# slytHErin: An Agile Framework for Encrypted Deep Neural Network Inference

Francesco Intoci*[1], Sinem Sav*[1], Apostolos Pyrgelis[1], Jean-Philippe Bossuat[2], Juan Ramón Troncoso-Pastoriza[2], and Jean-Pierre Hubaux[1,2]

[1] EPFL, Lausanne 1015, Switzerland
name.surname@epfl.ch
[2] Tune Insight SA, Lausanne 1015, Switzerland
name@tuneinsight.com

**Abstract.** Homomorphic encryption (HE), which allows computations on encrypted data, is an enabling technology for confidential cloud computing. One notable example is privacy-preserving Prediction-as-a-Service (PaaS), where machine-learning predictions are computed on encrypted data. However, developing HE-based solutions for encrypted PaaS is a tedious task which requires a careful design that predominantly depends on the deployment scenario and on leveraging the characteristics of modern HE schemes. Prior works on privacy-preserving PaaS focus solely on protecting the confidentiality of the client data uploaded to a remote model provider, e.g., a cloud offering a prediction API, and assume (or take advantage of the fact) that the model is held in plaintext. Furthermore, their aim is to either minimize the latency of the service by processing one sample at a time, or to maximize the number of samples processed per second, while processing a fixed (large) number of samples. In this work, we present slytHErin, an agile framework that enables privacy-preserving PaaS beyond the application scenarios considered in prior works. Thanks to its hybrid design leveraging HE and its multiparty variant (MHE), slytHErin enables novel PaaS scenarios by encrypting the data, the model or both. Moreover, slytHErin features a flexible input data packing approach that allows processing a batch of an arbitrary number of samples, and several computation optimizations that are model-and-setting-agnostic. slytHErin is implemented in Go and it allows end-users to perform encrypted PaaS on custom deep learning models comprising fully-connected, convolutional, and pooling layers, in a few lines of code and without having to worry about the cumbersome implementation and optimization concerns inherent to HE.

---

* These authors contributed equally to this work.

arXiv:2305.00690v1 [cs.CR] 1 May 2023

## 1   Introduction

With recent advances in deep learning, cloud service providers expose trained deep neural networks (DNNs) to end-users for prediction-as-a-service (PaaS) through their application programming interfaces (APIs) [4, 5, 7, 25, 56]. For instance, Amazon Forecast enables business analytics by performing forecasting on client time-series data [3] and Azure's Cognitive summarizes and classifies financial documents [6]. However, PaaS applications raise privacy concerns as both the user data (e.g., client time-series, text, or health data) and the machine learning model (due to intellectual property concerns), can be sensitive information, and cloud service providers must comply with privacy regulations such as CCPA [12], GDPR [20], and HIPAA [30]. Thus, it is now needed more than ever to protect the privacy of the data used in PaaS applications.

To enable privacy-preserving PaaS, various works propose performing encrypted DNN inference by employing homomorphic encryption (HE) schemes which allow computations directly on ciphertexts [8, 10, 11, 16, 17, 22, 29, 31, 33, 37, 38, 41, 42, 49]. However, to cope with the computational overhead introduced by HE operations and to account for the characteristics of modern HE schemes, e.g., their support of Single Instruction, Multiple Data (SIMD) operations, these works rely on various optimizations which are tailored to specific PaaS scenarios, the most common of which comprises a cleartext DNN model and encrypted data. As a result, existing HE-based works cannot support emerging scenarios, e.g., edge machine learning [2, 44, 55], that require outsourcing the prediction to the client (while protecting the model's intellectual property), or privacy-preserving federated learning where inference is performed on a model that is trained in encrypted form by multiple data providers [52, 53, 57]. Moreover, these works rely on data packing schemes adapted to specific DNN architectures and application requirements, aiming either to minimize prediction latency (typically by processing one sample at a time, e.g., for real-time analytics), or to maximize the number of samples processed per second (usually by collecting and then processing in parallel a large number of samples leveraging on SIMD capabilities).

In this work, we design slytHErin, an agile framework for encrypted DNN inference. Built on HE and its multiparty variant, our framework can be adapted to various and novel PaaS scenarios where: (i) the client's data is encrypted while the model is in cleartext, (ii) the client's data is in cleartext and the model is encrypted, and (iii) both the client's data and the model are encrypted. Moreover, slytHErin features application- and model-agnostic optimizations which make it suitable for various settings. For instance, slytHErin implements an intuitive and flexible packing scheme that efficiently enables SIMD operations for *arbitrary* batch sizes, and generic optimizations for encrypted matrix operations. We implement slytHErin in Go and provide the building blocks that enable the encrypted execution of any DNN model composed of fully-connected, convolutional, and pooling layers. Contrary to prior works, our implementation is not centered around a system model, specific assumptions, or DNN architectures, making it a versatile tool for securing different PaaS pipelines. Our evaluation shows that slytHErin achieves accuracy similar to performing

inference on cleartext data and/or models. Moreover, it yields an interesting trade-off between latency and throughput, and its overall performance is on par with that of the state-of-the-art HE-based inference solutions, while being more flexible than specialized solutions. Our implementation can be found on `https://github.com/ldsec/slytHErin`.

## 2    Related Work

Given the potential privacy issues that might arise in PaaS, a number of works that build encrypted PaaS frameworks have been proposed. These works rely on homomorphic encryption (HE) and/or multiparty computation (MPC) to protect the confidentiality of both the ML model and the client's evaluation data during prediction [8, 10, 11, 15, 16, 22, 29, 31, 33, 37, 38, 41, 42, 46, 48, 49].

**HE-based Solutions.** Cryptonets was the first work in this research direction that enabled DNN evaluation on encrypted data using an HE scheme [22]. Its overhead, in terms of latency, was later improved by Brutzkus et al. which proposed novel approaches to represent the input data [11]. Other works focus on improving the efficiency of encrypted matrix operations [32, 39] or on designing novel techniques for the encrypted evaluation of more complex ML models such as graph convolutional networks [47]. The latter has been used in downstream tasks such as human action recognition [34] achieving better latency than [11]. Other works develop compilers that ease the deployment of trained ML models with HE libraries, e.g., SEAL [54], HElib [27], or Palisade [50], for encrypted inference. Boemer et al. [8,9] build a graph compiler for SEAL that simplifies the use of a model trained with Tensorflow [1] or PyTorch [45] for encrypted PaaS. CHET, on the other hand, is a domain-specific optimizing compiler that allows the specification of tensor circuits suitable for HE-based DNN inference [18]. All of these works propose specific input data representations (packing) and optimizations for either latency or throughput for specific scenarios (e.g., featuring a cleartext model vs. encrypted data) and DNN architectures. Moreover, to cope with DNN non-linear operations that are not supported by HE schemes, e.g., activations, they either use interactions with the client [8], modify their functionality to low-degree polynomial functions [11, 18, 22, 34], or use polynomial approximations [13, 47].

