrypted gradients to the cloud server. Therefore, if the encryption scheme is CPA-secure, no bit of information on the data of the participants can be leaked.  $\square$

**Theorem 2 (Accuracy equivalence to asynchronous SGD).** *Our system in Figure 4, when  $I_k^{(\text{up})} = I_k^{(\text{down})} = [1, n_{pu}] \forall 1 \leq k \leq N$  and all ciphertexts are decrypted, functions as asynchronous SGD described in Section 2.2. Therefore, our system can achieve the same accuracy as that of asynchronous SGD.*

*Proof.* After decryption, the update rule of weight parameter becomes  $W_{\text{global}}^{(i)} := W_{\text{global}}^{(i)} - \alpha \cdot G^{(i)}$  for any  $1 \leq i \leq n_{pu}$  in which  $G^{(i)}$  is the gradient vector computed from data samples held by participant  $k$  (and the downloaded  $W_{\text{global}}$ ). Since the update rule is identical to (6) and each learning participant in our system functions as a replica (as in asynchronous SGD) when encryption is removed, the theorem follows.  $\square$

## 5 Instantiations of our system

In this section we use additively homomorphic encryption schemes to instantiate our system in Section 4. We use following schemes to show two instantiations of our system: LWE-based encryption (modern, potentially post-quantum), and Paillier encryption (classical, smaller key sizes, not post-quantum).

### 5.1 Using an LWE-based encryption

The mark  $\stackrel{\$}{\leftarrow}$  is for “sampling randomly from a discrete Gaussian” set, so that  $x \stackrel{\$}{\leftarrow} \mathbb{Z}_{(0,s)}$  means  $x$  appears with probability proportional to  $\exp(-\pi x^2/s^2)$ .

**LWE-based encryption.** We use an additively homomorphic variant of the public key encryption scheme in [13].

- **ParamGen**( $1^\lambda$ ): Fix  $q = q(\lambda) \in \mathbb{Z}^+$  and  $l \in \mathbb{Z}^+$ . Fix  $p \in \mathbb{Z}^+$  so that  $\gcd(p, q) = 1$ . Return  $pp = (q, l, p)$ .
- **KeyGen**( $1^\lambda, pp$ ): Take  $s = s(\lambda, pp) \in \mathbb{R}^+$  and  $n_{lwe} \in \mathbb{Z}^+$ . Take  $R, S \stackrel{\$}{\leftarrow} \mathbb{Z}_{(0,s)}^{n_{lwe} \times l}$ ,  $A \stackrel{\$}{\leftarrow} \mathbb{Z}_q^{n_{lwe} \times n_{lwe}}$ . Compute  $P = pR - AS \in \mathbb{Z}_q^{n_{lwe} \times l}$ . Return the public key  $pk = (A, P, n_{lwe}, s)$ , and the secret key  $sk = S$ .
- **Enc**( $pk, m \in \mathbb{Z}_p^{1 \times l}$ ): Take  $e_1, e_2 \stackrel{\$}{\leftarrow} \mathbb{Z}_{(0,s)}^{1 \times n_{lwe}}$ ,  $e_3 \stackrel{\$}{\leftarrow} \mathbb{Z}_{(0,s)}^{1 \times l}$ . Compute  $c_1 = e_1 A + p e_2 \in \mathbb{Z}_q^{1 \times n_{lwe}}$ ,  $c_2 = e_1 P + p e_3 + m \in \mathbb{Z}_q^{1 \times l}$ . Return  $c = (c_1, c_2)$ .
- **Dec**( $S, c = (c_1, c_2)$ ): Compute  $\bar{m} = c_1 S + c_2 \in \mathbb{Z}_q^{1 \times l}$ . Return  $m = \bar{m} \bmod p$ .
- **Add**( $c, c'$ ): For addition, compute and return  $c_{\text{add}} = c + c' \in \mathbb{Z}_q^{1 \times (n_{lwe} + l)}$ .

**Data encoding and encryption.** A real number  $a \in \mathbb{R}$  can be represented, with  $\text{prec}$  bits of precision, by an integer  $\lfloor a \cdot 2^{\text{prec}} \rfloor \in \mathbb{Z}$ . To realise the encryption  $\mathbf{E}(\cdot)$  in Figure 4, because both  $W_{\text{global}}^{(i)}$  and  $\alpha \cdot G^{(i)}$  are in the space  $\mathbb{R}^{L_i}$  where  $L_i$  is the length of the partition so that  $\sum_{i=1}^{n_{\text{pu}}} L_i = n_{\text{gd}}$ , it suffices to describe an encryption of a real vector  $r = (r^{(1)}, \dots, r^{(L_i)}) \in \mathbb{R}^{L_i}$ . The encryption is, for  $l = L_i$ ,

$$\mathbf{E}(r) = \text{lweEnc}_{pk} \left( \overbrace{\left( \underbrace{\lfloor r^{(1)} \cdot 2^{\text{prec}} \rfloor}_{\in \mathbb{Z}_p} \dots \underbrace{\lfloor r^{(L_i)} \cdot 2^{\text{prec}} \rfloor}_{\in \mathbb{Z}_p} \right)}^{\mathbb{Z}_p^{1 \times L_i}} \right). \quad (10)$$

