

# Security for Distributed Machine Learning based Software

Laurent Gomez<sup>1</sup>, Alberto Ibarrondo<sup>2</sup>, Marcus Wilhelm<sup>3</sup>, José Márquez<sup>4</sup>, and Patrick Duverger<sup>5</sup>

<sup>1</sup> SAP Global Security, SAP Security Research, France  
laurent.gomez@sap.com

<sup>2</sup> Eurecom, Sophia Antipolis, France  
alberto.ibarrondo@eurecom.com

<sup>3</sup> Hasso Plattner Institute, University Potsdam, Germany  
marcus.wilhelm@student.hpi.com

<sup>4</sup> Portfolio Strategy Technology Adoption, SAP SE, Germany  
jose.marquez@sap.com

<sup>5</sup> City of Antibes - Juan les Pins, France  
patrick.duverger@antibes.com

**Abstract.** Current developments in Enterprise Systems observe a paradigm shift, moving the needle from the backend to the edge sectors of those; by distributing data, decentralizing applications and integrating novel components seamlessly to the central systems. Distributively deployed AI capabilities will thrust this transition.

Several non-functional requirements arise along with these developments, security being at the center of the discussions. Bearing those requirements in mind, hereby we propose an approach to holistically protect distributed Deep Neural Network (DNN) based/enhanced software assets, i.e. confidentiality of their input & output data streams as well as safeguarding their Intellectual Property.

Making use of Fully Homomorphic Encryption (FHE), our approach enables the protection of Distributed Neural Networks, while processing encrypted data. On that respect we evaluate the feasibility of this solution on a Convolutional Neural Network (CNN) for image classification deployed on distributed infrastructures.

**Keywords:** Intellectual Property Protection, Fully Homomorphic Encryption, Neural Networks, Distributed Landscapes, Smart Cities

## 1 Introduction

### 1.1 Motivation

Until now, the backend (on-prem & cloud) deployments were considered as the single source of truth & unique point of access in regards of Enterprise Systems (ES). Nevertheless, a paradigm shift has been recently observed, by the deployment of ES assets towards the Edge sectors of the landscapes; by distributing data, decentralizing applications, de-abstracting technology and integrating edge components seamlessly to the central backend systems.

Capitalizing on recent advances on High Performance Computing along with the rising amounts of publicly available labeled data, Deep Neural Networks (DNN), as an implementation of AI, have and will revolutionize virtually every current application domain as well as enable novel ones like those on autonomous, predictive, resilient, self-managed, adaptive, and evolving applications.

Distributively deployed AI capabilities will thrust the above mentioned transition. As reported by Deloitte, “... *companies are incorporating artificial intelligence in particular, machine learning into their 'Internet of Things applications' and seeing capabilities grow, including improving operational efficiency and helping avoid unplanned downtime*” [28].

## 1.2 Problem Statement

The deployment of data processing capabilities throughout Distributed Enterprise Systems rises several security challenges related to the protection of input & output data [26] as well as of software assets.

In the specific context of distributed intelligence, DNN based/enhanced software will represent key investments in infrastructure, skills and governance, as well as in the acquisition of data and talents. The software industry is therefore in the direct need to safeguard these strategic investments by enforcing the protection of this new form of Intellectual Property.

Furthermore, on the wake of Data Protection (DP) regulations such as the EU-GDPR [26], Independent Software Vendors (ISVs) have the non-transferable requirement to comply with those.

Therefore, ISVs aim to protect both: data and the Intellectual Property of their AI-based software assets, deployed on potentially unsecure edge hardware & platforms [15].

## 1.3 State-of-the-Art

Security of Deep Neural Networks is a current research topic taking advantage of two major cryptographic approaches: variants of Fully Homomorphic Encryption/FHE [12] and Secure Multi-Party Computation/SMC [8]. While FHE techniques allow addition and multiplication on encrypted data, SMC enables arithmetic operations on data shared across multi-parties.

Several approaches can be found in the literature, at different phases of the development and deployment of DNNs.

*Secure Training* Secure DNN training has been addressed using FHE [16] and SMC [30], disregarding protection once the trained model is to be productively deployed. Other Machine Learning models such as linear and logistic regressions have also been trained in a secure way in [24]. In those approaches, confidentiality of training data is guaranteed, while runtime protection (i.e. input, model, output) is out of scope.

*Processing on Encrypted Data* At processing phase, SMC has led to cooperative solutions where several devices work together to obtain federated inferences [21], not supporting deployment of the trained DNN to trusted decentralized systems. DNN processing on FHE encrypted data is covered in CryptoNets [13], improved in [4] and [18]. More recently, in [2], the authors proposed a privacy-preserving framework for deep learning, making use of the SEAL [29] FHE library. While disclosure of data at runtime is prevented in these solutions, protection of DNN models remains out of the scope.

*Intellectual Property Protection of DNN Model* In [31], the authors tackles IP protection of DNN models through model watermarking. While infringement can be detected with this method, it can not be prevented. Furthermore, runtime protection of input, model and output are out of scope.

To the best of our knowledge, no other publication has holistically tackled the protection of both trained DNN models and data, targeting distributed untrusted systems.

#### 1.4 Data & Intellectual Property Protection for Deep Neural Networks

In this paper we propose a novel approach for the Intellectual Property Protection of DNN-based/enhanced software while enabling data protection at processing time, making use of concepts such as Fully Homomorphic Encryption (FHE).

Once trained, DNN model parameters (i.e. weights, biases) are encrypted homomorphically. The resulting (encrypted) DNN can be distributed across untrusted landscapes, preserving its IP while mitigating the risk of reverse engineering. At runtime, FHE-encrypted insights from encrypted input data are produced by the homomorphically encrypted DNN. Confidentiality of both trained DNN, input and output data will be therefore guaranteed.

Despite of recent improvements of FHE schemes [3], [5] and implementations [29, ?,?], homomorphic encryption remains computationally expensive. Hence it could represent a bottleneck having a negative impact on overall performance, and on the accuracy of encrypted DNNs outputs, handling encrypted inputs. In this paper, we therefore evaluate as well the overall performance (e.g. CPU, memory, disk usage) along with the accuracy of encrypted DNNs.

This paper is organized as follows: Section 2 details the fundamentals of our approach. Section 3 provides an overview of our solution. In Sections 4 and 5, we present the architecture and evaluation, concluding with an outlook in Section 6.

## 2 Fundamentals

### 2.1 Deep Neural Network

Figure 1 depicts a DNN with multiple layers. It is composed of  $L$  layers:

1. An **input layer**, the tensor of input data  $\mathbf{X}$
2.  $L - 1$  **hidden layers**, mathematical computations transforming  $\mathbf{X}$  sequentially.
3. An **output layer**, the tensor of output data  $\mathbf{Y}$ .

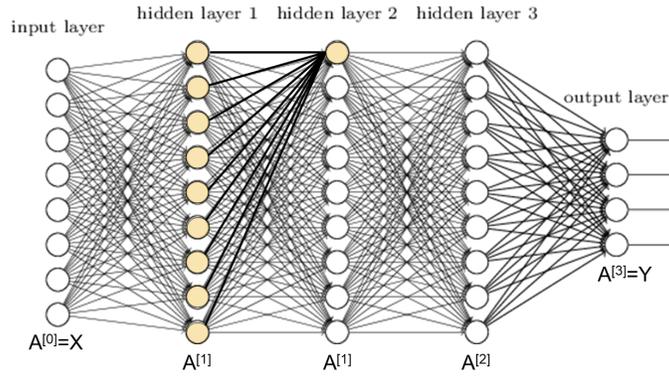


Fig. 1: Deep Neural Network [14]

We denote the output of layer  $i$  as a tensor  $\mathbf{A}^{[i]}$ , with  $\mathbf{A}^{[0]} = X$ , and  $\mathbf{A}^{[L]} = Y$ . Tensors can have different sizes and number of dimensions.