**Hybrid Approaches.** To ease the encrypted execution of non-linear functions, some works rely on hybrid approaches combining two-party computation with HE [31,33,37,48], or secret sharing with garbled circuits [42,46,49]. For instance, Liu et al. [37] utilize HE for matrix multiplications and garbled circuits for the non-linear activations. Juvekar et al. [33] employ HE for matrix-vector multiplication and convolution operations and garbled circuits for comparisons which are widely used in activation functions. Similarly, we provide a hybrid framework for privacy-preserving PaaS that supports a wide range of applications by relying on a multiparty variant of HE. Moreover, thanks to our generic data representation scheme and optimizations, our framework is agnostic of the DNN architecture

and parameters such as batch size, while achieving on par performance with the state-of-the-art.

## 3   Background

### 3.1   Homomorphic Encryption

Homomorphic encryption (HE) schemes enable the execution of arithmetic operations directly on ciphertexts, i.e., without requiring decryption; this makes them ideal candidates for privacy-preserving machine learning inference applications. In this work, we employ the Cheon-Kim-Kim-Song (CKKS) scheme [14], which is suitable for machine learning tasks as it enables approximate arithmetic over $\mathbb{C}^{\mathcal{N}/2}$ (hence, over real values as well). The ring $R_{Q_L} = \mathbb{Z}_{Q_L}[X]/(X^{\mathcal{N}} + 1)$ of dimension $\mathcal{N}$ with coefficients modulo $Q_L = \prod_{i=0}^{L} q_i$ defines the plaintext and ciphertext spaces, hence both plaintexts/ciphertexts are represented by polynomials of degree $\mathcal{N} - 1$ whose coefficients encode a vector of $\mathcal{N}/2$ values. The security of CKKS is based on the ring learning with errors problem [40]. CKKS supports the homomorphic evaluation of operations such as additions, multiplications, and rotations, and any operation is simultaneously performed on all encoded values, hence offering Single Instruction, Multiple Data (SIMD). Non-linear operations, e.g., comparisons, are supported via polynomial approximations, introducing a computation overhead versus accuracy tradeoff. CKKS is a leveled HE scheme, i.e., an $L$-depth circuit can be evaluated before the ciphertext is exhausted. Then, a costly procedure, called bootstrapping [21], is required to refresh the exhausted ciphertext and enable more operations on it. We refer to the traditional bootstrapping operation (performed by a single party) as *centralized bootstrapping*.

**Multiparty Homomorphic Encryption (MHE).** To make our framework adaptable to various PaaS scenarios (see Section 4), we also rely on a multiparty variant of the CKKS scheme [43]. In the multiparty homomorphic encryption (MHE) scheme, a set of parties (e.g., model-providers) collectively generate a public key while the corresponding secret key is secret-shared among them. This setting enables secure collaboration between $N$ parties, as parties use the collective public key to encrypt their inputs and perform joint operations on them using the MHE scheme. The result decryption by the client, however, requires the participation of all parties. Hence, this scheme ensures confidentiality under a passive adversary model with up to $N - 1$ collusions. Moreover, the multiparty CKKS scheme offers efficient multiparty computation protocols. For instance, it enables a collective bootstrapping operation, where the costly centralized bootstrapping which homomorphically evaluates the decryption and consumes many levels, is substituted by a lightweight one-round interactive protocol ($\mathsf{CBootstrap}(\cdot)$) which does not consume levels. Moreover, the scheme supports ($\mathsf{CKeySwitch}(\cdot)$), a collective key-switch operation which can change the encryption key of a ciphertext.

### 3.2   Deep Neural Networks

Deep neural networks (DNNs) are able to model complex non-linear relationships and find applicability in various domains such as computer vision. A DNN consists of multiple hidden layers between the input and output layers. Our framework enables the encrypted evaluation of DNNs comprising fully connected (FC), convolutional (Conv), and pooling (Pool) layers. We succinctly present the functionalities of these layer types:

**Fully Connected layer**: Given an input vector $\mathbf{x}$, a weight matrix $\mathbf{W}$ and a bias vector $\mathbf{b}$, a FC-layer computes $\mathbf{xW^T} + \mathbf{b}$.

**Convolutional layer**: Given an input tensor (e.g., an image) $\mathbf{X}$ with $c_i$ channels of dimensions $w \cdot h$ and a set of $c_o$ kernels $\mathbf{K}$ each made up of $c_i$ filters of size $f_w \cdot f_h$, a Conv-layer computes a tensor $\mathbf{O}$ with $c_o$ channels. Each channel $\mathbf{O}_i$ is computed as $\sum_{n=0}^{c_i} \mathbf{X}_n * \mathbf{K}_{i,n}$, with $\mathbf{X}_n$ the n-th channel of the input image, $\mathbf{K}_{i,n}$ the n-th filter of the i-th kernel, and $*$ the cross-correlation operator.

**Pooling layer**: It performs dimensionality reduction on the input. The most common types are SumPooling, AveragePooling, and MaxPooling, where the feature-map is the sum, average, and the maximum of the features in a region of the input, respectively. Max-Pooling requires non-linear operations, i.e., comparisons, which are non-trivial to implement under encryption, thus we only consider the first two types.

Each layer can be paired with an activation function which is evaluated on its output. The output of the DNN's last layer is the prediction result (output).

## 4   slytHErin Overview

Building on the CKKS HE scheme and its multiparty variant (see Section 3), we design a framework that is flexible for various encrypted PaaS scenarios (Figure 1). We first describe the involved entities before detailing slytHErin's objectives and workflow for each PaaS scenario.

- **Model-provider(s):** This entity (one or more) has trained an ML model and exposes it to end-users for queries (PaaS) through a prediction API hosted on a cloud service provider.
- **Client:** This entity is a user of the PaaS that inputs its own sensitive data which is evaluated on the model exposed by the model-provider. The client obtains the output of the PaaS process, i.e., the prediction.

We consider that the client and the model-provider are *honest-but-curious*, i.e., they follow the protocol specification, but they might try to infer information about each other's data. slytHErin's objective is to protect both the confidentiality of the client's and the model-provider's data. In particular, the model-provider should not learn any information about the client's evaluation data and the prediction result, whereas the client should not obtain any knowledge about the model beyond what can be inferred from the PaaS output.
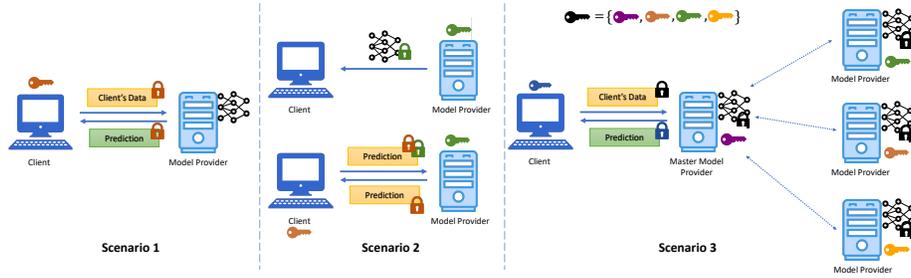
Fig. 1: Encrypted PaaS scenarios enabled by `slytHErin`. Encryption is depicted with a lock whose color is the same as the corresponding secret key. The black key (rightmost figure) corresponds to the model-providers' collective key. **Scenario 1:** The client sends its encrypted data to the model-provider that evaluates it on the plaintext model. **Scenario 2:** The encrypted model is sent to the client for evaluation on its cleartext data. **Scenario 3:** The client sends encrypted data to a cohort of model-providers that retain an encrypted model.

