

Multi-Party Replicated Secret Sharing over a Ring with Applications to Privacy-Preserving Machine Learning

Alessandro Baccarini, Marina Blanton, and Chen Yuan
Department of Computer Science and Engineering
University at Buffalo (SUNY), Buffalo, USA
{anbaccar, mblanton, chyuan}@buffalo.edu

Abstract

Secure multi-party computation has seen significant performance advances and increasing use in recent years. Techniques based on secret sharing offer attractive performance and are a popular choice for privacy-preserving machine learning applications. Traditional techniques operate over a field, while designing equivalent techniques for a ring \mathbb{Z}_{2^k} can boost performance. In this work we develop a suite of multi-party protocols for a ring in the honest majority setting starting from elementary operations to more complex with the goal of supporting general-purpose computation. We demonstrate that our techniques are substantially faster than their field-based equivalents and perform on par with or better than state-of-the-art techniques. We also evaluate our techniques on machine learning applications and show that they offer attractive performance for these applications.

1 Introduction

Secure multi-party computation has recently seen notable performance improvements that make privacy-preserving computation of increasingly complex functionalities on increasingly large data sets more practical than ever before. Recent significant interest in privacy-preserving machine learning (PPML) has brought to light secret sharing techniques which were often previously overlooked in the literature. Secret sharing (SS) offers superior performance for arithmetic operations such as matrix multiplications and has been extensively used for privacy-preserving neural network (NN) inference and training [62, 18, 52, 55, 17, 40, 21, 63, 31]. Because SS offers information-theoretic security, computation can proceed on short integers, aiding efficiency.

Traditionally performance of SS techniques has been measured in terms of two parameters: the number of interactive operations and the number of sequential interactive operations, or rounds. However, for some computations such as matrix multiplication local operations can dominate the overall cost. Traditional techniques such as Shamir SS [61] carry out computation on protected data over a field, most commonly set up as \mathbb{Z}_p with prime p . This makes frequent use of modulo reduction a necessity, increasing the cost of the computation. To improve performance and directly utilize native instructions of modern processors, researchers turned to computation over ring \mathbb{Z}_{2^k} [14, 7, 19, 23]. Unfortunately, Shamir SS – a popular and efficient choice for computation in the honest majority setting – cannot be used for computation over \mathbb{Z}_{2^k} and we must seek other options.

The honest majority setting, which assumes that only a minority of the parties carrying out the computation can be corrupt, offers great performance with reasonable trust assumptions, making a good performance-security trade-off. The techniques we are aware of in this setting which can perform computation over ring \mathbb{Z}_{2^k} for some k are limited to a fixed number of parties, most

commonly to 3 (see, e.g., [7, 45, 52, 18, 17]). This means that the techniques do not easily generalize to a larger number of participants, should there be a need to change the computation setup. This is the task we set to address in this work and generalize computation based on replicated secret sharing (RSS) to support more than $n = 3$ computational parties.

Our contributions can be summarized as follows:

- We design a comprehensive set of elementary building blocks for RSS over an arbitrary ring in the semi-honest setting. These building blocks include generating shares of pseudorandom integers and ring elements, multiplication, reconstructing a value from shares, multiplication followed by reconstruction as a single building block, denoted by `MulPub`, and inputting private values into computation. We optimize the solutions to lower communication complexity by relying on a pseudo-random function. This means that the techniques are computationally secure and they also come with formal security proofs. Our solutions are efficient and, for example, the cost of multiplication when instantiated with three parties matches custom results which apply to the three-party setting only [7, 62].
- We build on the techniques of [23] and [31] to develop higher-level protocols over \mathbb{Z}_{2^k} such as for random bit generation, comparisons, conversion between different ring sizes and more to enable general-purpose computation in this framework.
- We provide extensive benchmarks to evaluate performance of the developed techniques. We observe that when $n = 3$ our techniques are up to 450 times faster than their field-based counterparts for operations such as matrix multiplication and up to 24 times faster for comparisons. Incorporating recent advances in random bit generation yields even more promising results.
- We improve the techniques of [21] for securely evaluating quantized NNs and eliminate the need for fixed-point multiplication and large truncation, which enables us to use a significantly smaller ring.
- We also evaluate performance of our techniques on machine learning applications, namely, NN predictions, quantized NN inference, and support vector machine (SVM) classification. Similarly, our runtimes are significantly faster than similar field-based implementations and compare favorably to the state of the art.