Each layer  $\mathbf{A}^{[i]}$  depends on the mathematical computations performed at the previous layer  $\mathbf{A}^{[i-1]}$ . At each layer  $\mathbf{A}^{[i]}$ , two types of function can be computed:

- *Linear*: involving polynomial operations.
- *Non-linear*, involving non-linear operations, so called activation function, such as *max*, *exp*, *division*, ReLU, or Sigmoid.

**Linear Computation Layer** For the sake of clarity, we exemplify the inner linear computation with a Fully Connected (FC) layer, as depicted in Figure 2.

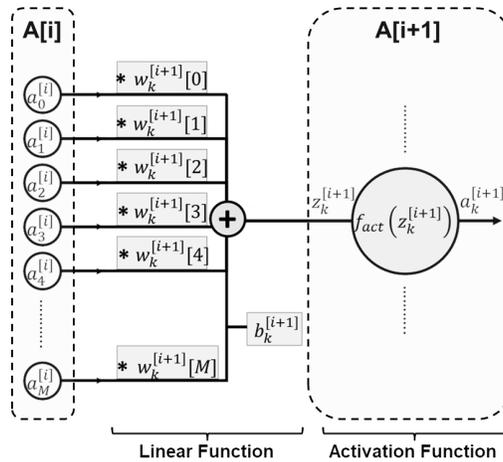


Fig. 2: Fully Connected layer with Activation Function [14]

A Fully Connected layer, noted  $\mathbf{A}^{[i]}$ , is composed of  $n$  parallel *neurons*, performing a  $\mathbb{R}^n \rightarrow \mathbb{R}^n$  transformation (see Figure 2). We define:

$\mathbf{a}^{[i]} = [a_0^{[i]} \dots a_k^{[i]} \dots a_N^{[i]}]^T$  as the output of layer  $\mathbf{A}^{[i]}$ ;

$\mathbf{z}^{[i]} = [z_0^{[i]} \dots z_k^{[i]} \dots z_N^{[i]}]^T$  as the linear output of layer  $\mathbf{A}^{[i]}$ ; ( $\mathbf{z}^{[i]} = \mathbf{a}^{[i]}$  if there is no activation function)

$\mathbf{b}^{[i]} = [b_0^{[i]} \dots b_k^{[i]} \dots b_N^{[i]}]^T$  as the bias for layer  $\mathbf{A}^{[i]}$ ;

$\mathbf{W}^{[i]} = [\mathbf{w}_0^{[i]} \dots \mathbf{w}_k^{[i]} \dots \mathbf{w}_N^{[i]}]^T$  as the weights for layer  $\mathbf{A}^{[i]}$ .

Neuron  $k$  performs a linear combination of the output of the previous layer  $\mathbf{a}^{[i-1]}$  multiplied by the weight vector  $\mathbf{w}_k^{[i]}$  and shifted with a bias scalar  $b_k^{[i]}$ , obtaining the linear combination  $z_k^{[i]}$ :

$$z_k^{[i]} = \left( \sum_{l=0}^M w_k^{[i][l]} * a_l^{[i-1]} \right) + b_k^{[i]} = \mathbf{w}_k^{[i]} * \mathbf{a}^{[i-1]} + b_k^{[i]} \quad [14] \quad (1)$$

Vectorizing the operations for all the neurons in layer  $\mathbf{A}^{[i]}$  we obtain the dense layer transformation:

$$\mathbf{z}^{[i]} = \mathbf{W}^{[i]} * \mathbf{a}^{[i-1]} + \mathbf{b}^{[i]} \quad [14] \quad (2)$$

where  $\mathbf{W}$  and  $\mathbf{b}$  are the parameters for layer  $\mathbf{A}^{[i]}$ .

**Activation Functions** Activation functions are the major source of non-linearity in DNNs. They are performed element-wise ( $\mathbb{R}^0 \rightarrow \mathbb{R}^0$ , thus easily vectorized), and are generally located after linear transformations such as Fully Connected layers.

$$a_k^{[i]} = f_{act} \left( z_k^{[i]} \right) \quad (3)$$

Several activation functions have been proposed in the literature but *Rectified Linear Unit (ReLU)* is currently considered as the most efficient activation function for DL. Several variants of ReLU exist, such as Leaky ReLU[23], ELU[7] or its differentiable version *Softplus*.

$$\begin{aligned} ReLU(z) &= z^+ = \max(0, z) \\ Softplus(z) &= \log(e^z + 1) \end{aligned} \quad [14] \quad (4)$$

## 2.2 Homomorphic Encryption

While preserving data privacy, Homomorphic Encryption (HE) schemes allow certain computations on ciphertext without revealing neither its inputs nor its internal states. Gentry [12] first proposed a Fully Homomorphic Encryption (FHE) scheme, which theoretically could compute any kind of arithmetic circuit, but is computationally intractable in practice. FHE evolved into more efficient schemes preserving addition and

multiplication over encrypted data, such as BGV [3], FV [11] or CKKS [5], allowing approximations of multiplicative inverse, exponential and logistic function, or discrete Fourier transformation. Similar to asymmetric encryption, a public-private key pair ( $pub, priv$ ) is generated.

**Definition 1.** An encryption scheme is called homomorphic over an operation  $\odot$  if it supports the following

$$\begin{aligned} Enc_{pub}(m) &= \langle m \rangle_{pub}, \forall m \in \mathcal{M} \\ \langle m_1 \odot m_2 \rangle_{pub} &= \langle m_1 \rangle_{pub} \odot \langle m_2 \rangle_{pub}, \forall m_1, m_2 \in \mathcal{M} \end{aligned}$$

where  $Enc_{pub}$  is the encryption algorithm and  $\mathcal{M}$  is the set of all possible messages.

**Definition 2.** Decryption is performed as follows

$$\begin{aligned} Enc_{pub}(m) &= \langle m \rangle_{pub}, \forall m \in \mathcal{M} \\ Dec_{priv}(\langle m \rangle_{pub}) &= m \end{aligned}$$

where  $Dec_{priv}$  is the decryption algorithm and  $\mathcal{M}$  is the set of all possible messages.

### 2.3 Challenges

Even though HE schemes seem theoretically promising, their usage comes with several drawbacks, particularly when applied to Deep Learning.

**Noise budget** In Gentry’s lattice-based HE schemes[12] and subsequent variants of it, ciphertexts contain a small term of random noise drawn from some probability distribution. While every operation performed on a ciphertext increases the noise of the resulting ciphertext, it is important to keep the noise below a certain threshold, because once the noise reaches that threshold, it is no longer possible to decrypt the ciphertext. To estimate the current magnitude of noise, a **noise budget** can be calculated, that starts as a positive integer, decreases with subsequent operations and reaches 0 exactly when the ciphertext becomes indecipherable. The noise budget is more strongly affected by multiplications as by additions.

