

Leia: A Lightweight Cryptographic Neural Network Inference System at the Edge

Xiaoning Liu[✉], Bang Wu, Xingliang Yuan[✉], *Member, IEEE*, and Xun Yi[✉]

Abstract—The advances in machine learning have revealed its great potential for emerging mobile applications such as face recognition and voice assistant. Models trained via a Neural Network (NN) can offer accurate and efficient inference services for mobile users. Unfortunately, the current deployment of such service encounters privacy concerns. Directly offloading the model to the mobile device violates model privacy of the model owner, while feeding user input to the service compromises user privacy. To address this issue, we propose Leia, a lightweight cryptographic NN inference system at the edge. Leia is designed from two mobile-friendly perspectives. First, it leverages the paradigm of edge computing wherein the inference procedure keeps the model closer to the mobile user to foster low latency service. Specifically, Leia's architecture consists of two non-colluding edge services to obviously perform NN inference on the encoded user data and model. Second, Leia's realization makes the judicious use of potentially constrained computational and communication resources in edge devices. We adapt the Binarized Neural Network (BNN), a trending flavor of NN with low inference overhead, and purely choose the lightweight secret sharing techniques to realize secure blocks of BNN. We implement Leia and deploy it on Raspberry Pi. Empirical evaluations on benchmark and medical datasets via various models demonstrate the practicality of Leia.

Index Terms—Secure computation, privacy-preserving mobile application, neural network inference, edge computing.

I. INTRODUCTION

RECENT flourishing of Machine Learning has promoted the Neural Network (NN) powered mobile applications such as face detection cameras and speech recognition assistants. An NN prediction service is typically deployed through two scenarios. One is the on-device NN prediction service, where the application downloads the pre-trained model from a company owning the model, and performs an inference task over user's data on mobile device [1]. Another scenario

relies on the cloud service provider (e.g., Google Cloud AI [2]), where user data and model are delivered to cloud who runs the NN inference task and sends back the prediction result to the mobile device. Unfortunately, both scenarios are troublesome due to the increasingly raised privacy issues. User's data contains sensitive information about their daily activities. Uploading such data to cloud in plaintext can put individual's privacy in danger [3], [4]. On the other hand, from the aspect of model owners, their models are valuable and often trained on proprietary data [4], [5]. The unauthorized exposure of the proprietary model and underlying confidential data inflicts severe commercial damages.

To alleviate the privacy issues, one approach seems plausible is to delegate the inference tasks with encrypted model and user data to a centralized cloud server. However, this approach relies on heavy cryptographic techniques like (fully) homomorphic encryption (HE). An efficient alternative tends to use secure multi-party computation (MPC) techniques with specialized designs that can execute NN inference over encrypted user input data and/or encrypted models. For example, Delphi [4] and MiniONN [5] proceed the inference tasks between the mobile device and the model owner, while continuous interaction is involved between them during secure computation. Namely, both parties have to be online and connected throughout the entire inference process. It is noteworthy that the above rigid operational confinement might not be always feasible in cellular networks. Other systems like XONN [6] and Quotient [7] employ constant-round secure 2-party computation protocols that result in less engagement of mobile device and model owner. But these systems consume large bandwidth costs due to the MPC techniques relied on, i.e., Garbled Circuits (GC), which may become the bottleneck at the edge. Later in Section VI, we demonstrate the considerable bandwidth saving by our design compared to GC based realization.

Our Contributions: In this paper, we propose and enable Leia, a cryptographic NN inference system executing at the edge. Leia takes the edge based architecture as a starting point, harnessing the novelty from system, cryptographic and machine learning domains. The combination endows Leia privacy assurance and seamless embracement of NN powered mobile applications.

Our first insight is to leverage edge computing, wherein the processing keeps the model closer to the mobile user so as to foster low latency service of NN inference [8], [9]. However, devising edge based architecture is non-trivial for Leia's scenario. Both models and user's data should fully be protected against edge nodes during inference. Moreover,

Manuscript received November 6, 2020; revised April 7, 2021 and November 11, 2021; accepted December 12, 2021. Date of publication December 24, 2021; date of current version January 14, 2022. This work was supported in part by the Australian Research Council (ARC) Discovery Projects (DP200103308, DP180103251, and DP190102835) and in part by the ARC Linkage Project LP160101766. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. Nele Mentens. (Corresponding author: Xingliang Yuan.)

Xiaoning Liu and Xun Yi are with the School of Computing Technologies, RMIT University, Melbourne, VIC 3001, Australia (e-mail: xiaoning.trust@gmail.com; xun.yi@rmit.edu.au).

Bang Wu and Xingliang Yuan are with the Faculty of Information Technology, Monash University, Clayton, VIC 3800, Australia (e-mail: bang.wu@monash.edu; xingliang.yuan@monash.edu).

This article has supplementary downloadable material available at <https://doi.org/10.1109/TIFS.2021.3138611>, provided by the authors.

Digital Object Identifier 10.1109/TIFS.2021.3138611

1556-6021 © 2021 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission.

See <https://www.ieee.org/publications/rights/index.html> for more information.

relaxing the constraint of the model owner and mobile device being online is expedient because of the dynamic network effects in cellular networks. To this end, we resort to the edge nodes as two non-colluding computational services to fulfill the above privacy objective and operational requirement simultaneously. The model owner encodes the model and sends the encoded one to edge nodes only once. After that, the user mobile device can submit the encoded input and obtain the encoded inference result. Within the process, edge nodes obliviously perform inference without further interacting with either the model owner or user mobile device.

The edge based paradigm entails judicious usage of computational and communication resources. We thus choose the relatively lightweight secret sharing techniques to offer Leia security guarantee. Yet, merely transforming the NN inference procedure into cryptographic operations does not necessarily achieve satisfied efficiency in communication and computation for mobile and edge devices. Instead, our second insight is the adoption of the Binarized Neural Network (BNN) [15], a special flavor of NN model with weights and activations are all confined to ± 1 . Because the small BNN model can drastically reduce the resource demand, and the beneath operations over binary values are more compatible with cryptographic primitives, it becomes our natural choice. Thereby, we subtly build Leia from ground up with secure layer functions, including the secure linear layers (secure convolutional layer SCONV and the secure fully connected layer SFC), the secure batch normalization function (SecBN), the secure binary activation function (SecBA), and the secure max pooling layer (SMP). They are highly customized for the BNN and securely realized based on secret sharing as the building blocks of Leia.

We implement and deploy Leia to Raspberry Pi. Our evaluation is comprehensively performed over benchmark datasets (i.e., MNIST and CIFAR-10) and multiple real-world medical datasets on breast cancer, diabetes, liver disease and thyroid. We apply eight different models to evaluate Leia's performance. Our results show that Leia can produce a 97% accurate prediction result by 4s in the edge environment for MNIST. For a 23-layer network on CIFAR-10, Leia achieves up to $22\times$ bandwidth saving compared to the prior art. In addition, all medical inference tasks are performed within 3.6s and require less than $50\mu s$ and 10ms on the user and model owner side.

Organization: The rest of the paper is organized as follows. Sec. II investigates related work. Sec. III introduces the preliminaries used in this paper. Sec. IV overviews the system architecture and threat model. Sec. V expatiates on the protocol designs. Empirical evaluation is given in Sec. VI. The paper is concluded in Sec. VII.

II. RELATED WORKS

A. Privacy-Preserving Neural Network Training and Inference

CryptoNets [10] adapts the leveled HE to perform privacy-preserving NN inference in an outsourced environment. Despite some optimizations have been employed, CryptoNets still suffers from intensive computational overheads due to heavy weight HE. Some other systems [4], [5], [12] consider a different setting, in which the client directly communicates

with the model owner for inference, but does not want to reveal its input. Gazelle [12] devises a 2PC secure NN inference framework combining the lattice-based HE and Yao's Garbled Circuits (GC). After that, Delphi [4] and MiniONN [5] are proposed with careful optimizations and achieve higher efficiency. Different from our work, the above systems require *continuous interactions* between the client and the model owner during secure computation. Our work, instead, leveraging edge computing paradigm, delegates the whole secure inference protocol to the edge and as such relaxes both parties from always being online.

SecureML [11] is the first to propose privacy-preserving training and inference system with tailored MPC protocols. The proposed system considers various learning problems including linear regression, logistic regression, and NN. Besides, the more MPC-friendly activation functions are subtly devised and realized with GC and designed Oblivious Transfer (OT) protocols. Quotient [7], XONN [6], and BANNERS [13] design privacy-preserving NN inference over quantized NN models with weights which are restricted within $\{-1, +1, 0\}$ and $\{-1, +1\}$, respectively. Such quantization allows for the conversion from arithmetic operations to Boolean operations, and thus are more compatible with the protocols realized with GC free-XOR optimizations [6], [7] and Boolean sharing techniques [13]. The GC-based designs [6], [7] usually require substantial bandwidth and thus are more suitable for the steady and high-throughput network. However, the escalated bandwidth can be prohibitive and the bottleneck at the edge and may incur additional charges by cellular network service providers. Meanwhile, the work BANNERS [13] is specially tailored under a three-party setting which requires a more complex deployment overhead at the edge. Our work subtly builds a privacy-assured collaborative inference protocol from ground up, harnessing the novelty from both cryptographic and machine learning literature. Our system originates a design where an all-binarized neural network inference procedure is securely carried out with customized layer functions under lightweight MPC primitives, and thus is particularly suitable for the application deployment for mobile devices. Empirical evidence shown in Section VI confirms that, for the equivalent functionalities, Leia's realizations introduce $30\text{-}2500\times$ less bandwidth costs than the GC-based realizations. To be more clear, we compare our work with prior art and summarize the main difference in Table I.

Meanwhile, recent years have drawn a growing interest in secure learning and inference protocols thwart the malicious parties. Notable works include a fully-decentralized linear model learning proposed in Helen [3], an end-to-end 3PC NN training and inference proposed in FALCON [16], 3PC inference framework over quantized NN proposed by Dalskov *et al.* [17], and the 3PC inference over BNN in BANNERS [13]. Leia's design lays a solid foundation for extending towards a mobile-friendly secure NN inference system against malicious edge nodes. We leave a careful study to our future work. Some prior works [13], [16], [18] also focus on lightweight NN inference based on secret sharing techniques. We note that these systems are specially designed for three-party secure protocols that require a more complicated practical deployment than Leia's two-server setting.

TABLE I
HIGH-LEVEL COMPARISON OF DIFFERENT CRYPTOGRAPHIC NN SYSTEMS

system	crypto. primitives	system model	client engagement	multi-inter. with client	comm. cost [◊]	plaintext accu. presrv.	model privacy	co-design
CryptoNets [10]	FHE	centralized server	homomorphic enc./dec.	✗	high	✗	✗	-
SecureML [11]	GC, SS	2 servers [†]	shares generation	✗	-	✗	✓	-
MiniONN [5]	HE, GC, SS	client-server*	build GC	✓	high	✓	✓	-
Gazelle [12]	HE, GC, SS	client-server*	homomorphic enc./dec., build GC	✓	low	✗	✓	-
Quotient [7]	GC, SS	2 servers [†]	shares generation	✗	-	✓	✓	ternarized NN
XONN [6]	GC	client-server*	build GC, OT-based computation	✓	low	✓	✓	BNN
BANNERS [13]	Replicated SS	3 servers	shares generation	✗	low	✓	✓	BNN
Delphi [4]	HE, GC, SS	client-server*	GC evaluation	✓	low	✓	✓	-
Chameleon [14]	GC, SS	2 servers [†]	shares generation	✗	medium	✓	✓	-
Leia (this work)	SS	2 servers [†]	shares generation	✗	low	✓	✓	BNN

* Client-server two-party computation.

[†] Outsourced computation to two non-colluding servers.

[◊] The communication cost is evaluated based on the bandwidth reported in their papers, and [7], [11] have not reported their inference bandwidth costs.

B. Secure Multi-Party Computation Framework

Privacy-preserving machine learning protocols can be carried out via generic MPC techniques, such as Garbled Circuits, Secret Sharing [19]–[21], and the frameworks mixed with multiple primitives [22]–[24]. Among these, a line falls into devising machine learning specified MPC frameworks, such as (1) the two-party frameworks proposed in Chameleon [14], EzPC [25], (2) the three-party frameworks proposed in ABY³ [23], SecureNN [18], FALCON [16], and the work [17], and (3) the four-party framework in Trident [24]. For performance consideration, recent privacy-preserving machine learning systems opt for specialized and optimized designs instead of direct adoption of generic MPC frameworks [3], [4], [6], [7].

III. BACKGROUND

In this section, we introduce the preliminaries underlying Leia's construction. Our notation is summarized in Table II.

A. Binarized Neural Networks

Binarized Neural Network [15] (BNN), i.e., neural networks with weights and possible activations restricted to ± 1 . It comprises a sequence of *linear* and *non-linear* layers, where the number of layers L indicates the depth of network. The BNN inference takes as input a tensor representing the task-specific raw data, and produces a prediction result based on a trained model tensor. The all binarized weights make BNN model significantly smaller than an equivalent network with high precision weights [26].

1) *Linear Layers*: The linear layers typically can be of two types: the fully connected layers (FC) and the convolutional layers (CONV). Both types can be formulated as the vector dot product ($\text{VDP}(\cdot, \cdot)$) between two vectors $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{w} \in \mathbb{R}^n$ as: $\text{VDP}(\mathbf{x}, \mathbf{w}) = \sum_{k=1}^n w_k \cdot x_k$, where w_k, x_k indicate the k -th element of vector \mathbf{x}, \mathbf{w} . Let the parameters c_{in} (and $c_o \times c_{in}$), n_{in} (and n_w), m_{in} (and n_w) denote the number of *channel*, the *width* and *height* of input (and model), respectively.