### 4.1   Scenario 1: Encrypted client data - Cleartext model

This is the traditional HE-based PaaS setting, where a client encrypts its data with its own public key and sends the ciphertext to the model-provider that stores its ML model in plaintext form. The model-provider evaluates its model on the client's encrypted data – without interacting with the client – and returns the encrypted prediction to the client. The client decrypts the ciphertext with its secret key and obtains the prediction result. In this scenario, the client's data confidentiality is ensured as its inputs are encrypted throughout the DNN evaluation and the model-provider does not learn the prediction result. The model confidentiality is protected as the model remains on the model-provider's side. Scenario 1 represents a typical PaaS setting, where a model-holder exposes a prediction service that receives sensitive data as inputs [11, 16, 22, 33]. For instance, imagine a health-care insurance provider that uses its customer data and trains a DNN that predicts the probability of patient re-admission to a hospital. The model is exposed through an API to clients (e.g., hospitals) who wish to obtain predictions about their own cohorts of patients. However, hospitals cannot share their patient data with third-parties due to ethical and data privacy requirements, hence, `slytHErin` could be an enabler for such a service as it ensures data confidentiality.

### 4.2   Scenario 2: Cleartext client data - Encrypted model

In this scenario, the model-provider outsources the computation of the prediction to the client. However, the model is an intellectual property that needs to be protected. Thus, the model-provider encrypts its model with its own public key and sends it to the client in encrypted form. The client evaluates the

encrypted model on its own (plaintext) data and obtains an encrypted prediction. Finally, the client sends the prediction ciphertext to the model provider, which obliviously decrypts the result and communicates it back to the client (Section 5.6). The client's data confidentiality is ensured as its evaluation data is never transferred and the model-provider does not learn the prediction result due to the oblivious decryption phase. The model confidentiality is protected as the model is encrypted with the model-provider's public key. Scenario 2 is suitable for applications that require outsourcing a trained model to the client side for predictions. For instance, this could be the case for model trading platforms that offer a *try-before-you-buy* option, where customers locally test the performance of an ML model on their data before purchasing it. Another relevant application is model outsourcing to edge devices [2, 44], e.g., mobile phones or smartwatches, that monitor their owners' activity and provide feedback to them through predictions, e.g., health recommendations or activity tracking [55]. We note that this is a novel PaaS scenario enabled by slytHErin.

### 4.3   Scenario 3: Encrypted client data - Encrypted model

In this scenario, we assume that the model-provider is represented by a cohort of $N$ nodes that have collectively trained a DNN on their joint data with a state-of-the-art encrypted collaborative learning framework [52, 53, 57]. For this, we rely on a multiparty variant of homomorphic encryption (MHE). In particular, the nodes (model-providers) generate a collective public key (black key in Figure 1, **Scenario 3**) whose corresponding secret key is secret-shared among them (colored keys in Figure 1, **Scenario 3**). We assume that the nodes collectively train a DNN on their data and retain it under encryption for PaaS to mitigate model-targeting attacks and protect its intellectual property. For this scenario, the client encrypts its evaluation data with the collective public key and a master node from the cohort performs the prediction (with both the model and the data encrypted) with the assistance of the other nodes for collective interactive operations (e.g., ciphertext refresh – CBootstrap($\cdot$), Sections 3.1 and 5.6). Finally, the ciphertext storing the prediction result is re-encrypted (i.e., CKeySwitch($\cdot$), Sections 3.1 and 5.6) under the public key of the client which decrypts it to obtain the prediction. In this case, both the model and the client data are encrypted with the cohort's collective public key, hence, their confidentiality is ensured as long as one of the cohort nodes is honest and does not participate in decryption. The confidentiality of the prediction output is protected, as only the client can decrypt it. Scenario 3 is suitable for PaaS applications after a model-provider outsources the model training procedure to a cohort of $N$ nodes that leverage on distributed learning techniques for improved efficiency or after a federation of $N$ model-providers, each with their own data, uses a state-of-the-art framework to train a collective ML model *under encryption* [19, 52, 53]. We note that previous works that focused on encrypted DNN inference do not support (or implement) inference on encrypted models or collaborative functionalities such as bootstrapping or re-encryption.

## 5   Cryptographic Building Blocks

We describe `slytHErin`'s underlying cryptographic building blocks that make it flexible and efficient for different encrypted PaaS scenarios (Section 4) and various DNN architectures. We first introduce the data packing approach adopted to encode/encrypt the input data (Section 5.1). Then, we describe the algorithms used to evaluate fully-connected, convolutional, and pooling layers under encryption in Sections 5.2 and 5.3, respectively. We also present several optimizations that `slytHErin` implements (Section 5.4) and how non-linear activation functions are evaluated (Section 5.5). Finally, we present the multiparty computation protocols which allow `slytHErin` to support novel PaaS scenarios (Section 5.6).

### 5.1   Input Data Packing

Modern homomorphic encryption schemes can encode (pack) a vector of values into one ciphertext, thus enabling SIMD operations via the parallel computation of a function on all ciphertext slots. Designing an efficient packing scheme is crucial, yet challenging, due to the costs of re-arranging the ciphertext slots via rotations. Prior work on encrypted DNN inference [11, 22, 31, 33, 34, 37, 48] designed efficient packing schemes but these are tailored to specific system models and assumptions (e.g., the client's availability for the evaluation of certain operations). `slytHErin` employs a simple yet generic data packing scheme that is agnostic of the encrypted PaaS scenario and also flexible in terms of batch size that results in optimized latency and throughput. Given a batch consisting of $n$ input samples each with $d$ features, a naive approach is to encrypt/encode each feature of an input sample separately, yielding an inefficient execution due to the high number of ciphertexts/plaintexts. To leverage on SIMD operations and enable efficient encrypted inference, we flatten the batch and encrypt/encode all values in a single ciphertext/plaintext. For an input sample represented by a tensor of size $h \times r \times c$ (where, e.g., for an image, $h$ is the number of channels, while $r$ and $c$ represent the size of the pixel matrix of each channel), we encrypt/encode a batch of size $n$ in a tensor of size $n \times h \times r \times c$ as follows: First, we row-flatten ($\mathsf{RowFlatten}(\cdot)$) each of the $n$ tensors, such that the batch-tensor is transformed into a matrix of size $n \times d$, with $d = h \times r \times c$. This is done by iterating through all the channels of the input, by row-flattening the corresponding 2D matrix, and by horizontally stacking their flattened representation. The $n \times d$ matrix is then transposed and row-flattened ($\mathsf{TensorFlatten}(\cdot)$), thus yielding a vector of size $m = d \times n$. Our packing scheme requires that $m \leq s$, where $s$ is the ciphertext capacity (i.e., $s = \mathcal{N}/2$ for CKKS) and if that is not possible, we employ block matrix arithmetic optimizations (see Section 5.4).

