

METEOR: Improved Secure 3-Party Neural Network Inference with Reducing Online Communication Costs

Ye Dong^{1,2}, Xiaojun Chen^{1,2,✉}, Weizhan Jing^{1,2}, Kaiyun Li^{1,2}, Weiping Wang^{1,2}

¹ Institute of Information Engineering, Chinese Academy of Sciences, Beijing, China

² School of Cyber Security, University of Chinese Academy of Sciences, Beijing, China
{dongye,chenxiaojun,jingweizhan,likaiyun,wangweiping}@iie.ac.cn

Abstract

Secure neural network inference has been a promising solution to private Deep-Learning-as-a-Service, which enables the service provider and user to execute neural network inference without revealing their private inputs. However, the expensive overhead of current schemes is still an obstacle when applied in real applications. In this work, we present METEOR, an online communication-efficient and fast secure 3-party computation neural network inference system against semi-honest adversary in honest-majority. The main contributions of METEOR are two-fold: i) We propose a new and improved 3-party secret sharing scheme stemming from the *linearity* of replicated secret sharing, and design efficient protocols for the basic cryptographic primitives, including linear operations, multiplication, most significant bit extraction, and multiplexer. ii) Furthermore, we build efficient and secure blocks for the widely used neural network operators such as Matrix Multiplication, ReLU, and Maxpool, along with exploiting several specific optimizations for better efficiency. Our total communication with the setup phase is a little larger than SecureNN (PoPETS’19) and FALCON (PoPETS’21), two state-of-the-art solutions, but the gap is not significant when the online phase must be optimized as a priority. Using METEOR, we perform extensive evaluations on various neural networks. Compared to SecureNN and FALCON, we reduce the online communication costs by up to 25.6× and 1.5×, and improve the running-time by at most 9.8× (resp. 8.1×) and 1.5× (resp. 2.1×) in LAN (resp. WAN) for the online inference.

CCS Concepts

• Security and privacy → Privacy-preserving protocols.

Keywords

Privacy, Security, Secret Sharing, Neural Network

1 Introduction

In the Deep-Learning-as-a-Service (DLaaS) paradigm, the service provider offers a trained neural network (NN), and a user calls a well-defined API for data analysis. Aiming to alleviate the privacy concerns associated with DLaaS [1, 3], existing works have introduced secure computation to enable *Secure Inference*. Secure inference exploits cryptographic primitives to ensure that the only information available for the user is the inference result, and nothing more is revealed to either party.

Secure inference protocols can provide high privacy protection, but the key concern is how to obtain privacy with satisfying efficiency. Note that different cryptographic tools offer their characteristics and trade-offs. In particular, fully homomorphic encryption (FHE)-based methods are efficient in communication but still limited by expensive computation burdens [30, 31, 35]. Garbled circuits [75] (GC)-based schemes only require a constant round of interactions but have a high communication overhead and are expensive for arithmetic operations [5, 61]. Secret sharing [66]-based approaches provide efficient arithmetic operations and support non-linear functions [48, 51, 52, 62, 72, 73] using much less communication, yet usually require interactions in proportion to the depth of Multiplication (MULT) gates. Among the secret sharing-based works, 2-out-of-3 replicated secret sharing-based secure 3-party computation (3PC) approaches [51, 73] have achieved significant improvements and gained much attention.

However, the online communication of replicated secret sharing is still the efficiency bottleneck even in the *semi-honest* model. The costly online communication might limit the users’ query throughput, especially in the WAN setting. Therefore, improving the online communication (and running-time) is challenging and a priority in real applications.

The online communication mainly stems from MULT, and we analyze the detailed costs in the following aspects: **i) Costs of Resharing:** Multiplying two ℓ -bit integers (2-MULT) generates 3-out-of-3 secret shared intermediate results. This requires interactive communication of ℓ bits per party for *resharing* 3-out-of-3 shares into 2-out-of-3 shares in 1 round for maintaining correctness and consistency; **ii) Costs of 2-MULT with Faithful Truncation:** When evaluating 2-MULT on two ℓ -bit fixed-point inputs, the parties need to truncate the product to avoid overflow. However, the best known protocol for 2-MULT equipped with faithful truncation needs an online communication of $\frac{4}{3}\ell$ bits per party in 1 round (incur $\frac{1}{3}\ell$ more bits than 2-MULT on integers); **iii) Costs of N -MULT:** N -MULT takes N ℓ -bit integers as inputs and multiplies them to produce the product. N -MULT plays an important role in extracting the most significant bit, but existing works achieve N -MULT by utilizing 2-MULT in a tree-manner. This incurs an online communication of $(N - 1)\ell$ bits in $\lceil \log_2 N \rceil$ rounds.

In this paper, we focus on improving the online communication costs of the latter two kinds of MULT gates. Although our method requires more costs for the setup phase (*i.e.*, communicating $(2^N - 1 - N)\ell$ bits in $\log_2 N$ rounds per party for N -MULT), our significant improvements in the online phase are beneficial in real applications and might be of independent interest. Formally, our techniques and contributions are as follows:

✉ Xiaojun Chen is the corresponding author.

Table 1: Online communication (Comm., in bits) and round complexity of SecureNN, FALCON, and METEOR. q is the smallest prime with $q \geq \ell$, and $k = \lceil \log_2 q \rceil$. For MSB Ext., we set $N = 3, 4$ for N -MULT. †: In SecureNN, MSB Ext. is equal to ShareConvert+Compute MSB (except Select Share), we summarize the corresponding communication and round costs here.

Framework	SecureNN		FALCON		METEOR		
Complexity	Comm.	Round	Comm.	Round	Comm.	Round	
2-MULT	Integers	$\frac{4}{3}\ell$	1	ℓ	1	ℓ	1
	Fixed-Point	$\frac{4}{3}\ell$	1	$\frac{4}{3}\ell$	1	ℓ	1
N -MULT	$(N-1)\frac{4}{3}\ell$	$\lceil \log_2 N \rceil + 1$	$(N-1)\ell$	$\lceil \log_2 N \rceil$	ℓ	1	
MSB Ext.	$\frac{8}{3}\ell k + \frac{14}{3}\ell^\dagger$	7^\dagger	$(2\ell+1)k + \ell$	$\lceil \log_2(\ell+1) \rceil + 3$	$(\frac{4}{3}\ell+1)k$	$\lceil \log_4(\ell+1) \rceil + 2$	
MUX	$\frac{4}{3}\ell$	1	$1 + \ell$	2	ℓ	1	

Our Techniques We propose METEOR, an online communication-efficient and fast secure neural network inference system. METEOR achieves its performance improvements via our improved 3PC protocols and specific optimizations for secure NN operators. Following previous works [72, 73], our 3PC protocols are secure against a *semi-honest* adversary in honest-majority. We build several primitives with a focus on online efficiency by exploiting a *function-dependent* but *input-independent* setup.

Our construction is similar to the sharing semantics of ABY2.0 [56], but exploits a different perspective from the *linearity* of replicated secret sharing [6, 51, 73]: For the MULT, we only need linear operations of replicated secret sharing to generate 2-out-of-3, instead of 3-out-of-3, secret shares in the online phase. This new perspective can accelerate the 2-MULT for fixed-point inputs and N -MULT for integer inputs, and bring several further optimizations for more complex primitives, such as most significant bit extraction (MSB Ext.) and multiplexer (MUX). Besides, our *linearity* perspective is more straightforward to be generalized to any kind of linear secret sharing. Detailed comparison is shown in § 3.1.

Contributions Formally, we have the following contributions:

- **Improved Secure 3-Party Computation:** We propose an improved 3PC secret sharing scheme ($\llbracket \cdot \rrbracket$ -sharing) and construct a set of basic cryptographic primitives, including linear operations (Lops), MULT, MSB Ext., MUX, and etc. Our primitives are more online communication- and round-efficient than that of SecureNN [72] and FALCON [73]. The detailed theoretical improvements are shown in Table 1.
- **Optimized Secure NN Operators:** Furthermore, we construct fast protocols for Matrix Multiplication (MatMul), ReLU function, and Maxpool (MP) based on our basic primitives with specific optimizations. Compared to SecureNN [72] and FALCON [73], we achieve 1.2-6 \times and 1.3-1.6 \times improvements in terms of communication costs. Meanwhile, we are approximately 1.8-20 \times and 1.5 faster for the secure online evaluation of NN operators, respectively.
- **Efficient Secure Inference:** In the end, we perform extensive secure inference experiments on various neural networks and datasets in both LAN and WAN settings: i) For single inference, we reduce the online communication by upto 25.6 \times and 1.5 \times , and improve the online running-time by at most 9.8 \times (resp. 8.1 \times) and 1.5 \times (resp. 2.1 \times) in LAN (resp. WAN) compared to SecureNN and FALCON, respectively. ii) For batch inference, METEOR is more scalable than FALCON. Specially, we improve the communication

Table 2: Notation table.

P_i	party i in 3PC
\mathbf{X}	uppercase bold letter denotes matrix
\mathbf{x}	lowercase bold letter denotes vector
x	lowercase letter denotes scalar
$x[i]$	the i th bit of x
$[\cdot]$	3-out-of-3 sharing
$\langle \cdot \rangle$	2-out-of-3 replicated secret sharing
$\llbracket \cdot \rrbracket$	our 3PC secret sharing
\mathbb{Z}_L	discrete ring modulo L with $L = 2^\ell$
\mathbb{F}_q	field modulo prime q
\mathcal{F}_f	the ideal functionality for $f(\cdot)$
ϵ_R	random sample
κ	the symmetric security parameter

and running-time by both around 1.5 \times in WAN. Our source code is available: <https://github.com/Ye-D/Meteor>.

Organization We present the background and preliminaries in § 2, and give a high-level overview of METEOR in § 3. We propose efficient protocols for the basic primitives in § 4 and justify their security in § 5. And in § 6, we construct the optimized secure NN operators. The experimental results are illustrated in § 7. We discuss related works in § 8 and conclude this work in § 9.

2 Background & Preliminaries

We introduce the background and preliminaries about neural network and 3PC replicated secret sharing in this section.

2.1 Notations

The main notations are summarized in Table 2.

2.2 Neural Network

The computational flow of a neural network is composed of multiple linear and non-linear layers. Each layer receives input and processes it to produce an output that serves as input to the next layer.

Linear Layers Typical linear layers in NN inference include Fully-Connected (FC), Convolution (CONV), and Batch Normalization (BatchNorm, only being linear layer in NN inference):

- **FC:** Given input vector $\mathbf{x} \in \mathbb{R}^{n \times 1}$, a FC layer generates the output $\mathbf{y} \in \mathbb{R}^{m \times 1}$ as $\mathbf{y} = \mathbf{W}\mathbf{x} + \mathbf{b}$, where $\mathbf{W} \in \mathbb{R}^{m \times n}$ is the weight matrix and $\mathbf{b} \in \mathbb{R}^{m \times 1}$ is the bias term. More generally, neural networks often take a batch of images as inputs $\mathbf{X}^{n \times |B|}$ ($|B|$ is the batchsize), thus the FC layer can be computed with matrix multiplication as $\mathbf{Y} = \mathbf{W}\mathbf{X} + \mathbf{B}$.

- **CONV:** The CONV layer computes the dot product of a small weight matrix (*filter*) and the neighborhood of an element of the input. The process is sliding each filter with a certain *stride*, and the size of filter is called *filter size*. For a generalized exposition on CONV, please refer to [72].
- **BatchNorm:** A BatchNorm layer is typically applied to shift its input x to amenable ranges. During the inference, the BatchNorm parameters γ and β are fixed, BatchNorm normalizes x as $\gamma \cdot x + \beta$.

Non-Linear Layers NN uses activation functions to model non-linear relationships between input and output. And Pool functions sometimes are applied.

- **Activation:** The activation functions are applied in element-wise. One of the most popular activation functions is ReLU function: $\text{ReLU}(x) = \max(0, x)$. Other activation functions include Sigmoid, Tanh, and etc [54];
- **Pool:** Pooling arranges inputs into several windows and aggregates elements of each window. Maxpool (resp. Avgpool) calculates the maximum (resp. average) for each window.

2.3 3PC Replicated Secret Sharing

Secret value $x \in \mathbb{Z}_{2^\ell}$ is shared by three random values $x_0, x_1, x_2 \in \mathbb{Z}_{2^\ell}$ with $x = x_0 + x_1 + x_2 \pmod{2^\ell}$ [6, 51, 73]. In 3-out-of-3 sharing ($[\cdot]$ -sharing), P_i has $[x]_i = x_i$. In replicated secret sharing (2-out-of-3, $\langle \cdot \rangle$ -sharing), P_i gets $\langle x \rangle_i = (x_i, x_{i+1})$. Without special declaration, we compute in \mathbb{Z}_{2^ℓ} and omit $\pmod{2^\ell}$ for brevity.