The FC layer takes as input a vector $\mathbf{x} \in \mathbb{R}^{n_{in}}$, applies a set of the weight vectors $\mathbf{w}_1, \dots, \mathbf{w}_{n_o} \in \{-1, +1\}^{n_{in}}$ and bias vector $\mathbf{bias} \in \mathbb{R}^{n_o}$, and outputs a vector $\mathbf{z} \in \mathbb{R}^{n_o}$. For $k \in [1, n_o]$, it repeatedly proceeds the $z_k = \text{VDP}(\mathbf{x}, \mathbf{w}_k) + \text{bias}_k$,

where z_k is the k -th element of the resulting vector \mathbf{z} . The result \mathbf{z} can be submitted to the non-linear operations.

The CONV layer is normally applied for image classification. It takes as input an image tensor $\mathbf{X} \in \mathbb{R}^{c_{in} \times n_{in} \times m_{in}}$, applies the weight tensor $\mathbf{W} \in \{-1, +1\}^{c_o \times c_{in} \times n_w \times m_w}$ and bias vector $\mathbf{bias} \in \mathbb{R}^{c_o}$, and outputs a feature map tensor $\mathbf{Z} \in \mathbb{R}^{c_o \times n_o \times m_o}$. For each input matrix $X \in \mathbb{R}^{n_{in} \times m_{in}}$, the CONV layer repeatedly moves the weight matrix $W \in \{-1, +1\}^{n_w \times m_w}$ as a sliding window, from left to right and top-down with given stride, until passing through the entire image matrix. Let \mathcal{T} be the total moves to slide the window. For τ -th move that $\tau \in \mathcal{T}$, it flattens X_τ inside the sliding window and W as vectors \mathbf{x}_τ and \mathbf{w} , and proceeds $z_\tau = \text{VDP}(\mathbf{x}_\tau, \mathbf{w}) + \text{bias}_\tau$ to obtain z_τ as one element in resulting matrix $Z \in \mathbb{R}^{n_o \times m_o}$. A further illustration is given in Section I of the supplementary materials.

2) *Binary Activation*: The non-linear binary activation function (BA) is attached on each neuron. It takes as input the real-valued activation $a \in \mathbb{R}$ outputted from previous operation, performs element-wise $\text{sign}(\cdot)$ function: $\text{sign}(a) = +1$, if $a \geq 0$; and $\text{sign}(a) = -1$, otherwise; and outputs the sign bit as the binarized activation $a \in \{-1, +1\}$.

3) *Batch Normalization*: Batch Normalization [27] (BN) is widely adopted in modern NN to regularize the model to avoid the activations growing too large to unstablize the model. The procedure of the BN function during inference performs as follows: (1) it performs element-wise normalization on each neuron's feature a via $\hat{a} = (a - \mu)/\delta$ where μ is the *running mean* and the non-zero δ is *running variance* the of training dataset; and (2) it applies the *scale* parameter $\gamma \in \mathbb{R}$ and the *shift* parameter $\beta \in \mathbb{R}$ to get the output $z = \gamma \cdot \hat{a} + \beta$. Parameters $\mu, \delta, \gamma, \beta$ are highly dependent on the training data and indicate the data distribution, thus have to be protected during secure computation.

4) *Max Pooling Layer*: Max pooling layer (MP) is normally applied straight after the CONV layer to downsample the image. It takes the matrices outputted from the CONV layer, and obtains the maximum value within a sliding window as an element of the output matrix.

B. Cryptographic Primitives

1) *Correlated Oblivious Transfer*: Correlated Oblivious Transfer (COT) [28] is a cryptographic primitive allowing

TABLE II
NOTATION USED IN OUR PAPER

Parameters	
$\mathbf{X}, X, \mathbf{x}$	Input image tensor, matrix, vector
$\mathbf{W}, W, \mathbf{w}$	Weight tensor, matrix, vector whose elements $\in \{-1, +1\}$
\mathbf{a}, a	Activation vector, an activation on a neuron
$\mu, \delta, \gamma, \beta$	Batch normalization parameters: the running mean, the running variance, the scale, the shift
ℓ, κ	Bit length
n, m	Vector length n , number of vectors m
General notation	
x_k, \mathbf{x}_k	The k -th element of vector \mathbf{x} , the k -th vector of set \mathbf{x}
x^b	Superscript b denotes encoded data $x^b \in \{0, 1\}$
i	Identifier of a party i that $i \in \{0, 1\}$
$\langle x \rangle_i^A$	Arithmetic shares of value x held by party i
$\llbracket x \rrbracket_i$	Boolean shares of value x held by party i
$\langle x \rangle_i^A \pm \langle y \rangle_i^A$	Addition/subtraction over arithmetic shares
$\langle x \rangle_i^A \cdot \langle y \rangle_i^A$	Multiplication over arithmetic shares
$\llbracket x \rrbracket_i + \llbracket y \rrbracket_i$	Bitwise XOR over boolean shares
$\llbracket x \rrbracket_i \cdot \llbracket y \rrbracket_i$	Bitwise AND over boolean shares
$\text{mtri}_i, \text{atri}_i$	Multiplication triples, Boolean AND triples held by party i

for secure two-party computation. It is one particular OT flavor with improved practicality. Given engaged two parties, a *sender* P_0 holding its input a pair of binary strings $m_0, m_1 \in \{0, 1\}^\ell$, and a *receiver* P_1 holding its input a choice bit $b \in \{0, 1\}$, a COT performs as follows. P_0 constructs and inputs a correlation function $f_\Delta(\cdot)$ to link m_0 and m_1 in a way that m_0 is a random value and $m_1 = f_\Delta(m_0)$. On input m_0, f_Δ from P_0 and b from P_1 , the COT outputs $m_b \in \{m_0, m_1\}$ to P_1 . It ensures that P_0 learns nothing about b , and P_1 learns nothing about m_{1-b} . Note that the function $f_\Delta(\cdot)$ is a correlation robust random oracle $H : \{0, 1\}^\ell \rightarrow \{0, 1\}^\ell$. We denote the above described COT functionality as $(\perp; m_b) \leftarrow \text{COT}(m_0, f_\Delta(\cdot); b)$. The n -times COT_ℓ (i.e., $n \times \text{COT}_\ell$) can be run in parallel, where each COT_ℓ is used for transferring the ℓ -bit strings.

2) *Arithmetic Sharing and Multiplication Triple*: Arithmetic sharing [20] additively shares an ℓ -bit secret value x in the ring \mathbb{Z}_{2^ℓ} as $\langle x \rangle_0^A + \langle x \rangle_1^A \equiv x \pmod{2^\ell}$. In this paper, all operations over Arithmetic shares are performed under mod 2^ℓ unless explicitly mentioned. Addition/subtraction over shares ($\langle z \rangle_i^A = \langle x \rangle_i^A \pm \langle y \rangle_i^A$), multiplication by a public value ($\langle z \rangle_i^A = \eta \cdot \langle x \rangle_i^A$) can be efficiently calculated by each party P_i ($i \in \{0, 1\}$) at local without any interaction. Multiplication over two shares ($\langle z \rangle_i^A = \langle x \rangle_i^A \cdot \langle y \rangle_i^A$) requires assistance with pre-computed multiplication triple (denoted as mtri_i), i.e., a type of Beaver's triple [29] in the format $\langle c \rangle_i^A = \langle a \rangle_i^A \cdot \langle b \rangle_i^A$. To multiplying two shares, each party P_i sets $\langle e \rangle_i^A = \langle x \rangle_i^A - \langle a \rangle_i^A$ and $\langle f \rangle_i^A = \langle y \rangle_i^A - \langle b \rangle_i^A$. The parties interact to reconstruct e and f . At the end, P_i sets $\langle z \rangle_i^A = i \cdot e \cdot f + f \cdot \langle a \rangle_i^A + e \cdot \langle b \rangle_i^A + \langle c \rangle_i^A$. Note that multiplication triples are data-independent and can be generated through cryptographic approaches [22], or by a third party [30].

3) *Boolean Sharing and Boolean AND Triple*: Boolean sharing [19], [22] can be viewed as the Arithmetic sharing over \mathbb{Z}_2 . It produces two Boolean shares $\llbracket x \rrbracket_0$ and $\llbracket x \rrbracket_1$ of a secret bit x between two parties P_0 and P_1 , respectively. Reconstruction of x is performed via $x = \llbracket x \rrbracket_0 \oplus \llbracket x \rrbracket_1$. The XOR operation (\oplus) over two Boolean shares is identical to the addition over Arithmetic shares in \mathbb{Z}_2 , i.e., each party P_i computes locally $\llbracket z \rrbracket_i = \llbracket x \rrbracket_i \oplus \llbracket y \rrbracket_i$. Meanwhile, the bitwise

AND operation (\wedge) over Boolean shares works similar to the multiplication as $\llbracket z \rrbracket = \llbracket x \rrbracket \wedge \llbracket y \rrbracket$. It is calculated with the assistance of pre-computed Boolean AND triple (denoted as atri_i) $\llbracket c \rrbracket = \llbracket a \rrbracket \wedge \llbracket b \rrbracket$.

IV. SYSTEM OVERVIEW

A. Architecture

1) *System Overview*: Fig. 1 illustrates Leia's system architecture, which involves three entities: the mobile device, the model owner, and the two distinct edge nodes S_0 and S_1 . The model owner obtains a customized NN model based on proprietary data and resorts to Leia to provide secure inference service for its users without revealing the model in cleartext. In practice, the model owner can be an ML-powered mobile application developing company (e.g., SnapML [31] and Apple Safari) or the enterprise providing ready-made intelligence for those applications (e.g., Amazon Rekognition [32], Google DeepMind Health [33]). The mobile device of the user collects user's data input and asks the mobile application for inference tasks without revealing the private user input. The two edge nodes can be deployed from separate edge computing service providers like Azure IoT Edge [34] and AWS Lambda@Edge [35].

2) *Workflow*: From a high level point of view, Leia's cryptographic NN inference service is operated as follows. The model owner holds a BNN model and deploys the encoded model tensor \mathbf{W} to the two edge nodes S_0 and S_1 , where \mathbf{W} is secret-shared (i.e., Boolean shares) into $\llbracket \mathbf{W} \rrbracket_0$ and $\llbracket \mathbf{W} \rrbracket_1$. Once the mobile device invokes an NN inference request, the raw input will be protected as a secret-shared (i.e., Arithmetic shares) tensor $\langle \mathbf{X} \rangle_0^A$ and $\langle \mathbf{X} \rangle_1^A$, and be submitted to corresponding edge nodes. After that, the two nodes run secure collaborative inference procedure (green box in Fig. 1), and send the encoded inference result vector $\langle \mathbf{z} \rangle_0^A$ and $\langle \mathbf{z} \rangle_1^A$ to the mobile device who can reconstruct the result \mathbf{z} to produce a label to classify the user's input. Leia's secure inference procedure consists of a series of essential computational blocks in BNN. A typical block in Leia assembles several layer functions, including the secure convolutional layer (SCONV) or secure fully connected layer (SFC), the secure batch normalization function (SecBN), the secure binary activation function (SecBA) and the secure max pooling layer (SMP). Note that, for efficient realization, we craft a secure normalized binary activation function (SecNBA) combining the SecBN and SecBA functions (see details in Section V-B.2).

3) *Usage Scenarios*: Before formalizing Leia's threat model, we demonstrate our usage scenarios in real-world applications.

a) *Face verification for payment authorization*: A typical use case can be a mobile banking application, which authorizes a legitimate user to do payment via the face verification. In a cleartext version, a user photo is taken by the mobile device and submitted to the application. The banking application evaluates the user photo under its NN model trained over its customer profiles and sends the validation result to the mobile device. By integrating Leia into the banking application, the entire verification procedure is delegated and securely proceeded at the edge. Both the photo containing user's biometric and home interior information and the bank's proprietary model are perfectly protected.

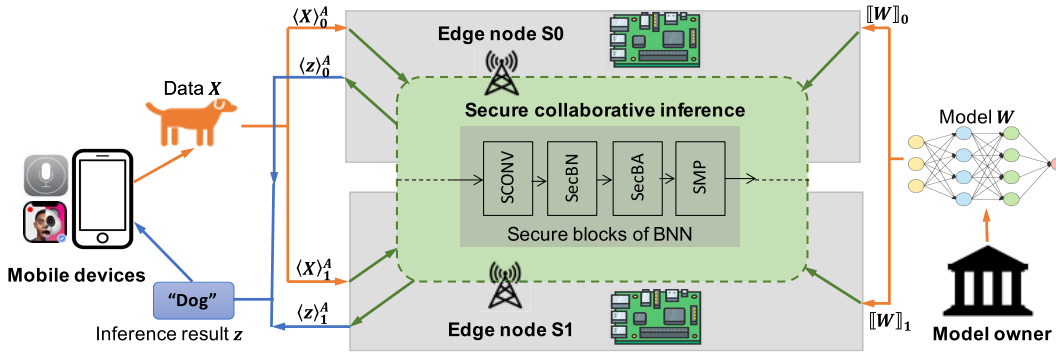


Fig. 1. System architecture.

b) Medical imaging on portable device: The usage of Leia includes but is not limited to the mobile devices. Instead, we note that Leia's deployment can be generalized to any resources-constrained devices, such as the IoT cameras, portable medical imaging devices. Herein we give an example to deploy Leia into a handheld medical imaging scanner, like the handheld CT scanner [36]. In this scenario, the CT scanner images the patient's chest as a CT scan. It aims to evaluate the patient's CT under the NN provided by some medical imaging service providers (the model owner), like Google DeepMind Health [33]. However, the CT scan indicating patient's health information is forbidden to share with the service provider in cleartext due to legislation. They thus can proceed the above evaluation via Leia. The scanner can install Leia and submits the protected patient's CT through Leia's API. Meanwhile, the edge nodes in Leia have already equipped with the protected NN supplied by the medical imaging service provider. Leia securely proceed the NN inference over encrypted CT scan and returns only the inference result to the scanner side. Neither the user of scanner or the service provider can learn the private data about each other.