### 5.2   Matrix Multiplication

To support the evaluation of fully-connected layers under encryption, `slytHErin` relies on the following matrix multiplication algorithm. Given two encrypted matrices, $\mathbf{A}$ and $\mathbf{W}$, where $\mathbf{A}$ is of size $n \times d$ and $\mathbf{W}$ of size $d \times h$, `slytHErin`

$$\begin{bmatrix} a_{00} & a_{01} & a_{02} \\ a_{10} & a_{11} & a_{21} \\ a_{20} & a_{21} & a_{22} \end{bmatrix} \times \begin{bmatrix} w_{00} & w_{01} & w_{02} \\ w_{10} & w_{11} & w_{21} \\ w_{20} & w_{21} & w_{22} \end{bmatrix} = \begin{bmatrix} b_{00} & b_{01} & b_{02} \\ b_{10} & b_{11} & b_{21} \\ b_{20} & b_{21} & b_{22} \end{bmatrix}$$

$$[w_{00}, w_{00}, w_{00}, w_{11}, w_{11}, w_{11}, w_{22}, w_{22}, w_{22}] \odot [a_{00}, a_{10}, a_{20}, a_{01}, a_{11}, a_{21}, a_{02}, a_{12}, a_{22}]$$
$$+[w_{10}, w_{10}, w_{10}, w_{21}, w_{21}, w_{21}, w_{02}, w_{02}, w_{02}] \odot [a_{01}, a_{11}, a_{21}, a_{02}, a_{12}, a_{22}, a_{00}, a_{10}, a_{20}]$$
$$+[w_{20}, w_{20}, w_{20}, w_{01}, w_{01}, w_{01}, w_{12}, w_{12}, w_{12}] \odot [a_{02}, a_{12}, a_{22}, a_{00}, a_{10}, a_{20}, a_{01}, a_{11}, a_{21}]$$
$$=[b_{00}, b_{10}, b_{20}, b_{01}, b_{11}, b_{21}, b_{02}, b_{12}, b_{22}]$$

Same format as input

Fig. 2: Multiplication of two matrices $\mathbf{A}$ and $\mathbf{W}$ of size $3 \times 3$.

implements their multiplication following the diagonal approach of [28]. First, $\mathbf{W}$ is represented by its *generalized diagonals* [28], where the element *i,j* of the diagonal is: $d_{i,j} = \mathbf{W}_{(i+j) \bmod d, j}$. Additionally, we replicate $n$ times the element $d_{i,j}$. The matrix multiplication, then, can be evaluated as follows:

$$\mathbf{A} \times \mathbf{W} = \sum_{i=1}^{d} \mathbf{d_i} \odot \mathsf{RotateCyclic}_{d \times i}(\mathsf{RowFlatten}(\mathbf{A}^T))$$

where $\mathsf{RotateCyclic}_k(\mathbf{v})$ represents a cyclic rotation of the values in $\mathbf{v}$ by $k$ positions to the left and $\odot$ represents the Hadamard product. Figure 2 represents a multiplication of two $3 \times 3$ matrices with this algorithm.

### 5.3  Convolutional and Pooling Layers

To evaluate convolutional layers under encryption, `slytHErin` represents the convolution operation as a matrix multiplication by expressing the filter as a *Toeplitz* matrix [23, 26]. For ease of presentation, consider a toy-example with a convolution between a single-channel input $\mathbf{I} \in R^{3 \times 3}$ and a filter $\mathbf{h} \in R^{2 \times 2}$ operating on the input with unitary stride and no padding. We can compute the convolution as: $\mathbf{O} = \mathsf{TensorFlatten}(\mathbf{h} * \mathbf{I})^T = \mathbf{h}' \times \mathbf{I}'$ where $\mathbf{I}' = \mathsf{TensorFlatten}(\mathbf{I})^T$ and $\mathbf{h}' = \mathcal{T}(\mathbf{h})$ for a function $\mathcal{T}$ that returns a Toeplitz matrix [23] as follows:

$$\mathbf{h}' = \begin{pmatrix} h_{1,1} & h_{1,2} & 0 & h_{2,1} & h_{2,4} & 0 & 0 & 0 & 0 \\ 0 & h_{1,1} & h_{1,2} & 0 & h_{2,1} & h_{2,4} & 0 & 0 & 0 \\ 0 & 0 & 0 & h_{1,1} & h_{1,2} & 0 & h_{2,1} & h_{2,4} & 0 \\ 0 & 0 & 0 & 0 & h_{1,1} & h_{1,2} & 0 & h_{2,1} & h_{2,4} \end{pmatrix}$$

Note that computing $\mathbf{O^T} = \mathbf{I'^T} \times \mathbf{h'^T}$ allows us to utilize the matrix multiplication algorithm and the input data packing protocol of Sections 5.2 and 5.1, respectively. Moreover, $\mathbf{O^T}$ is a valid input to any subsequent layer in the DNN

architecture, without requiring any re-packing, hence avoiding the cost of slot re-arrangement. slytHErin generalizes this method for convolutional layers with $k$ kernels, each with $m$ filters, and $n$ inputs with $m$ channels. slytHErin also supports SumPooling and AveragePooling layers: these are evaluated by treating them as convolutional layers, and employing the method previously described.

### 5.4   Optimizations

**Complex-Number Trick.** To optimize the input data packing scheme (Section 5.1), slytHErin employs the complex-number trick [51]: Since the CKKS plaintext space is $\mathbb{C}^{\mathcal{N}/2}$, we can leverage the imaginary part of complex numbers and pack (up to) two values in one plaintext slot. This allows us to effectively perform the multiplication and sum of two values with just one multiplication. As a toy example, let us consider the vectors: $\mathbf{a} = (a_1, \dots), \mathbf{b} = (b_1, \dots), \mathbf{c} = (c_1, \dots)$, and $\mathbf{d} = (d_1, \dots)$. To compute $\mathbf{a} \odot \mathbf{c} + \mathbf{b} \odot \mathbf{d} = (a_1 c_1 + b_1 d_1, \dots)$, we compress the first two and the two last vectors each into one vector with the following complex representation: $\mathbf{g} = (a_1 + i b_1, \dots)$, $\mathbf{h} = (c_1 - i d_1, \dots)$. Then, $\mathbf{g} \odot \mathbf{h} = (a_1 c_1 + b_1 d_1 + ie, \dots)$ for some value $e$, and the real part of the result can be extracted with complex conjugation, addition and constant multiplication. We apply this technique to the input matrix $\mathbf{A}$ and to the weight matrix $\mathbf{W}$. In particular, we embed pairs of adjacent columns of $\mathbf{A}$ into one column, i.e., column $k$ is paired with column $k + 1 \bmod d$, where $d$ is the number of columns, hence the entry $\mathbf{A}_{(k,j)}$ becomes $\mathbf{A}_{(k,j)} + i\mathbf{A}_{(k,j+1)}$. For $\mathbf{W}$, we compress the pairs of adjacent *diagonals* into one, padding with an extra 0-diagonal if the number of diagonals is odd. The newly packed matrix $\tilde{\mathbf{W}}$ has $\lceil \frac{d}{2} \rceil$ diagonals instead of $d$, reducing the complexity of the matrix multiplication algorithm by a factor of 2.