For vectors  $r, t \in \mathbb{R}^{L_i}$ , the decryption of

$$\mathbf{E}(r) + \mathbf{E}(-t) \in \mathbb{Z}_q^{1 \times (n_{lwe} + L_i)} \quad (11)$$

will yield, for all  $1 \leq j \leq L_i$ ,

$$\lfloor r^{(j)} \cdot 2^{\text{prec}} \rfloor - \lfloor t^{(j)} \cdot 2^{\text{prec}} \rfloor \in \mathbb{Z}_p \subset (-p/2, p/2] \quad (12)$$

and hence

$$u^{(j)} = \lfloor r^{(j)} \cdot 2^{\text{prec}} \rfloor - \lfloor t^{(j)} \cdot 2^{\text{prec}} \rfloor \in \mathbb{Z} \quad (13)$$

if  $p/2$  is large enough (see below). The subtraction  $r^{(j)} - t^{(j)} \in \mathbb{R}$  is computed via  $u^{(j)}/2^{\text{prec}} \in \mathbb{R}$ , so that finally  $r - t \in \mathbb{R}^L$  is obtained after decryption as desired. To get (13) from (12), it suffices that  $p/2 > 2 \cdot 2^{\text{prec}}$ , as via normalisation, we can assume  $-1 < r^{(j)}, t^{(j)} < 1$ . In general, to handle  $n_{\text{gradupd}}$  additive terms without overflow, it is necessary that  $p/2 > n_{\text{gradupd}} \cdot 2^{\text{prec}}$ , or equivalently,

$$p > n_{\text{gradupd}} \cdot 2^{\text{prec}+1}. \quad (14)$$

**Lemma 1 (Choosing parameters).** *When  $n_{lwe} \geq 3000$ ,  $s = 8$ , it is possible to set*

$$\log_2 q \approx \log_2 p + \log_2 n_{\text{gradupd}} + \log_2(167.9\sqrt{n_{lwe}} + 33.9) + 1$$

in which  $p$  satisfies (14), and  $n_{\text{gradupd}}$  is the number of gradient updates at each processing unit of the cloud server in Figure 4. For example, when  $n_{lwe} = 3000$ ,  $p = 2^{48} + 1$ ,  $n_{\text{gradupd}} = 2^{15}$ , it is possible to set  $q = 2^{77}$ .

It is worth noting that the parameters in Lemma 1 is very conservative, due to the consideration of unlikely large noises in the proof. Therefore,  $n_{\text{gradupd}}$  can be larger than stated. Indeed, we experimentally check with  $q = 2^{77}$ ,  $n_{lwe} = 3000$ ,  $p = 2^{48} + 1$ ,  $s = 8$  and confirms that  $n_{\text{gradupd}}$  can be twice (namely  $n_{\text{gradupd}} = 2 \cdot 2^{15}$ ) as stated in Lemma 1 without any decryption error.

**Theorem 3 (Increased communication factor, LWE-based).** *The communication between the server and participants of our system is*

$$\frac{n_{pu}n_{lwe} \log_2 q}{n_{gd} \cdot \text{prec}} + \frac{\log_2 q}{\text{prec}}$$

time of the communication of the corresponding Asynchronous SGD, in which  $(n_{lwe}, p, q)$  is parameters of the encryption scheme,  $n_{gd}$  is the number of gradient variables represented by `prec` bits.

*Proof.* In asynchronous SGD (Section 2.2), each replica sends  $n_{gd} = \sum_{i=1}^{n_{pu}} L_i$  gradients (each of `prec` bits) to the parameter server at each iteration, so that the communication cost for one iteration in bits is

$$\text{PlainBits} = n_{gd} \cdot \text{prec}$$

In our system, let us compute the ciphertext length that each participant sends to the cloud parameter server at each iteration. By (10), the ciphertext sent to processing unit  $\text{PU}_i$  is in  $\mathbb{Z}_q^{1 \times (n_{lwe} + L_i)}$  so that its length in bits is  $(n_{lwe} + L_i) \log_2 q$ . The length in bits of all ciphertexts sent to the parameter server from the corresponding learning participant at each interaction is at most

$$\text{EncryptedBits} = \sum_{i=1}^{n_{pu}} (n_{lwe} + L_i) \log_2 q = n_{pu}n_{lwe} \log_2 q + n_{gd} \log_2 q.$$

Therefore, the increased factor is

$$\frac{\text{EncryptedBits}}{\text{PlainBits}} = \frac{n_{pu}n_{lwe} \log_2 q}{n_{gd} \cdot \text{prec}} + \frac{\log_2 q}{\text{prec}}$$

ending the proof.  $\square$

## 5.2 Using Paillier encryption

**Paillier encryption.** With public key  $pk = n$  (a large positive integer), the encryption of an integer  $m \in \{0, \dots, n-1\}$  is  $\text{PaiEnc}_{pk}(m) = r^n(1+n)^m \bmod n^2$  in which  $r$  is chosen randomly from  $\{0, \dots, n-1\}$ . The encryption is additively homomorphic because the ciphertext product  $\text{PaiEnc}_{pk}(m_1)\text{PaiEnc}_{pk}(m_2) \bmod n^2$  becomes an encryption of  $m_1 + m_2 \bmod n$ . For decryption and CPA security, see the paper [16].