Because our techniques are based on RSS, it is expected that they will be used with a relatively small number of parties. This is similar to most efficient techniques based on Shamir SS (e.g., [16, 11]) which also rely on RSS for certain operations.

2 Related Work

Secret sharing [61, 10] is a popular choice for secure multi-party computation, and common options include Shamir SS [61], additive SS, and more recently RSS [36] for three parties. Computation over rings, and specifically \mathbb{Z}_{2^k} , has recently gained attention, and publications that use this setting include [14, 7, 45, 19, 25, 23, 30, 4, 39, 21]. We can distinguish between three-party techniques based on RSS such as [14, 7, 45, 25, 30, 4, 39]; multi-party techniques based on additive SS such as [19, 23], often for the setting with no honest majority; and ad-hoc techniques for three or four parties that utilize one or more types of rings with constructions for specific applications such as [38] and others.

The first category is the closest to this work and includes Sharemind [14], a well-developed framework for three-party computation with a single corruption using custom protocols; Araki et al. [7] who use three-party with a single corruption to support arithmetic or Boolean circuits; and several compilers from passively secure to actively secure protocols [45, 25, 30, 4]. Dalskov et al. [22]

also studied four-party computation with a single corruption. We are not aware of existing multi-party techniques with honest majority over a ring which extend beyond three parties or multi-party protocols based on RSS over a ring. While RSS is meaningful only for a small number of parties, we still find it desirable to support more participants and build additional techniques for this setting. For example, if our matrix multiplication protocol over a ring with three parties is 100 times faster than field-based computation, it will remain faster even if the work increases when the number of parties is larger than 3.

We rely on the results of Damgard et al. [23] for some of our protocols. While this work is for the SPDZ $_{\mathbb{Z}_2^k}$ framework [19] in the malicious setting with no honest majority, once we develop elementary building blocks, the structure of higher-level protocols can remain similar. Composite protocols such as comparison, conversion, and truncation require a large number of random bits. We leverage the edaBit protocol from [31] to efficiently generate sets of binary and arithmetic shared bits. Their technique improves upon the daBit technique [59]. Rabbit [48] builds on daBits [59] and edaBits [31] and developed an efficient n -party comparison protocol by relying on commutativity of addition over fields and rings. Their protocol offers significant improvement over [31] in most adversarial settings over a field, but remains comparable with a passively secure honest majority over a ring.

Literature on PPML is also related to this work. We distinguish between two-party solutions where one party holds the model and the other party holds the input on which the model is to be evaluated and between multi-party (typically, three-party) solutions. Publications from the first category include MiniONN [46] which studied NN evaluation using SS and homomorphic encryption; Gazelle [38] which combined homomorphic encryption with garbled circuits (GC) and additive SS; DELPHI [51] which improved upon these techniques; CryptFlow2 [57] which also built upon this work with a focus on deep NNs; and ABY2.0 [54] which expanded upon the ABY [28] framework with fundamental operations relying on Beaver circuit randomization [8]. Chameleon [58] incorporated GCs, the GMW protocol [50], and additive SS. MP2ML [12] introduced a hybrid MPC and homomorphic encryption framework based on ABY and nGraph-HE [13], respectively. SIRNN [56] provides semi-honest two-party protocols and improved approximations of several continuous functions.

Multi-party constructions provide protocols for training and prediction across multiple parties. SecureML [53] was one of the first publications to provide a two-server architecture for training NNs (as well as several other machine learning applications). ABY3 [52] combines techniques based on replicated and binary SS with and GCs in the three-party setting with honest majority. These techniques are further improved in Trident [18] and extended to the four-party setting. SecureNN [62] provides three- and four-party protocols for a variety of NN functions under the same security assumption as ABY3. Their protocols are asymmetric, where parties have dedicated roles in a computation. This work is improved upon with FALCON [63] by adding malicious security with honest majority and combining the techniques from SecureNN and ABY3.

ASTRA [17] is a three-party framework that uses SS over the ring \mathbb{Z}_2^k under both semi-honest and malicious security assumptions. Similar to SecureNN, protocols are asymmetric. BLAZE [55] builds upon this work in a similar setting. Abspoel et al. [5] applies the MP-SPDZ [39] framework for secure outsourced training of decision trees. Their system operates under the three-party, honest-majority assumption with RSS. The Manticore framework [15] provides support for real-number arithmetic under full-threshold semi-honest security. Dalskov et al. [21] were the first to address quantized NN inference using secure multi-party computation. Their system is built into MP-SPDZ and benchmarked on the MobileNets [35] network architecture. Keller et al. [40] conducts quantization-based training and inference with three parties and one semi-honest corruption.