B. Threat Model and Privacy Goal

Leia considers the following threat model: all entities are the *semi-honest* parties, and the two edge nodes are *non-colluding* computing devices. Specifically, each party will faithfully follow the prescribed secure inference protocol yet trying to deduce the information from the transcripts exchanged during the protocol execution. When corruption happens, a semi-honest adversary can compromise at most one of the edge nodes and either the mobile user or the model owner, while the other parties remain honest. We note that our considered semi-honest threat model is consistent with a great majority of prior works [4], [6], [7], [11], [14], [22], [30]), including the works under the two-server model [7], [11], [30].

Considering semi-honest mobile user, model owner, and the edge nodes makes sense and is practical in Leia's targeted NN inference applications. First, the engaged model owner and the edge nodes are from business-driven companies which do not willing to ruin their reputation and business models to behave in a malicious way and collusion. Moreover, we follow the existing works in the literature and consider the mobile user to be semi-honest. The primary objective is to protect the confidentiality of the valuable neural network model. For the

non-colluding assumption, we can regard them as from two distinct and well-established edge service providers (e.g., Microsoft Azure IoT Edge service [34] and Amazon Lambda@Edge service [35]), belonging to separated administrative domains and are hosted by the economical service providers to avoid collusion. It is worth noting that leveraging such a two non-colluding servers has become increasingly appealing in many industrial projects as well. Examples include Facebook's CrypTen [37] and Cape Privacy's TFEncrypted [38].

Leia guarantees both the privacy of mobile user's data and the model privacy. It hides both the user's data and the model values (i.e., the trained weights and coefficients) from being known by the edge nodes. Meanwhile, Leia is consistent with the security guarantees in prior neural network inference works [4], [5], [7], [11]. That is, the parameters of network architecture are considered as hyper-parameters already known by the edge nodes, including the number of layers, the sizes of weight matrices, and types of operations used in each layer. Such hyper-parameters are data independent and not proprietary since they are usually described in scientific and white papers. We are aware that a malicious user can exploit the inference service as a blackbox oracle to perform attacks to extract auxiliary information from prediction results. Like prior cryptographic inference systems [4]–[6], we emphasize that protecting against such attacks is a complementary problem beyond Leia's security scope [39]. Mitigation strategies can consider the adoption of differentially private training algorithms [40].

V. OUR PROPOSED DESIGN

A. Secure Linear Layers

We present in this section the secure realizations of linear layers, i.e., the secure convolutional layer (SCONV) and the secure fully connected layer (SFC). As mentioned above, they can be expressed as $VDP(\mathbf{x}, \mathbf{w}) + bias$ over n -dimensional layer input vector \mathbf{x} , weight vector \mathbf{w} , and the $bias \in \mathbb{R}$ attached on each neuron. Note that the hidden layer's input is the activation vector \mathbf{a} . All weights and activations are restricted as ± 1 in BNN, except the real-valued first layer input $\mathbf{x} \in \mathbb{R}^n$. In Leia, we encode $\mathbf{w}, \mathbf{a} \in \{+1, -1\}^n$ to $\mathbf{w}^b, \mathbf{a}^b \in \{1, 0\}^n$ based on the sign bits, i.e., $+1 \rightarrow 1$ and $-1 \rightarrow 0$. Here, we make an important observation from the machine learning literature [27] that the bias can be removed

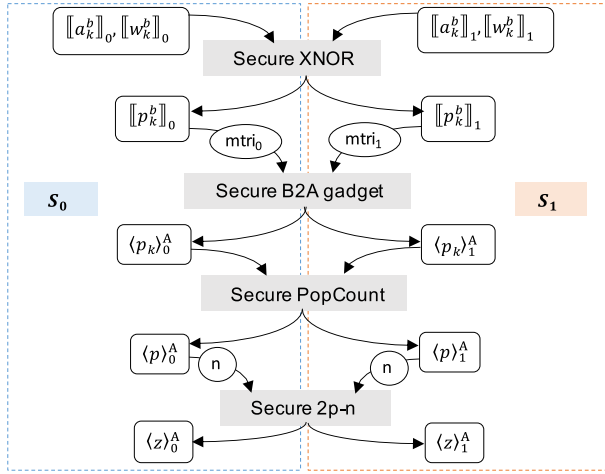


Fig. 2. An overview of the SecBVDP function.

if applying batch normalization, because the shift β in BN achieves the same effect as the bias. Well-known open source learning framework follows this treatment like PyTorch [41]. Likewise, we set the bias as 0 to avoid the involvement of real-valued bias and make our design more compatible with the MPC techniques.

To realize private linear transformation, our design carefully protects both mobile user's data (i.e., the input) and the BNN model (i.e., the weights) with lightweight secret sharing techniques (Arithmetic sharing and Boolean sharing). In particular, we introduce the secure Boolean-VDP function, and the secure Boolean-Arithmetic-VDP function which performs VDP over mixed share representations. They are the main building blocks to realize the linear layers.

1) *Secure Hidden Layer VDP*: The secure Boolean-VDP function (SecBVDP) computes VDP for the hidden layers. It takes as input a set of Boolean shared binary activation vector $\llbracket \mathbf{a}^b \rrbracket_i$ and weight vector $\llbracket \mathbf{w}^b \rrbracket_i$, and outputs the Arithmetic shares $\langle z \rangle_i^A$ of their Boolean-VDP result, where $i \in \{0, 1\}$ is the identifier of each edge node. The activation vector \mathbf{a}^b is the output of the binary activation function in BNN, which is naturally binarized. We note that the VDP operation on two plaintext binary vectors can be converted to a simpler XNOR-PopCount operation [15]. That is, for $\mathbf{a}, \mathbf{w} \in \{-1, +1\}$, the element-wise multiplication $p_k = a_k \cdot w_k$ is switched to bitwise-XNOR via $p_k^b = \text{XNOR}(a_k^b, w_k^b) = \neg(a_k^b \oplus w_k^b)$ when $\mathbf{a}^b, \mathbf{w}^b \in \{0, 1\}$, where $k \in [1, n]$. Meanwhile, the accumulation over all multiplication results $p = \sum_{k=1}^n p_k$ for $p_k \in \{-1, +1\}$ is converted to PopCount followed by a $2p-n$. That is, counting the number of "1"s in the resulting binary vector $\mathbf{p}^b = (p_1^b, p_2^b, \dots, p_n^b) \in \{0, 1\}$ as p and setting the result to $2p - n$.

Following this convention, Fig. 2 gives a high-level illustration of the SecBVDP function on secret-shared data. It consists of four atomic operations: the *secure XNOR*, the *secure B2A gadget*, the *secure PopCount*, and the *secure 2p-n*. The *secure XNOR* performs element-wise XNOR operation over every element of the shared binary input vector $\llbracket a_k^b \rrbracket$ and weight vector $\llbracket w_k^b \rrbracket$, and outputs the Boolean-shared XNOR results $\llbracket p_k^b \rrbracket$, where $k \in [1, n]$. Prior to the secure PopCount operation, the *secure B2A gadget* needs to be applied to every

Input: Boolean shares of binary activation vector $\mathbf{a}^b \in \{0, 1\}^n$, and binary weight vector $\mathbf{w}^b \in \{0, 1\}^n$.
Output: Arithmetic shares of Binary-VDP result $z = \text{VDP}(\mathbf{a}^b, \mathbf{w}^b)$.

Secure XNOR:

1) For each $k \in [1, n]$, S_i sets $\llbracket p_k^b \rrbracket_i = \llbracket a_k^b \rrbracket_i \oplus \llbracket w_k^b \rrbracket_i \oplus i$.

Secure B2A(\cdot) gadget:

2) S_0 and S_1 convert $\llbracket p_k^b \rrbracket \in \mathbb{Z}_2$ to $\langle p_k \rangle^A \in \mathbb{Z}_{2^\ell}$ as follows:

- S_0 sets two variables $\langle u_k \rangle_0^A = \llbracket p_k^b \rrbracket_0$, $\langle v_k \rangle_0^A = 0$;
- S_1 sets two variables $\langle u_k \rangle_1^A = 0$, $\langle v_k \rangle_1^A = \llbracket p_k^b \rrbracket_1$;
- S_0 and S_1 set $\langle p_k \rangle_i^A = \langle u_k \rangle_i^A + \langle v_k \rangle_i^A - 2 \cdot \langle u_k \rangle_i^A \cdot \langle v_k \rangle_i^A$.

Secure PopCount:

3) S_i counts the number of "1" via $\langle p \rangle_i^A = \sum_{k=1}^n \langle p_k \rangle_i^A$;

Secure 2p-n:

4) At the end, S_i sets $\langle z \rangle_i^A = 2\langle p \rangle_i^A - i \cdot n$.

Fig. 3. The secure Boolean-VDP function SecBVDP(\cdot, \cdot) based on XNOR-PopCount.

$\llbracket p_k^b \rrbracket$, converting from over \mathbb{Z}_2 to over \mathbb{Z}_{2^ℓ} . With the assist of pre-generated multiplication triples $\text{mtri}_0, \text{mtri}_1$, it outputs the Arithmetic-shared XNOR result $\langle p_k \rangle^A$. This is because that $\llbracket p_k^b \rrbracket$ is shared as $\llbracket p_k^b \rrbracket_0 + \llbracket p_k^b \rrbracket_1 \pmod{2}$, whereas the PopCount result is an aggregated integer that should be shared as Arithmetic shares $\langle p \rangle^A$. Naive sum $\langle p \rangle_i^A = \sum_{k=1}^n \llbracket p_k^b \rrbracket_i$ cannot correctly perform the modular addition over \mathbb{Z}_{2^ℓ} . Given an obvious example that $\mathbf{p}^b = (0, 0)$ with two "0"s and shared as $\llbracket \mathbf{p}^b \rrbracket_0 = (1, 1)$ and $\llbracket \mathbf{p}^b \rrbracket_1 = (1, 1)$, direct aggregation produces a wrong result "4" instead of the expected result "0". The *secure PopCount* then aggregates $\langle p_k \rangle^A$ to $\langle p \rangle^A$. At the end, the *secure 2p-n* is calculated with the system parameter n , i.e., the length of the vector.

Fig. 3 expatiates on the realization of the SecBVDP function corresponding to the above four operations. The *secure XNOR* is realized in step 1. Each edge node S_i locally calculates $\llbracket p_k^b \rrbracket_i = \llbracket a_k^b \rrbracket_i \oplus \llbracket w_k^b \rrbracket_i \oplus i$ to obtain its shared XNOR result. The *secure B2A gadget* is conducted in step 2. To do so, two variables are set to $u_k = \llbracket p_k^b \rrbracket_0$ and $v_k = \llbracket p_k^b \rrbracket_1$. S_0 and S_1 jointly perform the conversion to obtain their shares $\langle p_k \rangle_i^A$, following the expression $\langle p_k \rangle^A = \langle u_k + v_k - 2 \cdot u_k \cdot v_k \rangle^A \pmod{2^\ell}$. Thereafter, S_i locally computes the *secure PopCount* in step 3 and the *secure 2p-n* in step 4, and obtains the shared result $\langle z \rangle_i^A$ at the end.

2) *Secure First Layer VDP*: This subsequent section presents Leia's secure realization of the first layer, where the two inputs submitted to VDP calculation are the real-valued matrix of the user's data (e.g., image) and the binarized weight matrix.

a) *Common approach and its limitation*: To realize the first layer, a common way seems plausible is to protect both real-valued data matrix and binarized weight matrix via Arithmetic sharing, and then perform VDP on Arithmetic-shared data. However, protecting the weights as Arithmetic shares would substantially exaggerate the bandwidth to transmit the model and waste BNN's advancement. Besides, the multiplications over Arithmetic shares require the assistant of multiplication triples [29], and generating the disposable triples incurs intensive bandwidth costs that scale linearly with the number of triples. As the bandwidth is the bottleneck at the edge, overwhelming amount of bandwidth will decrease

Input: Arithmetic shares of integer input vector $\mathbf{x} \in \mathbb{Z}^n$, Boolean shares of binary weight vector $\mathbf{w}^b \in \{0, 1\}^n$.
Output: Arithmetic shares of result $z = \text{VDP}(\mathbf{x}, \mathbf{w}^b)$.
 The $n \times \text{COT}_\ell$ protocol with f_Δ :

- 1) The k -th COT_ℓ that $k \in [1, n]$ computes $\langle u_k \rangle^A$ of Eq. 1 as:
 - a) S_0 is the sender, and S_1 is the receiver;
 - b) S_0 sets f_Δ to Eq. 3, $m_0 = r_u \in_R \mathbb{Z}_{2^\ell}$; S_1 sets $\mathbf{b}_u = \llbracket w_k^b \rrbracket_1$;
 - c) S_0 and S_1 run $(\perp; m_{\mathbf{b}_u}) \leftarrow \text{COT}(m_0, f_\Delta(m_0); \mathbf{b}_u)$;
 - d) S_1 obtains $m_{\mathbf{b}_u}$ and sets $\langle u_k \rangle_1^A = m_{\mathbf{b}_u}$.
 - e) S_0 sets $\langle u_k \rangle_0^A = -r_u + \llbracket w_k^b \rrbracket_0 \cdot \langle x_k \rangle_0^A$;

The $n \times \text{COT}_\ell$ protocol with g_Δ :

- 2) The k -th COT_ℓ that $k \in [1, n]$ computes $\langle v_k \rangle^A$ of Eq. 2 as:
 - a) S_1 is the sender, and S_0 is the receiver;
 - b) S_1 sets g_Δ to Eq. 4, $m_0 = r_v \in_R \mathbb{Z}_{2^\ell}$; S_0 sets $\mathbf{b}_v = \llbracket w_k^b \rrbracket_0$;
 - c) S_1 and S_0 run $(\perp; m_{\mathbf{b}_v}) \leftarrow \text{COT}(m_0, g_\Delta(m_0); \mathbf{b}_v)$;
 - d) S_0 obtains $m_{\mathbf{b}_v}$ and sets $\langle v_k \rangle_0^A = m_{\mathbf{b}_v}$;
 - e) S_1 sets $\langle v_k \rangle_1^A = -r_v + \llbracket w_k^b \rrbracket_1 \cdot \langle x_k \rangle_1^A$.
- 3) At the end, S_i sets $\langle z \rangle_i^A = \sum_{k=1}^n (\langle u_k \rangle_i^A + \langle v_k \rangle_i^A)$.