Sharing and Reconstruction To achieve functionality $\mathcal{F}_{\text{SHARE}}^{\langle \cdot \rangle}$, secret owner samples random $x_1, x_2 \in_R \mathbb{Z}_{2^\ell}$, sets $x_0 = x - x_1 - x_2$, and sends $\langle x \rangle_i = (x_i, x_{i+1})$ to P_i . And to implement $\mathcal{F}_{\text{REC}}^{\langle \cdot \rangle}$, P_i sends x_{i+1} to P_{i-1} such that P_{i-1} reconstructs $x = x_0 + x_1 + x_2$ for $i \in \{0, 1, 2\}$. **Linear Operations** Let (c_1, c_2, c_3) be public constants, and $(\langle x \rangle, \langle y \rangle)$ be two secret-shared values. Then, $\langle c_1x + c_2y + c_3 \rangle$ can be computed as $(c_1x_0 + c_2y_0 + c_3, c_1x_1 + c_2y_1, c_1x_2 + c_2y_2)$ where P_i can compute its share locally. When $(c_1 = 1, c_2 = 1, c_3 = 0)$, we get $\langle x + y \rangle$.

Multiplication Functionality $\mathcal{F}_{\text{MULT}}^{\langle \cdot \rangle}$ multiplies two shared values $\langle x \rangle$ and $\langle y \rangle$, existing protocol achieves this as follows: i) First, P_i computes $z_i = x_i y_i + x_{i+1} y_i + x_i y_{i+1}$ locally such that z_i is $[\cdot]$ -shared. ii) Parties then perform *re-sharing* by letting P_i sends $z'_i = \alpha_i + z_i$ to P_{i-1} , where $\alpha_0 + \alpha_1 + \alpha_2 = 0$ (P_i can generate α_i in the setup phase as [6, 51, 73]). iii) Finally, $\{(z'_0, z'_1), (z'_1, z'_2), (z'_2, z'_0)\}$ form $\langle x \cdot y \rangle$.

In the case of $\ell > 1$ (e.g., $\ell = 64$) which support arithmetic operations (e.g., $+$, $-$, and \cdot), we refer to this type as *Arithmetic Sharing* and use notation $\langle \cdot \rangle$. *Boolean Sharing* ($\langle \cdot \rangle^2$) refers to $\ell = 1$ where $+$, $-$ and \cdot are respectively replaced by bit-wise \oplus and \wedge .

MSB Extraction The key step of comparing $\langle x \rangle \geq \langle y \rangle$ in two's complement representation is extracting the most significant bit of $\langle z \rangle = \langle x \rangle - \langle y \rangle$. General methods either re-interpret the arithmetic sharing as boolean sharing and evaluate an addition circuit on boolean shares to compute $\langle \text{msb}(z) \rangle^2$, or employ garbled circuits to extract the most significant bit. Recently, Wagh *et al.* proposed an efficient MSB Ext. method based on wrap function and bit decomposition in replicated secret sharing [73]. We follow their approach in METEOR but optimize the online efficiency with our improved secret sharing scheme. And we plan to improve other MSB Ext. methods [29, 49, 51] with our novel secret sharing for future work.

Fixed-Point Representation In secure NN inference, we need to encode floating-point numbers as integers in rings [51, 52, 73]. Given floating-point $x \in \mathbb{R}$, its encoding is as: $x = \lfloor 2^d \cdot x \rfloor \pmod{2^\ell}$, where it is usually $\ell = 64$ and $d = 13$ as [73]. In this way, we use $[0, 2^{\ell-1})$ to represent $x \in \mathbb{R}^+$, and $[2^{\ell-1}, 2^\ell)$ for negative values.

3 A High-Level Overview of METEOR

We first present an overview of our $[\cdot]$ -sharing semantics in § 3.1. Then, we show the design and threat model of METEOR in § 3.2.

3.1 Overview of $[\cdot]$ -Sharing Semantics

Costs Analysis of MULT Existing $\langle \cdot \rangle$ -sharing based 3PC approaches need ℓ bits per party in 1 round for evaluating 2-MULT with 2 integer inputs (*resharing*), and require more costs for 2-MULT with fixed-point inputs and N -MULT with integer inputs: i) When multiplying 2 fixed-point inputs x and y , the parties need to reveal $[z + r']$ and compute $(z + r')/2^d - \langle r \rangle$ for faithful truncation, where $z = xy$, $r = r'/2^d$, and (r, r') are in secret. As $[z + r']$ is of $[\cdot]$ -sharing, it needs $\frac{4}{3}\ell$ bits communication per party for mask-and-reveal in the online phase [51, 73]. ii) When multiplying N integers, parties need to perform *resharing* for each 2-MULT. This requires an online communication of $(N - 1)\ell$ bits per party in $\lceil \log_2 N \rceil$ rounds [73].

The main costs of MULT stem from **Resharing**-related operations. Inspired by this conclusion, we are wondering: *Will the efficiency (e.g., communication and running-time) be improved if we can maintain the $\langle \cdot \rangle$ -sharing format during the whole computation?*

Linearity of $\langle \cdot \rangle$ -Sharing From § 2.3, we notice the linear operations of $\langle \cdot \rangle$ -values lead to $\langle \cdot \rangle$ -shared results locally (no communication). This is true for two $\langle \cdot \rangle$ -shared inputs $(c_1 \langle x \rangle + c_2 \langle y \rangle + c_3)$, and can be easily generalized to three or more $\langle \cdot \rangle$ -shared inputs.

$[\cdot]$ -Sharing Inspired by the sharing semantics of [9, 10, 69] and with the *linearity* of $\langle \cdot \rangle$ -sharing in mind, we propose an improved 3PC secret sharing ($[\cdot]$ -sharing) as follows:

DEFINITION 1 ($[\cdot]$ -SHARING). *A value $x \in \mathbb{Z}_{2^\ell}$ is said to be $[\cdot]$ -shared among $\{P_0, P_1, P_2\}$ if there exists random ψ_x and m_x such that: i) ψ_x is $\langle \cdot \rangle$ -shared among $\{P_0, P_1, P_2\}$; ii) $m_x = x - \psi_x$ is known to all parties in clear. The share of P_i is $[\mathbf{x}]_i = (m_x, \langle \psi_x \rangle_i)$ for $i \in \{0, 1, 2\}$.*

For brevity, we use notations $\psi_{x_1 \dots x_n} = \psi_{x_1} \psi_{x_2} \dots \psi_{x_n}$ and $m_{x_1 \dots x_n} = m_{x_1} m_{x_2} \dots m_{x_n}$. Similarly, $[\cdot]^q$ -sharing is for $x \in \mathbb{F}_q$ and $[\cdot]^2$ -sharing is for $x \in \mathbb{Z}_2$, where we use modulo q in $[\cdot]^q$ -sharing and replace $+$, $-$ by \oplus and \cdot by \wedge in $[\cdot]^2$ -sharing.

With $[\cdot]$ -sharing, we can evaluate MULT by computing the relatively expensive multiplication of secret random $\langle \psi \rangle$ s in the setup phase, such that *the online phase only involves the linear operations of $\langle \cdot \rangle$ -sharing*. Taking 2-MULT(x, y) with integer inputs as an example, the parties compute $\langle \psi_z \rangle = \langle \psi_x \rangle \langle \psi_y \rangle$ in the setup phase, and compute m_z with linear operations and 1 round of revealing. $[\cdot]$ -sharing also needs ℓ bits per party, but gives the following benefits:

- For 2-MULT with fixed-point inputs, we only need ℓ bits per party in 1 round for mask-and-reveal in the online phase. This is because the intermediate results are in $\langle \cdot \rangle$ -shared fashion. Hence, we improve the communication by 1.3 \times ;
- For N -MULT with integer inputs, we only need linear operations of $\langle \cdot \rangle$ -sharing to generate $\langle \cdot \rangle$ -shared product of N integers in the online phase, since all multiplications

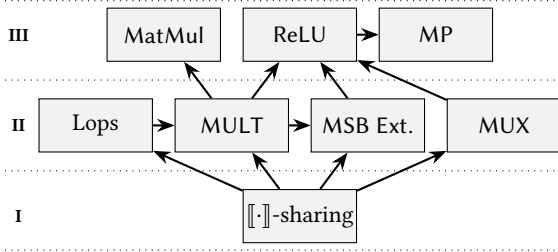


Figure 1: Dependency of protocols in METEOR.

among $\langle\psi\rangle$ s can be evaluated in the setup phase. Therefore, our approach needs an online communication of ℓ bits per party in 1 round, which is independent of N and first achieved in the regime of 3PC. Compared to prior methods, we improve the online communication by $N\times$ and rounds by $\lceil\log_2 N\rceil\times$.

What’s more, we propose efficient $\llbracket\cdot\rrbracket$ -sharing based protocols for other primitives and NN operators in respective § 4 and § 6.

Comparison to ABY2.0 [56]. Patra *et al.* has proposed similar sharing semantics to improve the online efficiency of 2PC [56] inspired by ASTRA [18] and [55], but our $\llbracket\cdot\rrbracket$ -sharing is different in the following aspects:

- **Beaver-Friendly v.s. Linearity:** ABY2.0 is inspired by Beaver triples [7] and reduces the communication by sharing the inputs in a Beaver-friendly format. However, $\llbracket\cdot\rrbracket$ -sharing stems from the *linearity* of replicated secret sharing. They might be equivalent in some settings (*e.g.*, 2PC), but our *linearity* perspective is more straightforward to be generalized to other linear secret sharing.
- **2PC v.s. 3PC:** For the setup phase, ABY2.0 exploits Oblivious Transfer (OT) [33] or HE [26] to generate correlated randomness, but we utilize the multiplication protocol of $\langle\cdot\rangle$ -sharing (free of OT or HE). Therefore, $\llbracket\cdot\rrbracket$ -sharing is more efficient in setup when *honest-majority in 3PC* is available.

3.2 Design of METEOR

Our METEOR, as depicted in Figure 1, consists of three layers:

- **I:** We first propose an improved 3-party secret sharing scheme ($\llbracket\cdot\rrbracket$ -sharing) inspired by the *linearity* of replicated secret sharing in 3PC.
- **II:** Secondly, we design efficient protocols for the most common basic cryptographic primitives, *i.e.*, Linear Operations (Lops), MULT, MSB Ext., and MUX.
- **III:** Thirdly, we build secure blocks for the widely used NN operators, such as MatMul, ReLU, and MP, together with specific optimizations to support fast secure NN inference.

Function-dependent but Input-independent Setup Following [72, 73], we also focus on the online efficiency. METEOR is cast into a *function-dependent* but *input-independent* setup phase, and an *input-dependent* online phase as [19, 56, 57]. In the setup phase, we generate the *function-dependent* but *input-independent* correlated randomness for a given function to improve the online efficiency. This setup is available and widely utilized in many applications.

Threat Model Following works [72, 73], METEOR resists semi-honest adversaries in honest-majority [47]. Namely, each party

follows the protocol, but may *individually* try to learn information about other inputs:

DEFINITION 2 (SEMI-HONEST SECURITY). Let Π be a three-party protocol running in real-world and $\mathcal{F} : (\{0, 1\}^n)^3 \rightarrow (\{0, 1\}^m)^3$ be the ideal randomized functionality. We say Π securely computes \mathcal{F} in presence of a single semi-honest adversary if for every corrupted party P_i ($i \in \{0, 1, 2\}$) and every input $\mathbf{x} \in (\{0, 1\}^n)^3$, there exists an efficient simulator \mathcal{S} such that:

$$\{\text{view}_{i,\Pi}(\mathbf{x}), \text{output}_{\Pi}(\mathbf{x})\} \stackrel{c}{\approx} \{\mathcal{S}(I, x_i, \mathcal{F}_i(\mathbf{x})), \mathcal{F}(\mathbf{x})\},$$

where $\text{view}_{i,\Pi}(\mathbf{x})$ is the view of P_i in the execution of Π on \mathbf{x} , $\text{output}_{\Pi}(\mathbf{x})$ is the output of all parties, and $\mathcal{F}_i(\mathbf{x})$ denotes the i th output of $\mathcal{F}(\mathbf{x})$.

4 Improved Secure 3-Party Computation

In this section, we present the detailed constructions of sharing and reconstruction (§ 4.1), linear operations (§ 4.2), MULT (§ 4.3), MSB Ext. (§ 4.4), and MUX (§ 4.5).

4.1 Sharing and Reconstruction

Sharing $\Pi_{\text{SHARE}}(\mathbf{x})$ achieves $\mathcal{F}_{\text{SHARE}}^{\llbracket\cdot\rrbracket}$ by enabling P_i (secret owner) to generate a $\llbracket\cdot\rrbracket$ -sharing of its x . In the setup phase, all parties together sample random $\langle\psi_x\rangle$ using existing $\mathcal{F}_{\text{RAND}}^{\langle\cdot\rangle}$ with P_i gets ψ_x in clear (*c.f.* Appendix A). In the online phase, P_i reveals $m_x = x - \psi_x$.

Reconstruction We describe our protocol $\Pi_{\text{REC}}(\llbracket\mathbf{x}\rrbracket)$ for $\mathcal{F}_{\text{REC}}^{\llbracket\cdot\rrbracket}$ that reconstructs \mathbf{x} as follows: Given $\llbracket\mathbf{x}\rrbracket$, parties invoke $\mathcal{F}_{\text{REC}}^{\langle\cdot\rangle}$ to reconstruct ψ_x and locally compute $\mathbf{x} = m_x + \psi_x$.