In order to cope with that challenge, encryption parameters can be adjusted accordingly to the required computation depth of an arithmetic circuit. In addition, Gentry introduced the so called bootstrapping procedure, which resets the noise budget of a ciphertext, but requires significant additional computational costs. Recently in [5], the authors proposed a optimized bootstrapping approach with improved performance.

**FHE libraries and APIs** As summarized in Table 1, multiple FHE libraries are available. Depending on the supported HE schemes, those libraries show noticeable difference on performance (e.g. computational, memory consumption), on supported operations type (e.g. addition, multiplication, negative, square, division), datatype (e.g. floating point, integer), and chipset infrastructure (e.g. CPU, GPU).

In addition, and regardless on their level of maturity and performance, HE libraries can be configured through several encryption parameters such as:

- Polynomial degree or modulus: which determines the available noise budget and strongly affects the performance.
- Plaintext modulus: which is mostly associated to the size of input data.
- Security parameter: which sets the reached level of security in bits of the cryptosystem (e.g. 128, 192, 256-bit security level).

Fine-tuning of those encryption parameters enables developers to optimize the performance of encryption and encrypted operations. The selection of the right encryption parameters depends on the size of the plaintext data, targeted accuracy loss or level of security.

Library	Language	Dependencies	License	Description
HElib[17]	C++ 11	NTL, GMP	Apache 2.0	Mature. Low level implementation, hard to use but complete. Ciphertext packing, integers, bootstrapping, multi-threading.
PALISADE[25]	C++ 11	None	Copyrighted	Many functionalities & multiple schemes. Well documented but fairly new. Ciphertext packing, integers, fractionals, bootstrapping, multi-threading.
SEAL[29]	C++ 17	None	Microsoft	Well documented and easy to use. Ciphertext packing, fractionals, automatic parameter selection and multi-threading. Latest version SEAL 3.0 - Oct'18
FHEW[10]	C++ 11	FFTW	GNU-GPLv2	NAND gate with ciphertext packing, works over bits.
TFHE[6]	C++ 11	FFTW	Apache 2.0	Binary gates at 20ms per gate. Works over bits. Bootstrapping is included in all operations.
cuFHE[9]	CUDA C++ 11	NVIDIA CUDA <i>arch</i> >= 6.0	MIT	Binary gates at 0.1ms per gate. Works over bits. Bootstrapping included in operations.

Table 1. FHE implementation libraries [14]

**Linear function support only** By construction, linear functions, composed of addition and multiplication operations, are seamlessly protected by FHE. But, non-linear activation functions such as ReLU or Sigmoid require approximation to be computed with FHE schemes.

The challenge lies on the transformation of activation functions into polynomial approximations supported by HE schemes. We elaborate more on approximation of activation functions in Section 3.2.

**Supported plaintext type** The vast majority of HE schemes allow operations on integers [17, 29], while others use booleans [6] or floating point numbers [5, 29]. In the case of integer supporting HE schemes, rational numbers can be approximated using fixed-point arithmetic by scaling with a scaling factor and rounding.

**Performance** FHE schemes are computationally expensive and memory consuming. In addition, ciphertexts are often significantly bigger than plaintexts and thus use more memory and disk space.

Even if in the past years the performance of FHE made it impractical, recent FHE schemes show promising throughput. New FHE libraries take also advantage of GPU acceleration.

In addition, modern implementations of HE schemes such as **HELib** [17], **SEAL** [29], or **PALISADE** [25] benefit from Single Instruction Multiple Data (SIMD), allowing multiple integers to be stored in a single ciphertext and vectorizing operations, which can accelerate certain applications significantly.

### 3 Approach

As introduced in Section 1.2, the delivery of DNN-enriched insights come at a cost. ISVs aim to guarantee data security, together with the IP protection of their DNN-based software assets, deployed on potentially unsecure edge hardware & platforms. In order to achieve those security objectives on DNN, we utilize FHE schemes to operate on ciphertext at runtime.

Consequently, secure training of DNN is out of scope of our approach as we focus on runtime execution. We assume that DNN training already preserves both data privacy & confidentiality, and the resulting trained model. Once a model is trained, as discussed in Section 2.1, we obtain a set of parameters for each DNN layer; i.e weights  $\mathbf{W}^{[i]}$  and biases  $\mathbf{b}^{[i]}$  for Fully Connected layers DNN are not solely made of FC layers, and in [14], we identified different type of linear operations parameters within DNN such as Batch Normalization[19] or Convolutional Layer[20]. Those parameters constitute the IP to be protected when deploying a DNN to distributed systems.

#### 3.1 Linear Computation Layer Protection

Our approach is agnostic from the type of layer. In [14], we detail the encryption of layers such as Convolutional Layer or Batch Normalization. For sake of simplicity, we exemplify the encryption of DNN layers parameters on FC layers. Since FC are simply a linear transformation on the previous layer's outputs, encryption is achieved straightforwardly as follows

$$\begin{aligned} \langle \mathbf{z}^{[i]} \rangle_{\text{pub}} &= \langle \mathbf{W}^{[i]} * \mathbf{a}^{[i-1]} + \mathbf{b}^{[i]} \rangle_{\text{pub}} \\ &= \langle \mathbf{W}^{[i]} \rangle_{\text{pub}} * \langle \mathbf{a}^{[i-1]} \rangle_{\text{pub}} + \langle \mathbf{b}^{[i]} \rangle_{\text{pub}} \end{aligned} \quad [14] \quad (5)$$

**Fully Connected Layer (FC)** Also known as Dense Layer, it is composed of  $N$  parallel *neurons*, performing a  $\mathbb{R}^1 \rightarrow \mathbb{R}^1$  transformation (Figure 1). We will define:

$\mathbf{a}^{[i]} = [a_0^{[i]} \dots a_k^{[i]} \dots a_N^{[i]}]^T$  as the output of layer  $i$ ;

$\mathbf{z}^{[i]} = [z_0^{[i]} \dots z_k^{[i]} \dots z_N^{[i]}]^T$  as the linear output of layer  $i$ ; ( $\mathbf{z}^{[i]} = \mathbf{a}^{[i]}$  if there is no activation function)

$\mathbf{b}^{[i]} = [b_0^{[i]} \dots b_k^{[i]} \dots b_N^{[i]}]^T$  as the bias of layer  $i$ ;

$\mathbf{W}^{[i]} = [\mathbf{w}_0^{[i]} \dots \mathbf{w}_k^{[i]} \dots \mathbf{w}_N^{[i]}]^T$  as the weights of layer  $i$ .