3 Preliminaries

3.1 Secure Multi-Party Computation

We consider a secure multi-party setting with n computational parties, out of which at most t can be corrupt. We work in the setting with honest majority, i.e., $t > n/2$ and semi-honest participants and use simulation-based security (see Appendix B for detail).

As customary with SS techniques, the set of computational parties does not have to coincide with (and can be formed independently from) the set of parties supplying inputs in the computation (input providers) and the set of parties receiving output of the computation (output recipients). Then if a computational party learns no output, the computation should reveal no information to that party. Consequently, if we wish to design a functionality that takes secret-shared input and produces shares of the output, any computational party should learn nothing from protocol execution.

3.2 Secret Sharing

A SS scheme allows one to produce shares of secret x such that access to a predefined number of shares reveals no information about x . In the context of secure multi-party computation, each of the n participants receives one or more shares x_i and in the case of (n, t) threshold SS schemes, possession of shares stored at any t or fewer parties reveals no information about x , while access to shares stored at $t + 1$ or more parties allows for reconstruction of x . Of particular importance are linear SS schemes, which have the property that a linear combination of secret shared values can be performed locally on the shares. Examples of linear SS schemes include additive SS with $x = \sum_i x_i$ (as used in Sharemind [14] with $n = 3$ and in SPDZ [26] with any n), Shamir SS which realizes (n, t) secret sharing with $t < n/2$ and represents a share as evaluation of a polynomial on a distinct point, and RSS discussed next.

3.3 Replicated Secret Sharing

Our techniques utilize RSS [36] which has an associated access structure Γ . An access structure is defined by qualified sets $Q \in \Gamma$, which are the sets of participants who are granted access, and the remaining sets of the participants are called unqualified sets. In the context of this work we only consider threshold structures in which any set of t or fewer participants is not authorized to learn information about private values (i.e., they form unqualified sets), while any $t + 1$ or more participants are able to jointly reconstruct the secret (and thus form qualified sets). RSS can be defined for any $n \geq 2$ and any $t < n$. To secret-share private x using RSS, we treat x as an element of a finite ring \mathcal{R} and additively split it into shares x_T such that $x = \sum_{T \in \mathcal{T}} x_T$ (in \mathcal{R}), where \mathcal{T} consists of all maximal unqualified set of Γ (i.e., all sets of t parties in our case). Then each party $p \in [1, n]$ stores shares x_T for all $T \in \mathcal{T}$ subject to $p \notin T$. In the general case of (n, t) -threshold RSS, the total number of shares is $\binom{n}{t}$ with $\binom{n-1}{t}$ shares stored by each party, which can become large as n and t grow. In what follows, we use notation $[x]$ to mean that (private) x is secret shared among the parties using RSS.

Example. In the $(4, 2)$ setting, \mathcal{T} consists of 6 sets $\mathcal{T} = \{\{1, 2\}, \{1, 3\}, \{1, 4\}, \{2, 3\}, \{2, 4\}, \{3, 4\}\}$ and thus there are 6 corresponding shares for every secret-shared x . Then party 1 stores shares $x_{\{2,3\}}, x_{\{2,4\}}, x_{\{3,4\}}$, party 2 stores $x_{\{1,3\}}, x_{\{1,4\}}, x_{\{3,4\}}$, etc.

The parties will need to perform computation on secret shared values. The first important property of RSS is that it is linear. For example, to add $[a]$ and $[b]$, party p computes $a_T + b_T$ (in \mathcal{R}) for each $T \in \mathcal{T}$ that p stores. A number of other operations, such as multiplications,

reconstructing a value from its shares, etc., are interactive. We consequently describe in Section 4 the way we realize these operations. An important optimization on which we rely is non-interactive evaluation of a pseudo-random function (PRF) using RSS in the computational (as opposed to information-theoretic) setting as proposed in [20]; see Section 4 for detail.

In what follows, we use the notation \leftarrow to denote output of randomized algorithms, while the notation $=$ refers to deterministic assignment.