4.2 Linear Operations

$\llbracket\cdot\rrbracket$ -sharing is linear in the sense that given $\llbracket\mathbf{x}\rrbracket$, $\llbracket\mathbf{y}\rrbracket$ and public constants c_1, c_2 , and c_3 , parties can compute $\llbracket\mathbf{z}\rrbracket = c_1 \cdot \llbracket\mathbf{x}\rrbracket + c_2 \cdot \llbracket\mathbf{y}\rrbracket + c_3$ by locally setting $(m_z = c_1 \cdot m_x + c_2 \cdot m_y + c_3, \langle\psi_z\rangle = c_1 \cdot \langle\psi_x\rangle + c_2 \cdot \langle\psi_y\rangle)$.

4.3 Multiplication

We consider 2-input multiplication (2-MULT) and N -input multiplication (N -MULT). The former is employed in secure FC/CONV, while the latter plays an important role in secure MSB Ext.

4.3.1 2-Input Multiplication We first consider the multiplication of two integers. Given the $\llbracket\cdot\rrbracket$ -shares of integers x and y , functionality $\mathcal{F}_{2\text{-MULT}}^{\llbracket\cdot\rrbracket}$ generates $\llbracket\mathbf{z}\rrbracket$ with $z = xy$. For z , we will need:

$$\begin{aligned} m_z &= z - \psi_z = xy - \psi_z \\ &= (m_x + \psi_x)(m_y + \psi_y) - \psi_z \\ &= m_x m_y + m_x \psi_y + m_y \psi_x + \psi_{xy} - \psi_z \end{aligned} \quad . \quad (1)$$

In the setup phase, parties compute the input-independent $\langle\psi_{xy}\rangle = \langle\psi_x\rangle\langle\psi_y\rangle$. And in the online phase, parties compute $\langle m_z \rangle$ locally and collaboratively reveal it. So the challenge is reduced to generate $\langle\psi_{xy}\rangle$ given $\langle\psi_x\rangle$ and $\langle\psi_y\rangle$. We leverage $\mathcal{F}_{\text{MULT}}^{\langle\cdot\rangle}$ to accomplish this task as § 2.3. The protocol is in Figure 2. $\Pi_{2\text{-MULT}}$ needs an online communication of ℓ bits per party in 1 round.

Fixed-Point Multiplication Extension As analyzed in § 3, we truncate the product (*i.e.*, $xy/2^d$ where x and y are in fixed-point) after each multiplication in secure NN inference. Existing faithful truncation method [51, 73] needs $\frac{4}{3}\ell$ bits per party in online phase.

Protocol $\Pi_{2\text{-MULT}}$

- Input:** $P_0, P_1,$ and P_2 hold $[\cdot]$ -shared $[\mathbf{x}]$ and $[\mathbf{y}]$.
Output: $[\mathbf{z}] = [\mathbf{x}] \cdot [\mathbf{y}]$.
- **Setup:**
 - (1) Parties generate $\langle \psi_z \rangle$ using functionality $\mathcal{F}_{\text{RAND}}^{(\cdot)}$.
 - (2) Parties execute $\langle \psi_{xy} \rangle = \mathcal{F}_{\text{MULT}}^{(\cdot)}(\langle \psi_x \rangle, \langle \psi_y \rangle)$.
 - **Online:**
 - (1) P_i locally computes $\langle m_z \rangle_i = m_{xy} + m_x \langle \psi_y \rangle_i + m_y \langle \psi_x \rangle_i + \langle \psi_{xy} \rangle_i - \langle \psi_z \rangle_i$.
 - (2) Parties exchange the shares of $\langle m_z \rangle$ to reconstruct m_z .
 - (3) P_i outputs $[\mathbf{z}]_i = (m_z, \langle \psi_z \rangle_i)$.

Figure 2: 2-Input Multiplication Protocol for Integers.

To reduce the costs, we propose online free faithful truncation at the same online communication as 2-MULT for integers: i) In the setup phase, parties generate $(\langle \psi_z \rangle, \langle \psi'_z \rangle)$ with $\psi_z = \psi'_z / 2^d$ using the optimized binary circuits [51]. ii) In the online phase, parties compute and reveal $\langle m'_z \rangle = xy - \langle \psi'_z \rangle$ as equation (1) and reveal m'_z . iii) Parties set $[\mathbf{z}] = (m'_z / 2^d, \langle \psi'_z \rangle)$ with $m'_z / 2^d + \psi_z = xy / 2^d$ holds. **Online Communication** The correctness and precision guarantees of our method are similar to prior works [51, 73], our main contributions here lie in the online communication improvements. As $\langle m'_z \rangle$ is in $\langle \cdot \rangle$ -sharing, our $\Pi_{2\text{-MULT}}$ with truncation for fixed-point inputs needs ℓ bits per party in 1 round for revealing it during the online phase, achieving 1.3 \times improvements.

4.3.2 N -Input Multiplication Functionality $\mathcal{F}_{N\text{-MULT}}^{[\cdot]}$ multiplies N integers for any positive constant N . From the fact that secret random $\langle \psi \rangle$ s are input-independent, we can multiply them in the setup phase. Therefore, we only need linear operations of $\langle \cdot \rangle$ -shared values to get the $\langle \cdot \rangle$ -shared product of N integers in the online phase. Taking (x_1, x_2, \dots, x_N) as inputs, we have $m_z = \prod_{t=1}^N x_t - \psi_z = \prod_{t=1}^N (m_{x_t} + \psi_{x_t}) - \psi_z = \sum_{\mathcal{T} \subseteq \{1, \dots, N\}} (\prod_{j \notin \mathcal{T}} m_{x_j} \cdot \prod_{k \in \mathcal{T}} \psi_{x_k}) - \psi_z$.

In the setup phase, parties can compute the input-independent $\langle \cdot \rangle$ -shares of $\{\prod_{k \in \mathcal{T}} \psi_{x_k}\}_{\mathcal{T} \subseteq \{1, \dots, N\}}$ exploiting $\mathcal{F}_{\text{MULT}}^{(\cdot)}$. In the online phase, the parties only need to reveal m_z . The details are in Figure 5. **Online Communication** The online communication remains just ℓ bits per party in 1 round independent of the fan-in. In contrast, previous $\langle \cdot \rangle$ -sharing based methods require $(N - 1)\ell$ bits per party in $\lceil \log_2 N \rceil$ rounds. In the setup phase, the above method requires $(2^N - 1 - N)\ell$ bits per party in $\lceil \log_2 N \rceil$ rounds. To balance the burden in the setup and online phases, we set $N = 3$ and 4 as [56].

4.4 Secure MSB Extraction

Given $[\mathbf{x}]$, functionality $\mathcal{F}_{\text{SecMSB}}^{[\cdot]}$ extracts $[\text{msb}(\mathbf{x})]^2$ securely. From $x = m_x + \psi_x$ for $[\mathbf{x}]$, we can write

$$\text{msb}(\mathbf{x}) = \text{msb}(m_x) \oplus \text{msb}(\psi_x) \oplus c, \quad (2)$$

where c is the carry bit of m_x and ψ_x modulo $\frac{1}{2}$ (ignoring their msb), which is formalized as $c = (2m_x + 2\psi_x \geq L) = (2\psi_x \geq L - 2m_x)$. Let $s = 2\psi_x$ and $b = L - 2m_x$, our key insights are as follows:

- i) ψ_x is independent of inputs, we can compute $[\text{msb}(\psi_x)]^2$ and $[\cdot]^q$ -shares of bits of s in the setup phase.
- ii) m_x and L are public in the online phase.

Protocol Π_{SecMSB}

- Input:** $P_0, P_1,$ and P_2 hold $[\mathbf{x}] = (m_x, \langle \psi_x \rangle)$.
Output: $[\text{msb}(\mathbf{x})]^2$.
- **Setup:**
 - (1) The parties call functionality $\mathcal{F}_{\text{PreMSB}}^{[\cdot]}$ to generate $[\text{msb}(\psi_x)]^2, \{[s[i]]^q\}_{i=1}^\ell$ with $s = 2\psi_x, ([\lambda]^2, [\lambda]^q)$ with $\lambda \in_R \mathbb{Z}_2$, and $[\zeta]^q$ with $\zeta \in_R \mathbb{F}_q^*$.
 - (2) The parties run the setup phase of $\mathcal{F}_{N\text{-MULT}}^{[\cdot]}$.
 - **Online:**
 - (1) Compute $b = L - 2m_x$.
 - (2) **for** $i \in [\ell, \ell - 1, \dots, 1]$, all parties compute in \mathbb{F}_q :
 - (3) $[u[i]]^q = (1 - 2[\lambda]^q) \cdot ([s[i]]^q - b[i])$.
 - (4) $[w[i]]^q = [s[i]]^q + b[i] - 2[s[i]]^q b[i]$.
 - (5) $[e[i]]^q = [u[i]]^q + 1 + \sum_{k=i+1}^\ell [w[k]]^q$.
 - (6) **end for**
 - (7) Compute $[d]^q = [\zeta]^q \cdot \prod_{i=1}^\ell [e[i]]^q \pmod{q}$ using $\mathcal{F}_{N\text{-MULT}}^{[\cdot]}$, and reveal d .
 - (8) Set $[c]^2 = \lambda' \oplus [\lambda]^2$, where $\lambda' = (d \neq 0)$.
 - (9) Output $[\text{msb}(\mathbf{x})]^2 = \text{msb}(m_x) \oplus [\text{msb}(\psi_x)]^2 \oplus [c]^2$.

Figure 3: Secure Most Significant Bit Extraction Protocol.
Protocol Π_{MUX}

- Input:** $P_0, P_1,$ and P_2 hold $[\cdot]$ -shared $[\mathbf{x}], [\mathbf{y}],$ and $[\mathbf{v}]^2$.
Output: $[\mathbf{z}]$ with $z = x$ if $v = 1$ and $z = y$ otherwise.
- **Setup:**
 - (1) Parties generate $\langle \psi \rangle$ using functionality $\mathcal{F}_{\text{RAND}}^{(\cdot)}$.
 - (2) Parties convert $\langle \psi_v^\ell \rangle = \mathcal{F}_{\text{Bit2A}}^{(\cdot)}(\langle \psi_v \rangle^2)$ for $[\mathbf{v}]^2$.
 - (3) Parties compute $[u] = [\mathbf{x}] - [\mathbf{y}]$.
 - (4) Parties invoke $\mathcal{F}_{\text{MULT}}^{(\cdot)}$ to compute $\langle \psi_u \psi_v^\ell \rangle = \langle \psi_u \rangle \cdot \langle \psi_v^\ell \rangle$.
 - **Online:**
 - (1) P_i locally computes $\langle m_{uv^\ell} \rangle_i = (1 - 2\langle \psi_v^\ell \rangle_i) m_u m_v^\ell + m_u \langle \psi_v^\ell \rangle_i + \langle \psi_u \rangle_i m_v^\ell + (1 - 2m_v^\ell) \langle \psi_u \psi_v^\ell \rangle_i - \langle \psi \rangle_i$.
 - (2) Parties reconstruct m_{uv^ℓ} and set $[u \cdot v^\ell] = (m_{uv^\ell}, \langle \psi \rangle)$.
 - (3) Parties locally compute and output $[\mathbf{z}] = [u \cdot v^\ell] + [\mathbf{y}]$.

Figure 4: Secure Multiplexer Protocol.

With our key insights in mind, we propose protocol Π_{SecMSB} as Figure 3. In the setup phase, we manage to generate $[\text{msb}(\psi_x)]^2, \{[s[i]]^q\}_{i=1}^\ell, [\lambda]^2, [\lambda]^q,$ and $[\zeta]^q$ using $\mathcal{F}_{\text{PreMSB}}^{[\cdot]}$. We construct protocol Π_{PreMSB} (c.f., Appendix § C) for $\mathcal{F}_{\text{PreMSB}}^{[\cdot]}$ based on [51, 73].

In the online phase, the challenge is computing c . Inspired by [73], we propose an optimized method as Figure 3. The key point is that $s < b$ if and only if there $\exists i \in [\ell, \dots, 1]$ subjected to $e[i] = 0$, which means $s[k] = b[k]$ for $\forall k > i$ and $(s[i] = 0, b[i] = 1)$. Otherwise, we have $s \geq b \Leftrightarrow e[i] \neq 0$ for $\forall i \in [\ell, \dots, 1]$.

Online Communication We use $\Pi_{N\text{-MULT}}$ with $N = 3, 4$. For the online phase: i) Steps 2-6 need $\ell \lceil \log_2 q \rceil$ bits per party in 1 round. ii) Step 7 needs $\approx (\frac{1}{3}\ell + 1) \lceil \log_2 q \rceil$ bits per party in $\lceil \log_4(\ell + 1) \rceil + 1$ rounds. iii) Steps 1, 8, 9 are locally. Therefore, the online phase needs $\approx (\frac{4}{3}\ell + 1) \lceil \log_2 q \rceil$ bits per party in $\lceil \log_4(\ell + 1) \rceil + 2$ rounds.