Fig. 4. The secure Boolean-Arithmetic-VDP function $\text{SecBAVDP}(\cdot, \cdot)$ based on COT.

the overall performance and introduce additional charges by cellular network service provider.

b) The secure Boolean-Arithmetic-VDP function: To minimize the overall bandwidth costs at the edge, we craft the secure Boolean-Arithmetic-VDP function (SecBAVDP) in Fig. 4, allowing for direct multiplication on mixed share representations. It takes as input the Boolean shares $\llbracket \mathbf{w}^b \rrbracket$ of the binary weight vector $\mathbf{w} \in \{0, 1\}^n$ and the Arithmetic shares $\langle \mathbf{x} \rangle^A$ of the real-valued input vector $\mathbf{x} \in \mathbb{Z}^n$. and outputs the Arithmetic-shared result of $\text{VDP}(\mathbf{x}, \mathbf{w}^b)$ as $\langle z \rangle^A$. Our observation is that the element-wise multiplication $\langle z_k \rangle^A = \llbracket w_k^b \rrbracket \cdot \langle x_k \rangle^A$ can be expressed as

$$\begin{aligned} \langle z_k \rangle^A &= \langle u_k \rangle^A + \langle v_k \rangle^A; \\ \langle u_k \rangle^A &= \langle (\llbracket w_k^b \rrbracket_0 \oplus \llbracket w_k^b \rrbracket_1) \cdot \langle x_k \rangle_0^A \rangle^A; \end{aligned} \quad (1)$$

$$\langle v_k \rangle^A = \langle (\llbracket w_k^b \rrbracket_0 \oplus \llbracket w_k^b \rrbracket_1) \cdot \langle x_k \rangle_1^A \rangle^A. \quad (2)$$

Eq. 1 and Eq. 2 then can be efficiently calculated by two customized COT_ℓ protocols corresponding to the correlation functions f_Δ and g_Δ , respectively.

The first COT_ℓ protocol with f_Δ calculates $\langle u_k \rangle^A$, where S_0 acts as the sender and S_1 acts as the receiver. By treating $\llbracket w_k^b \rrbracket_1$ as the choice bit \mathbf{b}_u , the logic can be expressed as

$$\begin{aligned} \langle u_k \rangle^A &= (1 - \mathbf{b}_u) \cdot (\llbracket w_k^b \rrbracket_0 \cdot \langle x_k \rangle_0^A) + \mathbf{b}_u \cdot (\neg \llbracket w_k^b \rrbracket_0 \cdot \langle x_k \rangle_0^A) \\ &= \underbrace{(\llbracket w_k^b \rrbracket_0 \cdot \langle x_k \rangle_0^A)}_{S_0 \text{ at local}} + \underbrace{\mathbf{b}_u \cdot (\neg \llbracket w_k^b \rrbracket_0 \cdot \langle x_k \rangle_0^A - \llbracket w_k^b \rrbracket_0 \cdot \langle x_k \rangle_0^A)}_{\text{COT}_\ell \text{ with } f_\Delta \text{ and a choice bit } \mathbf{b}_u}. \end{aligned}$$

The former part of above formula is calculated by S_0 at local, while the latter part is performed by COT_ℓ with correlation function f_Δ . In particular, S_0 sets the correlation function as

$$f_\Delta(s) = s + (\neg \llbracket w_k^b \rrbracket_0 \cdot \langle x_k \rangle_0^A - \llbracket w_k^b \rrbracket_0 \cdot \langle x_k \rangle_0^A), \quad (3)$$

where s is any input from S_0 . Once invoked, S_0 sets his input m_0 as a random number r_u , while S_1 sets his input to \mathbf{b}_u . Then S_0 and S_1 run $(\perp; m_{\mathbf{b}_u}) \leftarrow \text{COT}_\ell(m_0, f_\Delta(m_0); \mathbf{b}_u)$.

Meanwhile, S_0 computes $\langle u_k \rangle_0^A = -r_u + \llbracket w_k^b \rrbracket_0 \cdot \langle x_k \rangle_0^A$. The reconstructed $\langle u_k \rangle^A$ is equivalent to Eq. 1.

Similarly, the second COT_ℓ protocol with g_Δ treats $\llbracket w_k^b \rrbracket_0$ as the choice bit \mathbf{b}_v and computes $\langle v_k \rangle^A$ as

$$\begin{aligned} \langle v_k \rangle^A &= \underbrace{(\llbracket w_k^b \rrbracket_1 \cdot \langle x_k \rangle_1^A)}_{S_1 \text{ at local}} \\ &\quad + \underbrace{\mathbf{b}_v \cdot (\neg \llbracket w_k^b \rrbracket_1 \cdot \langle x_k \rangle_1^A - \llbracket w_k^b \rrbracket_1 \cdot \langle x_k \rangle_1^A)}_{\text{COT}_\ell \text{ with } g_\Delta \text{ and a choice bit } \mathbf{b}_v}, \end{aligned}$$

where S_1 acts as the sender and S_0 acts as the receiver. In particular, S_1 sets the correlation function as

$$g_\Delta(s) = s + (\neg \llbracket w_k^b \rrbracket_1 \cdot \langle x_k \rangle_1^A - \llbracket w_k^b \rrbracket_1 \cdot \langle x_k \rangle_1^A), \quad (4)$$

where s is any input from S_1 . Once invoked, S_1 sets his input m_0 to a random number r_v , while S_0 sets his input to \mathbf{b}_v . S_1 and S_0 run $(\perp; m_{\mathbf{b}_v}) \leftarrow \text{COT}_\ell(m_0, g_\Delta(m_0); \mathbf{b}_v)$. S_1 sets $\langle v_k \rangle_1^A = -r_v + \llbracket w_k^b \rrbracket_1 \cdot \langle x_k \rangle_1^A$. This $\langle v_k \rangle^A$ is equivalent to Eq. 2. At the end, S_0 and S_1 locally aggregate $\sum_{k=1}^n (\langle u_k \rangle_i^A + \langle v_k \rangle_i^A)$ as his shared result $\langle z \rangle_i^A$. Note that our COT-based SecBAVDP requires to $2n$ calls of COT_ℓ with $2n(\ell + \lambda)$ bits and computing $2n \cdot 3$ hashing, where λ is the security parameter (128 in our work), while the OT-based approach requires $4n$ calls of OT with $4n\ell(\ell + \lambda)$ bits.

c) The secure first layer VDP function: To perform the secure first layer VDP function (Sec1VDP), we encode the weight vector \mathbf{w} as a tuple $(+\mathbf{w}^b, -\mathbf{w}^b)$, where $+\mathbf{w}^b, -\mathbf{w}^b \in \{0, 1\}^n$. That is, when an element $w = +1$, the corresponding tuple is $^+w \leftarrow 1$ and $^-w \leftarrow 0$; while when $w = -1$, it is encoded as $^+w \leftarrow 0$ and $^-w \leftarrow 1$. The Sec1VDP function takes as input the Arithmetic shares of integer input vector $\langle \mathbf{x} \rangle^A$, Boolean shares of binary weight vectors $\llbracket +\mathbf{w}^b \rrbracket, \llbracket -\mathbf{w}^b \rrbracket$, and outputs the Arithmetic shares of feature $\langle z \rangle^A$. Given the two edge nodes and the SecBAVDP function, the $\text{Sec1VDP}(\langle \mathbf{x} \rangle^A, \llbracket +\mathbf{w}^b \rrbracket, \llbracket -\mathbf{w}^b \rrbracket)$ proceeds as follows:

- 1) S_0 and S_1 run to get $\langle +z \rangle_i^A \leftarrow \text{SecBAVDP}(\llbracket +\mathbf{w}^b \rrbracket, \langle \mathbf{x} \rangle^A)$.
- 2) S_0 and S_1 run to get $\langle -z \rangle_i^A \leftarrow \text{SecBAVDP}(\llbracket -\mathbf{w}^b \rrbracket, \langle \mathbf{x} \rangle^A)$.
- 3) S_i locally computes $\langle z \rangle_i^A = \langle +z \rangle_i^A - \langle -z \rangle_i^A$.

B. Secure Batch Normalization and Binary Activation

1) Common Approach and Its Limitation: Batch normalization and binary activation are usually applied as a combination on each linear layer, following the linear transformation. Apart from the output layer, the prediction results are the output from the batch normalization without binary activation. At a high level, the combination of such two functions proceeds the functionality via

$$a^b = \text{sign}(\epsilon_1 \cdot a + \epsilon_2), \quad (5)$$

$$\epsilon_1 = \frac{\gamma}{\delta}, \epsilon_2 = \beta - \frac{\gamma \mu}{\delta} \quad (6)$$

where ϵ_1, ϵ_2 are preprocessed parameters derived from the trained BN parameters μ, δ, γ and β . During our model training procedure over plaintext, we observe that the trained ϵ_1, ϵ_2 are real-valued, where their integer parts before radix point are usually very small (i.e., 0 or 1) and the fractional

parts can last for a few digits (e.g., 10 digits). To handle the real-valued numbers in secure computation, a common way is to scale the ϵ_1, ϵ_2 to integers with a certain precision factor 2^q , followed by a ring conversion applied on the secret-shared a . Such a conversion is normally scaling up the $\langle a \rangle^A \in \mathbb{Z}_{2^\ell}$ to $\langle a' \rangle^A \in \mathbb{Z}_{2^\kappa}$ where $\kappa > \ell + q$. After sharing the ϵ_1, ϵ_2 in \mathbb{Z}_{2^κ} as $\langle \epsilon_1 \rangle^A, \langle \epsilon_2 \rangle^A \in \mathbb{Z}_{2^\kappa}$, the computation $\langle y \rangle^A = \langle \epsilon_1 \rangle^A \cdot \langle a' \rangle^A + \langle \epsilon_2 \rangle^A$ can be securely carried out over \mathbb{Z}_{2^κ} . And the $\text{sign}(\langle y \rangle^A)$ can be securely realized via a most significant bit (MSB) extraction.

The limitations of such a common way are two-fold. First, an additional ring conversion operation has to be applied on each neuron leading to higher computational costs. Second, the enlarged ring \mathbb{Z}_{2^κ} leads to a more complicated bitwise MSB extraction. In general, this MSB extraction operation follows the bit extraction protocol in [30] that performs non-local operations on the bit string of $\langle y \rangle^A \in \mathbb{Z}_{2^\kappa}$. It requires the interactions between the two edge nodes with the complexity scaling linearly with the length of the bit string (i.e., κ). So, a larger ring size leads to heavier bandwidth costs which is undesired at the edge. To address the above challenges, we propose two secure functions with special treatments: 1) the secure normalized binary activation function (SecNBA) for the first layer and hidden layers; and 2) the secure batch normalization function (SecBN) in the output layer.

2) *Secure Normalized Binary Activation*: The SecNBA function combines the secure batch normalization and the secure binary activation. We observe that its functionality defined as Eq. 5 can be transformed to

$$a^b = \text{sign}(\epsilon_1) \cdot \text{sign}(a + \epsilon) = \text{XNOR}(\zeta, z); \quad (7)$$

$$z = \text{sign}(y), y = a + \epsilon; \quad (8)$$

$$\zeta = \text{sign}(\epsilon_1), \epsilon = \frac{\epsilon_2}{\epsilon_1} = \frac{\delta\beta}{\gamma} - \mu. \quad (9)$$

Such a transformation results in a much simpler problem without the ring conversion. Through our careful examination, we identify that $\epsilon = \epsilon_2/\epsilon_1$ is real-valued number with large integer part. We thus quantize ϵ directly as integer and share it in \mathbb{Z}_{2^ℓ} , so as to circumvent the conversion between different rings. It is noteworthy that, a similar transformation has been also proposed in XONN [6] and BANNERS [13]. However, their designs do not consider the case that the parameter ϵ_1 is negative value. In fact, during our plaintext model training, we observed that ϵ_1 in batch normalization can be a negative value, and thus its sign bit $\zeta = \text{sign}(\epsilon_1)$ can not be ignored. Moreover, the multiplication between the sign bits ζ and z can be carried out via XNOR operation. Since ζ and ϵ are independent of inference input, they can be pre-generated by the model owner.