Neuron  $k$  performs a linear combination of the output of the previous layer  $\mathbf{a}^{[i-1]}$  multiplied by the weight vector  $\mathbf{w}_k^{[i]}$  and shifted with a bias scalar  $b_k^{[i]}$ , obtaining the linear combination  $z_k^{[i]}$ :

$$z_k^{[i]} = \left( \sum_{l=0}^M w_k^{[i][l]} * a_l^{[i-1]} \right) + b_k^{[i]} = \mathbf{w}_k^{[i]} * \mathbf{a}^{[i-1]} + b_k^{[i]} \quad [14] \quad (6)$$

Vectorizing the operations for all the neurons in layer  $i$  we obtain the dense layer transformation:

$$\mathbf{z}^{[i]} = \mathbf{W}^{[i]} * \mathbf{a}^{[i-1]} + \mathbf{b}^{[i]} \quad [14] \quad (7)$$

*Protecting FC layer* Since FC is a linear layer, it can be directly computed in the encrypted domain using additions and multiplications. Vectorization is achieved straightforwardly:

$$\begin{aligned} \langle \mathbf{z}^{[i]} \rangle_{\text{pub}} &\equiv \langle \mathbf{W}^{[i]} * \mathbf{a}^{[i-1]} + \mathbf{b}^{[i]} \rangle_{\text{pub}} \\ &= \langle \mathbf{W}^{[i]} \rangle_{\text{pub}} * \langle \mathbf{a}^{[i-1]} \rangle_{\text{pub}} + \langle \mathbf{b}^{[i]} \rangle_{\text{pub}} \end{aligned} \quad [14] \quad (8)$$

$$\begin{aligned} \langle \mathbf{z}^{[i]} \rangle_{\text{pub}} &\equiv \langle \mathbf{W}^{[i]} * \mathbf{a}^{[i-1]} + \mathbf{b}^{[i]} \rangle_{\text{pub}} \\ &= \langle \mathbf{W}^{[i]} \rangle_{\text{pub}} * \langle \mathbf{a}^{[i-1]} \rangle_{\text{pub}} + \langle \mathbf{b}^{[i]} \rangle_{\text{pub}} \end{aligned} \quad [14] \quad (9)$$

$$\langle \mathbf{a}_k^{[i]} \rangle_{\text{pub}} \equiv \langle \mathbf{f}_{\text{approxact}}(\mathbf{z}_k^{[i]}) \rangle_{\text{pub}} \quad [14] \quad (10)$$

**Convolutional Layer (Conv)** Conv layers constitute a key improvement for image recognition and classification using NNs. The  $\mathbb{R}^{2^3} \rightarrow \mathbb{R}^{2^3}$  linear transformation involved is **spatial convolution**, where a 2D  $s * s$  filter (a.k.a. *kernel*) is multiplied to the 2D input image in subsets (*patches*) with size  $s * s$  and in defined steps (**strides**), then added up and then shifted by a bias (see Figure 3). For input data with several channels or *maps* (e.g.: RGB counts as 3 channels), the filter is applied to the same patch of each map and then added up into a single value of the output image (cumulative sum across maps). A map in Conv layers is the equivalent of a neuron in FC layers. We define:

- $\mathbf{A}_k^{[i]}$  as the map  $k$  of layer  $i$ ;
- $\mathbf{Z}_k^{[i]}$  as the linear output of map  $k$  of layer  $i$ ;
- ( $\mathbf{Z}_k^{[i]} = \mathbf{A}_k^{[i]}$  in absence of activation function)
- $b_k^{[i]}$  as the bias value for map  $k$  in layer  $i$
- $\mathbf{W}_k^{[i]}$  as the  $s * s$  filter/kernel for map  $k$ .

This operation can be vectorized by smartly replicating data [27]. The linear transformation can be expressed as:

$$\mathbf{Z}_k^{[i]} = \left( \sum_{m=0}^{M \text{ maps}} \mathbf{A}_m^{[i-1]} \oplus \mathbf{W}_k^{[i]} \right) + b_k^{[i]} \quad [14] \quad (11)$$

*Protecting Convolutional Layers* Convolution operation can be decomposed in a series of vectorized sums and multiplications over patches of size  $s * s$  :

$$\begin{aligned} \langle \mathbf{Z}_k^{[i]} \rangle_{\text{pub}} &= \left\langle \left( \sum_{m=0}^{M \text{ maps}} \mathbf{A}_m^{[i-1]} \oplus \mathbf{W}_k^{[i]} \right) + b_k^{[i]} \right\rangle_{\text{pub}} = \\ &= \sum_{m=0}^{M \text{ maps}} \langle \mathbf{A}_m^{[i-1]} \oplus \mathbf{W}_k^{[i]} \rangle_{\text{pub}} + \langle b_k^{[i]} \rangle_{\text{pub}} = \\ &= \left\{ \sum_{m=0}^M \langle \mathbf{A}_m^{[i-1]} [j] \rangle_{\text{pub}} * \langle \mathbf{W}_k^{[i]} \rangle_{\text{pub}} \right\}_{[s*s]} + \langle b_k^{[i]} \rangle_{\text{pub}} \end{aligned} \quad [14] \quad (12)$$

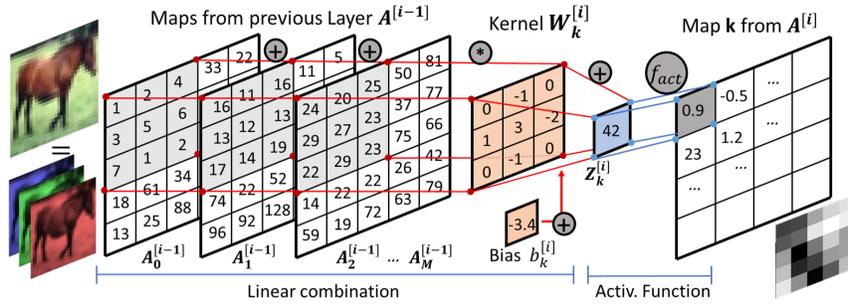


Fig. 3: Conv layer with activation for map  $k$  [14]

**Pooling Layer** This layer reduces the input size by using a packing function. Most commonly used functions are **max** and **mean**. Similarly to convolutional layers, pooling layers apply their packing function to patches (subsets) of the image with size  $s * s$  at strides(steps) of a defined number of pixels, as depicted in Figure 4.

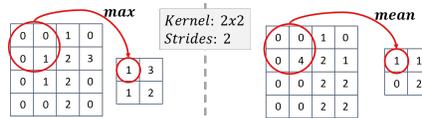


Fig. 4: Max and Mean packing for Pooling layers [14]

*Protecting Pooling layer* Max can be approximated by the sum of all the values in each patch of size  $s * s$ , which is equivalent to scaled *mean* pooling. *Mean* pooling can be scaled (sum of values) or standard (multiplying by  $1/N$ ). By employing a flattened input, pooling becomes easily vectorized.

### Other Techniques

- **Batch Normalization (BN)** reduces of the range of input values by 'normalizing' across data batches: subtracting mean and dividing by standard deviation. BN also allows finer tuning using trained parameters  $\beta$  and  $\gamma$  ( $\epsilon$  is a small constant used for numerical stability).

$$a_k^{[i+1]} = BN_{\gamma,\beta}(a_k^{[i]}) = \gamma * \frac{a_k^{[i]} - E[a_k^{[i]}}{\sqrt{Var[a_k^{[i]}] + \epsilon}} + \beta \quad [14] \quad (13)$$

**Protection of BN** is achieved by treating division as the inverse of a multiplication.

$$\begin{aligned} \langle a_k^{[i+1]} \rangle_{\text{pub}} &= \langle \gamma \rangle_{\text{pub}} * \left( \langle a_k^{[i]} \rangle_{\text{pub}} - \langle E[a_k^{[i]}] \rangle_{\text{pub}} \right) \\ &* \left\langle \frac{1}{\sqrt{Var[a_k^{[i]}] + \epsilon}} \right\rangle_{\text{pub}} + \langle \beta \rangle_{\text{pub}} \end{aligned} \quad [14] \quad (14)$$

- **Dropout and Data Augmentation** only affect training procedure. They don't require protection.
- **Residual Block** is an aggregation of layers where the input is added unaltered at the end of the block, thus allowing the layers to learn incremental ('residual') modifications (Figure 5).

$$\mathbf{A}^{[i]} = \mathbf{A}^{[i-1]} + ResBlock(\mathbf{A}^{[i-1]}) \quad (15)$$

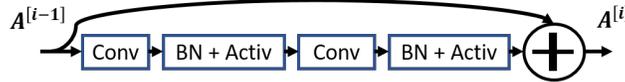


Fig. 5: Example of a possible Residual Block [14]

**Protection of ResBlock** is achieved by protecting the sum and the layers inside ResBlock:

$$\langle \mathbf{A}^{[i]} \rangle_{\text{pub}} = \langle \mathbf{A}^{[i-1]} \rangle_{\text{pub}} + \langle ResBlock(\mathbf{A}^{[i-1]}) \rangle_{\text{pub}} \quad [14] \quad (16)$$

### 3.2 Activation Function Protection

Due to their innate non-linearity, activation functions need to be approximated with polynomials to be encrypted with FHE. Several approaches have been elaborated in the literature. In [22] and [13], the authors proposed to use a square function as activation function. The last layer, a sigmoid activation function, is only applied during training.