**Block Matrix Arithmetic.** When the size of the input batch exceeds the ciphertext capacity, slytHErin employs block-matrix arithmetic [52]. The input matrix $\mathbf{A}$ of size $n \times d$, is represented as a block-matrix $\bar{\mathbf{A}}$ of size $q \times p$, i.e., a matrix consisting of *blocks* (or sub-matrices) of size $\frac{n}{q} \times \frac{d}{p}$ for some divisors $q$ and $p$ of $n$ and $d$, respectively. Similarly, the weight matrix $\mathbf{W}$ of size $d \times h$ is partitioned to enable the multiplication $\bar{\mathbf{O}} = \bar{\mathbf{A}} \times \bar{\mathbf{W}}$ under two constraints: (i) $\bar{\mathbf{W}}$ must have $p$ row partitions, and (ii) every inner block $\mathbf{W}_{k,j}$ must be compatible for matrix multiplication with the inner blocks $\mathbf{A}_{i,k}$. $\bar{\mathbf{O}}$ is a block-matrix of size $n \times h$ with $q$ row partitions and $m$ column partitions (and $m$ the number of column partitions of $\bar{\mathbf{W}}$). Each block $\mathbf{O}_{i,j}$ is computed as: $\mathbf{O}_{i,j} = \sum_{k=1}^{p} \mathbf{A}_{i,k} \mathbf{W}_{k,j}$. Hence, by choosing suitable partitions, each matrix inner block is small enough to be encrypted/encoded independently following the input data packing and the generalized-diagonals approach described earlier (Sections 5.1 and 5.2). Figure 3 represents the encryption of matrix $\mathbf{A}$ with $2 \times 2$ partitioning. Then, the matrix multiplication between two large matrices is evaluated as a series of sums and multiplications between these smaller blocks. Given a model to evaluate (i.e., the dimensions of its layers), the number of input features, and a set of CKKS parameters, slytHErin follows a heuristic-based approach to automatically find the best batch size and partition strategy. In more detail, slytHErin explores
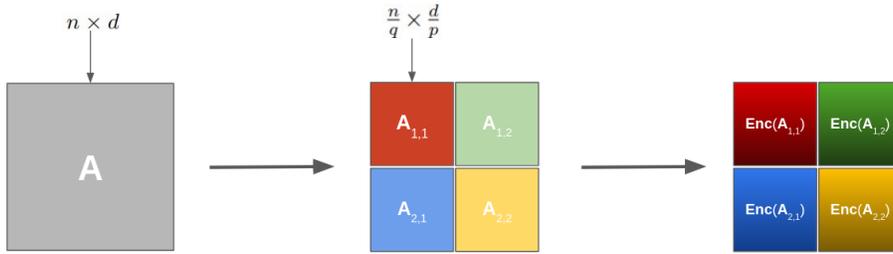
Fig. 3: Partitioning of an input matrix $\mathbf{A}$ in a $2 \times 2$ block matrix.

the space of possible splits, starting from divisors of the number of samples (if provided by the user) or divisors of the features dimension, and picks the split sequence and batch size that minimize the overall complexity of the pipeline in terms of homomorphic operations (i.e., it minimizes the number of homomorphic multiplications required to evaluate the model), thus optimizing throughput. In any case, the user can also declare a customized batch size which overrides the optimized batch size, and let `slytHErin` operate with a sub-optimal block matrix representation. An advantage of the block matrix arithmetic approach is that it is amenable to parallelization: Given $q \times p \times m$ threads, the matrix multiplication between two blocks $\mathbf{A}_{i,k}$ and $\mathbf{W}_{k,j}$ can be delegated to each thread, while using $q \times m$ of them to combine the individual results. Moreover, for a given set of cryptographic parameters and the corresponding evaluation keys, the client does not need to regenerate the keys for the evaluation of arbitrary size matrices, which is a computationally intensive task.

### 5.5 Non-Linear Operations

As non-polynomial functions, e.g., comparisons, are not computable under HE, some works modify common activation functions (e.g., ReLU) with simple polynomial functions [22] (e.g., $x^2$), or use polynomial approximations [53]. `slytHErin` employs the second approach and relies on *Chebychev interpolants* to approximate any Lipschitz continuous function on any finite real interval.

### 5.6 Multiparty Computation Protocols

We remind that `slytHErin` relies on CKKS and its multiparty variant (MHE) which enables interactive functionalities such as CBootstrap($\cdot$) for collective bootstrapping and CKeySwitch($\cdot$) for collective key-switching. The latter enables changing the encryption key of a ciphertext. In **Scenario 3**, the model-providers rely on these functionalities to refresh the ciphertexts noise and to change the encryption key of the prediction result, so that only the client can decrypt it.

We also design and implement an oblivious decryption protocol ObvDec($\cdot$), for **Scenario 2** (Section 4.2). In this protocol, the client masks its prediction result (encrypted under the model provider secret key) with an encryption of 0

under an ephemeral secret key, and sends the result to the model provider, which can remove one layer of encryption from the result (by invoking the decryption procedure of CKKS), without exposing the underlying plaintext. The result is finally sent to the client that unmasks it.

## 6  Experimental Evaluation

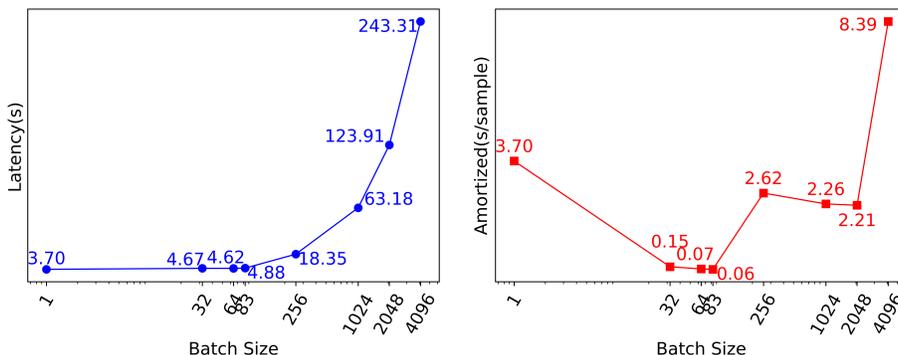### 6.1  Implementation and Experimental Setup

We implemented `slytHErin` in Go [24], using Lattigo as the cryptographic library [35]. Our implementation is modular, reusable, and easy to adapt to several PaaS applications. Detailed documentation can be found along with our source code on `https://github.com/ldsec/slytHErin`. We evaluate `slytHErin` using the following DNN architectures:

- **NN5**: A 5-layer convolutional neural network described in [22] for which we replace the square activation function with a degree 2 Chebyshev approximation of Softplus.
- **NN20**: A 20-layer DNN composed of convolutional and fully connected layers described in [16] ($\sim$754K model parameters) for which we replace the activation functions with a degree-63 approximation of SiLU and train it with the MSE loss function.
- **NN50**: Similar to **NN20** but comprising 50 layers ($\sim$1M model parameters [16]).

We use the **MNIST** dataset [36] for encrypted image classification, as it is the de-facto benchmark dataset used in prior work for privacy-preserving inference tasks [8, 9, 11, 16, 17, 22, 33]. All models were trained from scratch, achieving similar accuracy to the original works (and with minimal accuracy loss in the encrypted inference, none for **NN5**, approximately $\sim$0.13% for **NN20**, and $\sim$2% for **NN50**). The CKKS parameters are configured to achieve 128-bit security. For the multiparty interactive protocols, we deploy `slytHErin` on a local cluster with an average network delay of 20ms and 1Gbps bandwidth. All experiments were executed on machines running Ubuntu 22.04, with 12-core Intel Xeon E5-2680 2.5 GHz CPUs and 256GB RAM DDR4. The results are averaged over 3-5 runs.