Given above equations, we present details of the secure realization of the SecNBA function. It takes as input the Arithmetic shares of the feature $\langle a \rangle^A$ that outputted from the linear transformation, the shares of two preprocessed parameters $\langle \epsilon \rangle^A$ and $\llbracket \zeta \rrbracket$, and outputs the Boolean shares of binary normalized activation $\llbracket a^b \rrbracket$. As summarized in Fig. 5, we decompose the computation as three atomic operations at a high-level, i.e., the *secure y*, the *secure sign* (secure MSB gadget + secure MSB to sign), the *secure XNOR*. The *secure y* takes as input the shares of the preprocessed parameter $\langle \epsilon \rangle^A$ and feature $\langle a \rangle^A$ and produces the shares $\langle y \rangle^A$ defined in Eq. 8.

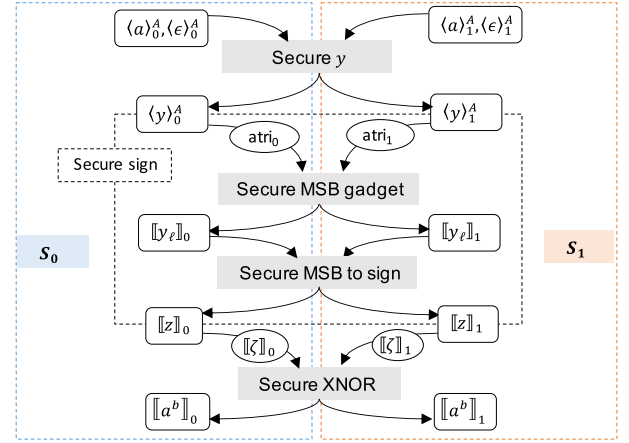


Fig. 5. An overview of the SecNBA function.

Input: Arithmetic shares of integer feature $a \in \mathbb{Z}$, Arithmetic shares of param. $\epsilon \in \mathbb{Z}$, Boolean shares of param. $\zeta \in \{0, 1\}$.
Output: Boolean shares of binarized activation $a^b \in \{0, 1\}$.
Secure y:
 1) S_i calculates $\langle y \rangle_i^A = \langle a \rangle_i^A + \langle \epsilon \rangle_i^A$.
Secure sign(\cdot):
 2) **Secure MSB(\cdot) gadget:** S_0 and S_1 run $\llbracket y_\ell \rrbracket_i \leftarrow \text{MSB}(\langle y \rangle_i^A)$.
 3) **Secure MSB to sign:** S_i sets $\llbracket z \rrbracket_i = \llbracket y_\ell \rrbracket_i \oplus i$.
Secure XNOR:
 4) S_i sets $\llbracket a^b \rrbracket_i = \llbracket z \rrbracket_i \oplus \llbracket \zeta \rrbracket_i \oplus i$.

Fig. 6. The secure normalized binary activation function SecNBA(\cdot, \cdot, \cdot).

The *secure sign* consists of the secure MSB gadget and the secure MSB to sign operations to extract the Boolean-shared sign bit $\llbracket z \rrbracket$. The secure MSB gadget extracts the shared MSB $\llbracket y_\ell \rrbracket$ of the input $\langle y \rangle^A$ with the assist of already generated Boolean AND triples $\text{atri}_0, \text{atri}_1$. Then $\llbracket y_\ell \rrbracket$ is converted to the shared sign bit $\llbracket z \rrbracket$ through the secure MSB to sign. The last *secure XNOR* produces the shared binary normalized activation $\llbracket a^b \rrbracket$ based on Eq. 7, given the shares of the sign bit $\llbracket z \rrbracket$ and preprocessed $\llbracket \zeta \rrbracket$.

Given above atomic operations, Fig. 6 details the corresponding realization of the SecNBA function. The *secure y* is realized in step 1, where each edge node S_i locally computes $\langle y \rangle_i^A$ by adding $\langle \epsilon \rangle_i^A$ to $\langle a \rangle_i^A$. Steps 2, 3 realize the *secure sign*. In step 2, the edge nodes S_0 and S_1 jointly execute the secure MSB gadget to obtain their shares of MSB, i.e., $\llbracket y_\ell \rrbracket_i$. This gadget employs the bit extraction protocol in [30], which is able to efficiently extract the MSB of the Arithmetic-shared values and produce a Boolean-shared MSB. The MSB is 0 of non-negative values (including 0) and 1 of the negative values, which is exactly the one's complement of a given sign bit. Then in step 3 (the secure MSB to sign), S_i performs logical negation on each $\llbracket y_\ell \rrbracket_i$ at local to obtain the shared sign bit $\llbracket z \rrbracket_i$. The *secure XNOR* is realized in step 4, where S_i conducts local XNOR operation over $\llbracket z \rrbracket_i$ and $\llbracket \zeta \rrbracket_i$, and finally gets its share of binary normalized activation $\llbracket a^b \rrbracket_i$. Details of the secure MSB gadget is given in Appendix.

Input: Boolean shares of m -number of n -dimensional binary activation vectors $\mathbf{a}_1^b, \dots, \mathbf{a}_m^b \in \{0, 1\}^n$, where the dimension n matches the size of pooling window.

Output: Boolean shares of m -number of binary pooling result $z_1^b, \dots, z_m^b \in \{0, 1\}$.

- 1) For $t \in [1, m]$, S_0 and S_1 compute the shares of maximum element in $[\mathbf{a}_t^b]_i = ([a_{t,1}^b]_i, \dots, [a_{t,n}^b]_i)$ based on Eq. 10:
 - a) For $k \in [1, n]$, S_i sets $[a_k^b]_i = [a_{t,k}^b]_i \oplus i$;
 - b) S_i sets variable $[c_t^b]_i = [a_{t,1}^b]_i$;
 - c) For $k \in [2, n]$, S_0 and S_1 set $[c_t^b]_i = [c_t^b]_i \wedge [a_{t,k}^b]_i$.
 - d) S_i sets $[z_t^b]_i = [c_t^b]_i \oplus i$.

Fig. 7. The secure max pooling function SecMP(\cdot).

3) *Secure Batch Normalization for Output Layer:* The secure batch normalization function (SecBN) is applied right after the secure linear transformation (SecBVDP) in the output layer. It takes as input the Arithmetic shares of activation $\langle a \rangle^A \in \mathbb{Z}_{2^\kappa}$ outputted from SecBVDP, and shares of two parameters $\langle \epsilon_1 \rangle^A, \langle \epsilon_2 \rangle^A \in \mathbb{Z}_{2^\kappa}$, and outputs the Arithmetic shares of the normalized activation $\langle z \rangle^A \in \mathbb{Z}_{2^\kappa}$. Here, the B2A gadget, i.e., the step 2 in the SecBVDP, performs the conversion from over \mathbb{Z}_2 to \mathbb{Z}_{2^κ} , and thus the output of SecBVDP is a feature already secret-shared in \mathbb{Z}_{2^κ} . Note that this ring conversion operation will not affect the following binary activation, as the output of SecBN is the shared inference result. The parameters ϵ_1, ϵ_2 are already enlarged during preprocessing, and secret shared in \mathbb{Z}_{2^κ} .

Given the two edge nodes, the pre-generated $\text{mtri}_i \in \mathbb{Z}_{2^\kappa}$, the SecBN($\langle a \rangle^A, \langle \epsilon_1 \rangle^A, \langle \epsilon_2 \rangle^A$) proceeds as follows: S_0 and S_1 compute $\langle z \rangle_i^A = \langle \epsilon_1 \rangle_i^A \cdot \langle a \rangle^A + \langle \epsilon_2 \rangle_i^A$ to obtain their shares of normalized activation.

C. Secure Binary Max Pooling Layer

The secure binary max pooling layer (SMP) is used to obtain the maximum values among the secret-shared binary activations within a certain sliding window. The number of n binary activations within the window can be denoted as a n -dimensional binary activation vector $\mathbf{a}^b = (a_1^b, \dots, a_n^b) \in \{0, 1\}^n$. We assume that there are overall m -number of vectors as $\mathbf{a}_1^b, \dots, \mathbf{a}_m^b$. We observe that its functionality over plaintext \mathbf{a}^b can be realized as the bitwise-OR operation on all bits of the vectors, i.e., the maximum value is $z^b = a_1^b \vee a_2^b \vee \dots \vee a_n^b$, so as to find if \mathbf{a}^b constitutes with any “1” bit. However, the key takeaway of a secure realization is to achieve obliviousness, i.e., for every step in a certain computation, both edge nodes have to proceed equivalent operations. Through carefully examination, we transform the logic to

$$z^b = \neg(\neg a_1^b \wedge \neg a_2^b \wedge \dots \wedge \neg a_n^b) \quad (10)$$

which is more compatible with our secret sharing based realization.

With this philosophy in mind, we present in Fig. 7 the proposed secure max pooling function (i.e., SecMP) design specialized for the two edge nodes case as the main building block of the MP layer. It takes as input the Boolean shares of m -number of n -dimensional binary activation vectors $[\mathbf{a}_t^b]$, determines the maximum element for each vector, and outputs their Boolean shares as the pooling result $[z_t^b]$, where

$t \in [1, m]$. For each $[\mathbf{a}_t^b]$, S_0 and S_1 securely proceeds Eq. 10 via the steps below. In step 1.a, S_i securely realizes $\neg a_{t,k}^b$ by XORing the share $[a_{t,k}^b]_i$ to its identifier i , where $k \in [1, n]$. In steps 1.b and 1.c, S_0 and S_1 iteratively perform AND over all of its shares, i.e., $[c_t^b]_i = [a_{t,1}^b]_i \wedge \dots \wedge [a_{t,n}^b]_i$. In step 1.d, S_i sets its output share $[z_t^b]_i$ as its identifier i XOR with $[c_t^b]_i$. We emphasize that all operations performed by the two nodes are identical, and thus endowing our SecMP function obliviousness.

D. Secure BNN Inference Protocol

Given above layer functions, we now describe our secure BNN inference protocol ϕ . It comprises two phases: the *preprocessing phase* performed by each entity individually, and the *secure inference phase* jointly carried out by the two non-colluding edge nodes.

1) *Preprocessing Phase:* During the preprocessing phase, the mobile user converts its task-specific raw input to a tensor $\mathbf{X} \in \mathbb{Z}^{c_{in} \times n_{in} \times m_{in}}$ and deploys the corresponding Arithmetic-shared tensors to S_0 and S_1 , respectively. Once received, they partition and flatten the shared tensor into a set of vectors, and the size of each vector equals to the sliding window size for the ease of subsequent VDP operations. Let \mathcal{T}^1 be the total moves to slide the window for the first layer. After flattening, S_0 and S_1 hold vectors $\langle (\mathbf{x})_{c,\tau} \rangle^A$, where $c \in [1, c_{in}]$ denotes the input channel, and $\tau \in [1, \mathcal{T}^1]$ indicates the τ -th sliding window. They also prepare triples during vacant time.

The model owner holds an L -layer BNN model. Each layer $l \in [1, L]$ is formed from a set of binarized weight vectors $(\mathbf{w}^l)_k$, i.e., a number of k vectors in the layer l . For the CONV and MP, $k = c_o^l$ (the number of output channels) and the length $|\mathbf{w}| = n_w \times n_w$. For FC, $k = n^l$ (the neurons of current layer) and $|\mathbf{w}| = n_{l-1}$ (the neurons of previous layer). For the first layer, the weight vectors are encoded as tuples $(+1 \rightarrow (1, 0)$ and $-1 \rightarrow (0, 1))$ as input of the Sec1VDP function, denoted as $[(+\mathbf{w}^l)_k]$, $[(-\mathbf{w}^l)_k]$. For hidden layers, the weight vectors are encoded based on the sign $(+1 \rightarrow 1$ and $-1 \rightarrow 0)$, denoted as $[(\mathbf{w}^l)_k]$. For BN, the model owner computes $\epsilon_1, \epsilon_2, \zeta$ according to Eq. 6 and Eq. 8, and generates shares $[(\zeta)_k^l]$, $\langle (\epsilon)_k^l \rangle^A \in \mathbb{Z}_{2^\ell}$, $\langle (\epsilon_1)_k^l \rangle^A, \langle (\epsilon_2)_k^l \rangle^A \in \mathbb{Z}_{2^\kappa}$. To this end, the model owner deploys all shares to the corresponding edge nodes S_0 and S_1 for coordinate processing.

2) *Secure Inference Phase:* Fig. 8 depicts the secure inference phase of an L -layer BNN. The inputs are the shares generated during preprocessing, including the shared user input vectors from the mobile user, the shared weight vectors and parameters from the model owner. The outputs are the last layer’s activations mapping a certain classification label after reconstruction. As NN can have distinct architectures assembled with layer functions, we present a typical one for demonstration purpose. It comprises the first SCONV layer, followed by an SMP layer, and $L - 2$ number of the SFC layers. For each linear layer, the SecNBA function is applied after the linear transformation. And the SecBN function is applied to the output layer.

For the first SCONV layer, S_0 and S_1 repeatedly execute the Sec1VDP function on shared weight and user’s data, inside every τ -th sliding window. The outputs are then summed across input channels as a set of features $\langle (a)_{k,\tau}^1 \rangle^A$ for each

Let $k \in [1, c_o^l]$, $\tau \in [1, T^l]$ for SCONV, SMP; $k \in [1, n^l]$ for SFC.
 First SCONV layer, $l = 1$:

- 1) S_0, S_1 run $\langle (a)_{k,\tau}^1 \rangle_i^A = \sum_{c=1}^{cin} \text{Sec1VDP}(\langle (w)_{k,\tau}^1 \rangle_i^A, \langle (x)_{c,\tau}^1 \rangle_i^A)$.
- 2) S_0, S_1 run $\langle (a)_{k,\tau}^1 \rangle_i \leftarrow \text{SecNBA}(\langle (a)_{k,\tau}^1 \rangle_i^A, \langle (\epsilon)_{k,\tau}^1 \rangle_i^A, \langle (\zeta)_{k,\tau}^1 \rangle_i^A)$.