Chabanne et al. used Taylor polynomials around  $x = 0$ , studying performance based on the polynomial degree [4]. In [18], Hesamifard et al. approximate instead the derivative of the function and then integrate to obtain their approximation.

Regardless on the approximation technique, we denote activation function  $f_{act}()$  approximation as

$$\mathbf{f}_{act}() \approx \mathbf{f}_{approxact}() \quad [14] \quad (17)$$

By construction, we have

$$\begin{aligned} \langle \mathbf{a}_k^{[i]} \rangle_{\text{pub}} &= \langle \mathbf{f}_{act}(\mathbf{z}_k^{[i]}) \rangle_{\text{pub}} \\ &\equiv \langle \mathbf{f}_{approxact}(\mathbf{z}_k^{[i]}) \rangle_{\text{pub}} \end{aligned} \quad [14] \quad (18)$$

- *Rectifier Linear Unit (ReLU)* is currently considered as the most efficient activation function for DL. Several variants have been proposed, such as Leaky ReLU[23], ELU[7] or its differentiable version *Softplus*.

$$\begin{aligned} ReLU(z) &= z^+ = \max(0, z) \\ Softplus(z) &= \log(e^z + 1) \end{aligned} \quad [14] \quad (19)$$

- *Sigmoid  $\sigma$*  The classical activation function. Its efficiency has been debated in the DL community.

$$Sigmoid(z) = \sigma(z) = \frac{1}{1 + e^{-z}} \quad [14] \quad (20)$$

- *Hyperbolic Tangent (tanh)* is currently being used in the industry because it is easier to train than ReLU: it avoids having any inactive neurons and it keeps the sign of the input.

$$tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \quad [14] \quad (21)$$

*Protecting Activation functions* Due to its innate non-linearity, they need to be approximated with polynomials. [13] proposed using only  $\sigma(z)$  approximating it with a square function. [4] used Taylor polynomials around  $x = 0$ , studying performance based on the polynomial degree. [18] approximate instead the derivative of the function and then integrate to obtain their approximation. One alternative would be to use Chebyshev polynomials.

## 4 Architecture

In this section we outline the architecture of our IP protection system, as depicted in Figure 6.

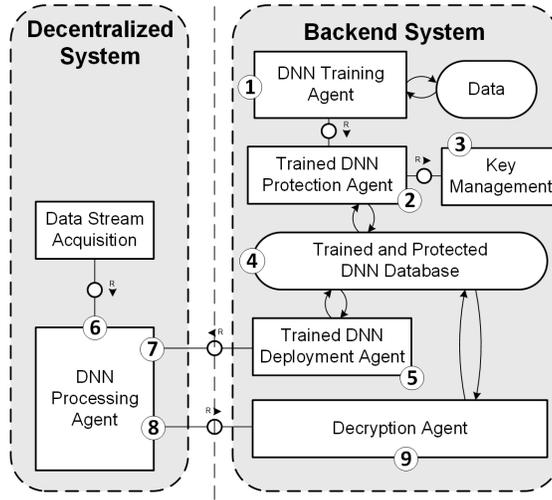


Fig. 6: Activity Diagram in our solution [14]

#### 4.1 Encryption of trained DNN

At backend-level, a DNN is trained by the *DNN Training Agent*, ①. Training outcome (NN architecture and parameters) is pushed to the *Trained DNN Protection Agent*, ③. Alternatively, an already trained DNN can be imported directly into the *Protection Agent*. The *DNN Protection Agent* generates a Fully Homomorphic key pair from the *Key Generator* component, ②. The DNN is then encrypted and stored together with its homomorphic key pair in the *Trained and Protected DNN Database*, ④.

#### 4.2 Deployment of trained and protected DNN

At the deployment phase, the *Trained DNN Deployment Agent* deploys the DNN on distributed systems, together with its public key, ⑤.

#### 4.3 DNN processing

On the distributed system, data is collected by a *Data Stream Acquisition* component, ⑦, and forwarded to the *DNN Processing Agent*, ⑥. Input layer does not involve any computation, and therefore can be seamlessly FHE encrypted as follows

$$\mathbf{X} \xrightarrow{\text{encryption}} \text{Enc}_{\text{pub}}(X) = \langle X \rangle_{\text{pub}} \quad [14] \quad (22)$$

Encrypted inferences are sent to the *Decryption Agent*, ⑧, for their decryption using the private key associated to the DNN, ⑨. FHE encryption propagates across the DNN layers, from input to output layer. By construction, output layer is encrypted homomorphically.

IP of the DNN, together with the computed inferences, is protected from any disclosure on the distributed system throughout the entire process.

The decryption of the last layer's output  $\mathbf{Y}$  is done with the private encryption key  $priv$ , as in standard asymmetric encryption schemes:

$$\langle \mathbf{A}^{[L]} \rangle_{pub} \xrightarrow{decryption} Dec_{priv} \left( \langle \langle \mathbf{A}^{[L]} \rangle_{pub} \rangle \right) = \mathbf{Y} \quad [14] \quad (23)$$

#### 4.4 Sequential Processes

**Encryption of Trained NN** Once a Neural Network is trained or imported, we encrypt all its parameters, using the Protected NN DataBase to store it and handle Homomorphic Keys (Figure 7).

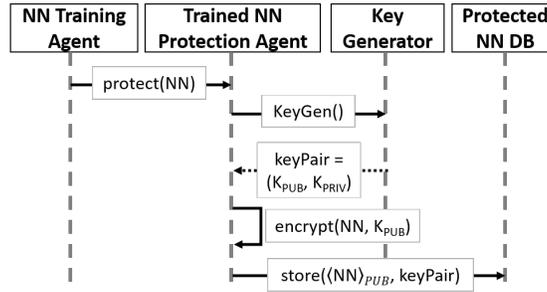


Fig. 7: Sequence diagram of Trained NN Encryption [14]

**Deploy Trained and Protected NN** The newly trained and protected deep neural network is deployed on the decentralized systems, including:

1. Network architecture;
2. Network model: Encrypted parameters;
3. Public encryption key.

**Encrypted Inference** On the decentralized system, data is collected and injected into the deployed NN. We must encrypt  $\mathbf{A}^{(0)} = \mathbf{X}$  with the public encryption key associated to the deployed NN (Figure 8).

**Inference Decryption** Encrypted inferences are sent to backend, together with an identifier of the NN used for the inference. The inference is homomorphically decrypted using the mapping private decryption key (Figure 9).