4.5 Multiplexer

Given $([\mathbf{x}], [\mathbf{y}], [\mathbf{v}]^2)$, functionality $\mathcal{F}_{\text{MUX}}^{[\cdot]}$ outputs $[\mathbf{z}] = [\mathbf{x}]$ if $v = 1$, and $[\mathbf{z}] = [\mathbf{y}]$ otherwise. This is $[\mathbf{z}] = ([\mathbf{x}] - [\mathbf{y}]) \cdot [\mathbf{v}]^2 + [\mathbf{y}]$.

Let $\llbracket \mathbf{u} \rrbracket = \llbracket \mathbf{x} \rrbracket - \llbracket \mathbf{y} \rrbracket$, the challenge is computing $\llbracket \mathbf{u} \rrbracket \cdot \llbracket \mathbf{v} \rrbracket^2$. A trivial solution is converting $\llbracket \mathbf{v} \rrbracket^2$ to $\llbracket \mathbf{v} \rrbracket$ by [29, 63] and computing $\llbracket \mathbf{u} \rrbracket \cdot \llbracket \mathbf{v} \rrbracket$ with $1 + \ell$ bits per party in 2 rounds. To reduce costs, we propose protocol Π_{MUX} . Denote the value of bit v in \mathbb{Z}_{2^ℓ} as v^ℓ . For $\llbracket \mathbf{v} \rrbracket^2 = (m_v, \langle \psi_v \rangle^2)$, $v^\ell = (m_v \oplus \psi_v)^\ell = m_v^\ell + \psi_v^\ell - 2m_v^\ell \psi_v^\ell$. Thus, for $\llbracket \mathbf{u} \rrbracket \cdot \llbracket \mathbf{v} \rrbracket^2$ we have

$$\begin{aligned} \mathbf{u} \cdot \mathbf{v}^\ell &= (m_u + \psi_u) \cdot (m_v^\ell + \psi_v^\ell - 2m_v^\ell \psi_v^\ell) \\ &= m_u m_v^\ell + m_u \psi_v^\ell - 2m_u m_v^\ell \psi_v^\ell + \psi_u m_v^\ell + \psi_u \psi_v^\ell - 2\psi_u m_v^\ell \psi_v^\ell \\ &= (1 - 2\psi_v^\ell) m_u m_v^\ell + m_u \psi_v^\ell + \psi_u m_v^\ell + (1 - 2m_v^\ell) \psi_u \psi_v^\ell \end{aligned} \quad (3)$$

In the setup phase, we compute $\langle \psi_v^\ell \rangle$ from $\langle \psi_v \rangle^2$ via $\mathcal{F}_{\text{Bit2A}}^{\langle \cdot \rangle}$ (c.f., Appendix § D) and compute $\langle \psi_u \psi_v^\ell \rangle$ using $\mathcal{F}_{\text{MULT}}^{\langle \cdot \rangle}$. In online phase, parties compute and reveal m_{uv^ℓ} , set $\llbracket \mathbf{u} \cdot \mathbf{v}^\ell \rrbracket = (m_{uv^\ell}, \langle \psi \rangle)$, and output $\llbracket \mathbf{z} \rrbracket = \llbracket \mathbf{u} \cdot \mathbf{v}^\ell \rrbracket + \llbracket \mathbf{y} \rrbracket$ as shown in Figure 4.

Online Communication Π_{MUX} needs ℓ bits per party in 1 round.

5 Security Analysis

Theorem 1 captures the security of our protocols, and the full proof is given in Appendix G.

THEOREM 1. *In the hybrid model, our protocols securely realize the functionalities $\mathcal{F}_{\text{SHARE}}^{\llbracket \cdot \rrbracket}$, $\mathcal{F}_{\text{REC}}^{\llbracket \cdot \rrbracket}$, $\mathcal{F}_{\text{Lops}}^{\llbracket \cdot \rrbracket}$, $\mathcal{F}_{2\text{-MULT}}^{\llbracket \cdot \rrbracket}$, $\mathcal{F}_{N\text{-MULT}}^{\llbracket \cdot \rrbracket}$, $\mathcal{F}_{\text{SecMSB}}^{\llbracket \cdot \rrbracket}$ and $\mathcal{F}_{\text{MUX}}^{\llbracket \cdot \rrbracket}$ against a semi-honest adversary \mathcal{A} , who corrupts no more than one party.*

6 Optimized Secure NN Operators

In § 6.1, we show the secure evaluation of FC, CONV, BatchNorm. In § 6.2, we construct secure ReLU. We give private MP and ReLU-MP equivalent switching in § 6.3. The full protocols and high-level buildings of other secure NN operators are shown in Appendix E.

6.1 Secure Matrix Multiplication

Protocol $\Pi_{2\text{-MULT}}$ can be easily vectorized to MatMul. Given $\llbracket \mathbf{X} \rrbracket = (m_X, \langle \psi_X \rangle)$ with dimension $m \times n$ and $\llbracket \mathbf{Y} \rrbracket = (m_Y, \langle \psi_Y \rangle)$ with dimension $n \times o$: i) In the setup phase, parties execute $\mathcal{F}_{\text{MULT}}^{\langle \cdot \rangle}$ to compute $\langle \psi_{XY} \rangle = \langle \psi_X \rangle \cdot \langle \psi_Y \rangle$. ii) In the online phase, parties locally compute $\langle m_Z \rangle = m_X \cdot m_Y + m_X \cdot \langle \psi_Y \rangle + \langle \psi_X \rangle \cdot m_Y + \langle \psi_{XY} \rangle - \langle \psi_Z \rangle$, reconstruct m_Z , and set $\llbracket \mathbf{Z} \rrbracket = (m_Z, \langle \psi_Z \rangle)$. We need mol bits (independent of n) in 1 round for the online phase. The full protocol is shown in Figure 8. The security of MatMul follows in $\mathcal{F}_{2\text{-MULT}}^{\llbracket \cdot \rrbracket}$ -hybrid model.

Secure FC & CONV We can leverage MatMul to achieve secure FC. For CONV, we reshape the input and filter to express convolution as MatMul for subsequent secure evaluation. To support fixed-point truncation, the parties can generate $(\langle \psi_Z \rangle, \langle \psi'_Z \rangle)$, and perform faithful truncation similarly as § 4.3.1 but in vectorization.

Fusing CONV & BatchNorm BatchNorm often goes after CONV, we can fuse them into one for better efficiency [4]. Suppose the trained parameters for BatchNorm and CONV are $(\mathbf{W}_{\text{BN}}, \mathbf{b}_{\text{BN}})$ and $(\mathbf{W}_{\text{CONV}}, \mathbf{b}_{\text{CONV}})$, we can replace them by a single CONV with $(\mathbf{W} = \mathbf{W}_{\text{BN}} \cdot \mathbf{W}_{\text{CONV}}, \mathbf{b} = \mathbf{W}_{\text{BN}} \cdot \mathbf{b}_{\text{CONV}} + \mathbf{b}_{\text{BN}})$. Therefore, we can compute both layers together at the same costs as secure CONV.

6.2 Secure ReLU

The activation function considered in this work is the rectified linear unit (ReLU). Taking x as input, $\mathcal{F}_{\text{ReLU}}^{\llbracket \cdot \rrbracket}$ returns x if $x \geq 0$, and 0 otherwise. To achieve $\mathcal{F}_{\text{ReLU}}^{\llbracket \cdot \rrbracket}$ securely, it suffices to first extract $\llbracket \text{msb}(x) \rrbracket^2$ using $\mathcal{F}_{\text{SecMSB}}^{\llbracket \cdot \rrbracket}$, and then execute $\mathcal{F}_{\text{MUX}}^{\llbracket \cdot \rrbracket}(\llbracket x \rrbracket, \llbracket y \rrbracket, \llbracket \text{msb}(x) \rrbracket^2 \oplus 1)$ with $y = 0$. The details are shown in Figure 9.

Our method needs an online communication of $(\frac{4}{3}\ell + 1)\lceil \log_2 q \rceil + \ell$ bits per party in $\lceil \log_4(\ell + 1) \rceil + 3$ rounds. And the security of Π_{ReLU} is easily to see in $(\mathcal{F}_{\text{SecMSB}}^{\llbracket \cdot \rrbracket}, \mathcal{F}_{\text{MUX}}^{\llbracket \cdot \rrbracket})$ -hybrid model.

6.3 Secure Maxpool

Given $\llbracket \cdot \rrbracket$ -shared vector $\mathbf{x} = (x_1, x_2, \dots, x_n)$ of size- n , the goal of functionality $\mathcal{F}_{\text{MP}}^{\llbracket \cdot \rrbracket}$ is to compute the maximum value among the n elements. $\mathcal{F}_{\text{MP}}^{\llbracket \cdot \rrbracket}$ can be implemented on top of ReLU. The key point is that the parties update $\llbracket \text{max} \rrbracket = \llbracket x_i \rrbracket$ if and only if $\text{ReLU}(\llbracket \text{max} \rrbracket - \llbracket x_i \rrbracket) = 0$ ($\Leftrightarrow \text{max} < x_i$). The full protocol is shown in Figure 10.

Furthermore, benefiting from binary sort on the inputs and small amounts of bookkeeping [73], the online phase needs approximately $(n-1)((\frac{4}{3}\ell + 1)\lceil \log_2 q \rceil + \ell)$ bits per party in $\lceil \log_2 n \rceil (\lceil \log_4(\ell + 1) \rceil + 3)$ rounds, and the security follows in the $\mathcal{F}_{\text{ReLU}}^{\llbracket \cdot \rrbracket}$ -hybrid model.

6.3.1 ReLU-MP Equivalent Switching MP is usually applied after ReLU, but they are commutative operators in NN inference: $\text{ReLU}(\text{MP}(\cdot)) = \text{MP}(\text{ReLU}(\cdot))$.

There is no significant performance difference of the alternation in cleartext, but $\text{ReLU}(\text{MP}(\cdot))$ is much more efficient than $\text{MP}(\text{ReLU}(\cdot))$ in MPC since the former reduces the number of ReLU operations significantly [43]. We thus evaluate MP before ReLU.

7 Evaluations

In this section, we present the system implementation and the detailed experimental results. The online costs of micro NN operators are in § 7.1. And in § 7.2, we present the results for single inference. In § 7.3, we show running-time and scalability for batch inference.

Experimental Details We implement METEOR on top of FALCON in C++ and run our experiments on Intel(R) Xeon(R) Silver 4314 CPU @ 2.40 GHz with 500 GB RAM in both LAN and WAN with a single thread. For LAN, our bandwidth is about 1 GB/s and round trip time (rtt) is about 1 ms. For WAN, our bandwidth is about 40 MB/s and rtt is about 70 ms. For fair comparisons, we re-run SecureNN and FALCON in our settings with semi-honest security.

Optimizations Following FALCON [73], we focus on online efficiency and do not take setup costs into account¹. As Softmax does not impact inference results, we omit it as [73]. We use Eigen library [2] for fast MatMul, uint64 for \mathbb{Z}_{2^ℓ} , $d = 13$, and $q = 67$. All experiments are executed 10 times, and we record the average².

Datasets & Neural Networks We select 2 standard benchmarking datasets: MNIST [44] and CIFAR-10 [42], and 6 standard network architectures: 3 from the secure ML community (Network-A [52], B [62], and C [48]) and 3 from the ML community (LeNet [45], AlexNet [41], and VGG16 [67]) for extensive experiments.

¹FALCON also does not take setup costs into consideration in experiments.

²The original results of communication in FALCON are incorrect due to the parallel implementation bugs and we have checked it with the authors. We re-run FALCON with a single thread. For details, please refer to <https://github.com/snwagh/falcon-public>.

Table 3: Online costs of NN operators of SecureNN, FALCON, and METEOR. For MatMul, the inputs are of size $m \times n$ and $n \times o$. And CONV is with input $m \times m$, c input channels, o output channels, and filter of $f \times f$. ReLU is computed in element-wise with size n . And MP is with $m \times m$ inputs, c input channels, and $f \times f$ window. Communication is in MB and Running-Time is in seconds.

Operator	Size	Comm.			Time (LAN)			Time (WAN)		
		SecureNN	FALCON	METEOR	SecureNN	FALCON	METEOR	SecureNN	FALCON	METEOR
MatMul $_{m,n,o}$	(784, 128, 10)	0.563	0.084	0.063	0.065	0.003	0.003	0.146	0.103	0.073
	(128, 500, 100)	0.642	0.137	0.102	0.258	0.007	0.012	0.346	0.118	0.076
CONV $_{m,c,o,f}$	(28, 1, 20, 5)	0.110	0.123	0.092	0.019	0.005	0.004	0.122	0.103	0.099
	(8, 16, 50, 5)	0.143	0.008	0.006	0.019	0.002	0.002	0.091	0.102	0.066
ReLU $_n$	128 \times 128	3.845	2.376	1.556	0.205	0.035	0.040	1.076	0.623	0.433
	576 \times 20	2.703	1.670	1.094	0.148	0.025	0.030	0.815	0.566	0.457
MP $_{m,c,f}$	(24, 20, 2)	2.143	1.287	0.821	0.123	0.038	0.031	1.412	1.355	0.929
	(8, 50, 4)	0.744	0.483	0.285	0.139	0.118	0.079	5.523	6.409	4.323

Table 4: Online costs of single inference for NN-A, -B, -C, and LeNet on MNIST of SecureNN, FALCON, and METEOR. Communication is in MB and Running-Time is in seconds.