SMP layer, $l = 2$:

- 3) S_0, S_1 run $\langle (a)_{k,\tau}^2 \rangle_i \leftarrow \text{SecMP}(\langle (a)_{k,\tau}^1 \rangle_i^A)$.

Remaining SFC layers, $l \in [3, L]$:

- 4) S_i flattens $\langle (a)_{k,\tau}^{l-1} \rangle_i = \langle (a)_{k,\tau}^{l-1} \rangle_{i,1}, \dots, \langle (a)_{k,\tau}^{l-1} \rangle_{i,n^{l-1}}$.
- 5) S_0, S_1 run $\langle (a)_{k,\tau}^l \rangle_i \leftarrow \text{SecBVDP}(\langle (a)_{k,\tau}^{l-1} \rangle_i^A, \langle (w)_{k,\tau}^l \rangle_i^A)$, where $(a)_{k,\tau}^{l-1}, (w)_{k,\tau}^l \in \{0, 1\}^{n^{l-1}}$.
- 6) If $l \neq L$, S_0, S_1 run $\langle (a)_{k,\tau}^l \rangle_i \leftarrow \text{SecNBA}(\langle (a)_{k,\tau}^l \rangle_i^A, \langle (\epsilon)_{k,\tau}^l \rangle_i^A, \langle (\zeta)_{k,\tau}^l \rangle_i^A)$.
- 7) Else, S_0, S_1 run $\langle (z)_{k,\tau}^L \rangle_i \leftarrow \text{SecBN}(\langle (a)_{k,\tau}^L \rangle_i^A, \langle (\epsilon)_{k,\tau}^L \rangle_i^A, \langle (\zeta)_{k,\tau}^L \rangle_i^A)$.
- 8) S_i outputs $\langle (z)_{k,\tau}^L \rangle_i^A$.

Fig. 8. Secure BNN inference phase of protocol ϕ .

output channel k , and submitted to the **SecNBA** function to obtain the shares of normalized binary activations. Afterwards, the **SMP** layer run the **SecMP** function to down sample the previous layer's activation vector into a single activation within each τ -th pooling window. The remaining $L - 2$ layers are the **SFC** layers. S_0 and S_1 firstly flatten the feature map outputted from the **SMP** layer as an one channel vector $\langle (a)_{k,\tau}^{l-1} \rangle_i$ across all c_o^{l-1} channels. Thereafter, for the l -th **SFC** layer, S_0 and S_1 execute the **SecBVDP** function on $\langle (w)_{k,\tau}^l \rangle_i^A$ and $\langle (a)_{k,\tau}^{l-1} \rangle_i^A$, and obtain features $\langle (a)_{k,\tau}^l \rangle_i^A$, where n_l is the number of neurons of the current layer and $k \in [1, n^l]$. For every feature, the **SecNBA** function is applied to obtain the activations $\langle (a)_{k,\tau}^l \rangle_i^A$. Likewise, the **SecBN** function is applied to the output layer L to obtain the results $\langle (z)_{k,\tau}^L \rangle_i^A$, where $k \in [1, n^L]$. To this end, S_0 and S_1 obtain the shares of inference result mapping a certain classification label. They then send back the shares of result to the mobile user who can reconstruct to get the prediction. The remark of complexity is provided in Section II of the supplementary materials.

a) Security guarantees: For our secure BNN inference protocol ϕ , we define security based on the *Universally Composable* (UC) security framework [42]. Under a general protocol composition operation (*universal composition*), the security of ϕ is preserved. Given a semi-honest *admissible adversary* \mathcal{A} who can compromise at most one of the two non-colluding edge nodes S_0, S_1 and either the mobile user or the model owner. This reflects on the property that S_0, S_1 are non-colluding servers, i.e., if S_0 is compromised by \mathcal{A} , S_1 acts honestly; vice versa. Leia's protocol follows the security of the Arithmetic sharing [20], Boolean sharing [19] and COT [28]. Leia properly protects the user data, model, Beaver's triples, and intermediate results outputted from layer functions as secret shares in \mathbb{Z}_{2^t} and \mathbb{Z}_2 . Given above, we argue that ϕ UC-realizes an ideal functionality \mathcal{F} against \mathcal{A} . The security captures the property that the only data learned by any compromised parties are their inputs and outputs from ϕ , but nothing about the data of the remaining honest parties.

b) Differences from prior art: We emphasis that Leia and XONN are different regarding the system models, application scenarios, the utilized cryptographic tools, and the designs of the secure layer functions.

Leia's overarching goal is to design secure NN inference system amiable for the recourse-constrained mobile devices.

Such resources encompass the hardware designs with limited computational power, stringent energy consumption, and more importantly, the unstable cellular network environment. To embrace the above rigid operational demands, we leverage edge computing to *fully delegate* our system to the edge devices, and as such the mobile devices do not need to always stay online during the secure inference. In particular, all secure computations are executed amongst the co-located edge nodes including the interactions. In comparison, XONN focuses on the scenario where an *interactive protocol* is executed between the client and the server. That is, XONN requires the client have *symmetric computational capabilities* to the server and always engaging in the whole secure inference with continuous interactions, which is not applicable to deploy in the dynamic cellular network and constrained resources devices. Meanwhile, Leia's edge-aided system architecture facilitates the model owner *dynamically* fine-tuning its service, where the neural network model can be regularly updated without republishing the mobile application, which in contrast is not enabled in XONN.

Despite the different system models, the realizations of the secure layer functions in XONN and Leia are entirely different. At a high-level, XONN directly resorts to the generic two-party secure computation framework, i.e., Yao's GC (with optimizations) and OT, to securely realize each layer function involved in BNN. In comparison, Leia crafts and realizes the secure layer functions fully with lightweight cryptographic tools, i.e., the Secret Sharing techniques and COT; wherein each proposed building block underpinning the secure layer functions are carefully designed to be suitable for the edge computing paradigm. The GC based approaches require substantial network resources and typically introduce larger latency than the secret-sharing based realizations [18], [22], [43]. Moreover, our designed COT for the secure binary-integer vector dot product saves half of the communication rounds to the oblivious condition addition based function, and this saving could be very significant enhancement due to the massive multiplication operations in processing NN inference.

We observe that our work [44] is concurrent and independent with prior art FALCON [16] and BANNERS [13], yet their designs are fundamentally different from ours. Both FALCON and BANNERS focus on the honest-majority malicious security under the three-server setting. Their secure protocols are built upon the 2-out-of-3 replicated secret sharing, while Leia is built upon additive secret sharing techniques [19], [20] and correlated oblivious transfer [28]. Besides, the implementation and real-world deployment of three-server protocols are more complex, whereas our secure and lightweight two-server inference protocol is advantageous in implementation and practical deployment to the edge devices.

VI. PERFORMANCE EVALUATION

A. Implementation and Setup

We implement a prototype of Leia in Java. All experiments are executed on two Raspberry Pi devices to simulate the edge environment. The devices are Raspberry Pi 4 Model B running Raspbian Linux 10 (buster) and equipped with Quad core Cortex-A72 (ARM v8) 64-bit SoC @ 1.5GHz processor, 4GB RAM, and gigabit Ethernet. Consist with prior



Fig. 9. Deployment on Raspberry Pi with power meter.

TABLE III
BANDWIDTH COST OF ATOMIC LAYER FUNCTIONS (IN KB)

# inputs	SecBVDP		SecBAVDP		SecNBA	SecBN	SecMP
	3×3	5×5	3×3	5×5			
Leia	0.3	0.8	16.0	17.3	0.2	0.03	0.01
GC-baseline	22.1	24.4	468.0	1257.4	25.9	78.3	20.9
savings	73×	30×	29×	130×	146×	2610×	2090×

art [5], [6], we evaluate Leia in the LAN setting. We use FlexSC [45] for the Extended OTs [28] and implement our designed COT protocol (i.e., the SecBAVDP function). Regarding Arithmetic sharing, we set the size of the ring as $\mathbb{Z}_{2^{32}}$ for the first layer and output layer, and $\mathbb{Z}_{2^{16}}$ for the remaining hidden layers. The reported measurements make use of the MNIST and CIFAR-10 datasets, i.e., the two commonly-used classification benchmarks in prior work [6], [7]. We evaluate Leia on a variety of different BNN models, where the models M1 and M2 are trained on MNIST, and the models C1 and C2 are trained on CIFAR-10. To demonstrate Leia's practicability in real-world applications, we further evaluate Leia on four medical datasets, i.e., breast cancer [46], diabetes [47], liver disease [48], and thyroid [49] on the models D1, D2, D3, D4, respectively. The details of our adopted model architectures can be found in Section III of the supplementary materials. For model training, we use PyTorch backend with standard BNN training algorithm [15]. We further use COOWOO power meter [50] to evaluate the energy consumption of Leia when deploying the real-world medical applications. Fig. 9 demonstrates our deployment.

B. Evaluation

1) *Microbenchmarks*: We present performance benchmarks of secure layer functions as the basic building blocks used for secure BNN inference. For demonstration purpose, we choose 3×3 and 5×5 sliding windows to show the performance of the SecBVDP and the SecBAVDP functions, i.e., the secure VDP operations over 9-dimensional vectors and 25-dimensional vectors, respectively. These two window sizes are common-used and adapted to our CONV layer. Likewise, we employ the 2×2 window to demonstrate the performance of the SecMP function.

We summarize the computational cost of the proposed secure layer functions in Table IV. The time consumption of the SecBAVDP function consists of two parts: 1) the constant initialization cost of the COT protocol ($\sim 3s$); and

TABLE IV
TIME COST OF ATOMIC LAYER FUNCTIONS (IN S)

# inputs	SecBVDP		SecBAVDP		SecNBA	SecBN	SecMP
	3×3	5×5	3×3	5×5			
10^3	0.5	0.5	4.8	4.9	0.6	0.4	0.06
10^4	3.1	3.4	20.4	22.8	5.1	2.6	0.4
10^5	29.4	31.6	198.3	217.9	50.1	23.0	3.1

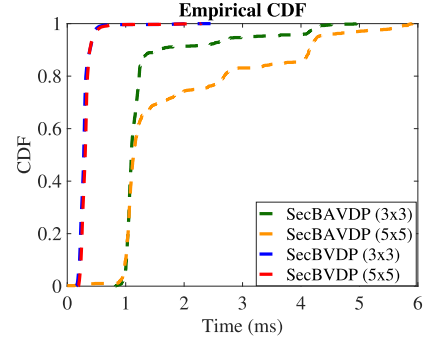


Fig. 10. Unit time of linear functions.

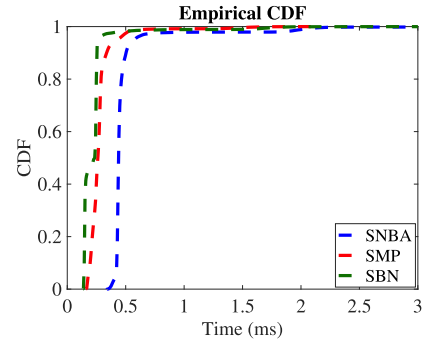


Fig. 11. Unit time of nonlinear functions.

2) the time to compute VDP over mixed share representations raising linearly with the number of calls. For the rest functions, their latencies ascend linearly in the growth of the number of executions yet with slight fluctuations. Besides, we grasp 10K executions of each secure layer function, and utilize the empirical cumulative distribution function (ECDF) to shed light on the distribution of their unit execution time. Fig. 10 depicts the distribution of unit run time of the secure linear functions, i.e., the SecBVDP and SecBAVDP functions. For overwhelming amount of executions, the unit executions of the SecBVDP function with 3×3 and 5×5 windows can be completed within 0.5ms. Besides, the unit execution time of the SecBAVDP function without the aforementioned constant COT initialization cost. As shown, more than 90% executions take 1ms and 4.5ms for 3×3 and 5×5 windows, respectively. Fig. 11 exhibits the time costs of single executions of the nonlinear functions, i.e., the SecNBA, SecBN, SecMP (with 2×2 window) functions. All three functions can be done within 1ms.

a) *Comparison with GC-based realization*: The communication costs of the secure layer functions are reported in Table III. We implement and evaluate the baseline based on GC with its free-XOR and half-AND optimizations, which

TABLE V
PERFORMANCE OF THE SCONV FUNCTION OF HIDDEN LAYERS

model	input	kernel	feature	stride, padding	#SecBAVDP	time (s)	comm. (MB)
M2	$16 \times 12 \times 12$	$16 \times 16 \times 5 \times 5$	$16 \times 8 \times 8$	1, -	16×1024	0.9	12.5
C1	$16 \times 32 \times 32$	$16 \times 16 \times 3 \times 3$	$16 \times 32 \times 32$	1, 0	16×16384	6.2	63.3
C2	$16 \times 32 \times 32$	$32 \times 16 \times 3 \times 3$	$32 \times 32 \times 32$	1, 0	16×32768	14.5	126.6

input: $c_{in} \times n_{in} \times m_{in}$; kernel: $c_o \times c_{in} \times n_w \times n_w$; output: $c_o \times n_o \times m_o$.

realizes the equivalent functionality for each of the secure layer function. In general, Leia's realizations require $30-79\times$, and $150-2500\times$ less communication for the linear and non-linear functions than the corresponding GC-based realizations. In detail, for the secure linear functions, the communication of Leia is $73\times$ and $30\times$ less for the SecBVDP function, and $29\times$ and $130\times$ less for the SecBAVDP function, over 3×3 and 5×5 windows respectively. For the non-linear functions, Leia achieves $146\times$, $2610\times$, and $2090\times$ bandwidth savings of the SecNBA, SecBN and SecMP costs compared with GC-based realizations. The reported results testify that the prior constructions relying on GC [6], [7], [11], [12] require a network environment with high bandwidth. They might not be applicable for our considered application scenario, i.e., the secure inference deployed at the edge with limited network conditions.