4 Basic Protocols

Recall that RSS enjoys the linear property. In addition to adding secret-shared values, we use the ability to add/subtract known integers to a secret-shared value $[a]$ and multiply a secret-shared value $[a]$ by a known integer. Addition $[a] + b$ converts b to $[b]$ without using randomness (e.g., we could set one share to b and the remaining shares to 0 to maintain $\sum_{T \in \mathcal{T}} b_T = b$). Multiplication $c = [a] \cdot b$ sets $c_T = a_T \cdot b$ (in \mathcal{R}) $\forall T \in \mathcal{T}$.

For convenience and without loss of generality, we let $n = 2t + 1$. When $n > 2t + 1$, $2t + 1$ parties can carry out the computation on a reduced set of shares in such a way that there is no need to involve the remaining parties in the computation.

4.1 Random Number Generation

We will be using two types of random number generation, which we discuss here.

4.1.1 PRG

Invocation of $[a_1], [a_2], \dots \leftarrow \text{PRG}([s])$ is realized by independently executing a PRG algorithm on each share of s without interaction between the parties. Because the output of $\text{PRG}([s])$ is private, we expect it to produce a sequence of secret-shared values (represented as ring elements). Furthermore, in our construction we only call the PRG to obtain random (secret-shared) ring elements. This means that calling $\text{PRG}(s_T)$ to produce pseudo-random a_T will result in $\text{PRG}([s])$ generating $[a]$, where a is pseudo-random as well because $a = \sum_{T \in \mathcal{T}} a_T$ (in \mathcal{R}). This is similar to evaluating a PRF on a secret-shared key in the RSS setting without interaction in [20].

$\text{PRG}(s_T)$ can be realized internally using any suitable algorithm, as long as it is consistent among the computational parties. For example, because of the speed of AES encryption on modern processors, one might implement $\text{PRG}(s_T) = \text{PRF}(s_T, 0) \parallel \text{PRF}(s_T, 1) \parallel \dots$, where $\text{PRF} : \mathcal{R} \times \{0, 1\}^k \rightarrow \mathcal{R}$ is a PRF instantiated with AES.

Let $G = \text{PRG}([s])$. When the output of G is not consumed all at once, we use notation $G.\text{next}$ to retrieve the next (secret-shared) element from G . Similarly, if $G_T = \text{PRG}(s_T)$, notation $G_T.\text{next}$ refers to the next pseudo-random share output by G_T .

4.1.2 PRandR

$[a] \leftarrow \text{PRandR}()$ computes a secret-shared random element of ring \mathcal{R} . We implement this function by executing $\text{PRG}([k]).\text{next}$, where k is a system-wide key. The key k is set up at the system initialization time (in the form of secret shares) and does not change throughout program execution.

4.2 Multiplication

Multiplication $[c] \leftarrow [a] \cdot [b]$ is realized using the fact that $[a] \cdot [b] = \sum_{T_1, T_2 \in \mathcal{T}} a_{T_1} \cdot b_{T_2}$ (in \mathcal{R}). Note that for any (T_1, T_2) pair, there will be a party holding shares T_1 and T_2 , and thus performing this

operation involves local multiplication and addition over different choices of T_1, T_2 . More formally, let mapping $\rho : \mathcal{T} \times \mathcal{T} \rightarrow [1, n]$ denote a function that for each pair $(T_1, T_2) \in \mathcal{T}^2$ dedicates a party $p \in [1, n]$ responsible for computing the product $a_{T_1} \cdot b_{T_2}$ (clearly, p must possess shares T_1 and T_2). For performance reasons, we also desire that ρ distributes the load among the parties as fairly as possible.

As a result of this (local) computation, the parties hold additive shares of the product $a \cdot b = c$, which needs to be converted to RSS for consecutive computation. This conversion was realized in early publications [49, 9] by having each party create replicated secret shares of their result and distribute each share to the parties entitled to knowing it (i.e., party p receives shares from each party for each $T \in \mathcal{T}$ subject to $p \notin T$). This results in each participant creating $\binom{n}{t}$ shares and sending $\binom{n-1}{t}$ of them to each party. Consequentially, each participant adds the values received for share T and stores the sum as c_T , for each T in its possession.