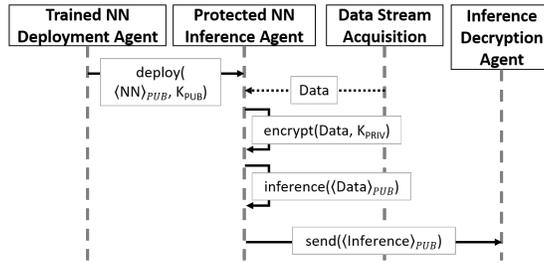


Fig. 8: Sequence diagram of inference processing [14]

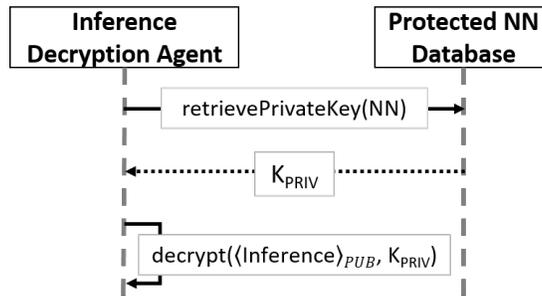


Fig. 9: Sequence diagram of inference decryption [14]

## 5 Evaluation

As detailed in Section 2.3, FHE introduces additional computational costs at each step of the DNN life-cycle. In this section, we evaluate performance overhead from computation time, memory load and disk usage perspectives at DNN model and processing encryption and output decryption.

### 5.1 Hardware Setup

As *backend*, we use a NVIDIA DGX-1<sup>6</sup> server, empowered with 8 Tesla V100 GPUs. This machine is theoretically not resource-constrained (computation & memory). We reasonably neglect the impact of the performance overhead introduced by FHE on DNN trained model encryption and output decryption.

We deploy and execute our encrypted DNN on a NVIDIA Jetson-TX2<sup>7</sup>. Powered by NVIDIA Pascal architecture, this platform embeds 256 CUDA cores, CPU HMP Dual Denver 22 MB L2 + Quad ARM® A572 MB L2, and 8 GB of memory. This platform gets closer to the hardware configuration of a Distributed Enterprise System.

<sup>6</sup> <https://www.nvidia.com/en-us/data-center/dgx-1/>

<sup>7</sup> <https://www.nvidia.com/en-us/autonomous-machines/embedded-systems-dev-kits-modules/>

## 5.2 Software Setup

*DNN Model* As demonstrated in Section 3, our approach is fully agnostic from NN topology, or implementation. For the sake of our evaluation, involving several modifications to the NN model, we choose a simple CNN classifier<sup>8</sup>, implemented with the Keras library<sup>9</sup>. Two datasets have been used in our experiment: CIFAR10<sup>10</sup>, for image classification, and MNIST<sup>11</sup> for handwritten digits classification.

As depicted in Figure 10, we distinguish two main parts in this CNN: a *feature extractor* and a *classifier*. The feature extractor reduces the amount of information from the input image, into a set of high level and more manageable features. This step facilitates the subsequent classification of the input data.

Composed of four layers,  $[FC \rightarrow ReLU \rightarrow FC \rightarrow Softmax]$ , the classifier categorizes the input data according to the extracted features, and outputs discrete probability distribution over 10 classes of objects.

As reference point, we evaluate key performance figures at model training and processing time without encryption. Once trained, the size of the CNN plaintext model is 9.6Mb. On Jetson TX2, single unencrypted image classification is computed on average in 89.1ms.

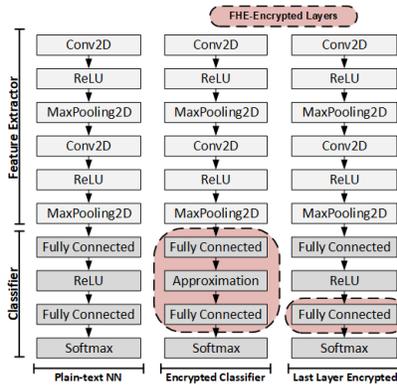


Fig. 10: Keras Convolutional Neural Network.

*FHE library* As introduced in section 2.3, several libraries are available for FHE. We use SEAL [29] C library from Microsoft Research running on CPU. This choice is motivated by the library's performance, support of multiple schemes such as BGV [3], stability, and documentation. The use of SEAL, implemented in C++, with the Keras

<sup>8</sup> [https://github.com/keras-team/keras/blob/master/examples/cifar10\\_cnn.py](https://github.com/keras-team/keras/blob/master/examples/cifar10_cnn.py)

<sup>9</sup> <https://keras.io>

<sup>10</sup> <https://www.cs.toronto.edu/~kriz/cifar.html>

<sup>11</sup> <http://yann.lecun.com/exdb/mnist/>

Python library requires some engineering efforts. To enable both fast performance of the native C++ library and rapid prototyping using Python, we use Cython<sup>12</sup>.

We conduct our evaluation with the BGV scheme [3], utilizing the integer encoding with SIMD support. To handle the floating-point DNN parameters, we use fixed-point arithmetic with a fixed scaling factor, similarly to CryptoNets[13]. This has no noticeable impact on the classification accuracy, if a suitable scaling factor is applied. The SIMD operations allow for optimized performance through vectorization.

### 5.3 Linearization

We tackle the problem of linearization of the ReLU functions following approaches: we approximate it with a modified square function, and we skip activation function. The modified square function  $x^2 + 2x$  (see Figure 13) is derived from the ReLU approximation proposed in [4]. In order to optimize the computation of that function on ciphertexts, we used simpler coefficients.

In order to evaluate the impact of these approaches, we trained the CNN on the CIFAR10 and MNIST datasets, replacing the last ReLU activation. Depicted in Figure 11 and Figure 12, we report the accuracy loss. Both approximations have merely a minor impact on the output classification accuracy.

Skipping the last activation function shows good results on this simple CNN, but we do not want to generalize to any other DNN or dataset.

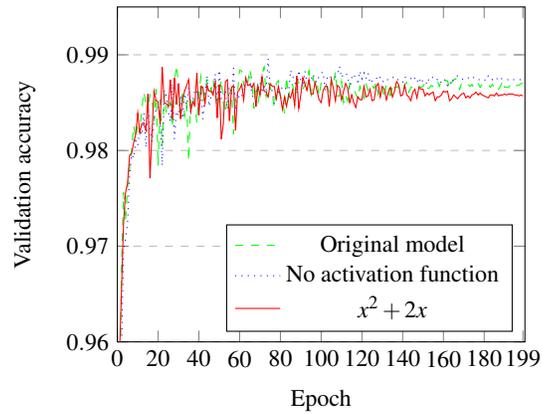


Fig. 11: Classification Accuracy with ReLU Approximation - MNIST Dataset.

### 5.4 Experimentation Results

**Model & Data Protection** Intellectual Property-wise, we consider the feature extractor as of minor importance, as CNNs generally use state of the art feature extractor. The IP

<sup>12</sup> <https://cython.org/>

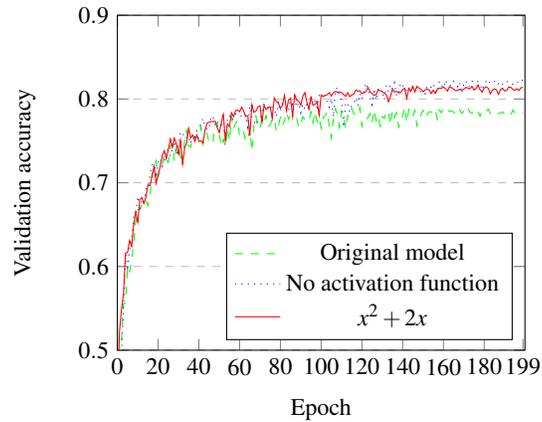


Fig. 12: Classification Accuracy with ReLU Approximation - CIFAR10 Dataset.