Framework	NN-A		NN-B		NN-C		LeNet		
	Comm.	Time	Comm.	Time	Comm.	Time	Comm.	Time	
LAN	SecureNN	0.700	0.175	1.352	0.236	2.954	0.416	6.314	0.861
	FALCON	0.041	0.035	0.168	0.031	1.641	0.128	2.415	0.160
	METEOR	0.027	0.024	0.111	0.024	1.066	0.104	1.568	0.143
WAN	SecureNN	0.700	5.223	1.352	6.734	2.954	9.666	6.314	14.892
	FALCON	0.041	1.544	0.168	1.191	1.641	5.517	2.415	6.759
	METEOR	0.027	1.056	0.111	0.836	1.066	3.049	1.568	3.176

7.1 Online Costs of Micro Benchmarks

We present the online costs of the NN operators, including MatMul, CONV, ReLU, and MP, in Table 3. Our improvements are as follows: **Communication Improvements** For the linear operators (*i.e.*, MatMul and CONV), we improve the communication costs by approximately 1.2-6 \times and 1.3 \times over SecureNN and FALCON, respectively. This improvement arises from our free truncation technique. For ReLU and MP, we achieve respective 2.5 \times and 1.6 \times communication improvements compared with SecureNN and FALCON due to our communication efficient Π_{SecMSB} and Π_{MUX} .

Running-Time Improvements In the LAN setting, we improve the running-time by approximately 2-20 \times for linear operators and 4 \times for non-linear operators compared with SecureNN. Meanwhile, we achieve comparable running-time in comparison to FALCON in LAN. In the WAN setting, we are approximately 1.8 \times and 1.5 \times faster than SecureNN and FALCON, respectively.

7.2 Online Costs of Single Inference

The online costs of single inference are shown as follows.

Evaluation on MNIST We perform experiments on NN with MNIST as SecureNN [72] and FALCON [73], and the results are illustrated in Table 4. Compared to SecureNN, we improve the communication by upto 25.6 \times (11.1 \times on average), and the running-time by upto 9.8 \times (6.8 \times on average) and 8.1 \times (5.2 \times on average) in respective LAN and WAN. Besides, we reduce the communication by 1.5 \times on average, and are upto 1.5 \times (1.3 \times on average) and 2.1 \times (1.7 \times on average) faster than FALCON in LAN and WAN.

Evaluation on CIFAR10 As illustrated in Table 5, we evaluate METEOR on NN with CIFAR10 to demonstrate our improvements.

Table 5: Online costs of single inference on CIFAR10. Communication is in MB and Running-Time is in seconds.

Framework	AlexNet		VGG16		
	Comm.	Time	Comm.	Time	
LAN	FALCON	4.075	0.480	44.844	3.428
	METEOR	2.562	0.429	29.424	4.240
WAN	FALCON	4.075	13.522	44.844	38.559
	METEOR	2.562	8.997	29.424	27.174

For communication, we reduce the costs by 1.6 \times on average. For inference time, we achieve comparable efficiency in LAN and $\approx 1.5\times$ improvements in WAN over FALCON. Note in this case, our improvements in LAN are not as significant as that in WAN. The reason is that in LAN, the overall time is more restricted by the computation burden since there is enough bandwidth and a small rtt [73]. Therefore, the improvements in communication gain limited running-time benefits when the NN is computationally expensive. However, as MPC protocols are more likely to be executed in WAN, our improvements are meaningful in practical applications.

The improvements mainly stem from two aspects: i) our efficient MatMul improves the online communication of linear layers by 1.3 \times in comparison to FALCON. ii) More importantly, OUR online communication- and round-efficient protocols Π_{SecMSB} and Π_{MUX} improve the online efficiency of secure ReLU and MP functions.

7.3 Online Costs of Batch Inference

In this section, we measure the amortized online running-time for batch inference and show our improvements in scalability.

Amortized Running-Time Table 6 shows the running-time of METEOR over a batch of 128 images on AlexNet and VGG16 in LAN and WAN settings. For AlexNet, the amortized time for per-image

Table 6: Online Running-Time in seconds of METEOR with a batch of $|B| = 1$ and 128 images of CIFAR-10.

$ B $	LAN		WAN	
	1	128	1	128
AlexNet	0.429	18.649	8.997	203.304
VGG16	4.240	282.682	27.174	2404.509

drops from 0.429 s to 0.146 s (2.9 \times improvements) in LAN and from 8.997 s to 1.588 s (5.7 \times improvements) in WAN. While for VGG16, we achieve 1.9 \times and 1.4 \times time reduction for single image inference by batch processing in respective LAN and WAN. The improvement mainly comes from batch processing amortizes the computation and latency costs for each image.

Scalability Evaluation To present our scalability improvements against FALCON, we further measure the communication and running-time of METEOR on AlexNet and VGG16 with different batchsize as Figure 11 in Appendix F. We have the following findings: i) Given the batchsize $|B|$, we improve the online communication costs by $\approx 1.5\times$ in comparison to FALCON, which is consistent with the analysis for single inference in § 7.2. ii) For the online inference running-time, we achieve comparable efficiency in LAN but $\approx 1.5\times$ improvements in WAN compared to FALCON, which is not unexpected; after all, the communication improvements have mere gains to running-time in LAN as analyzed in § 7.2. Also, the scalability improvement is primarily due to our proposed efficient protocols.

8 Related Work

Secure NN inference using MPC has gained much attention recently. In the earlier stage, privacy-preserving machine learning mainly focused on traditional machine learning models such as linear regression [12, 27, 28, 65], logistic regression [68], decision trees [37], k-means clustering [14, 34], and SVM [71, 76].

In the area of two-party computation (2PC), CryptoNets [30] was one of the earliest works to use homomorphic encryption for secure NN inference, CryptoDL [31] developed approximate and low-degree polynomials to implement non-linear functions for efficiency improvements over CryptoNets. Mohassel *et al.* proposed SecurML [52] in the two-server setting with secret sharing and GC. Meanwhile, Liu *et al.* designed fast matrix multiplications protocols in MiniONN [48]. DeepSecure [64] uses GC to develop a privacy-preserving deep learning prediction framework, and GAZELLE [35] combines techniques from HE and MPC to achieve fast private inference. EzPC [17] is a ABY-based [26] framework, and its follow-up works [32, 58–60] focus on improving performance.

In order to solve the performance bottleneck of 2PC, recent works introduce a third party to assist computations. Chameleon [62] used the same technique as in [48] to complete the matrix multiplication operations but employed a semi-honest third-party to generate correlated randomness for multiplication triplets in offline. In order to solve the computational bottleneck incurred by the garbled circuits [6], both SecureNN [72] and Cryptflow [43] constructed novel protocols for non-linear functions such as ReLU and Max-pool that completely avoid the use of GC with the help of a third-party. What’s more, schemes based on 3PC replicated secret sharing also provide better overall efficiency [51, 73]. FALCON is one of the most efficient methods and can evaluate large NN such as VGG16 and AlexNet. Additionally, works such as BLAZE [57], Fantastic

Four [23], FLASH [15], and Tetrad [40], and [24, 46] are evaluated in MPC ($N \geq 3$), but these works focus on resisting malicious adversaries and require much more costs. And more works focus on dishonest-majority settings [13, 21, 25, 36]. ASTRA [18] and SWIFT [39] proposed 3PC sharing based on 2PC additive sharing and constructed protocols against malicious adversaries. METEOR is based on 3PC replicated secret sharing and semi-honestly secure, our protocols are more succinct and easier to follow. ScionFL [10] extended ABY2.0 to multi-party, and its Inner-Product protocol needs 2 elements per party in 2 rounds, which is 2 \times more expensive than ours. And our method is independent of ScionFL.

There are some other works combining quantized NN with MPC [5, 22, 61]. Riazi *et al.* proposed [61], where the weights and activations are in ± 1 , and they used GC and Oblivious Transfer (OT) to provide constant round private inference. QUOTIENT [5] was proposed to realize the secure computation of ternarized NN, where the weights are in $\{-1, 0, 1\}$. The author converts the ternarized multiplication into two binary multiplications and completes them based on OT. And other functions are all processed by GC. Therefore, prior private binary (ternarized) NN inference schemes suffer from the enormous communication costs introduced by GC, and they are even slower than secret sharing-based approaches for floating-point NN. Recently, some works proposed to utilize hardware, such as GPU, to accelerate the computation of MPC [20, 50, 53, 70, 74]. Specially, GForce is a 2PC inference framework, it proposed stochastic rounding and truncation layers to fuse (de)quantization between non-linear/linear layers for better efficiency, and a suite of GPU-friendly protocols for common operations. CryptGPU is a 3PC training and inference framework, it embedded cryptographic operations of discrete secret-shared values into floating-point operations to exploit existing CUDA kernels, and proposed several optimizations to softmax. These works mainly improved computation efficiency, while METEOR focuses on reducing communication.

9 Conclusion & Future Work

In METEOR, we propose an improved 3PC secret sharing scheme from the *linearity* of replicated secret sharing and construct secure blocks for secure NN inference. Extensive evaluations also present our improvements. For future work, we are willing to improve other MSB Ext. methods [29, 49, 51] with our novel secret sharing and the setup communication costs for better efficiency.

Acknowledgments

This work is supported by The National Key Research and Development Program of China No. 2020YFB1006100 and the Strategic Priority Research Program of Chinese Academy of Sciences, Grant No. XDC02040400.

References

- [1] 1996. *The Health Insurance Portability and Accountability Act of 1996 (HIPAA)*. <https://www.hhs.gov/hipaa/index.html>
- [2] 2011. Eigen library. <https://eigen.tuxfamily.org/>
- [3] 2016. *Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (GDPR)*. <https://gdpr-info.eu/>
- [4] 2022. Fusing Convolution and Batch Norm using Custom Function. https://pytorch.org/tutorials/intermediate/custom_function_conv_bn_tutorial.html