### 6.2  Empirical Results

We first demonstrate how `slytHErin` supports different batch sizes by evaluating **NN5** on **Scenario 1** (Section 6.2.1). We also compare `slytHErin` with prior work on private PaaS as **NN5** is the predominantly used benchmark. Then, in Section 6.2.2, we evaluate **NN20** on **Scenario 3** to discuss `slytHErin`'s scalability aspects with the number of model-providers. Finally, we demonstrate `slytHErin`'s application and model agility by evaluating the more complex model **NN50** in all scenarios of Section 4 (Section 6.2.3).

(a) slytHErin's latency for **NN5** and different batch sizes.

(b) slytHErin's amortized runtime for **NN5** and different batch sizes.

Fig. 4: slytHErin's amortized runtime with different batch sizes.

**6.2.1   Elastic Data Packing.** We demonstrate the benefits of our packing approach (Section 5.1), by benchmarking **NN5** [22] in the traditional PaaS setting (**Scenario 1**) for various batch sizes. For this experiment, slytHErin heuristically estimates the optimal batch size for throughput at 83, as described in Section 5.4; this is experimentally confirmed by observing Figures 4a and 4b. In particular, Figure 4a shows slytHErin's latency for varying batch sizes up to 4,096 in semi-log scale. We observe a linear increase in latency after the optimal size. This is expected, as slytHErin automatically splits batch sizes larger than the optimal size into sub-batches of optimal size, and processes them sequentially. Figure 4b shows the amortized runtime of slytHErin for variable batch sizes: We observe that a batch of size 83 is indeed the optimal point which minimizes the amortized runtime (or maximizes the throughput). Finally, we compare slytHErin's performance with related works that evaluate **NN5** in the same application scenario with polynomial activation functions (thus, we exclude Gazelle [33] which relies on Garbled Circuits). Table 1 shows that slytHErin's performance is on par with or better than previous works, while providing enhanced flexibility in terms of batch size. The approach followed by CryptoNets and inspired works [8, 17, 22] allows them to achieve a good throughput by processing large batches of data items (up to $\mathcal{N}/2$), but their runtime is independent of the batch size (hence, it will not decrease for smaller batches as per Table 1). Conversely, the approach followed by LoLa [11] achieves low latency for a single sample, but cannot amortize the runtime when processing multiple samples. With slytHErin, the end-user can define its custom batch size without a major impact on performance.

---

[3] While slytHErin and related works [11, 17, 22] employ similar hardware for testing, we note that nGraph-HE2 [8] employs compiler optimizations and a more performant hardware with 376GB of RAM and 112 cores.

| | Latency (s) | | |
|---|---|---|---|
| Framework | Batch size = 1 | Batch size = 83 | Batch size = 4,096 |
| CryptoNets [22] | 250 | 250 | 250 |
| Faster CryptoNets [17] | 39.1 | 3,245 | 160,153 |
| LoLa [11] | 2.2 | 182.6 | 8,951 |
| nGraph-HE2 [8][3] | 2.05 | 2.05 | 2.05 |
| `slytHErin` | 3.7 | 4.08 | 243.4 |

Table 1: Latency comparison between `slytHErin` and prior encrypted frameworks for the evaluation of **NN5** and various batch sizes.

| # of Parties | Latency (s) | Throughput (samples/s) |
|---|---|---|
| **3** | 245.58 ($\pm$0.50) | 1.19 |
| **5** | 238.15 ($\pm$4.12) | 1.22 |
| **10** | 278.19 ($\pm$9.11) | 1.05 |
| **20** | 354.17 ($\pm$10.66) | 0.82 |

Table 2: `slytHErin`'s performance for **NN20** on **Scenario 3** (Section 4.3) with increasing number of parties (model-providers).
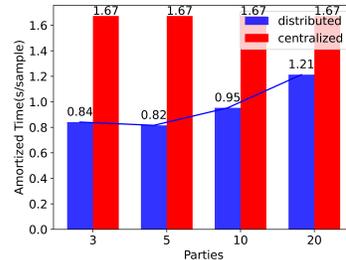


Fig. 5: Benchmarking decentralized vs. centralized bootstrapping on encrypted **NN20** for variable number of parties (model-providers).

**6.2.2 Interactive MPC Protocols.** We evaluate **NN20** in **Scenario 3** where the model is trained and retained under encryption by multiple parties using a privacy-preserving collaborative training framework [19, 52, 53] (Section 4.3). Note that this scenario requires the use of the collective bootstrapping protocol $\mathsf{CBootstrap}(\cdot)$, thus, it is not supported by prior encrypted inference frameworks. Table 2 shows `slytHErin`'s latency and throughput for increasing number of parties, while in Figure 5 we compare `slytHErin`'s amortized runtime when employing $\mathsf{CBootstrap}(\cdot)$ versus the centralized bootstrapping. Overall, we observe a linear increase (decrease) in `slytHErin`'s latency (throughput) as the number of parties increases. We note that the $\mathsf{CBootstrap}(\cdot)$ operation is executed in an asynchronous fashion by the master model provider, i.e., the protocol is initiated concurrently with all the model providers and the output is generated as soon as the last party provides its share. For this reason, we can even experience lower latency when increasing the number of parties by a limited amount (3 vs. 5), as the protocol becomes particularly sensible to the network conditions. In any case, the benefits of employing $\mathsf{CBootstrap}(\cdot)$ over the centralized version (when possible) are evident, as the former enables refreshing the ciphertext

| | Latency(s) | Amortized (s/sample) | Throughput (samples/s) | Avg. latency/layer (s) |
|---|---|---|---|---|
| **Plaintext model Encrypted data (Scenario 1)** | 2,496.83 | 4.26 | 0.234 | 48.95 |
| **Encrypted model Plaintext data (Scenario 2)** | 2,699.75 | 4.62 | 0.216 | 52.93 |
| **Encrypted model Encrypted data (Scenario 3)** | 613.52 | 2.09 | 0.476 | 12.02 |

Table 3: `slytHErin`'s performance for **NN50** in **Scenarios 1**, **2**, and **3** (Section 4). For **Scenario 3**, the number of model-providers is $N=3$.

noise with an efficient interactive protocol, rather than with a computationally expensive homomorphic circuit.

**6.2.3   Application and Model Agility.** Finally, we demonstrate the high degree of flexibility offered by `slytHErin`, both in terms of variety of enabled use-cases and supported architectures, by evaluating a more complex model on all the scenarios described in Section 4. In particular, we benchmark `slytHErin` with **NN50** and a batch of 585 samples on: (i) **Scenario 1** with encrypted data and a plaintext model (Section 4.1), (ii) **Scenario 2** with an encrypted model and plaintext data (Section 4.2), and (iii) **Scenario 3** where the encrypted model is kept by $N=3$ model-providers and encrypted data. Note that evaluating **NN50** in **Scenarios 1** and **2** requires the invocation of the centralized bootstrapping operation, that is not supported by most of the related works [8, 11, 22, 33].