More recent work, e.g., [7] and others traded information-theoretic security (in the presence of secure channels) for communication efficiency by having the parties generate shared (pseudo-)random values. We pursue this direction as well. However, if this idea is applied naively, it results in unnecessarily high overhead. In particular, if we instruct each party p to generate all shares for its secret, some shares will be known to more than t participants and thus do not contribute to secrecy. Instead, our solution eliminates shares that p does not possess and thus do not contribute to secrecy. Thus, our construction utilizes key material consistent with the setup of the RSS scheme. In particular, we use the same key setup as in PRSS, where k_T is known by all $p \notin T$. Then when a party needs to generate a pseudo-random share associated of its value for share T , the party will draw it from the PRG seeded with k_T .

We, however, note that multiple participants may need to draw from the PRG seeded with k_T to produce shares of their values, and it is generally not safe to use the same secret to protect multiple values, which is also the case in our application. Instead, multiple elements might be drawn from the PRG (seeded with k_T) to protect different values, and consistent use of the PRG with each seed can be setup by the participants ahead of time so that that information is public knowledge.

In addition to mapping ρ , our multiplication protocol uses another mapping $\chi : [1, n] \rightarrow \mathcal{T}$, which for each party p specifies the share T (subject to $p \notin T$) that p communicates (with all other shares of p 's value being produced as pseudo-random elements). As before, we desire to choose the values of $\chi(p)$ as to evenly distribute the load and communication.

The above intuition leads us to the optimized n -party multiplication protocol given as Protocol 1. After computing its private value $v^{(p)}$ according to ρ , each party p distributes it into $\binom{n-1}{t}$ additive shares (one of which is communicated while others are computed using PRGs). Afterwards, each party sets its c_T as a sum of $t+1$ shares (computed or received) of values $v^{(p')}$ for each party p' entitled to shares c_T . This matches the fact that each share a_T of secret a is maintained by $t+1$ parties. Correctness is achieved by ensuring that in Protocol 1 two different participants p and p' with access to shares T consistently associate the values that they draw from G_T with shares belonging to different parties by always processing the values in the increasing order of participants' IDs. Preparation of the shares in Protocol 1 is done on lines 10–16, where a participant either masks its share with a pseudo-random value because it is used by another party or forms its own shares and the value to be transmitted.

In this protocol, each party on average sends t ring elements and draws $\binom{n-1}{t} - 1 + (n-1)\binom{n-2}{t} - t$ pseudo-random ring elements (which is $(t+1)(\binom{n-1}{t} - 1)$ when $n = 2t+1$).¹The latter can be explained by using $\binom{n-1}{t} - 1$ pseudo-random shares for its value being re-shared and $\binom{n-2}{t}$ shares

¹It is possible to distribute the load evenly among the parties by appropriately setting the χ function.

Protocol 1 $[c] \leftarrow [a] \cdot [b]$

```
// pre-distributed values are  $[k]$  and public maps  $\rho$  and  $\chi$ 
// define  $G_T = \text{PRG}(k_T)$ 
1: each  $p \in [1, n]$  does the following
2:   let  $S_p = \{T \in \mathcal{T} \mid p \notin T\}$ ;
3:    $v^{(p)} = \sum_{T_1, T_2 \in \mathcal{T}, \rho(T_1, T_2) = p} a_{T_1} b_{T_2}$  (in  $\mathcal{R}$ );
4:   for  $T \in S_p$  do
5:      $c_T = 0$ ;
6:      $v_{\chi(p)}^{(p)} = v^{(p)}$ ;
7:   end for
8:   for  $p' \in [1, n]$  in order do
9:     for  $T \in S_p$  do
10:      if  $(p' \neq p) \wedge (p' \notin T) \wedge (\chi(p') \neq T)$  then
11:         $c_T = c_T + G_T.\text{next}$  (in  $\mathcal{R}$ );
12:      else if  $(p' = p) \wedge (\chi(p) \neq T)$  then
13:         $z = G_T.\text{next}$ ;
14:         $c_T = c_T + z$  (in  $\mathcal{R}$ );
15:         $v_{\chi(p)}^{(p)} = v_{\chi(p)}^{(p)} - z$  (in  $\mathcal{R}$ );
16:      end if
17:    end for
18:  end for
19:  send  $v_{\chi(p)}^{(p)}$  to each  $p' \notin \chi(p)$  (other than itself);
20:  for  $p' \in [1, n]$  such that  $p \notin \chi(p')$  do
21:    after receiving  $v_{\chi(p')}^{(p')}$  from  $p'$ , set  $c_{\chi(p')} = c_{\chi(p')} + v_{\chi(p')}^{(p')}$  (in  $\mathcal{R}$ );
22:  end for
23:   $c_{\chi(p)} = c_{\chi(p)} + v_{\chi(p)}^{(p)}$  (in  $\mathcal{R}$ );
24: return  $[c]$ ;
```