- [5] Nitin Agrawal, Ali Shahin Shamsabadi, Matt J Kusner, and Adrià Gascón. 2019. QUOTIENT: two-party secure neural network training and prediction. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*. 1231–1247.
- [6] Toshinori Araki, Jun Furukawa, Yehuda Lindell, Ariel Nof, and Kazuma Ohara. 2016. High-throughput semi-honest secure three-party computation with an honest majority. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*. 805–817.
- [7] Donald Beaver. 1991. Efficient multiparty protocols using circuit randomization. In *Annual International Cryptology Conference*. Springer, 420–432.
- [8] Mihir Bellare, Ran Canetti, and Hugo Krawczyk. 1996. Pseudorandom functions revisited: The cascade construction and its concrete security. In *Proceedings of 37th Conference on Foundations of Computer Science*. IEEE, 514–523.
- [9] Aner Ben-Efraim, Michael Nielsen, and Eran Omri. 2019. Turbospeedz: Double your online SPDZ! Improving SPDZ using function dependent preprocessing. In *Applied Cryptography and Network Security: 17th International Conference, ACNS 2019, Bogota, Colombia, June 5–7, 2019, Proceedings 17*. Springer, 530–549.
- [10] Yaniv Ben-Itzhak, Helen Möllering, Benny Pinkas, Thomas Schneider, Ajith Suresh, Oleksandr Tkachenko, Shay Vargaftik, Christian Weinert, Hossein Yalame, and Avishay Yanai. 2022. ScionFL: Secure Quantized Aggregation for Federated Learning. *arXiv preprint arXiv:2210.07376* (2022).
- [11] Dan Bogdanov, Sven Laur, and Jan Willems. 2008. Sharemind: A framework for fast privacy-preserving computations. In *European Symposium on Research in Computer Security*. Springer, 192–206.
- [12] Raphael Bost, Raluca Ada Popa, Stephen Tu, and Shafi Goldwasser. 2015. Machine learning classification over encrypted data. In *Symposium on Network and Distributed System Security (NDSS)*.
- [13] Lennart Braun, Daniel Demmler, Thomas Schneider, and Oleksandr Tkachenko. 2022. MOTION—A Framework for Mixed-Protocol Multi-Party Computation. *ACM Transactions on Privacy and Security* 25, 2 (2022), 1–35.
- [14] Paul Bunn and Rafail Ostrovsky. 2007. Secure two-party k-means clustering. In *Proceedings of the 14th ACM conference on Computer and communications security*. 486–497.
- [15] Megha Byali, Harsh Chaudhari, Arpita Patra, and Ajith Suresh. 2020. FLASH: Fast and Robust Framework for Privacy-preserving Machine Learning. *Proc. Priv. Enhancing Technol.* 2020, 2 (2020), 459–480.
- [16] Ran Canetti. 1998. Security and Composition of Multi-party Cryptographic Protocols. *Cryptology ePrint Archive, Paper 1998/018*. <https://eprint.iacr.org/1998/018> <https://eprint.iacr.org/1998/018>
- [17] Nishanth Chandran, Divya Gupta, Aseem Rastogi, Rahul Sharma, and Shardul Tripathi. 2019. EzPC: programmable and efficient secure two-party computation for machine learning. In *2019 IEEE European Symposium on Security and Privacy (EuroS&P)*. IEEE, 496–511.
- [18] Harsh Chaudhari, Ashish Choudhury, Arpita Patra, and Ajith Suresh. 2019. Astra: High throughput 3pc over rings with application to secure prediction. In *Proceedings of the 2019 ACM SIGSAC Conference on Cloud Computing Security Workshop*. 81–92.
- [19] Harsh Chaudhari, Rahul Rachuri, and Ajith Suresh. 2019. Trident: Efficient 4pc framework for privacy preserving machine learning. *arXiv preprint arXiv:1912.02631* (2019).
- [20] Zheng Chen, Feng Zhang, Amelie Chi Zhou, Jidong Zhai, Chenyang Zhang, and Xiaoyong Du. 2020. ParSecureML: An efficient parallel secure machine learning framework on GPUs. In *49th International Conference on Parallel Processing*. 1–11.
- [21] Ronald Cramer, Ivan Damgård, Daniel Escudero, Peter Scholl, and Chaoping Xing. 2018. SPDZ^k: efficient MPC mod 2^k for dishonest majority. In *Annual International Cryptology Conference*. Springer, 769–798.
- [22] Anders Dalskov, Daniel Escudero, and Marcel Keller. 2020. Secure evaluation of quantized neural networks. *Proceedings on Privacy Enhancing Technologies* 2020, 4 (2020), 355–375.
- [23] Anders Dalskov, Daniel Escudero, and Marcel Keller. 2021. Fantastic four: Honest-majority four-party secure computation with malicious security. In *30th {USENIX} Security Symposium ({USENIX} Security 21)*.
- [24] Anders Dalskov, Daniel Escudero, and Ariel Nof. 2022. Fast Fully Secure Multi-Party Computation over Any Ring with Two-Thirds Honest Majority. *Cryptology ePrint Archive* (2022).
- [25] Ivan Damgård, Daniel Escudero, Tore Frederiksen, Marcel Keller, Peter Scholl, and Nikolaj Volgushev. 2019. New primitives for actively-secure MPC over rings with applications to private machine learning. In *2019 IEEE Symposium on Security and Privacy (SP)*. IEEE, 1102–1120.
- [26] Daniel Demmler, Thomas Schneider, and Michael Zohner. 2015. ABY-A framework for efficient mixed-protocol secure two-party computation.. In *NDSS*.
- [27] Wenliang Du, Mikhail J Atallah, et al. 2001. Privacy-Preserving Cooperative Scientific Computations.. In *csfw*, Vol. 1. Citeseer, 273.
- [28] Wenliang Du, Yungshiang S Han, and Shigang Chen. 2004. Privacy-preserving multivariate statistical analysis: Linear regression and classification. In *Proceedings of the 2004 SIAM international conference on data mining*. SIAM, 222–233.
- [29] Daniel Escudero, Satrajit Ghosh, Marcel Keller, Rahul Rachuri, and Peter Scholl. 2020. Improved primitives for MPC over mixed arithmetic-binary circuits. In *Annual International Cryptology conference*. Springer, 823–852.
- [30] Ran Gilad-Bachrach, Nathan Dowlin, Kim Laine, Kristin Lauter, Michael Naehrig, and John Wernsing. 2016. Cryptonets: Applying neural networks to encrypted data with high throughput and accuracy. In *International Conference on Machine Learning*. 201–210.
- [31] Ehsan Hesamifard, Hassan Takabi, and Mehdi Ghasemi. 2017. Cryptodl: Deep neural networks over encrypted data. *arXiv preprint arXiv:1711.05189* (2017).
- [32] Zhicong Huang, Wenjie Lu, Cheng Hong, and Jiansheng Ding. 2022. Cheetah: Lean and Fast Secure Two-Party Deep Neural Network Inference. *Cryptology ePrint Archive, Report 2022/207*. <https://ia.cr/2022/207>.
- [33] Yuval Ishai, Joe Kilian, Kobbi Nissim, and Erez Petrank. 2003. Extending oblivious transfers efficiently. In *Annual International Cryptology Conference*. Springer, 145–161.
- [34] Geetha Jagannathan and Rebecca N Wright. 2005. Privacy-preserving distributed k-means clustering over arbitrarily partitioned data. In *Proceedings of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining*. 593–599.
- [35] Chiraag Juvekar, Vinod Vaikuntanathan, and Anantha Chandrakasan. 2018. {GAZELLE}: A low latency framework for secure neural network inference. In *27th {USENIX} Security Symposium ({USENIX} Security 18)*. 1651–1669.
- [36] Marcel Keller. 2020. MP-SPDZ: A versatile framework for multi-party computation. In *Proceedings of the 2020 ACM SIGSAC conference on computer and communications security*. 1575–1590.
- [37] Agnes Kiss, Masoud Naderpour, Jian Liu, N Asokan, and Thomas Schneider. 2019. Sok: Modular and efficient private decision tree evaluation. *Proceedings on Privacy Enhancing Technologies* 2019, 2 (2019), 187–208.
- [38] Brian Knott, Shobha Venkataraman, Awni Hannun, Shubho Sengupta, Mark Ibrahim, and Laurens van der Maaten. 2021. Crypten: Secure multi-party computation meets machine learning. *Advances in Neural Information Processing Systems* 34 (2021), 4961–4973.
- [39] Nishat Koti, Mahak Pancholi, Arpita Patra, and Ajith Suresh. 2021. {SWIFT}: Super-fast and Robust Privacy-Preserving Machine Learning. In *30th {USENIX} Security Symposium ({USENIX} Security 21)*.
- [40] Nishat Koti, Arpita Patra, Rahul Rachuri, and Ajith Suresh. 2021. Tetrad: Actively Secure 4PC for Secure Training and Inference. *arXiv preprint arXiv:2106.02850* (2021).
- [41] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. ImageNet classification with deep convolutional neural networks. *Advances in neural information processing systems* 25 (2012), 1097–1105.
- [42] V. Nair Krizhevsky and G. Hinton. 2014. The CIFAR-10 dataset.
- [43] Nishant Kumar, Mayank Rathee, Nishanth Chandran, Divya Gupta, Aseem Rastogi, and Rahul Sharma. 2019. Cryptflow: Secure tensorflow inference. *arXiv preprint arXiv:1909.07814* (2019).
- [44] Yann LeCun. 2017. MNIST database. <http://yann.lecun.com/exdb/mnist/>
- [45] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. 1998. Gradient-based learning applied to document recognition. *Proc. IEEE* 86, 11 (1998), 2278–2324.
- [46] Yi Li and Wei Xu. 2019. PrivPy: General and scalable privacy-preserving data mining. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 1299–1307.
- [47] Yehuda Lindell and Benny Pinkas. 2009. A proof of security of Yao’s protocol for two-party computation. *Journal of cryptography* 22, 2 (2009), 161–188.
- [48] Jian Liu, Mika Juuti, Yao Lu, and Nadarajah Asokan. 2017. Oblivious neural network predictions via minion transformations. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*. 619–631.
- [49] Eleftheria Makri, Dragos Rotaru, Frederik Vercauteren, and Sameer Wagh. 2021. Rabbit: Efficient Comparison for Secure Multi-Party Computation. In *International Conference on Financial Cryptography and Data Security*. Springer, 249–270.
- [50] Pratyush Mishra, Ryan Lehmkuhl, Akshayaram Srinivasan, Wenting Zheng, and Raluca Ada Popa. 2020. Delphi: A cryptographic inference service for neural networks. In *29th {USENIX} Security Symposium ({USENIX} Security 20)*. 2505–2522.
- [51] Payman Mohassel and Peter Rindal. 2018. ABY3: A mixed protocol framework for machine learning. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security*. 35–52.
- [52] Payman Mohassel and Yupeng Zhang. 2017. Secureml: A system for scalable privacy-preserving machine learning. In *2017 IEEE Symposium on Security and Privacy (SP)*. IEEE, 19–38.
- [53] Lucien KL Ng and Sherman SM Chow. 2021. GForce: GPU-Friendly Oblivious and Rapid Neural Network Inference. In *30th USENIX Security Symposium (USENIX Security 21)*. 2147–2164.
- [54] Chigozie Nwankpa, Winifred Ijomah, Anthony Gachagan, and Stephen Marshall. 2018. Activation functions: Comparison of trends in practice and research for deep learning. *arXiv preprint arXiv:1811.03378* (2018).

- [55] Satsuya Ohata and Koji Nuida. 2020. Communication-efficient (client-aided) secure two-party protocols and its application. In *Financial Cryptography and Data Security: 24th International Conference, FC 2020, Kota Kinabalu, Malaysia, February 10–14, 2020 Revised Selected Papers*. Springer, 369–385.
- [56] Arpita Patra, Thomas Schneider, Ajith Suresh, and Hossein Yalame. 2021. {ABY2.0}: Improved {Mixed-Protocol} Secure {Two-Party} Computation. In *30th USENIX Security Symposium (USENIX Security 21)*. 2165–2182.
- [57] Arpita Patra and Ajith Suresh. 2020. BLAZE: blazing fast privacy-preserving machine learning. *arXiv preprint arXiv:2005.09042* (2020).
- [58] Deevashwer Rathee, Anwesh Bhattacharya, Rahul Sharma, Divya Gupta, Nishanth Chandran, and Aseem Rastogi. 2022. SECFLOAT: Accurate Floating-Point meets Secure 2-Party Computation. In *2022 IEEE Symposium on Security and Privacy (SP)*. IEEE, 576–595.
- [59] Deevashwer Rathee, Mayank Rathee, Rahul Kranti Kiran Goli, Divya Gupta, Rahul Sharma, Nishanth Chandran, and Aseem Rastogi. 2021. SIRNN: A Math Library for Secure RNN Inference. *arXiv preprint arXiv:2105.04236* (2021).
- [60] Deevashwer Rathee, Mayank Rathee, Nishant Kumar, Nishanth Chandran, Divya Gupta, Aseem Rastogi, and Rahul Sharma. 2020. CryptFlow2: Practical 2-Party Secure Inference. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3372297.3417274>
- [61] M Sadegh Riazi, Mohammad Samragh, Hao Chen, Kim Laine, Kristin Lauter, and Farinaz Koushanfar. 2019. {XONN}: XNOR-based Oblivious Deep Neural Network Inference. In *28th {USENIX} Security Symposium ({USENIX} Security 19)*. 1501–1518.
- [62] M Sadegh Riazi, Christian Weinert, Oleksandr Tkachenko, Ebrahim M Songhori, Thomas Schneider, and Farinaz Koushanfar. 2018. Chameleon: A hybrid secure computation framework for machine learning applications. In *Proceedings of the 2018 on Asia Conference on Computer and Communications Security*. 707–721.
- [63] Dragos Rotaru and Tim Wood. 2019. Marbled circuits: Mixing arithmetic and boolean circuits with active security. In *International Conference on Cryptology in India*. Springer, 227–249.
- [64] Bitu Darvish Rouhani, M Sadegh Riazi, and Farinaz Koushanfar. 2018. Deepsecure: Scalable provably-secure deep learning. In *Proceedings of the 55th Annual Design Automation Conference*. 1–6.
- [65] Ashish P Sanil, Alan F Karr, Xiaodong Lin, and Jerome P Reiter. 2004. Privacy preserving regression modelling via distributed computation. In *Proceedings of 10th ACM SIGKDD international conference on Knowledge discovery and data mining*. 677–682.
- [66] Adi Shamir. 1979. How to share a secret. *Commun. ACM* 22, 11 (1979), 612–613.
- [67] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).
- [68] Aleksandra B Slavkovic, Yuval Nardi, and Matthew M Tibbits. 2007. "Secure" Logistic Regression of Horizontally and Vertically Partitioned Distributed Databases. In *Seventh IEEE International Conference on Data Mining Workshops*. IEEE, 723–728.
- [69] Ajith Suresh. 2021. Mpcleague: robust MPC platform for privacy-preserving machine learning. *arXiv preprint arXiv:2112.13338* (2021).
- [70] Sijun Tan, Brian Knott, Yuan Tian, and David J Wu. 2021. CRYPTGPU: Fast Privacy-Preserving Machine Learning on the GPU. *arXiv preprint arXiv:2104.10949* (2021).
- [71] Jaideep Vaidya, Hwanjo Yu, and Xiaoqian Jiang. 2008. Privacy-preserving SVM classification. *Knowledge and Information Systems* 14, 2 (2008), 161–178.
- [72] Sameer Wagh, Divya Gupta, and Nishanth Chandran. 2019. Securren: 3-party secure computation for neural network training. *Proceedings on Privacy Enhancing Technologies* 2019, 3 (2019), 26–49.
- [73] Sameer Wagh, Shruti Tople, Fabrice Benhamouda, Eyal Kushilevitz, Prateek Mittal, and Tal Rabin. 2020. FALCON: Honest-Majority Maliciously Secure Framework for Private Deep Learning. *arXiv preprint arXiv:2004.02229* (2020).
- [74] Jean-Luc Watson, Sameer Wagh, and Raluca Ada Popa. 2022. Piranha: A GPU Platform for Secure Computation. *Cryptology ePrint Archive* (2022).
- [75] Andrew Chi-Chih Yao. 1986. How to generate and exchange secrets. In *27th Annual Symposium on Foundations of Computer Science (sfcs 1986)*. IEEE, 162–167.
- [76] Hwanjo Yu, Jaideep Vaidya, and Xiaoqian Jiang. 2006. Privacy-preserving svm classification on vertically partitioned data. In *Pacific-asia conference on knowledge discovery and data mining*. Springer, 647–656.