**Packing data in encryption.** As the plaintext space has  $\log_2 n \geq 2048$  bits, we can pack many non-negative integers  $I_1, \dots, I_t$  each of `prec` bits into each Paillier plaintext as follows.

$$\text{PaiEnc}_{pk} \left( \overbrace{\left( \underbrace{[I_1 0_{\text{pad}}]}_{\text{prec} + \text{pad} \text{ bits}} \cdots \underbrace{[I_t 0_{\text{pad}}]}_{\text{prec} + \text{pad} \text{ bits}} \right)}^{\lfloor \log_2 n \rfloor \text{ bits}} \right)$$

in which  $0_{\text{pad}}$  is the zero padding of `pad` bits, which helps preventing overflows in ciphertext additions. Typically,  $\text{pad} \approx \log_2 n_{\text{gradupd}}$  as we need  $n_{\text{gradupd}}$  additions of ciphertexts. Moreover, as the number of plaintext bits must be less than  $\log_2 n$ , it is necessary that  $t(\text{prec} + \text{pad}) \leq \log_2 n$ . Therefore,

$$t = \left\lfloor \frac{\lfloor \log_2 n \rfloor}{\text{prec} + \text{pad}} \right\rfloor$$

which is the upper-bound of packing  $\text{prec}$ -bit integers into one Paillier plaintext. To handle both negative and positive integers, we can use the bijection  $z \in [0, 2^{\text{prec}} - 1] \mapsto z - \lfloor z/2^{\text{prec}} \rfloor \cdot 2^{\text{prec}}$ .

As a real number  $0 \leq r < 1$  can be represented as an integer of form  $\lfloor r \cdot 2^{\text{prec}} \rfloor$ , the above packing method can be used to encrypt around  $\lfloor \log_2 n \rfloor / (\text{prec} + \text{pad})$  real numbers in the range  $[0, 1)$  with precision  $\text{prec}$ , tolerating around  $2^{\text{pad}}$  ciphertext additions.

To realise the encryption  $\mathbf{E}(\cdot)$  in Figure 4, it suffices to describe an encryption of a real vector  $r = (r^{(1)}, \dots, r^{(L_i)}) \in \mathbb{R}^{L_i}$  because both  $W_{\text{global}}^{(i)}$  and  $\alpha \cdot G^{(i)} \in \mathbb{R}^{L_i}$  in which  $\sum_{i=1}^{n_{pu}} L_i = n_{gd}$ . The encryption  $\mathbf{E}(r)$  consists of approximately  $\lceil L_i/t \rceil$  Paillier ciphertexts:

$$\text{PaiEnc}_{pk} \left( \overbrace{\underbrace{\lfloor r^{(1)} \cdot 2^{\text{prec}} \rfloor 0_{\text{pad}} \cdots \lfloor r^{(t)} \cdot 2^{\text{prec}} \rfloor 0_{\text{pad}}}_{\text{prec}+\text{pad bits}}}^{\lfloor \log_2 n \rfloor \text{ bits}}, \dots, \text{PaiEnc}_{pk} \left( \overbrace{\underbrace{\lfloor r^{(L_i-t+1)} \cdot 2^{\text{prec}} \rfloor 0_{\text{pad}} \cdots \lfloor r^{(L_i)} \cdot 2^{\text{prec}} \rfloor 0_{\text{pad}}}_{\text{prec}+\text{pad bits}}}^{\lfloor \log_2 n \rfloor \text{ bits}} \right) \right).$$

**Theorem 4 (Increased communication factor, Paillier-based).** *The communication between the server and participants of our system is*

$$2 \left( 1 + \frac{\text{pad}}{\text{prec}} \right)$$

*times of the communication of the corresponding Downpour SGD, in which  $\text{pad}$  is the number of 0 paddings to tolerate  $2^{\text{pad}}$  additions (equal to the number of gradient updates on the server), and  $\text{prec}$  is the bit precision of numbers.*

*Proof.* In our system in Figure 4, each participant needs to encrypt and send at most  $\sum_{i=1}^{n_{pu}} L_i = n_{gd}$  gradients to the server, in which  $L_i$  is the length of  $G^{(i)}$ . Therefore, with the above packing method for Paillier encryption, the number of Paillier ciphertexts sent from each participant is around

$$\sum_{i=1}^{n_{pu}} \frac{L_i}{t} \approx \frac{n_{gd}}{t} \approx \frac{n_{gd}(\text{prec} + \text{pad})}{\lfloor \log_2 n \rfloor}.$$

Since each Paillier ciphertext is of  $2\lfloor \log_2 n \rfloor$  bits, the number of bits for each communication is around

$$\text{EncryptedBits} \approx 2\lfloor \log_2 n \rfloor \cdot \frac{n_{gd}}{t} \approx 2n_{gd}(\text{prec} + \text{pad}).$$