that it has in common with any other party except the t values that it receives with a symmetric communication pattern. (Recall that each party maintains $\binom{n-1}{t}$ shares of a secret and has $\binom{n-2}{t}$ shares in common with any other party). When the communication pattern is not symmetric, the overall amount of work and communication remains unchanged, but it may be distributed differently. Thus, we refer to the average work and communication in that case.

Compared to other results, the three-party version of our protocol matches communication of recent multiplication from [7], which is available only for three parties and improves on communication of Sharemind's three-party multiplication from [41] by a factor of 2. For multi-party multiplication it can be desirable to use a different communication pattern when a king reconstructs a protected value and communicates it to others (as in, e.g., [24]) which scales better as n grows. However, our version has lower communication when $n = 3$, uses fewer rounds, and n cannot be large with RSS.

Example. With three parties, we could have party 1 (in possession of shares $\{2\}$ and $\{3\}$) compute (and add) products $a_{\{2\}}b_{\{2\}}$, $a_{\{2\}}b_{\{3\}}$, and $a_{\{3\}}b_{\{2\}}$, party 2 (in possession of shares $\{1\}$ and $\{3\}$) compute products $a_{\{3\}}b_{\{3\}}$, $a_{\{1\}}b_{\{3\}}$, and $a_{\{3\}}b_{\{1\}}$, and party 3 (in possession of shares $\{1\}$ and $\{3\}$) compute products $a_{\{1\}}b_{\{1\}}$, $a_{\{1\}}b_{\{2\}}$, and $a_{\{2\}}b_{\{1\}}$. This defines mapping ρ . Also let $\chi(1) = \{2\}$, $\chi(2) = \{3\}$, and $\chi(3) = \{1\}$. This, for example, means that when party 1 divides its computed value $v^{(1)}$ into shares $v_{\{2\}}^{(1)}$ and $v_{\{3\}}^{(1)}$, the latter is computed using a PRG, while the

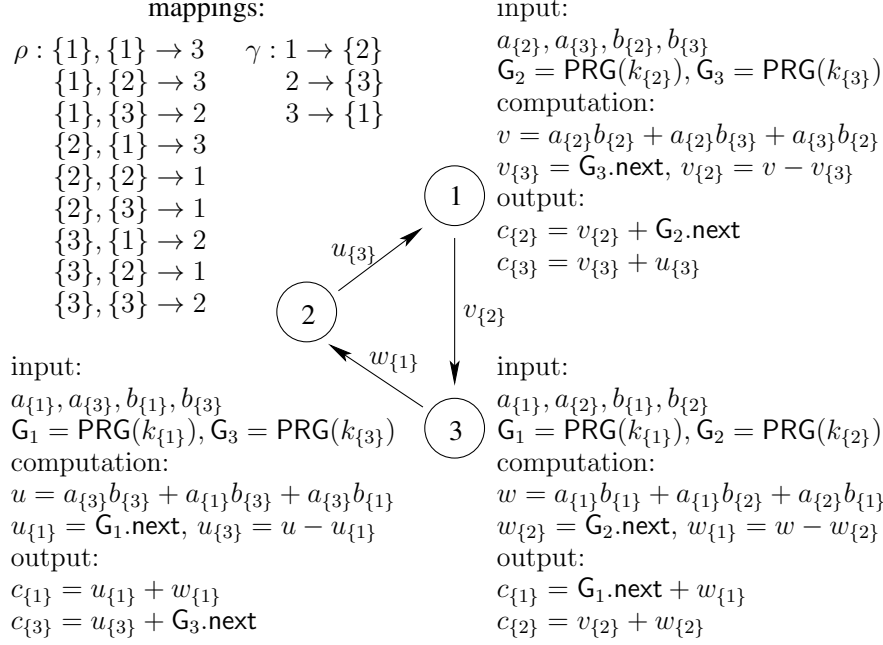


Figure 1: Sample three-party multiplication $[a] \cdot [b]$; arithmetic is in \mathcal{R} .