In particular, for the COT-based SecBAVDP function, we emphasize that the adoption of such regime saves the overall bandwidth consumption at a system level. Such retrenchment includes the cost of protecting each weight element as 32-bit shares in $\mathbb{Z}_{2^{32}}$ to a tuple of 1-bit shares in \mathbb{Z}_2 , and the cost of generation of multiplication triples in $\mathbb{Z}_{2^{32}}$. As shown by the empirical result, the GC-based realizations produce $30\times$ and $73\times$ bandwidth consumptions higher than the Leia's realizations for 9-dimensional and 25-dimensional vectors, respectively. We further report the bandwidth costs of the realizations based on multiplication triples as 270KB and 790KB, amounting to one magnitude larger than Leia's bandwidth.

2) *Linear Layers*: We report the performance of secure linear transformations (i.e., SCONV and SFC) below, which comprise the majority of Leia's overall inference overhead.

Table VI and Table V benchmark the performance of the SCONV layer function as the first layer and the hidden layer, respectively. The reported results are in line with our specified network architectures of M2, C1, and C2. As they consist plenty of convolutional hidden layers, we choose to show the performance of their second layers (the most complicated hidden layers) for the ease of demonstration. Note that, the M1, D1, D2, D3, D4 networks consist of only fully connected layers. The complexity of the SCONV layer function is determined by a set of parameters: 1) the number of input channels c_{in} and output channels c_o ; 2) the dimensions of input image; 3) the kernel size (i.e., the sliding window size s), stride, and padding regime. These parameters directly reflect on the number of invocations of SecBAVDP/SecBVDP as shown. The key takeaway here is our runtime optimization of batch processing to amortize the overhead of executing SecBAVDP/SecBVDP. In detail, we flatten the input matrices across multiple channels yet within the same sliding window as a single vector, and conduct SecBAVDP/SecBVDP

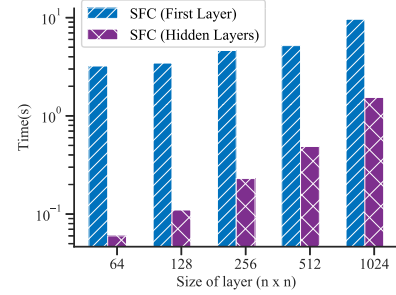


Fig. 12. Time cost of the SFC layer.

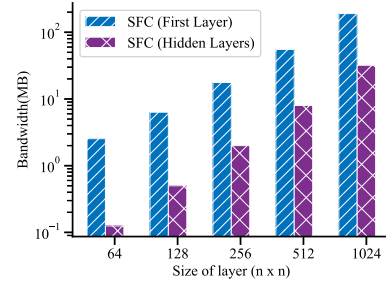


Fig. 13. Comm. cost of the SFC layer.

over it in a batch. We take as an example the complexity of C1's hidden SCONV layer reported in Table V. It is proceeded in the batch integrating with 16-channel input matrices. As a result, the calls of SecBVDP (3×3 window) are reduced from 230400 to 16384, speeding up the time from 68s to 6.2s accordingly.

Fig. 12 and Fig. 13 depict the computational and communication overheads of the SFC layer function as the first layer and hidden layers, respectively. They are evaluated over a series of $n \times n$ fully connected layers, i.e., both the input and weight are n -dimensional vectors. Followed by the growth of n , the time of the hidden SFC layer ascends linearly attributed to our batch processing optimization, while the bandwidth ascends quadratically with the growth of dimension n . For the first SFC layer, the computational overhead is primarily dominated by the constant COT initialization time, and the bandwidth grows with the layer size.

3) *Leia's Protocol on MNIST and CIFAR-10*: We evaluate Leia's cryptographic inference protocol on MNIST and CIFAR-10 datasets. Table VIII summarizes the overall performance. We overview the network architectures here, and present more details in Section III of the supplementary materials. The online phase of Leia is executed at the edge. The networks M1 and M2 for MNIST dataset are relatively simple, and Leia can produce high-quality prediction results within 4s

TABLE VI
PERFORMANCE OF THE SCONV FUNCTION OF FIRST LAYER

model	input	kernel	feature	stride, padding	#SecBAVDP	time (s)	comm. (MB)
M2	$1 \times 28 \times 28$	$16 \times 1 \times 5 \times 5$	$16 \times 24 \times 24$	1, -	1×18432	30	310
C1/C2	$3 \times 32 \times 32$	$16 \times 3 \times 3 \times 3$	$16 \times 32 \times 32$	1, 0	3×32768	46	490

input: $c_{in} \times n_{in} \times m_{in}$; kernel: $c_o \times c_{in} \times n_w \times n_w$; output: $c_o \times n_o \times m_o$.

TABLE VII
PERFORMANCE SUMMARY OF THE MEDICAL APPLICATIONS

network	time (s) ^a edge	comm. (MB) edge	time (μ s) mobile user	time (ms) model owner	accuracy Leia	accuracy plaintext
D1	3.15	0.57	42.1	2	98.23%	97.37%
D2	3.19	0.65	24.2	1.5	74.14%	80.17%
D3	3.22	1.06	24.8	2.4	78.45%	80.17%
D4	3.64	3.67	47.4	9.7	92.04%	93.64%

^a time at edge includes ~ 3 s OT initialization time.

TABLE VIII
PERFORMANCE SUMMARY OF THE BENCHMARKING NETWORKS

dataset	network	time (s) ^a edge	comm. (MB) edge	time (ms) ^b mobile user	time (s) ^c model owner	accuracy Leia	accuracy plaintext	layers
MNIST	M1	4.0	19.7	0.4	0.05	97.0%	97.0%	3SFC, 2SecNBA, 1SecBN
	M2	37.4	328.1	0.4	0.5	99.12%	99.12%	2SCONV, 2SFC, 2SMP, 3SecNBA, 1SecBN
CIFAR-10	C1	123.1	919.4	1.1	5.1	71.68%	69.03%	9SCONV, 1SFC, 3SMP, 9SecNBA, 1SecBN ^d
	C2	199	1829.9	1.2	15.7	81.0%	77.88%	

^a Time at edge includes ~ 3 s OT initialization time.

^b Cost of generating shares of an image during preprocessing.

^c One-time cost of generating shares of the model during preprocessing.

^d C1 and C2 have same layers but with different weight size.

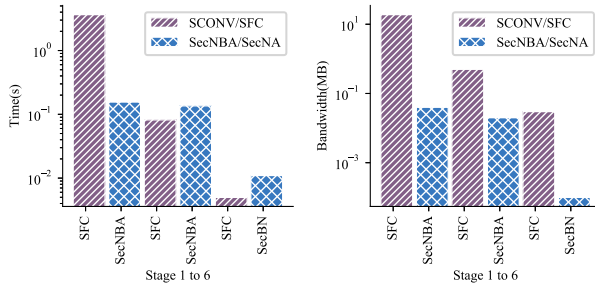


Fig. 14. Performance breakdown of M1. Left: time cost. Right: bandwidth cost.

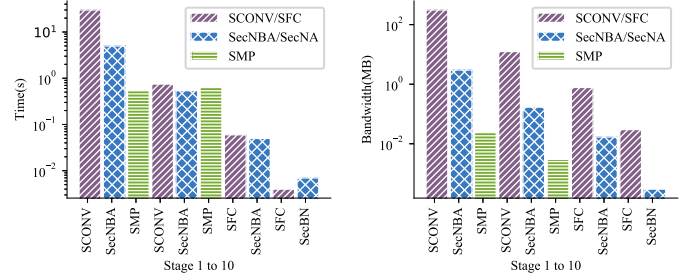


Fig. 15. Performance breakdown of M2. Left: time cost. Right: bandwidth cost.

and 37.4s, respectively. The more complex C1 and C2 for CIFAR-10 dataset involve 13 layers (23 stages), and their executions require about 2min and 3.3min respectively. The workload of the mobile user is light, which confirms that Leia is amiable to the resource limited portable devices. The one-time overhead of the model owner is determined by the model size. Such cost does not aggravate workload on the model owner, as generating shares of the most complicated network C2 can be completed within 15.7s.

a) Performance breakdown: To gain a more comprehensive understanding of resource consumption, we demonstrate the performance breakdown of each network. Fig. 14 and Fig. 15 show the time cost (left figure) and bandwidth cost (right figure) for each stage of M1 (6 stages) and M2 (10 stages) on MNIST dataset, respectively. Since C1 and C2

share the same architecture (different weight size), Fig. 16 reports the time (top figure) and bandwidth (bottom figure) for each stage of C1 (23 stages) for demonstration purpose. As seen, the first layer occupies most of the resources, and the linear functions usually require more workload than the non-linear functions.

b) Accuracy: The effectiveness is demonstrated in Table VIII via the accuracy comparison between Leia's prediction results and the plaintext's results. For the M1 and M2 networks evaluated on MNIST dataset, Leia's prediction results are accurate as the plaintext (i.e., 97% and 99%, respectively). Besides, Leia achieves the accuracy of 69% and 81% for the C1 and C2 networks evaluated on CIFAR-10 dataset, amounting to slight accuracy impacts compared with the plaintext results. Such variations can be attributed to the quantization of batch normalization parameters,

Input: Arithmetic shares of integer feature $y \in \mathbb{Z}$.
Output: Boolean shares of MSB $y_\ell \in \{0, 1\}$.

- 1) S_i decomposes $\langle y \rangle_i^A$ to a bit string $\langle y_1 \rangle_i^A, \dots, \langle y_\ell \rangle_i^A$;
- 2) For each $k \in [1, \ell]$:
 S_0 sets $\llbracket u_k \rrbracket_0 = \langle y_k \rangle_0^A$, $\llbracket v_k \rrbracket_0 = 0$, $\llbracket t_k \rrbracket_0 = \langle y_k \rangle_0^A$;
 S_1 sets $\llbracket u_k \rrbracket_1 = 0$, $\llbracket v_k \rrbracket_1 = \langle y_k \rangle_1^A$, $\llbracket t_k \rrbracket_1 = \langle y_k \rangle_1^A$;
 S_0 and S_1 set $\llbracket d_k \rrbracket_i = \llbracket u_k \rrbracket_i \wedge \llbracket v_k \rrbracket_i$ in a batch;
- 3) S_i sets variable $\llbracket c_1 \rrbracket_i = \llbracket d_1 \rrbracket_i$;
- 4) For $k \in [2, \ell - 1]$:
 S_i sets $\llbracket d_k \rrbracket_i = \llbracket d_k \rrbracket_i \oplus i$;
 S_0 and S_1 set $\llbracket e_k \rrbracket_i = \llbracket t_k \rrbracket_i \wedge \llbracket c_{k-1} \rrbracket_i \oplus i$;
 S_0 and S_1 set $\llbracket c_k \rrbracket_i = \llbracket e_k \rrbracket_i \wedge \llbracket d_k \rrbracket_i \oplus i$;
- 5) S_i sets the MSB to $\llbracket y_\ell \rrbracket_i = \llbracket t_\ell \rrbracket_i \oplus \llbracket c_{\ell-1} \rrbracket_i$.

Fig. 17. The secure MSB(\cdot) gadget.

in BNN, including the binarized linear layers, the commonly used non-linear binary activation function (the sign function), and the max pooling over binarized weights. These essential secure computational blocks can be scaled to support more neural networks. For example, our secure sign function can be used as a building block of secure ReLU function and secure max pooling over integers. Moreover, our proposed secure inference system can be deployed to other inference tasks, ranging from medical image segmentation, object detection, to natural language processing.

VII. CONCLUSION

In this paper, we propose Leia, a lightweight cryptographic NN inference system at the edge. Leia resorts to the edge based architecture, to foster a low-latency service and relax the constraint of the model owner and mobile device being online. To cater for the operational needs of edge environment, Leia is co-designed with the advancement from both machine learning and cryptographic areas. With the highly customized secure layer functions on binarized neural network, Leia enables an oblivious inference service guaranteeing both user and model privacy. Comprehensive empirical validation on benchmark and medical datasets demonstrates Leia practical and applicable for the real-world scenarios.

APPENDIX

Fig. 17 presents the secure MSB(\cdot) gadget. It follows the bit extraction protocol in [30], which is able to efficiently extract the MSB of the Arithmetic-shared values and produce a Boolean-shared MSB. The protocol takes as input the Arithmetic shared of integer feature $\langle y \rangle^A \in \mathbb{Z}^{2^\ell}$, extracts the ℓ -th bit as MSB of y , and outputs its Boolean shares $\llbracket y_\ell \rrbracket$. Let $y = u + v \pmod{2^\ell}$. The idea is that the difference between the sum of bit strings of u, v and the bitwise-XOR of the bit strings of u, v, y is equal to the carry bits c_1, \dots, c_ℓ . This can be realized by an ℓ -bit ripple carry logic, where every carry bit is calculated by a full adder and propagated to the next full adder, and finally the MSB of y is outputted by the ℓ -th full adder.