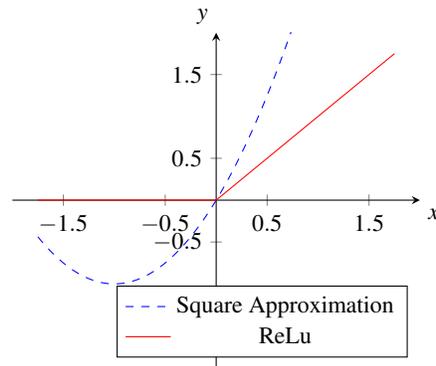


Fig. 13: ReLU Approximation as Square Function

of the model rather lies in the parameters, weights and bias, of the trained classifier. For that reason, we encrypt the classifier only, as a first step towards full model encryption, as depicted in Figure 10. To better understand the impact of computation depth, we also complete our evaluation with the encryption of the last FC layer only.

Confidentiality-wise, we evaluate the impact of extracted features encryption by comparing processing performance on an encrypted model with plaintext and encrypted feature extractor outputs.

As depicted in Figure 10, we evaluate our approach on three modified versions of the model:

- Last FC Layer Encrypted
- Full Classifier Encrypted with no Activation Function
- Full Classified Encrypted with our Modified Square Activation Function

Confidentiality-wise, we evaluate the impact of extracted features encryption by comparing processing performance on an encrypted model with plaintext and encrypted feature extractor outputs.

In order to optimize our approach, we omit the *Softmax* layer within the classifier. This layer does not have any influence on the classification results, as *Softmax* layer is mostly required at training phase, to normalize network outputs probability distribution, for more consistent loss calculations.

The overall experiment as described in section 4 has been applied 5 times on each model. We report average evaluation metrics for each step: model encryption, processing encryption and decryption.

**DNN Model Encryption** Each trained CNN model is encrypted on DGX-1’s CPU. In Table 2, we depict the resource consumption average on the following metrics:

- Time to Compute: Time to encrypt the model.
- Model Size: Size of resulting encrypted model.
- Memory Load: Overall memory usage for model encryption.

We target three security levels: 128, 192, and 256-bits. For each of those, we optimize SEAL parameters as introduced in section 2.3, maximizing performance, and minimizing leftover noise budget. Note that the security levels can have a counter-intuitive effect on performance, where for instance 192-bit security level might be faster than 128-bit security level. This can be explained by the fact that 128-bit security level offers more (unnecessary) noise budget, depending on the choice of FHE scheme parameters (e.g. plaintext modulus, polynomial degree). Therefore, we target a remaining noise budget as close as possible to zero.

Compared to the plaintext model size (9.6Mb), encrypted model size increases by a factor of 8,22 in the best case, up to 1173,33.

	Achieved Security Level (bits)								
	128	192	256	128	192	256	128	192	256
	<i>Full Classifier - <math>x^2 + 2x</math></i>			<i>Full Classifier - No Act.</i>			<i>Last Layer</i>		
<b>Time to Compute (s)</b>	256.7	191.5	212.6	86.9	78.0	96.0	3.4	3.4	7.3
<b>Model Size</b>	11G	6.4G	11G	2.5G	4.4G	2.5G	79Mb	79Mb	158Mb
<b>Memory (Mb)</b>	2257.9	1112.4	2389.5	557.9	459.2	566.1	300.6	301.4	324.2

Table 2. Model Encryption

**DNN Processing Encryption** The three encrypted CNN models deployed on Jetson-TX2 for CPU based encrypted processing. At this stage, we evaluate the following metrics

- Time to compute: Processing time for an encrypted classification.
- Memory: Memory usage for encrypted classification.

- **Remaining Noise Budget:** At the end of processing encryption, we evaluate the remaining noise budget, which determines if additional encryption operations could be performed on the output vector.

In Table 3 and Table 4, we depict the performance of encrypted processing with plaintext and encrypted previous layer outputs. We study the impact of confidentiality preservation of the preceding layer outputs. SEAL library supports secure computation over plaintext and ciphertext producing ciphertext. As a consequence, output of the last MaxPooling2D layer can be processed in FHE-encrypted Fully Connected layer. Secure computation between plaintext and ciphertext has a lower impact on performance.

We observe a slight performance improvement on time to compute and memory between 128 and 192-bit security level. This is due to the FHE parameters optimization as described in Section 5.4, where initial noise budget is oversized for 128-bit security level, which has a direct impact to performance.

Experiment results show that, depending on the level of achieved security, and targeted scenario, we can achieve at best encrypted classification in 2.1s (for 128 level security and only one layer encrypted). In the worst case, with encrypted input, full classifier encrypted with a modified square function as activation layer, 5627s (93 mins) is required for a single classification.

	Achieved Security Level (bits)								
	128	192	256	128	192	256	128	192	256
	<i>Full Classifier - <math>x^2 + 2x</math></i>			<i>Full Classifier - No Act.</i>			<i>Last Layer</i>		
<b>Time to Compute (s)</b>	287.2	174.6	1221.3	43.7	32.9	90.2	2.1	2.1	4.5
<b>Memory Load (Mb)</b>	4683.1	2924.3	4899.2	1162.5	869.2	2342.3	53.9	54.0	117.7
<b>Remaining Noise Budget</b>	221.6	80.2	88.8	91.4	23.8	16.2	58.8	20.2	60.2

Table 3. Runtime Encryption with Plaintext Input

	Achieved Security Level (bits)								
	128	192	256	128	192	256	128	192	256
	<i>Full Classifier - <math>x^2 + 2x</math></i>			<i>Full Classifier - No Act.</i>			<i>Last Layer</i>		
<b>Time to Compute (s)</b>	2902.6	1048.6	5627	367.1	272.2	835.9	19.2	19.2	40.2
<b>Memory Load (Mb)</b>	5047.8	5809	4906	2314.8	1733.7	4644.5	119.5	118.5	253
<b>Remaining Noise Budget</b>	206.0	65.0	76.0	79.0	11.0	3.0	45.0	10.0	52.0

Table 4. Runtime Encryption with Encrypted Input

**Decryption** Following our approach, encrypted output are decrypted by the backend, on DGX-1. We therefore consider decryption as not computationally expensive, compared to encryption. Results are available in Table 5.

	Achieved Security Level (bits)								
	128	192	256	128	192	256	128	192	256
	<i>Full Classifier - <math>x^2 + 2x</math></i>			<i>Full Classifier - No Act.</i>			<i>Last Layer</i>		
<b>Time to Compute (s)</b>	2.9	1.7	3.2	0.6	0.6	1.0	0.2	0.1	0.2
<b>Memory Load (Mb)</b>	963.8	397.4	2062.5	123.4	73.4	267.1	17.8	17.8	38.7

Table 5. Decryption - Performance

## 6 Conclusion

In this paper, we discuss and evaluate a holistic approach for the protection of distributed Deep Neural Network (DNN) enhanced software assets, i.e. confidentiality of their input & output data streams as well as safeguarding their Intellectual Property. On that matter, we take advantage of Fully Homomorphic Encryption (FHE). We evaluate the feasibility of this solution on a Convolutional Neural Network (CNN) for image classification.