On the other hand, note that each replica in Asynchronous SGD needs to send  $n_{gd}$  gradients each of  $\text{prec}$  bits to the server, the communication cost in bits is  $\text{PlainBits} = n_{gd} \cdot \text{prec}$ . Therefore, the increased factor is

$$\frac{\text{EncryptedBits}}{\text{PlainBits}} = \frac{2n_{gd}(\text{prec} + \text{pad})}{n_{gd} \cdot \text{prec}} = 2 \left( 1 + \frac{\text{pad}}{\text{prec}} \right)$$

as claimed. □

## 6 Concrete evaluations with LWE-based encryption

We take  $n_{lwe} = 3000$ ,  $s = 8$ ,  $p = 2^{48} + 1$ ,  $n_{\text{gradupd}} = 2^{15}$ , and  $q = 2^{77}$  following Lemma 1. These parameters for  $(n_{lwe}, s, q)$  conservatively ensure that the LWE assumption has at least 128-bit security according to recent attacks [4, 7, 13, 14]. We will take  $n_{pu} \in \{1, 10\}$ , depending on the number of gradients.

### 6.1 Increased factors in communication

Let us consider multiple number of gradient parameters  $n_{gd}$ :

**Table 3.** LWE-based encryption scheme ( $n_{lwe} = 3000$ ,  $s = 8$ ,  $p = 2^{48} + 1$ , and  $q = 2^{77}$ ) running times with various numbers of gradients.

		Number of gradients $n_{gd}$							
		20000	27882	50000	52650	100000	200000	300000	402250
Processing time using 1 thread of computation									
<b>Enc. via (10)</b>	msec	164.4	213.1	364.6	454.8	732.9	1395.1	2031.6	2797.4
<b>Dec. of (10)</b>	msec	143.6	194.5	354.6	412.6	666.1	1355.0	1996.5	2738.0
<b>Add of ciphertexts via (11)</b>	usec	51	72.3	137.7	91.2	211.5	441.7	481.0	698.4
Processing time using 20 threads of computation									
<b>Enc. via (10)</b>	msec	30.2	50.8	93.2	70.3	188.9	415.8	587.8	875.4
<b>Dec. of (10)</b>	msec	31.9	44.7	85.8	78.5	196.0	460.5	644.9	820.9
<b>Add of ciphertexts via (11)</b>	usec	9.5	11.1	15.8	14.0	19.5	38.7	41.1	56.4

- $n_{gd} = 109386$ : this number of gradient parameters is used with the dataset MNIST [3]. Specifically, consider an MLP with form 784 (input) – 128 (hidden) – 64 (hidden) – 10 (output). The number of gradients for this network is  $(784 + 1)128 + (128 + 1)64 + (64 + 1)10 = 109386$ . We will consider cases where real numbers are represented by 32, so that  $\text{prec} = 32$ . Theorem 3 tells us that the increased communication factor between our system and the related Downpour SGD is

$$\frac{n_{pu}n_{lwe} \log_2 q}{n_{gd} \cdot \text{prec}} + \frac{\log_2 q}{\text{prec}} = \frac{1 \cdot 3000 \cdot 77}{109386 \cdot 32} + \frac{77}{32} \approx 2.47.$$

- $n_{gd} = 402250$ : this is used in [19] with the dataset SVHN [15]. The increased communication factor becomes

$$\frac{n_{pu}n_{lwe} \log_2 q}{n_{gd} \cdot \text{prec}} + \frac{\log_2 q}{\text{prec}} = \frac{1 \cdot 3000 \cdot 77}{402250 \cdot 32} + \frac{77}{32} \approx 2.42.$$

- $n_{gd} = 42 \cdot 10^6$ : this number of gradient parameters is used in [8] for speech data. As the number of gradients is large, consider  $n_{pu} = 10$ . The increased communication factor becomes

$$\frac{n_{pu}n_{lwe} \log_2 q}{n_{gd} \cdot \text{prec}} + \frac{\log_2 q}{\text{prec}} = \frac{10 \cdot 3000 \cdot 77}{42 \cdot 10^6 \cdot 32} + \frac{77}{32} \approx 2.4.$$

**Our system with Paillier encryption.** We take  $\log_2 n = 3072$ ,  $\text{pad} = 15$ , so that the increased communication factor becomes, via Theorem 4,

$$\frac{\text{EncryptedBits}}{\text{PlainBits}} = 2 \left( 1 + \frac{\text{pad}}{\text{prec}} \right) = 2 \left( 1 + \frac{15}{32} \right) \approx 2.93$$

which are independent of  $n_{gd}$ .