Table 3 shows the performance results for all scenarios. First, we note that by leveraging on our data packing approach and processing multiple samples in a SIMD fashion, `slytHErin` achieves reasonable runtime given the complexity of the **NN50** model (**Scenario 1**). For reference, the original work by Chillotti et al. achieves at best an amortized time of 37.69s/sample and a throughput of 0.02samples/s. Then, we also observe that `slytHErin`'s generic optimizations enable the efficient evaluation of encrypted models: Evaluating **NN50** under encryption on **Scenario 2**, which involves a matrix multiplication, addition, polynomial activation, and centralized bootstrapping operations, is only ∼7% slower than evaluating a plaintext model evaluation (c.f. **Scenario 1**). `slytHErin` achieves the best performance results on **Scenario 3** thanks to its support for interactive multiparty protocols such as collective bootstrapping (Section 6.2.2). Overall, we remark that `slytHErin` is the first framework for encrypted inference that can support all these application scenarios.

## 7   Conclusion

In this work, we presented `slytHErin`, an agile framework for privacy-preserving deep neural network inference using homomorphic encryption. Thanks to our hybrid design that leverages on HE and its multiparty variant, and generic setting-agnostic optimizations, `slytHErin` can support various and novel scenarios for encrypted inference featuring untrusted model providers and clients. These scenarios include: (i) the client sending encrypted data to an untrusted model-provider for inference, (ii) the model-provider sending an encrypted model to a client for local inference (without the need of mutual trust between them), and (iii) the client sending the encrypted data to a cohort of model-providers holding an encrypted model. Thus, `slytHErin` extends the applicability of privacy-preserving PaaS beyond previous works. Moreover, with our intuitive and flexible input data packing scheme, `slytHErin` can be adapted to various deep neural network architectures and can accommodate diverse application requirements, being able to process an arbitrary number of samples without incurring major performance loss. Our experimental results show that the simplicity of our packing approach and the agility of our framework does not harm its performance as it is on par with, and occasionally better than, state-of-the-art related works, while introducing an increased degree of flexibility over previous works.

## References

1. Abadi, M., et al.: TensorFlow: Large-scale machine learning on heterogeneous systems (2015), `http://tensorflow.org/`, software available from tensorflow.org
2. Almeida, M., Laskaridis, S., Venieris, S.I., Leontiadis, I., Lane, N.D.: Dyno: Dynamic onloading of deep neural networks from cloud to device. ACM Trans. Embed. Comput. Syst. **21**(6) (oct 2022). https://doi.org/10.1145/3510831, `https://doi.org/10.1145/3510831`
3. Amazon Forecast. `https://aws.amazon.com/forecast/` (2023), (Accessed: 2023-01-01)
4. Machine Learning on AWS. `https://aws.amazon.com/machine-learning/` (2023), (Accessed: 2023-01-01)
5. Azure Machine Learning. `https://azure.microsoft.com/en-us/products/machine-learning/` (2023), (Accessed: 2023-01-01)
6. Microsoft Azure Cognitive Service. `https://learn.microsoft.com/en-us/azure/cognitive-services/language-service/` (2023), (Accessed: 2023-01-01)
7. Machine Learning made beautifully simple for everyone . `https://bigml.com/` (2023), (Accessed: 2023-01-01)
8. Boemer, F., Costache, A., Cammarota, R., Wierzynski, C.: nGraph-HE2: A high-throughput framework for neural network inference on encrypted data. In: ACM WAHC (2019)
9. Boemer, F., Lao, Y., Wierzynski, C.: ngraph-he: A graph compiler for deep learning on homomorphically encrypted data. CoRR **abs/1810.10121** (2018), `http://arxiv.org/abs/1810.10121`
10. Boura, C., Gama, N., Georgieva, M.: Chimera: a unified framework for b/fv, tfhe and heaan fully homomorphic encryption and predictions for deep learning. IACR Cryptol. ePrint Arch. **2018**, 758 (2018)

11. Brutzkus, A., Gilad-Bachrach, R., Elisha, O.: Low latency privacy preserving inference. In: International Conference on Machine Learning. pp. 812–821. PMLR (2019)
12. California Consumer Privacy Act (CCPA). `https://www.oag.ca.gov/privacy/ccpa` (2023), (Accessed: 2023-01-01)
13. Chabanne, H., de Wargny, A., Milgram, J., Morel, C., Prouff, E.: Privacy-preserving classification on deep neural network. IACR Cryptol. ePrint Arch. **2017**, 35 (2017)
14. Cheon, J.H., Kim, A., Kim, M., Song, Y.: Homomorphic encryption for arithmetic of approximate numbers. In: ASIACRYPT (2017)
15. Chillotti, I., Gama, N., Georgieva, M., Izabachène, M.: Tfhe: Fast fully homomorphic encryption over the torus. Journal of Cryptology **33**(1), 34–91 (2020). https://doi.org/10.1007/s00145-019-09319-x
16. Chillotti, I., Joye, M., Paillier, P.: Programmable bootstrapping enables efficient homomorphic inference of deep neural networks. Cryptology ePrint Archive, Paper 2021/091 (2021)
17. Chou, E., Beal, J., Levy, D., Yeung, S., Haque, A., Fei-Fei, L.: Faster cryptonets: Leveraging sparsity for real-world encrypted inference. CoRR **abs/1811.09953** (2018), `http://arxiv.org/abs/1811.09953`
18. Dathathri, R., Saarikivi, O., Chen, H., Laine, K., Lauter, K., Maleki, S., Musuvathi, M., Mytkowicz, T.: Chet: An optimizing compiler for fully-homomorphic neural-network inferencing. In: Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation. p. 142–156. PLDI 2019, Association for Computing Machinery, New York, NY, USA (2019). https://doi.org/10.1145/3314221.3314628, `https://doi.org/10.1145/3314221.3314628`
19. Froelicher, D., Troncoso-Pastoriza, J.R., Pyrgelis, A., Sav, S., Sousa, J.S., Bossuat, J.P., Hubaux, J.P.: Scalable privacy-preserving distributed learning. PETS (2021)
20. The EU General Data Protection Regulation. `https://gdpr-info.eu/` (2023), (Accessed: 2023-01-01)
21. Gentry, C.: Fully homomorphic encryption using ideal lattices. In: Proceedings of the Forty-First Annual ACM Symposium on Theory of Computing. p. 169–178. STOC '09, Association for Computing Machinery, New York, NY, USA (2009). https://doi.org/10.1145/1536414.1536440, `https://doi.org/10.1145/1536414.1536440`
22. Gilad-Bachrach, R., Dowlin, N., Laine, K., Lauter, K., Naehrig, M., Wernsing, J.: Cryptonets: Applying neural networks to encrypted data with high throughput and accuracy. In: ICML (2016)
23. Gnacik, M., Łapa, K.: Using toeplitz matrices to obtain 2d convolution (10 2022). https://doi.org/10.21203/rs.3.rs-2195496/v1
24. Go Programming Language. `https://golang.org` (2023), (Accessed: 2023-01-01)
25. AI and machine learning products. `https://cloud.google.com/products/ai` (2023), (Accessed: 2023-01-01)
26. Gray, R.M.: Toeplitz and circulant matrices: A review. Foundations and Trends® in Communications and Information Theory **2**(3), 155–239 (2006). https://doi.org/10.1561/0100000006, `http://dx.doi.org/10.1561/0100000006`
27. Halevi, S., Shoup, V.: HElib - An Implementation of homomorphic encryption. `https://github.com/shaih/HElib/` (2014), (Accessed: 2023-01-01)
28. Halevi, S., Shoup, V.: Algorithms in helib. In: Annual International Cryptology Conference (CRYPTO). Springer (2014)