Our evaluation on NVIDIA DGX-1 and Jetson-TX2 shows promising results on the CNN image classifier. Firstly, the impact of activation function approximation is negligible, with almost no accuracy loss on output classification probability. Most of the overhead is introduced at processing time, affecting computation time & memory consumption. Performances vary from 2.1s for an encrypted classification, with only 53.9Mb consumed memory, up to 1h33m with almost 5Gb of consumed memory. This requires a balancing between expected classification throughput, targeted security level and encryption depth of the model. Currently this approach would be unrealistic for the protection of DNN-based real-time analytics. Still, the Industry calls for numerous scenarios – such as predictive maintenance – matching the current performance of our approach.

As future work, we aim to improve the performance of our approach by different means: following the constant evolution of FHE, such as with the recent CKKS scheme[5], acceleration of FHE libraries on GPU based infrastructure or optimized vectorized operations on FHE encrypted data[1]. In addition, we foresee a deployment of our solution into a Smart City scenario for risk prevention in public spaces; while expanding our approach to different types of DNNs, and complete encryption of CNNs, including the feature extraction layers.

## References

1. BADAWI, A. A., CHAO, J., LIN, J., MUN, C. F., JIE, S. J., TAN, B. H. M., NAN, X., AUNG, K. M. M., AND CHANDRASEKHAR, V. R. The alexnet moment for homomorphic encryption: Hcnn, the first homomorphic CNN on encrypted data with gpus. *CoRR abs/1811.00778* (2018).
2. BOEMER, F., RATNER, E., AND LENDASSE, A. Parameter-free image segmentation with SLIC. *Neurocomputing* 277 (2018), 228–236.
3. BRAKERSKI, Z., GENTRY, C., AND VAIKUNTANATHAN, V. Fully homomorphic encryption without bootstrapping. Cryptology ePrint Archive, Report 2011/277, 2011.

4. CHABANNE, H., DE WARGNY, A., MILGRAM, J., MOREL, C., AND PROUFF, E. Privacy-preserving classification on deep neural network. *IACR Cryptology ePrint Archive 2017* (2017), 35.
5. CHEON, J. H., HAN, K., KIM, A., KIM, M., AND SONG, Y. Bootstrapping for approximate homomorphic encryption. *IACR Cryptology ePrint Archive 2018* (2018), 153.
6. CHILLOTTI, I., GAMA, N., GEORGIEVA, M., AND IZABACHÈNE, M. Tfhe: Fast fully homomorphic encryption over the torus. *Cryptology ePrint Archive, Report 2018/421*, 2018. <https://eprint.iacr.org/2018/421>.
7. CLEVERT, D.-A., UNTERTHINER, T., AND HOCHREITER, S. Fast and accurate deep network learning by exponential linear units (elus). *arXiv preprint arXiv:1511.07289* (2015).
8. CRAMER, R., DAMGÅRD, I. B., ET AL. *Secure multiparty computation*. Cambridge University Press, 2015.
9. DAI, W., AND SUNAR, B. cuhe: A homomorphic encryption accelerator library. In *BalkanCryptSec* (2015), vol. 9540 of *Lecture Notes in Computer Science*, Springer, pp. 169–186.
10. DUCAS, L., AND MICCIANCIO, D. Fhew: bootstrapping homomorphic encryption in less than a second. In *Annual International Conference on the Theory and Applications of Cryptographic Techniques* (2015), Springer, pp. 617–640.
11. FAN, J., AND VERCAUTEREN, F. Somewhat practical fully homomorphic encryption. *Cryptology ePrint Archive, Report 2012/144*, 2012.
12. GENTRY, C. *A Fully Homomorphic Encryption Scheme*. PhD thesis, Stanford University, Stanford, CA, USA, 2009. AAI3382729.
13. GILAD-BACHRACH, R., DOWLIN, N., LAINE, K., LAUTER, K., NAEHRIG, M., AND WERNING, J. Cryptonets: Applying neural networks to encrypted data with high throughput and accuracy. 201–210.
14. GOMEZ, L., IBARRONDO, A., MÁRQUEZ, J., AND DUVERGER, P. Intellectual property protection for distributed neural networks - towards confidentiality of data, model, and inference. In *Proceedings of the 15th International Joint Conference on e-Business and Telecommunications, ICETE 2018 - Volume 2: SECRYPT, Porto, Portugal, July 26-28, 2018*. (2018), P. Samarati and M. S. Obaidat, Eds., SciTePress, pp. 313–320.
15. GOODFELLOW, I. Security and privacy of machine learning. RSA Conference, 2018.
16. GRAEPEL, T., LAUTER, K., AND NAEHRIG, M. MI confidential: Machine learning on encrypted data, 2012.
17. HALEVI, S., AND SHOUP, V. Algorithms in helib. In *International cryptology conference* (2014), Springer, pp. 554–571.
18. HESAMIFARD, E., TAKABI, H., AND GHASEMI, M. Cryptodl: Deep neural networks over encrypted data. *CoRR* (2017).
19. IOFFE, S., AND SZEGEDY, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning* (2015), pp. 448–456.
20. KRIZHEVSKY, A., SUTSKEVER, I., AND HINTON, G. E. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems* (2012), pp. 1097–1105.
21. LIU, J., JUUTI, M., LU, Y., AND ASOKAN, N. Oblivious neural network predictions via minionn transformations. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security* (2017), ACM, pp. 619–631.
22. LIVNI, R., SHALEV-SHWARTZ, S., AND SHAMIR, O. On the computational efficiency of training neural networks. In *Advances in Neural Information Processing Systems 27*, Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, Eds. Curran Associates, Inc., 2014, pp. 855–863.
23. MAAS, A. L., HANNUN, A. Y., AND NG, A. Y. Rectifier nonlinearities improve neural network acoustic models. In *Proc. icml* (2013), vol. 30, p. 3.

24. MOHASSEL, P., AND ZHANG, Y. Secureml: A system for scalable privacy-preserving machine learning. In *Security and Privacy (SP), 2017 IEEE Symposium on (2017)*, IEEE, pp. 19–38.
25. PALISADE. The palisade lattice cryptography library, 2018.
26. PARLIAMENT, E., AND COUNCIL. General data protection regulation, 2016.
27. REN, J. S., AND XU, L. On vectorization of deep convolutional neural networks for vision tasks. In *AAAI (2015)*, pp. 1840–1846.
28. SCHATSKY, D., KUMAR, N., AND BUMB, S. Intelligent IoT, Bringing the power of AI to the Internet of Things. Deloitte Insights, 2017.
29. Simple Encrypted Arithmetic Library (release 3.1.0). <https://github.com/Microsoft/SEAL>, Dec. 2018. Microsoft Research, Redmond, WA.
30. SHOKRI, R., AND SHMATIKOV, V. Privacy-preserving deep learning. In *Proceedings of the 22nd ACM SIGSAC conference on computer and communications security (2015)*, ACM, pp. 1310–1321.
31. UCHIDA, Y., NAGAI, Y., SAKAZAWA, S., AND SATOH, S. Embedding watermarks into deep neural networks. In *Proceedings of the 2017 ACM on International Conference on Multimedia Retrieval (2017)*, ACM, pp. 269–277.