Operation	Rounds	(3, 1) setting		(n, t) setting	
		Comm	Crypto ops	Comm	Crypto ops
PRG($[s]$).next and PRandR()	0	0	2	0	$\binom{n-1}{t}$
Mul($[a], [b]$)	1	1	2	t	$(t+1) \left(\binom{n-1}{t} - 1 \right)$
Open($[a]$)	1	1	0	t	0
MulPub($[a], [b]$)	1	2	2	$n-1$	$t \binom{n-1}{t}$
DotProd($\langle [a^1], \dots, [a^N] \rangle, \langle [b^1], \dots, [b^N] \rangle$)	1	1	2	t	$(t+1) \left(\binom{n-1}{t} - 1 \right)$

Table 1: Performance of basic RSS operations with computation and communication per party.

former is being sent to party 3 (i.e., the other party entitled to have that share). An illustration of the multiplication protocol with these mappings in the three-party setting is given in Figure 1.

We state security of multiplication as follows, with its proof available in Appendix B:

Theorem 1. *Multiplication $[c] \leftarrow [a] \cdot [b]$ is secure according to definition 1 in the (n, t) setting with $t = (n-1)/2$ in the presence of secure communication channels and assuming PRG is a pseudo-random generator.*

The computation associated with multiplication can be generalized to compute the dot-product of two secret-shared vectors $\text{DotProd}(\langle [a^1], \dots, [a^N] \rangle, \langle [b^1], \dots, [b^N] \rangle)$, or evaluate any other multivariate polynomial of degree 2, using the same communication and the same number of cryptographic operations as in multiplication. For that purpose, we only need to change the computation in step 3 of the multiplication protocol. For example, for DotProd , we modify step 3 to compute $v^{(p)} = \sum_{T_1, T_2 \in \mathcal{T}, \rho(T_1, T_2) = p} \sum_{i=1}^N a_{T_1}^i b_{T_2}^i$ (in \mathcal{R}), while the rest of the steps remain unchanged.

Table 1 shows performance of these and other basic protocols for the general (n, t) and the (3,1) settings. Communication is measured as the number of ring elements sent by each party and computation is the number of cryptographic operations (i.e., retrieval of the next pseudo-random element using a PRG) per party.

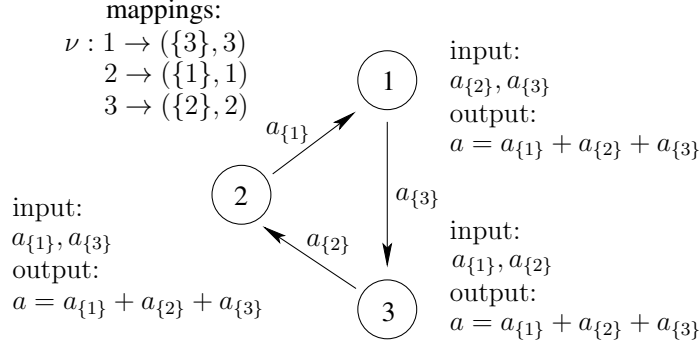


Figure 2: Sample three-party $\text{Open}([a])$; arithmetic is in \mathcal{R} .

4.3 Revealing Private Values

4.3.1 Open

Reconstruction of a secret shared value $a = \text{Open}([a])$ amounts to communicating missing shares to each party so that the value could be reconstructed locally from all shares. Recall that there are $\binom{n}{t}$ total shares and each party holds $\binom{n-1}{t}$ of them. Thus, during this operation each party is to receive $d = \binom{n}{t} - \binom{n-1}{t}$ missing shares.

Our next observation is that when n is not small, the value of d will exceed n and transmitting d messages to each party is not needed. Since the value is reconstructed as the sum of all shares, it is sufficient to communicate sums of shares instead of the individual shares themselves. Recall that $[a]$ can be reconstructed by $t + 1$ parties. This means that it is sufficient for a participant to receive one element (i.e., a sum of the necessary shares) from t other parties.

As before, we would like to balance the load between the parties and ideally have each party transmit the same amount of data. This means that we instruct each party to send information to t other parties according to another agreed upon mapping $\nu : [1, n] \rightarrow (\mathcal{T}, [1, n])^d$. For each party p , this mapping will specify which of p 's shares should be communicated to which other party. The mapping ν will then define computation associated with this operation: each p computes $\sum_{T, \nu(p) \subseteq T, p' \in T} a_T$ (in \mathcal{R}) for each $p' \neq p$ present in the mapping and sends the result to p' .