29. Hesamifard, E., Takabi, H., Ghasemi, M., Wright, R.: Privacy-preserving machine learning as a service. PETS (2018)
30. Centers for Medicare & Medicaid Services. The Health Insurance Portability and Accountability Act of 1996 (HIPAA). https://www.cms.gov/Regulations-and-Guidance/Administrative-Simplification/HIPAA-ACA/PrivacyandSecurityInformation (2023), (Accessed: 2023-01-01)
31. Huang, Z., jie Lu, W., Hong, C., Ding, J.: Cheetah: Lean and fast secure Two-Party deep neural network inference. In: 31st USENIX Security Symposium (2022)
32. Jiang, X., Kim, M., Lauter, K., Song, Y.: Secure outsourced matrix computation and application to neural networks. In: Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security. p. 1209–1222. CCS '18, Association for Computing Machinery, New York, NY, USA (2018). https://doi.org/10.1145/3243734.3243837, https://doi.org/10.1145/3243734.3243837
33. Juvekar, C., Vaikuntanathan, V., Chandrakasan, A.: Gazelle: A low latency framework for secure neural network inference. USENIX Security (2018)
34. Kim, M., Jiang, X., Lauter, K., Ismayilzada, E., Shams, S.: Secure human action recognition by encrypted neural network inference. Nature Communications **13**(1), 4799 (Aug 2022), https://doi.org/10.1038/s41467-022-32168-5
35. Lattigo: A library for lattice-based homomorphic encryption in go. https://github.com/ldsec/lattigo (2023), (Accessed: 2023-01-01)
36. LeCun, Y., Cortes, C.: MNIST handwritten digit database (2010)
37. Liu, J., Juuti, M., Lu, Y., Asokan, N.: Oblivious neural network predictions via MiniONN transformations. In: ACM CCS (2017)
38. Lloret-Talavera, G., Jorda, M., Servat, H., Boemer, F., Chauhan, C., Tomishima, S., Shah, N., Peña, A.: Enabling homomorphically encrypted inference for large dnn models. IEEE Transactions on Computers **PP**, 1–1 (04 2021). https://doi.org/10.1109/TC.2021.3076123
39. Lu, W.j., Sakuma, J.: More practical privacy-preserving machine learning as a service via efficient secure matrix multiplication. In: Proceedings of the 6th Workshop on Encrypted Computing & Applied Homomorphic Cryptography. p. 25–36. WAHC '18, Association for Computing Machinery, New York, NY, USA (2018)
40. Lyubashevsky, V., Peikert, C., Regev, O.: On ideal lattices and learning with errors over rings. Journal of the ACM (JACM) **60**(6), 1–35 (2013)
41. Meftah, S., Tan, B.H.M., Mun, C.F., Aung, K.M.M., Veeravalli, B., Chandrasekhar, V.: Doren: Toward efficient deep convolutional neural networks with fully homomorphic encryption. IEEE Transactions on Information Forensics and Security **16**, 3740–3752 (2021). https://doi.org/10.1109/TIFS.2021.3090959
42. Mishra, P., Lehmkuhl, R., Srinivasan, A., Zheng, W., Popa, R.A.: Delphi: A cryptographic inference service for neural networks. In: USENIX Security (2020)
43. Mouchet, C., Troncoso-pastoriza, J.R., Bossuat, J.P., Hubaux, J.P.: Multiparty homomorphic encryption from ring-learning-with-errors. PETS (2021)
44. Murshed, M.G.S., Murphy, C., Hou, D., Khan, N., Ananthanarayanan, G., Hussain, F.: Machine learning at the network edge: A survey. ACM Comput. Surv. **54**(8) (oct 2021). https://doi.org/10.1145/3469029, https://doi.org/10.1145/3469029
45. Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., Lerer, A.: Automatic differentiation in pytorch. In: 31st Conference on Neural Information Processing Systems (NIPS 2017) (2017)
46. Patra, A., Suresh, A.: Blaze: Blazing fast privacy-preserving machine learning. In: NDSS (2020)

47. Ran, R., Wang, W., Gang, Q., Yin, J., Xu, N., Wen, W.: CryptoGCN: Fast and scalable homomorphically encrypted graph convolutional network inference. In: Oh, A.H., Agarwal, A., Belgrave, D., Cho, K. (eds.) Advances in Neural Information Processing Systems (2022), `https://openreview.net/forum?id=VeQBBm1MmTZ`

48. Rathee, D., Rathee, M., Kumar, N., Chandran, N., Gupta, D., Rastogi, A., Sharma, R.: Cryptflow2: Practical 2-party secure inference. In: ACM CCS. p. 325–342 (2020)

49. Riazi, M.S., Samragh, M., Chen, H., Laine, K., Lauter, K.E., Koushanfar, F.: Xonn: Xnor-based oblivious deep neural network inference. In: USENIX Security (2019)

50. Rohloff, K.: The PALISADE Lattice Cryptography Library. `https://git.njit.edu/palisade/PALISADE` (2018)

51. Sav, S., Bossuat, J.P., Troncoso-Pastoriza, J.R., Claassen, M., Hubaux, J.P.: Privacy-preserving federated neural network learning for disease-associated cell classification. Patterns **3**(5) (2022). https://doi.org/10.1016/j.patter.2022.100487, `https://doi.org/10.1016/j.patter.2022.100487`

52. Sav, S., Diaa, A., Pyrgelis, A., Bossuat, J.P., Hubaux, J.P.: Privacy-preserving federated recurrent neural networks. CoRR **abs/2207.13947** (2022), `https://arxiv.org/abs/2207.13947`

53. Sav, S., Pyrgelis, A., Troncoso-Pastoriza, J.R., Froelicher, D., Bossuat, J.P., Sousa, J.S., Hubaux, J.P.: Poseidon: Privacy-preserving federated neural network learning. In: Network and Distributed System Security Symposium (NDSS) (2021)

54. Microsoft SEAL (release 3.3). `https://github.com/Microsoft/SEAL` (2023), (Accessed: 2023-01-01)

55. Sim, S.H., Paranjpe, T., Roberts, N., Zhao, M.: Exploring edge machine learning-based stress prediction using wearable devices. In: 2022 21st IEEE International Conference on Machine Learning and Applications (ICMLA). pp. 1266–1273 (2022). https://doi.org/10.1109/ICMLA55696.2022.00203

56. Watson Machine Learning . `https://cloud.ibm.com/catalog/services/watson-machine-learning` (2023), (Accessed: 2023-01-01)

57. Xu, G., Han, X., Xu, S., Zhang, T., Li, H., Huang, X., Deng, R.H.: Hercules: Boosting the performance of privacy-preserving federated learning. IEEE Transactions on Dependable and Secure Computing pp. 1–18 (2022). https://doi.org/10.1109/TDSC.2022.3218793