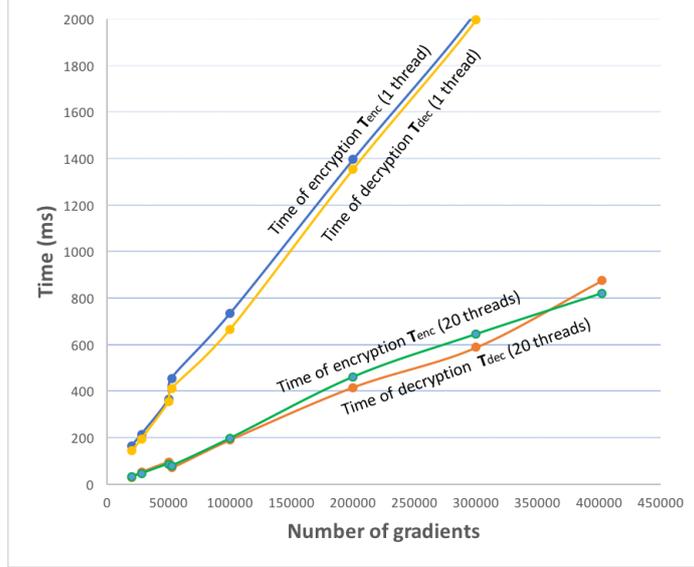
## 6.2 Estimating the computational costs

To estimate the running time of our system, we use the following formulas

$$\mathbf{T}_{\text{ours, one run}}^{(i)} = \mathbf{T}_{\text{original, one run}}^{(i)} + \mathbf{T}_{\text{enc}}^{(i)} + \mathbf{T}_{\text{dec}}^{(i)} + \mathbf{T}_{\text{upload}}^{(i)} + \mathbf{T}_{\text{download}}^{(i)} + \mathbf{T}_{\text{add}}, \quad (15)$$

$$\mathbf{T}_{\text{our system}} = \sum_{i=1}^N n_{\text{upload/download}}^{(i)} \times \mathbf{T}_{\text{ours, one run}}^{(i)} \quad (16)$$

where, in (15),  $\mathbf{T}_{\text{ours, one run}}^{(i)}$  is the running time of the participant  $i$  when doing the following: waiting for the addition of ciphertexts at the server ( $\mathbf{T}_{\text{add}}$ ); download the added ciphertext from the server ( $\mathbf{T}_{\text{download}}^{(i)}$ ); decrypted the downloaded ciphertext ( $\mathbf{T}_{\text{dec}}^{(i)}$ ); training with the downloaded



**Fig. 5.** Computational costs of the LWE-based encryption and decryption when  $n_{lwe} = 3000, s = 8, q = 2^{77}$ .

weight ( $\mathbf{T}_{\text{original, one run}}^{(i)}$ ); encrypting the resulting gradients ( $\mathbf{T}_{\text{enc}}^{(i)}$ ) and sending back that ciphertext to the server ( $\mathbf{T}_{\text{dec}}^{(i)}$ ). The total running time of our system  $\mathbf{T}_{\text{our system}}$  will be the sum of all  $N$  participants' running times, multiplied by the number  $n_{\text{repeat}}$  of repetitions, as expressed in (16).

**Environment.** Our codes for homomorphic encryption are in C++, and the benchmarks are over an Xeon CPU E5-2660 v3@ 2.60GHz server. To estimate the communication speed between each learning participants and the server, we assume a 1 Gbps network. To measure the running time of MLP, we use the Theano library over Cuda-8.0 and GPU Tesla K40m.

**Time of encryption, decryption and addition.** Table 3 gives the time for encryption, decryption and addition depending on the number of gradients  $n_{gd}$ , using  $n_{lwe} = 3000, s = 8, p = 2^{48} + 1$ , and  $q = 2^{77}$ . Figure 5 depicts the time of encryption and decryption using 1 thread and 20 threads of computation respectively.

**Multilayer Perceptron (MLP).** Consider an MLP with form 784 (input) – 128 (hidden) – 64 (hidden) – 10 (output). The number of gradients for this network is  $(784+1)128+(128+1)64+(64+1)10 = 109386$ . For 32-bit precision, these gradients are around  $109386 \times 32 / (8 \times 10^6) \approx 0.437$  MB in plain. The ciphertext of these gradients is of  $0.437 \times 2.47 \approx 1.0$  MB as computed in Section 6.1 which can be sent in around 0.008 seconds (= 8 ms) via an 1 Gbps communication channel. The original running time of this MLP with one data batch of size 50 (MNIST images) is  $\mathbf{T}_{\text{original, one run}}^{(i)} \approx 4.6$  (ms). Therefore,

$$\begin{aligned} \mathbf{T}_{\text{ours, one run}}^{(i)} &= \mathbf{T}_{\text{original, one run}}^{(i)} + \mathbf{T}_{\text{enc}}^{(i)} + \mathbf{T}_{\text{dec}}^{(i)} + \mathbf{T}_{\text{upload}}^{(i)} + \mathbf{T}_{\text{download}}^{(i)} + \mathbf{T}_{\text{add}} \\ &\approx 4.6 + 188.9 + 196.0 + 8 + 8 + 9.5/10^3 \text{ (ms)} \\ &\approx 405.5 \text{ (ms)} \end{aligned}$$

Suppose that there are totally  $2 \times 10^4$  times of upload and download from all training participants to the server, namely  $\sum_{i=1}^N n_{\text{upload/download}}^{(i)} = 2 \times 10^4$ . For simplicity suppose  $\mathbf{T}_{\text{ours, one run}}^{(i)}$  is the same for all  $1 \leq i \leq N$ . In such case, the running time of our system can be estimated as

$$\mathbf{T}_{\text{our system}} = \mathbf{T}_{\text{ours, one run}}^{(i)} \times 2 \times 10^4 = 405.5 \times 2 \times 10^4 \text{ (ms)} \approx 2.25 \text{ (hours)}$$

to achieve around 97.14% accuracy over the testing set of  $10^4$  images.