A Non-Interactive $\langle \cdot \rangle$ -Randomness Generation

Functionality $\mathcal{F}_{\text{RAND}}^{\langle \cdot \rangle}$ enables parties sample a secret random number in $\langle \cdot \rangle$ -shared fashion using a pseudorandom function PRF [8]. Protocol Π_{RAND} achieves $\mathcal{F}_{\text{RAND}}^{\langle \cdot \rangle}$ as follows: Each pair of parties (i.e., P_i and P_j) maintain a pre-set random shared-key $k_{ij} \in_R \{0, 1\}^\kappa$, i.e., P_0 and P_1 share key k_{01} . Then, P_i and P_{i+1} generate $r_i \leftarrow \text{PRF}(k_{i,i+1}, \text{cnt})$ where cnt is a counter incremented for $i \in \{0, 1, 2\}$.

Protocol $\Pi_{N\text{-MULT}}$

Input: P_0, P_1 , and P_2 hold $[\cdot]$ -shared $\{\llbracket x_i \rrbracket\}_{i=1}^N$.

Output: $\llbracket z \rrbracket = \Pi_{i=1}^N \llbracket x_i \rrbracket$.

• **Setup:**

- (1) Parties generate $\langle \psi_z \rangle$ using functionality $\mathcal{F}_{\text{RAND}}^{\langle \cdot \rangle}$.
- (2) For $\mathcal{T} \subseteq \{1, \dots, N\}$, parties securely compute $\langle \Pi_{k \in \mathcal{T}} \psi_{x_k} \rangle = \Pi_{k \in \mathcal{T}} \langle \psi_{x_k} \rangle$ using $\mathcal{F}_{\text{MULT}}^{\langle \cdot \rangle}$ in a tree-manner.

• **Online:**

- (1) P_i locally computes $\langle m_z \rangle_i = \sum_{\mathcal{T} \subseteq \{1, \dots, N\}} (\Pi_{j \notin \mathcal{T}} m_{x_j} \cdot \langle \Pi_{k \in \mathcal{T}} \psi_{x_k} \rangle_i) - \langle \psi_z \rangle_i$.
- (2) Parties exchange the shares of $\langle m_z \rangle$ to reconstruct m_z .
- (3) P_i outputs $\llbracket z \rrbracket_i = (m_z, \langle \psi_z \rangle_i)$.

Figure 5: N -Input Multiplication Protocol for Integers.

Finally, P_i sets $\langle r \rangle_i = (r_i, r_{i+1})$ [6, 51]. Also, we can let P_i (secret owner) obtain r in clear by letting P_i keep $k_{i-1, i+1}$ and generate r_{i-1} locally.

As the pre-set shared-keys are secure (i.e., $\kappa = 128$), the security of Π_{RAND} is easy to see. Similarly, parties can generate shared random bit $\langle r \rangle^2$ in \mathbb{Z}_2 and $\langle r \rangle^q$ in \mathbb{F}_q .

B Protocol $\Pi_{N\text{-MULT}}$

Figure 5 show the N -input multiplication protocol for integers.

C Input-Independent Randomness Generation

Following FALCON [73], we generate the input-independent randomness for protocol Π_{SecMSB} in the setup phase as Figure 6: i) We first accomplish the bit decomposition of ψ_x in $\langle \cdot \rangle$ -sharing using $\mathcal{F}_{\text{BitDec}}^{\langle \cdot \rangle}$ [51], so that it is trivially to extract $\langle \text{msb}(\psi_x) \rangle^2$ and $\{\langle s[i] \rangle^2\}_{i=1}^\ell$ ($s = 2\psi_x$). ii) Then, we reshare $\langle \text{msb}(\psi_x) \rangle^2$ as $\llbracket \text{msb}(\psi_x) \rrbracket^2$. iii) Next, we convert $\{\langle s[i] \rangle^2\}_{i=1}^\ell$ to $\langle \cdot \rangle^q$ -shares in \mathbb{F}_q leveraging $\mathcal{F}_{\text{Bit2A}}^{\langle \cdot \rangle}$, and reshare $\{\langle s[i] \rangle^q\}_{i=1}^\ell$ as $\{\llbracket s[i] \rrbracket^q\}_{i=1}^\ell$. iv) Finally, we generate $\llbracket \zeta \rrbracket^q \in \mathbb{F}_q^*$ based on $\mathcal{F}_{\text{Prep}}^{\langle \cdot \rangle}$ from [73].

Protocol Π_{PreMSB} is a little expensive but practical since it can be executed in the setup phase. The correctness is guaranteed and security follows the analysis in [51, 73], and we omit it for brevity.

D Bit to Arithmetic Conversion for $\langle \cdot \rangle$ -Sharing

We first introduce functionality $\mathcal{F}_{\text{Bit2A}}^{\langle \cdot \rangle}$, which generates the arithmetic sharing of a bit $v \in \{0, 1\}$, given $\langle v \rangle^2 = (v_0, v_1, v_2)$. We utilize protocol $\Pi_{\text{Bit2A}}^{\langle \cdot \rangle}$ to achieve $\mathcal{F}_{\text{Bit2A}}^{\langle \cdot \rangle}$ as follows. Denote the value in \mathbb{Z}_{2^ℓ} of bit v as v^ℓ . Then we have

$$\begin{aligned} v^\ell &= v_0 \oplus v_1 \oplus v_2 \\ &= v_0^\ell + v_1^\ell + v_2^\ell - 2v_0^\ell v_1^\ell - 2v_1^\ell v_2^\ell - 2v_2^\ell v_0^\ell + 4v_0^\ell v_1^\ell v_2^\ell. \end{aligned} \quad (4)$$

For equation (4), only the last term $v_0^\ell v_1^\ell v_2^\ell$ requires communication and other terms $(v_i^\ell, v_i^\ell v_{i+1}^\ell)_{i=0}^2$ can be computed by P_i locally. For $v_0^\ell v_1^\ell v_2^\ell$, we firstly let P_0 compute $v_0^\ell v_1^\ell$. Then, parties can invoke Du-Atallah protocol [11, 62] to compute $v_0^\ell v_1^\ell v_2^\ell$ securely. Finally, we get the $[\cdot]$ -shared $[v_0^\ell v_1^\ell v_2^\ell]$ (step 1)-3) in Figure 7. Also, the other terms form $[\cdot]$ -sharing (e.g., for $[v_0^\ell]$, $[v_0^\ell]_0 = v_0^\ell$, $[v_0^\ell]_1 = [v_0^\ell]_2 = 0$),

Protocol Π_{PreMSB}

- Input:** $P_0, P_1,$ and P_2 hold $\langle \psi_x \rangle$.
Output: $\llbracket \text{msb}(\psi_x) \rrbracket^2, \{\llbracket s[i] \rrbracket^q\}_{i=1}^\ell$ with $s = 2\psi_x, (\llbracket \lambda \rrbracket^2, \llbracket \lambda \rrbracket^q)$ with $\lambda \in_R \mathbb{Z}_2$, and $\llbracket \zeta \rrbracket^q$ with $\zeta \in_R \mathbb{F}_q^*$.
- (1) Perform bit decomposition $\mathcal{F}_{\text{BitDec}}^{(\cdot)}$ from [51] to get $\langle \psi_x \rangle \rightarrow \langle \psi_x \rangle^2$.
 - (2) Sample random bits $\langle \psi_{\psi_x} \rangle^2, \langle \lambda \rangle^2,$ and $\langle \psi_\lambda \rangle^2,$ random values $\langle \psi_\lambda \rangle^q, \langle \psi_\zeta \rangle^q,$ and $\{\langle \psi_{s[i]} \rangle^q\}_{i=1}^\ell$ using $\mathcal{F}_{\text{RAND}}^{(\cdot)}$.
 - (3) Reveal $\langle m_{\psi_x} \rangle^2 = \langle \psi_x \rangle^2 [\ell] \oplus \langle \psi_{\psi_x} \rangle^2$ using $\mathcal{F}_{\text{REC}}^{(\cdot)}$, where $\langle \psi_x \rangle^2 [\ell]$ is the $\langle \cdot \rangle$ -shares of the msb of ψ_x . And set $\llbracket \text{msb}(\psi_x) \rrbracket^2 = (m_{\psi_x}, \langle \psi_{\psi_x} \rangle^2)$.
 - (4) Compute $\langle s \rangle^2 = (\langle \psi_x \rangle^2 \ll 1)$.
 - (5) Use $\mathcal{F}_{\text{Bit2A}}^{(\cdot)}$ for each bit of s to get $\langle s[i] \rangle^2 \rightarrow \langle s[i] \rangle^q$ with $i \in \{1, 2, \dots, \ell\}$ and $\langle \lambda \rangle^2 \rightarrow \langle \lambda \rangle^q$.
 - (6) Reveal $\langle m_{s[i]} \rangle^q = \langle s[i] \rangle^q - \langle \psi_{s[i]} \rangle^q \pmod{q}$ exploiting $\mathcal{F}_{\text{REC}}^{(\cdot)}$, and set $\llbracket s[i] \rrbracket^q = (m_{s[i]}, \langle \psi_{s[i]} \rangle^q), i \in \{1, 2, \dots, \ell\}$.
 - (7) Reveal $\langle m_\lambda \rangle^2 = \langle \lambda \rangle^2 \oplus \langle \psi_\lambda \rangle^2$ leveraging $\mathcal{F}_{\text{REC}}^{(\cdot)}$, and set $\llbracket \lambda \rrbracket^2 = (m_\lambda, \langle \psi_\lambda \rangle^2)$.
 - (8) Reveal $\langle m_\lambda \rangle^q = \langle \lambda \rangle^q - \langle \psi_\lambda \rangle^q \pmod{q}$ via $\mathcal{F}_{\text{REC}}^{(\cdot)}$, and set $\llbracket \lambda \rrbracket^q = (m_\lambda, \langle \psi_\lambda \rangle^q)$.
 - (9) Generate $\langle \zeta \rangle^q \in \mathbb{F}_q^*$ using $\mathcal{F}_{\text{Prep}}^{(\cdot)}$ from [73].
 - (10) Reveal $\langle m_\zeta \rangle^q = \langle \zeta \rangle^q - \langle \psi_\zeta \rangle^q \pmod{q}$ by $\mathcal{F}_{\text{REC}}^{(\cdot)}$, and set $\llbracket \zeta \rrbracket^q = (m_\zeta, \langle \psi_\zeta \rangle^q)$.
 - (11) Outputs $\llbracket \text{msb}(\psi_x) \rrbracket^2, \{\llbracket s[i] \rrbracket^q\}_{i=1}^\ell, (\llbracket \lambda \rrbracket^2, \llbracket \lambda \rrbracket^q),$ and $\llbracket \zeta \rrbracket^q$.

Figure 6: Secure MSB Extraction Pre-processing Protocol.
Protocol $\Pi_{\text{Bit2A}}^{(\cdot)}$

- Input:** $P_0, P_1,$ and P_2 hold $\langle v \rangle^2$ with $v \in \{0, 1\}$.
Output: $\langle v^\ell \rangle$.
- (1) P_2 samples $a_0, a_1 \in_R \mathbb{Z}_{2^\ell}$ and sends a_i to P_i for $i \in \{0, 1\}$.
 - (2) P_0 computes $v_0^\ell v_1^\ell + a_0$ and sends it to P_1 ; P_1 computes $v_2^\ell + a_1$ and sends it to P_0 .
 - (3) P_0 computes $[v_0^\ell v_1^\ell v_2^\ell]_0 = -a_0(v_2^\ell + a_1)$, P_1 computes $[v_0^\ell v_1^\ell v_2^\ell]_1 = v_2^\ell(v_0^\ell v_1^\ell + a_0)$, and P_2 computes $[v_0^\ell v_1^\ell v_2^\ell]_2 = a_0 a_1$.
 - (4) P_i locally computes $[v^\ell]_i = v_i^\ell - 2v_i^\ell v_{i+1}^\ell - [v_0^\ell v_1^\ell v_2^\ell]_i$.
 - (5) Parties reshare $[v^\ell]$ as $\langle v^\ell \rangle$.