REFERENCES

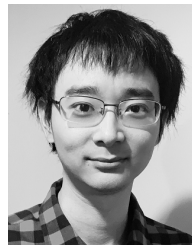
- [1] *Machine Learning and the Future of Mobile App Development*. Accessed: Feb. 13, 2019. [Online]. Available: <https://heartbeat.fritz.ai/machine-learning-and-the-future-of-mobile-app-development-13dd2a5d33>
- [2] *Google Cloud AI*. Accessed: 2021. [Online]. Available: <https://cloud.google.com/products/ai/>

- [3] W. Zheng, R. A. Popa, J. E. Gonzalez, and I. Stoica, "Helen: Maliciously secure cooperative learning for linear models," in *Proc. IEEE Symp. Secur. Privacy (SP)*, May 2019, pp. 724–738.
- [4] P. Mishra, R. Lehmkuhl, A. Srinivasan, W. Zheng, and R. A. Popa, "Delphi: A cryptographic inference system for neural networks," in *Proc. 29th USENIX Secur.*, Nov. 2020, pp. 2505–2522.
- [5] J. Liu, M. Juuti, Y. Lu, and N. Asokan, "Oblivious neural network predictions via MiniONN transformations," in *Proc. ACM SIGSAC Conf. Comput. Commun. Secur.*, Oct. 2017, pp. 619–631.
- [6] M. S. Riaz, M. Samragh, H. Chen, K. Laine, K. Lauter, and F. Koushanfar, "XONN: Xnor-based oblivious deep neural network inference," in *Proc. 28th USENIX Secur.*, 2019, pp. 1501–1518.
- [7] N. Agrawal, A. S. Shamsabadi, M. J. Kusner, and A. Gascón, "QUOTIENT: Two-party secure neural network training and prediction," in *Proc. ACM SIGSAC Conf. Comput. Commun. Secur.*, Nov. 2019, pp. 1231–1247.
- [8] *AI at the Edge: The Next Frontier of the Internet of Things*. Accessed: 2018. [Online]. Available: <https://iotbusinessnews.com/download/white-papers/AVNET-ai-at-the-edge-whitepaper.pdf>
- [9] L. Zhou, M. H. Samavatian, A. Bacha, S. Majumdar, and R. Teodorescu, "Adaptive parallel execution of deep neural networks on heterogeneous edge devices," in *Proc. 4th ACM/IEEE Symp. Edge Comput.*, Nov. 2019, pp. 195–208.
- [10] R. Gilad-Bachrach, N. Dowlin, K. Laine, K. Lauter, M. Naehrig, and J. Wernsing, "CryptoNets: Applying neural networks to encrypted data with high throughput and accuracy," in *Proc. ICML*, 2016, pp. 201–210.
- [11] P. Mohassel and Y. Zhang, "SecureML: A system for scalable privacy-preserving machine learning," in *Proc. IEEE Symp. Secur. Privacy (SP)*, May 2017, pp. 19–38.
- [12] C. Juvekar, V. Vaikuntanathan, and A. Chandrakasan, "Gazelle: A low latency framework for secure neural network inference," in *Proc. 27th USENIX Secur.*, 2018, pp. 1651–1669.
- [13] A. Ibarrondo, H. Chabanne, and M. Önen, "Banners: Binarized neural networks with replicated secret sharing," in *Proc. ACM Workshop Inf. Hiding Multimedia Secur.*, Jun. 2021, pp. 63–74.
- [14] M. S. Riaz, C. Weinert, O. Tkachenko, E. M. Songhori, T. Schneider, and F. Koushanfar, "Chameleon: A hybrid secure computation framework for machine learning applications," in *Proc. AsiaCCS*, 2018, pp. 707–721.
- [15] M. Courbariaux, I. Hubara, D. Soudry, R. El-Yaniv, and Y. Bengio, "Binarized neural networks: Training deep neural networks with weights and activations constrained to +1 or -1," 2016, *arXiv:1602.02830*.
- [16] S. Wagh, S. Tople, F. Benhamouda, E. Kushilevitz, P. Mittal, and T. Rabin, "FALCON: Honest-majority maliciously secure framework for private deep learning," in *Proc. Privacy Enhancing Technol.*, no. 1, 2021, pp. 188–208.
- [17] A. Dalskov, D. Escudero, and M. Keller, "Secure evaluation of quantized neural networks," in *Proc. Privacy Enhancing Technol.*, vol. 4, 2020, pp. 355–375.
- [18] S. Wagh, D. Gupta, and N. Chandran, "SecureNN: 3-party secure computation for neural network training," in *Proc. PETS*, 2019, pp. 26–49.
- [19] O. Goldreich, S. Micali, and A. Wigderson, "How to play any mental game, or a completeness theorem for protocols with honest majority," in *Proc. STOC*, 1987, pp. 307–328.
- [20] M. Atallah, M. Bykova, J. Li, K. Frikken, and M. Topkara, "Private collaborative forecasting and benchmarking," in *Proc. ACM Workshop Privacy Electron. Soc. (WPES)*, 2004.
- [21] M. Blanton, A. Kang, and C. Yuan, "Improved building blocks for secure multi-party computation based on secret sharing with honest majority," in *Proc. ACNS*. Cham, Switzerland: Springer, 2020, pp. 377–397.
- [22] D. Demmler, T. Schneider, and M. Zohner, "ABY—A framework for efficient mixed-protocol secure two-party computation," in *Proc. Netw. Distrib. Syst. Secur. Symp.*, 2015, p. 1.
- [23] P. Mohassel and P. Rindal, "ABY³: A mixed protocol framework for machine learning," in *Proc. ACM CCS*, 2018, pp. 35–52.
- [24] R. Rachuri and A. Suresh, "Trident: Efficient 4PC framework for privacy preserving machine learning," 2020, *arXiv:1912.02631*.
- [25] N. Chandran, D. Gupta, A. Rastogi, R. Sharma, and S. Tripathi, "EzPC: Programmable and efficient secure two-party computation for machine learning," in *Proc. IEEE Eur. Symp. Secur. Privacy (EuroS&P)*, 2019, pp. 496–511.
- [26] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi, "XNOR-Net: ImageNet classification using binary convolutional neural networks," in *Proc. ECCV*, 2016, pp. 525–542.
- [27] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," 2015, *arXiv:1502.03167*.

- [28] G. Asharov, Y. Lindell, T. Schneider, and M. Zohner, "More efficient oblivious transfer and extensions for faster secure computation," in *Proc. ACM SIGSAC Conf. Comput. Commun. Secur.*, 2013, pp. 535–548.
- [29] D. Beaver, "Efficient multiparty protocols using circuit randomization," in *Proc. Crypto*, 1991, pp. 420–432.
- [30] Y. Zheng, H. Duan, and C. Wang, "Towards secure and efficient outsourcing of machine learning classification," in *Proc. ESORICS Cham*, Switzerland: Springer, 2019, pp. 22–40.
- [31] *SnapML for Snapchat Lens Studio*. Accessed: 2021. [Online]. Available: <https://lensstudio.snapchat.com/guides/machine-learning/ml-overview/>
- [32] *Amazon Rekognition*. Accessed: 2021. [Online]. Available: <https://aws.amazon.com/rekognition>
- [33] (2020). *Google DeepMind Health*. [Online]. Available: <https://deepmind.com/blog/announcements/deepmind-health-joins-google-health>
- [34] *Azure IoT Edge*. Accessed: 2021. [Online]. Available: <https://azure.microsoft.com/en-au/services/iot-edge/>
- [35] *LambdaEdge*. Accessed: 2021. [Online]. Available: <https://aws.amazon.com/lambda/edge/>
- [36] D. D. B. Bates *et al.*, "Use of a portable computed tomography scanner for chest imaging of COVID-19 patients in the urgent care at a tertiary cancer center," *Emergency Radiol.*, vol. 27, no. 6, pp. 597–600, Dec. 2020.
- [37] B. Knott, S. Venkataraman, A. Hannun, S. Sengupta, M. Ibrahim, and L. van der Maaten, "CrypTen: Secure multi-party computation meets machine learning," in *Proc. NeurIPS Workshop Privacy-Preserving Mach. Learn.*, 2020, p. 1.
- [38] Cape Privacy. (2020). *TF Encrypted: Encrypted Deep Learning in TensorFlow*. [Online]. Available: <https://tf-encrypted.io/>
- [39] Y. Lindell and B. Pinkas, "Privacy preserving data mining," *J. Cryptol.*, vol. 15, no. 3, pp. 36–54, 2002.
- [40] L. Yu, L. Liu, C. Pu, M. E. Gursory, and S. Truex, "Differentially private model publishing for deep learning," in *Proc. IEEE Symp. Secur. Privacy (SP)*, May 2019, pp. 332–349.
- [41] *Pytorch for DenseNet*. Accessed: 2021. [Online]. Available: <https://github.com/pytorch/vision/blob/master/torchvision/models/densenet.py>
- [42] R. Canetti, "Universally composable security: A new paradigm for cryptographic protocols," in *Proc. 42nd IEEE Symp. Found. Comput. Sci.*, 2001, pp. 136–145.
- [43] A. Patra, T. Schneider, A. Suresh, and H. Yalame, "ABY2.0: Improved mixed-protocol secure two-party computation," in *Proc. USENIX Secur.*, vol. 21, 2020, pp. 2165–2182.
- [44] X. Liu, B. Wu, X. Yuan, and X. Yi, "Leia: A lightweight cryptographic neural network inference system at the edge," *IACR Cryptol. ePrint Arch.*, vol. 2020, p. 463, Apr. 2020.
- [45] X. Wang. (2018). *Flexsc*. [Online]. Available: <https://github.com/wangxiao1254/FlexSC>
- [46] *Breast Cancer*. Accessed: Sep. 25, 2016. [Online]. Available: <https://www.kaggle.com/uciml/breast-cancer-wisconsin-data/>
- [47] *Diabetes*. Accessed: Oct. 7, 2016. [Online]. Available: <https://www.kaggle.com/uciml/pima-indians-diabetes-database>
- [48] *Liver Disease*. Accessed: Sep. 21, 2017. [Online]. Available: <https://www.kaggle.com/uciml/indian-liver-patient-records>
- [49] *Thyroid*. Accessed: Jan. 1, 1987. [Online]. Available: <https://archive.ics.uci.edu/ml/datasets/Thyroid+Disease>
- [50] *COOWOO USB Digital Power Meter Tester*. [Online]. Available: <http://www.coowootech.com/tools.html>
- [51] R. Balestriero and R. G. Baraniuk, "Mad max: Affine spline insights into deep learning," *Proc. IEEE*, vol. 109, no. 5, pp. 704–727, 2020.
- [52] H. Yang, L. Duan, Y. Chen, and H. Li, "BSQ: Exploring bit-level sparsity for mixed-precision neural network quantization," in *Proc. ICLR*, 2021, p. 1.



Xiaoning Liu received the B.E. degree in computer science and technology from Henan University in 2012 and the M.S. degree in computer science from the City University of Hong Kong in 2013. She is currently a Ph.D. candidate with the School of Computing Technologies, RMIT University. Her research pivots on data privacy and security related to machine learning, data mining, cloud computing, and digital health, with the current focus on designing practical secure computation systems for neural networks powered applications.



Bang Wu received the B.S. degree from the Nanjing University of Aeronautics and Astronautics in 2015 and the M.S. degree from The University of Melbourne in 2018, both in electrical engineering. He is currently pursuing the Ph.D. degree with the Faculty of Information Technology, Monash University, Australia. His research interests include security and privacy in graph neural network systems.



Xingliang Yuan (Member, IEEE) is currently a Senior Lecturer (aka U.S. Associate Professor) at the Department of Software Systems and Cybersecurity, Faculty of Information Technology, Monash University, Australia. His research has been supported by the Australian Research Council, CSIRO Data61, and the Oceania Cyber Security Centre. His research interests include data security and privacy, secure networked systems, machine learning security and privacy, and confidential computing. In the past few years, his work has appeared in prestigious venues in cybersecurity, computer networks, and distributed systems, such as ACM CCS, NDSS, IEEE INFOCOM, IEEE TRANSACTIONS ON DEPENDABLE AND SECURE COMPUTING, IEEE TRANSACTIONS ON INFORMATION FORENSICS AND SECURITY, and IEEE TRANSACTIONS ON PARALLEL AND DISTRIBUTED SYSTEMS. He was the recipient of the Dean's Award for Excellence in Research by an Early Career Researcher at Monash Faculty of IT in 2020. He received the Best Paper Award in the European Symposium on Research in Computer Security (ESORICS) 2021, the IEEE Conference on Dependable and Secure Computing (IDSC) 2019, and the IEEE International Conference on Mobility, Sensing and Networking (MSN) 2015.



Xun Yi is currently a Professor with the School of Computing Technologies, RMIT University, Australia. He has published more than 150 research papers in international journals, such as the IEEE TRANSACTIONS ON COMPUTERS, IEEE TRANSACTIONS ON PARALLEL AND DISTRIBUTED SYSTEMS, IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS, IEEE TRANSACTIONS ON DEPENDABLE AND SECURE COMPUTING, IEEE TRANSACTIONS ON INFORMATION FORENSICS AND SECURITY, IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS, IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY, IEEE COMMUNICATIONS LETTERS, *IET Electronic Letters*, and conference proceedings. His research interests include applied cryptography, computer and networks security, mobile and wireless communication security, and privacy-preserving data mining. He has ever undertaken program committee members for more than 30 international conferences. Recently, he has led a few of Australia Research Council (ARC) Discovery Projects. Since 2014, he has been an Associate Editor of the IEEE TRANSACTIONS ON DEPENDABLE AND SECURE COMPUTING.