## 7 Conclusion

We show that, sharing gradients even partially over a parameter cloud server as in [19] may leak information on data. We then propose a new system that utilises additively homomorphic encryption to protect the gradients against the curious server. In addition to privacy-preserving, our system enjoys the property of not declining the accuracy of deep learning.

## Acknowledgment

This work is partially supported by Japan Science and Technology Agency CREST #JPMJCR168A.

## References

1. Machine Learning Coursera’s course. <https://www.coursera.org/learn/machine-learning>.
2. Stanford Deep Learning Tutorial. <http://deeplearning.stanford.edu>.
3. The MNIST dataset. <http://yann.lecun.com/exdb/mnist/>.
4. Y. Aono, X. Boyen, L. T. Phong, and L. Wang. Key-private proxy re-encryption under LWE. In G. Paul and S. Vaudenay, editors, *Progress in Cryptology - INDOCRYPT 2013. Proceedings*, volume 8250 of *Lecture Notes in Computer Science*, pages 1–18. Springer, 2013.
5. W. Banaszczyk. New bounds in some transference theorems in the geometry of numbers. *Mathematische Annalen*, 296(1):625–635, 1993.
6. W. Banaszczyk. Inequalities for convex bodies and polar reciprocal lattices in  $\mathbb{R}^n$ . *Discrete & Computational Geometry*, 13(1):217–231, 1995.
7. I. Chillotti, N. Gama, M. Georgieva, and M. Izabachène. Faster fully homomorphic encryption: Bootstrapping in less than 0.1 seconds. In *Advances in Cryptology - ASIACRYPT 2016 - 22nd International Conference on the Theory and Application of Cryptology and Information Security, Proceedings, Part I*, pages 3–33, 2016.
8. J. Dean, G. Corrado, R. Monga, K. Chen, M. Devin, Q. V. Le, M. Z. Mao, M. Ranzato, A. W. Senior, P. A. Tucker, K. Yang, and A. Y. Ng. Large scale distributed deep networks. In P. L. Bartlett, F. C. N. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *26th Annual Conference on Neural Information Processing Systems 2012.*, pages 1232–1240, 2012.
9. J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *J. Mach. Learn. Res.*, 12:2121–2159, July 2011.
10. R. Gilad-Bachrach, N. Dowlin, K. Laine, K. E. Lauter, M. Naehrig, and J. Wernsing. Cryptonets: Applying neural networks to encrypted data with high throughput and accuracy. In M. Balcan and K. Q. Weinberger, editors, *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016*, volume 48 of *JMLR Workshop and Conference Proceedings*, pages 201–210. JMLR.org, 2016.
11. O. Goldreich. *The Foundations of Cryptography - Volume 2, Basic Applications*. Cambridge University Press, 2004.
12. B. Hitaj, G. Ateniese, and F. Pérez-Cruz. Deep models under the GAN: information leakage from collaborative deep learning. *CoRR*, abs/1702.07464, 2017.
13. R. Lindner and C. Peikert. Better key sizes (and attacks) for LWE-based encryption. In A. Kiayias, editor, *Topics in Cryptology - CT-RSA 2011*, volume 6558 of *Lecture Notes in Computer Science*, pages 319–339. Springer, 2011.
14. M. Liu and P. Q. Nguyen. Solving BDD by enumeration: An update. In E. Dawson, editor, *Topics in Cryptology - CT-RSA 2013. Proceedings*, volume 7779 of *Lecture Notes in Computer Science*, pages 293–309. Springer, 2013.
15. Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning 2011*, 2011.
16. P. Paillier. Public-key cryptosystems based on composite degree residuosity classes. In J. Stern, editor, *Advances in Cryptology - EUROCRYPT ’99, Proceeding*, volume 1592 of *Lecture Notes in Computer Science*, pages 223–238. Springer, 1999.
17. L. T. Phong, Y. Aono, T. Hayashi, L. Wang, and S. Moriai. Privacy-preserving deep learning: Revisited and enhanced. In *Applications and Techniques in Information Security - 8th International Conference, ATIS 2017, Proceedings*, pages 100–110, 2017.
18. B. Recht, C. Ré, S. J. Wright, and F. Niu. Hogwild: A lock-free approach to parallelizing stochastic gradient descent. In J. Shawe-Taylor, R. S. Zemel, P. L. Bartlett, F. C. N. Pereira, and K. Q. Weinberger, editors, *NIPS 2011*, pages 693–701, 2011.
19. R. Shokri and V. Shmatikov. Privacy-preserving deep learning. In I. Ray, N. Li, and C. Kruegel, editors, *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security, 2015*, pages 1310–1321. ACM, 2015.

## A CPA security of the LWE-based homomorphic encryption scheme

We first recall the definition of CPA security.