Figure 7: Bit2A Conversion Protocol for $\langle \cdot \rangle$ -Sharing.

we hence can compute $[v^\ell]$ locally. In the end, parties can reshare $[v^\ell]$ as $\langle v^\ell \rangle$. The details are in Figure 7.

Protocol $\Pi_{\text{Bit2A}}^{(\cdot)}$ needs a communication of $\frac{7}{3}\ell$ bits per party in 3 rounds, and the security is fully analyzed in [11].

D.1 Bit to Arithmetic Conversion for $\llbracket \cdot \rrbracket$ -Sharing

The goal of functionality $\mathcal{F}_{\text{Bit2A}}^{[\cdot]}$ is to generate the arithmetic sharing of a given secret bit $\llbracket v \rrbracket^2 = (m_v, \langle \psi_v \rangle^2)$. Although we do not this

Protocol Π_{MM}

- Input:** $P_0, P_1,$ and P_2 hold $\llbracket \cdot \rrbracket$ -shared matrices $\llbracket \mathbf{X} \rrbracket$ and $\llbracket \mathbf{Y} \rrbracket$.
Output: $\llbracket \mathbf{Z} \rrbracket = \llbracket \mathbf{X} \rrbracket \cdot \llbracket \mathbf{Y} \rrbracket$.
- **Setup:**
- (1) Parties generate $\langle \psi_Z \rangle$ using functionality $\mathcal{F}_{\text{RAND}}^{(\cdot)}$.
 - (2) Parties execute $\langle \psi_{XY} \rangle = \mathcal{F}_{\text{MULT}}^{(\cdot)}(\langle \psi_X \rangle, \langle \psi_Y \rangle)$.
- **Online:**
- (1) P_i locally computes $\langle m_Z \rangle_i = m_{XY} + m_X \langle \psi_Y \rangle_i + m_Y \langle \psi_X \rangle_i + \langle \psi_{XY} \rangle_i - \langle \psi_Z \rangle_i$.
 - (2) Parties exchange the shares of $\langle m_Z \rangle$ to reconstruct m_Z .
 - (3) P_i sets and outputs $\llbracket \mathbf{Z} \rrbracket_i = (m_Z, \langle \psi_Z \rangle_i)$.

Figure 8: Secure Matrix Multiplication Protocol for Integers.
Protocol Π_{ReLU}

- Input:** $P_0, P_1,$ and P_2 hold $\llbracket \cdot \rrbracket$ -shared $\llbracket \mathbf{x} \rrbracket$.
Output: $\llbracket \mathbf{z} \rrbracket = \text{ReLU}(\llbracket \mathbf{x} \rrbracket)$.
- **Setup:**
- (1) Parties execute the setup phases of $\mathcal{F}_{\text{SecMSB}}^{[\cdot]}$ and $\mathcal{F}_{\text{MUX}}^{[\cdot]}$.
- **Online:**
- (1) Parties compute $\llbracket \text{msb}(\mathbf{x}) \rrbracket^2$ using $\mathcal{F}_{\text{SecMSB}}^{[\cdot]}$.
 - (2) Output $\llbracket \mathbf{z} \rrbracket = \mathcal{F}_{\text{MUX}}^{[\cdot]}(\llbracket \mathbf{x} \rrbracket, \llbracket \mathbf{y} \rrbracket, \llbracket \text{msb}(\mathbf{x}) \rrbracket^2 \oplus 1)$ with $\mathbf{y} = 0$.

Figure 9: Secure ReLU Protocol.
Protocol Π_{MP}

- Input:** $P_0, P_1,$ and P_2 hold $\llbracket \cdot \rrbracket$ -shared vector $\llbracket \mathbf{x} \rrbracket = (x_1, x_2, \dots, x_n)$.
Output: $\llbracket \text{max} \rrbracket$ with $\text{max} = \text{Max}(x_1, x_2, \dots, x_n)$.
- **Setup:**
- (1) Parties generate $\langle \psi_z \rangle$ using functionality $\mathcal{F}_{\text{RAND}}^{(\cdot)}$.
- **Online:**
- (1) Parties set $\llbracket \text{max} \rrbracket = \llbracket x_1 \rrbracket$.
 - (2) **for** $i = 2, 3, \dots, n$, **do**
 - (3) Parties compute $\llbracket \text{max} \rrbracket = \text{ReLU}(\llbracket \text{max} \rrbracket - \llbracket x_i \rrbracket) + \llbracket x_i \rrbracket$.
 - (4) **end for**
 - (5) Parties output $\llbracket \text{max} \rrbracket$.

Figure 10: Secure Maxpool Protocol.

functionality in this work, $\mathcal{F}_{\text{Bit2A}}^{[\cdot]}$ might be of independent interests and be helpful in other task such as privacy-preserving e-voting.

Given $\llbracket v \rrbracket^2 = (m_v, \langle \psi_v \rangle^2)$, we have

$$v^\ell = (m_v \oplus \psi_v)^\ell = m_v^\ell + \psi_v^\ell - 2m_v^\ell \psi_v^\ell. \quad (5)$$

In the setup phase, parties interactively generate the $\langle \cdot \rangle$ -shares of value ψ_v^ℓ using $\mathcal{F}_{\text{Bit2A}}^{(\cdot)}$ and sample random values $\langle \psi'_v \rangle$ in \mathbb{Z}_{2^ℓ} using $\mathcal{F}_{\text{RAND}}^{(\cdot)}$. In the online phase, parties locally compute $\langle v^\ell \rangle = m_v^\ell + \langle \psi'_v \rangle^\ell - 2m_v^\ell \langle \psi'_v \rangle^\ell$ and then reshare $\langle v^\ell \rangle$ as $\llbracket v^\ell \rrbracket$. Concretely, parties compute and reveal $\langle m'_v \rangle = \langle v^\ell \rangle - \langle \psi'_v \rangle^\ell$, and set $\llbracket v^\ell \rrbracket = (m'_v, \langle \psi'_v \rangle^\ell)$. **Online Communication** In the online phase, the parties reconstruct $\langle m'_v \rangle$ with a communication of ℓ bits per party in 1 round.

E Secure Buildings of NN operators

Figure 8-10 show the secure protocols for MatMul, ReLU, and MP.

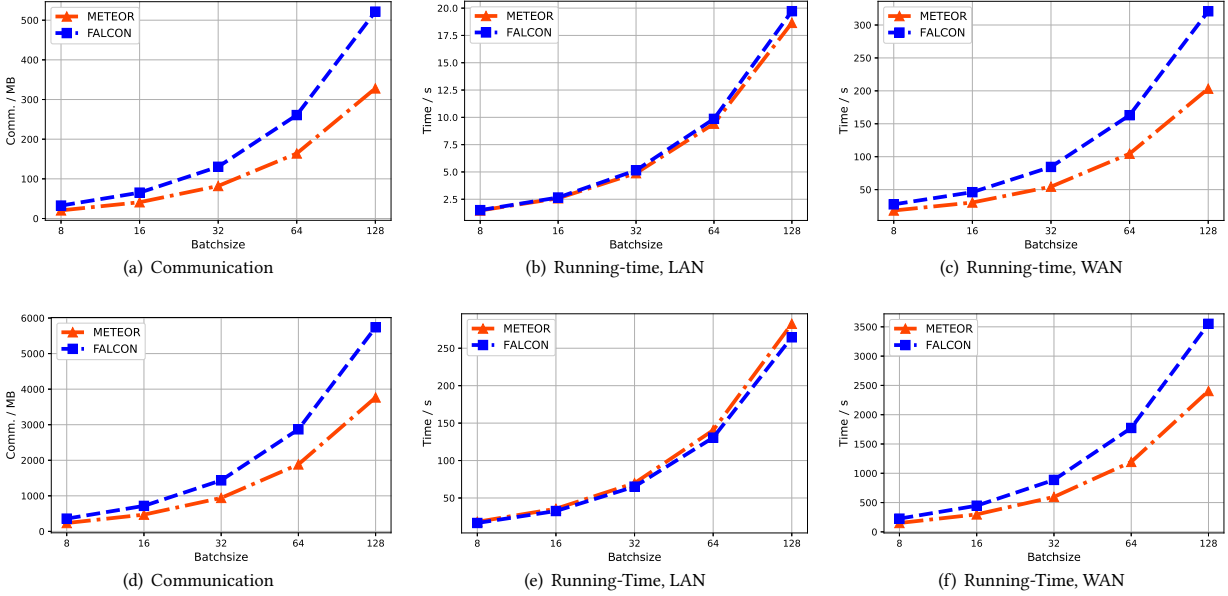


Figure 11: Online costs of batch inference of FALCON and METEOR for AlexNet and VGG16 on CIFAR-10, where 11(a)-11(c) are for AlexNet and 11(d)-11(f) are for VGG16. Communication is in MB and Running-Time is in seconds.

Secure Buildings of Other NN Operators: In addition, our basic primitives can support other secure NN operators: i) For other Activation functions (*i.e.*, Sigmoid and Tanh), we approximate them by piece-wise continuous polynomials following previous works [48, 52, 56]. The polynomials can be expressed as basic operations (*e.g.*, MULT & MSB Ext.), and thus we can construct their secure protocols on the top of the basic primitives. ii) Avgpool is much simpler than Maxpool. Note that the poolsize is in plaintext, parties thus can compute sum of their respective shares and truncate the sum by the poolsize to get the approximate average [38, 70]. As we do not employ these operators in METEOR, we omit their detailed protocols for brevity.

F Online costs of Batch Inference

The online costs of batch inference are illustrated in Figure 11.

G Security Proof

PROOF. Let the semi-honest adversary \mathcal{A} corrupt no more than one party, we now present the steps of the ideal-world adversary (simulator) \mathcal{S} for \mathcal{A} in the stand-alone model with security under sequential composition [16]. Our simulator \mathcal{S} for individual protocol is constructed as follows:

Security for Π_{SHARE} : For the instances where \mathcal{A} is the owner of the secret value x , \mathcal{S} has to do nothing since \mathcal{A} is not receiving any messages. \mathcal{S} receives m_x from \mathcal{A} on behalf of honest parties. For the instances where one honest party is the owner, \mathcal{S} sets $x = 0$ and follows protocol honestly.

Security for Π_{REC} : To reconstruct a value x , \mathcal{S} is given the output x , which is the output of \mathcal{A} . Using x and shares corresponding to honest parties, \mathcal{S} computes the shares corresponding to \mathcal{A} and sends this to \mathcal{A} on behalf of honest parties. \mathcal{S} sends the shares of honest parties to \mathcal{A} on behalf of honest parties.

Security for Π_{LOPS} : There is nothing to simulate as the protocol Π_{LOPS} is non-interactive.

Security for $\Pi_{2\text{-MULT}}$ & $\Pi_{N\text{-MULT}}$: For the setup phase, we consider the multiplication of $\langle \cdot \rangle$ -sharing as an ideal functionality $\mathcal{F}_{\text{MULT}}^{\langle \cdot \rangle}$ which multiplies the randomness. Since we make only black-box access to $\mathcal{F}_{\text{MULT}}^{\langle \cdot \rangle}$, the simulation for the same follows from the security of the underlying primitive used to instantiate $\mathcal{F}_{\text{MULT}}^{\langle \cdot \rangle}$ [6]. During the online phase, \mathcal{S} follows the step honestly using the data obtained from the corresponding setup phase.

Security for Π_{SecMSB} : For the setup phase, we invoke $\mathcal{F}_{\text{PreMSB}}$ in a black-box manner as FALCON [73]. Therefore, the simulation for the same follows from the security analyzed in [73]. For the online phase, we make black-box access to $\mathcal{F}_{\text{LOPS}}^{\llbracket \cdot \rrbracket}$, $\mathcal{F}_{2\text{-MULT}}^{\llbracket \cdot \rrbracket}$, and $\mathcal{F}_{N\text{-MULT}}^{\llbracket \cdot \rrbracket}$. To simulate the revealed d , \mathcal{S} samples a random number $r \in_R \mathbb{F}_q$, shares r as $\llbracket r \rrbracket^q$, and sends \mathcal{A} 's share to it. The security of protocol Π_{SecMSB} follows in the $(\mathcal{F}_{\text{PreMSB}}^{\llbracket \cdot \rrbracket}, \mathcal{F}_{2\text{-MULT}}^{\llbracket \cdot \rrbracket}, \mathcal{F}_{N\text{-MULT}}^{\llbracket \cdot \rrbracket})$ -hybrid model.

Security for Π_{MUX} : The setup phase is constructed directly on $\mathcal{F}_{\text{Bit2A}}^{\langle \cdot \rangle}$ and $\mathcal{F}_{\text{MULT}}^{\langle \cdot \rangle}$, thus the security is easily to see in $(\mathcal{F}_{\text{Bit2A}}^{\langle \cdot \rangle}, \mathcal{F}_{\text{MULT}}^{\langle \cdot \rangle})$ -hybrid model. For the online phase, \mathcal{S} follows the steps honestly using the data obtained from the corresponding setup phase.

This concludes the proof. \square