**Definition 2. (CPA security)** *With respect to a PHE scheme as in Definition 1, consider the following game between an adversary  $\mathcal{A}$  and a challenger:*

- **Setup.** *The challenger creates  $pp$  and key pairs  $(pk, sk)$ . Then  $pp$  and  $pk$  are given to  $\mathcal{A}$ .*
- **Challenge.**  *$\mathcal{A}$  chooses two plaintexts  $m_0, m_1$  of the same length, then submits them to the challenger, who in turn takes  $b \in \{0, 1\}$  randomly and computes  $C^* = \text{Enc}(pk, m_b)$ . The challenge ciphertext  $C^*$  is returned to  $\mathcal{A}$ , who then produces a bit  $b'$ .*

*A PHE scheme is secure against chosen plaintext attacks (CPA-secure) if the advantage*

$$\mathbf{Adv}_{\mathcal{A}}^{\text{cpa}}(\lambda) = \left| \Pr[b' = b] - \frac{1}{2} \right|$$

*is negligible in  $\lambda$ .*

Related to the decision LWE assumption  $\text{LWE}(n_{lwe}, s, q)$ , where  $n_{lwe}, s, q$  is the parameter for security, consider matrix  $A \xleftarrow{\$} \mathbb{Z}_q^{m \times n_{lwe}}$ , vectors  $r \xleftarrow{\$} \mathbb{Z}_q^{m \times 1}$ ,  $x \xleftarrow{\$} \mathbb{Z}_{(0,s)}^{n_{lwe} \times 1}$ ,  $e \xleftarrow{\$} \mathbb{Z}_{(0,s)}^{m \times 1}$ . Then vector  $Ax + e$  is computed over  $\mathbb{Z}_q$ . Define the following advantage of a poly-time probabilistic algorithm  $\mathcal{D}$ :

$$\mathbf{Adv}_{\mathcal{D}}^{\text{LWE}(n_{lwe}, s, q)}(\lambda) = \left| \Pr[\mathcal{D}(A, Ax + e) \rightarrow 1] - \Pr[\mathcal{D}(A, r) \rightarrow 1] \right|.$$

The LWE assumption asserts that  $\mathbf{Adv}_{\mathcal{D}}^{\text{LWE}(n_{lwe}, s, q)}(\lambda)$  is negligible as a function of  $\lambda$ .

*Claim.* The additively homomorphic encryption scheme in Section 5.1 is CPA-secure under the LWE assumption. Specifically, for any poly-time adversary  $\mathcal{A}$ , there is an algorithm  $\mathcal{D}$  of essentially the same running time such that

$$\mathbf{Adv}_{\mathcal{A}}^{\text{cpa}}(\lambda) \leq (l + 1) \cdot \mathbf{Adv}_{\mathcal{D}}^{\text{LWE}(n_{lwe}, s, q)}(\lambda).$$

*Proof.* The proof is standard and follows [13] closely. First, we modify the original game in Definition 2 by providing  $\mathcal{A}$  with a random  $P \xleftarrow{\$} \mathbb{Z}_q^{n_{lwe} \times l}$ . Namely  $P = pR - AS \in \mathbb{Z}_q^{n_{lwe} \times l}$  is turned to random. This is indistinguishable to  $\mathcal{A}$  thanks to the LWE assumption with secret vectors as  $l$  columns of  $S$ . More precisely, we need the condition  $\gcd(p, q) = 1$  to reduce  $P = pR - AS \in \mathbb{Z}_q^{n_{lwe} \times l}$  to the LWE form. Indeed,  $p^{-1}P = R + (-p^{-1}A)S \in \mathbb{Z}_q^{n_{lwe} \times l}$ . As  $A$  is random,  $A' = -p^{-1}A \in \mathbb{Z}_q^{n_{lwe} \times n_{lwe}}$  is also random. Therefore,  $P' = p^{-1}P \in \mathbb{Z}_q^{n_{lwe} \times l}$  is random under the LWE assumption which in turn means  $P$  is random as claimed.

Second, the challenge ciphertext  $c^* = e_1[A|P] + p[e_2|e_3] + [0|m_b]$  is turned to random. This relies on the LWE assumption with secret vector  $e_1$ , while also basing on the first modification so that  $[A|P]$  is a random matrix. Here, the condition  $\gcd(p, q) = 1$  is also necessary as above. Thus  $b$  is perfectly hidden after this change. The factor  $l + 1$  is due to  $l$  uses of LWE in changing  $P$  and 1 use in changing  $c^*$ .

## B Proof of Lemma 1

We will use following lemmas, whose proofs can be derived from [5, 6]. Below  $\langle \cdot, \cdot \rangle$  stands for inner product. Writing  $\|\mathbb{Z}_{(0,s)}^n\|$  is a short hand for taking a vector from the discrete Gaussian distribution of deviation  $s$  and computing its Euclidean norm.

**Lemma 2.** Let  $c \geq 1$  and  $C = c \cdot \exp(\frac{1-c^2}{2})$ . Then for any real  $s > 0$  and any integer  $n \geq 1$ , we have

$$\Pr \left[ \|\mathbb{Z}_{(0,s)}^n\| \geq \frac{c \cdot s \sqrt{n}}{\sqrt{2\pi}} \right] \leq C^n.$$

**Lemma 3.** For any real  $s > 0$  and  $T > 0$ , and any  $x \in \mathbb{R}^n$ , we have

$$\Pr \left[ \|\langle x, \mathbb{Z}_{(0,s)}^n \rangle\| \geq Ts\|x\| \right] < 2 \exp(-\pi T^2).$$

*Proof (Proof of Claim 1).* Suppose  $c_{\text{add}} = \sum_{i=1}^{n_{\text{add}}} c^{(i)} = \sum_{i=1}^{n_{\text{add}}} (c_1^{(i)}, c_2^{(i)})$  is the additions of  $n_{\text{add}}$  ciphertexts. In Claim 1,  $n_{\text{add}} = n_{\text{gradupd}}$ . Its decryption is

$$\begin{aligned} \bar{m}_{\text{add}} &= \sum_{i=1}^{n_{\text{add}}} \left[ c_1^{(i)} S + c_2^{(i)} \right] \\ &= \sum_{i=1}^{n_{\text{add}}} \left[ (e_1^{(i)} A + p e_2^{(i)}) S + e_1^{(i)} P + p e_3^{(i)} + m^{(i)} \right] \\ &= \sum_{i=1}^{n_{\text{add}}} \left[ e_1^{(i)} (AS + P) + p(e_2^{(i)} S + e_3^{(i)}) + m^{(i)} \right] \\ &= \sum_{i=1}^{n_{\text{add}}} \left[ p(e_1^{(i)} R + e_2^{(i)} S + e_3^{(i)}) + m^{(i)} \right] \in \mathbb{Z}_q^{1 \times l} \end{aligned}$$

Therefore, the accumulative noise after  $n_{\text{add}}$  ciphertext additions is

$$\text{Noise} = p \sum_{i=1}^{n_{\text{add}}} (e_1^{(i)} R + e_2^{(i)} S + e_3^{(i)}) \in \mathbb{Z}_q^{1 \times l}$$

The term  $e_1^{(i)} R + e_2^{(i)} S + e_3^{(i)} \in \mathbb{Z}_q^{1 \times l}$  has  $l$  components, each of which can be written as the inner product of two vectors of form

$$\begin{aligned} e &= (e_1^{(i)}, e_2^{(i)}, e_3^{(i)}) \\ x &= (r, \quad r', \mathbf{010}_{1 \times l}) \end{aligned}$$

where,

- vectors  $e_1^{(i)} \stackrel{\mathbb{g}}{\leftarrow} \mathbb{Z}_{(0,s)}^{1 \times n_{lwe}}$ ,  $e_2^{(i)} \stackrel{\mathbb{g}}{\leftarrow} \mathbb{Z}_{(0,s)}^{1 \times n_{lwe}}$ , and  $e_3^{(i)} \stackrel{\mathbb{g}}{\leftarrow} \mathbb{Z}_{(0,s)}^{1 \times l}$
- vectors  $r, r' \stackrel{\mathbb{g}}{\leftarrow} \mathbb{Z}_{(0,s)}^{1 \times n_{lwe}}$  represent corresponding columns in matrices  $R, S$ ; and  $\mathbf{010}_{1 \times l}$  stands for a vector of length  $l$  with all 0's except one 1.

We have

$$\begin{aligned} e &\in \mathbb{Z}_{(0,s)}^{1 \times (2n_{lwe} + l)} \\ \|x\| &\leq \|(r, r')\| + 1. \end{aligned}$$

Applying Lemma 2 for the above vector  $(r, r')$  of length  $2n_{lwe}$ , with overwhelming probability of

$$1 - C^{2n_{lwe}} (\geq 1 - 2^{-80} \text{ for our choices of parameters})$$

we have

$$\|x\| \leq \frac{c \cdot s \sqrt{2n_{lwe}}}{\sqrt{2\pi}} + 1.$$

We now use Lemma 3 with vectors  $x$  and  $e$ . Let  $\rho$  be the error per message symbol in decryption, we set  $2 \exp(-\pi T^2) = \rho$ , so  $T = \sqrt{\ln(2/\rho)}/\sqrt{\pi}$ . The bound on the inner product  $\langle x, e \rangle$  becomes  $Ts\|x\|$ , which is not greater than

$$U = \frac{s\sqrt{\ln(2/\rho)}}{\sqrt{\pi}} \left( \frac{c \cdot s\sqrt{2n_{lwe}}}{\sqrt{2\pi}} + 1 \right)$$

so that each component in  $\text{Noise} \in \mathbb{Z}_q^{1 \times l}$  will be less than

$$B = pn_{\text{add}} \cdot U.$$

For correctness, it suffices to set  $B = q/2$ , namely  $q = 2B$ . This means

$$\log_2 q = \log_2 p + \log_2 n_{\text{add}} + \log_2(U) + 1.$$

We take  $c = 1.1$  so that  $C^{2n_{lwe}} < 2^{-80}$  when  $n_{lwe} \geq 3000$ . Also take  $\rho = 2^{-80}$ ,  $s = 8$  so that

$$U = 167.9106\sqrt{n_{lwe}} + 33.8197$$

and the proof follows. □