Similar to other SS frameworks, simply opening the shares of a maintains security of the computation (in the sense that no information about private values is revealed beyond the opened value a). This is because we maintain that at the end of each operation secret-shared values are represented using random shares. In particular, it is clear that the result of $\text{PRG}([s]).\text{next}$ and $\text{PRandR}()$ produces random shares; shares are properly re-randomized during multiplication of $[a]$ and $[b]$, and shares of $[a] + [b]$ and $[a] - [b]$ are random if the shares of $[a]$ and $[b]$ are random themselves.

Example. With $n = 3$, we could have $\nu(1) = (\{3\}, 3)$, $\nu(2) = (\{1\}, 1)$, and $\nu(3) = (\{2\}, 2)$, which corresponds to $\nu(p) = (\{p-1\}, p-1)$ (where $p-1 = 3$ for $p = 1$), which corresponds to the communication pattern in Figure 2.

4.3.2 MulPub

Functionality $c = \text{MulPub}([a], [b])$ refers to multiplying two secret-shared $[a]$ and $[b]$ and opening their product c . The reason why we are discussing this functionality is because in the past this operation could be implemented more efficiently than multiplication followed by an opening in alternative SS frameworks (e.g., see [16]), and we pursue a similar direction here. In the protocol we present here, MulPub is realized using a single round without increasing communication cost.

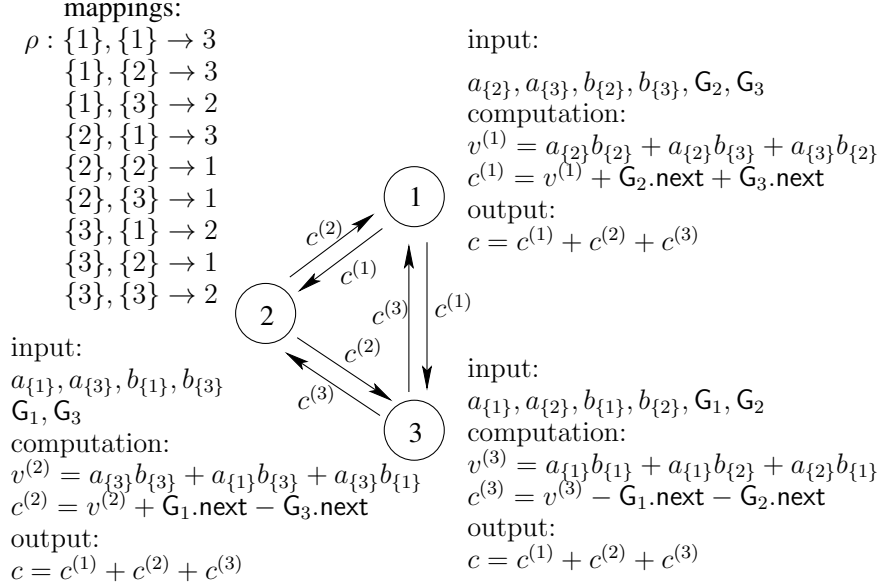


Figure 3: Sample three-party $\text{MulPub}([a], [b])$; arithmetic is in \mathcal{R} .

Executing multiplication followed by `Open` would double the number of rounds.

In multiplication, after computing a product, each locally computed value is no longer random and needs to be re-randomized prior to opening it. In our RSS setting, we realize this by relying on pseudo-random values locally computed by the parties. In particular, similar to other building blocks, we associate a secret key k_T with each $T \in \mathcal{T}$ (i.e., this is the same key shares used with `PRandR()` and multiplication) and use pseudo-random values $\mathbf{G}_T.\text{next}$ to protect the share of the product that each party locally computes, prior to that party revealing its randomized value to all others. To ensure that the product reconstructed by the parties is correct, we need to make sure that the sum of all blinding pseudo-random values equals to 0. In the three-party case, this can be accomplished by adding some pseudo-random values and subtracting others, as illustrated in Figure 3. With larger n and t we must be careful to draw new elements from each PRG to ensure that values released by different parties are protected using proper randomness without reusing them. This is similar to the logic used in multiplication. Then to realize this logic and ensure that all blinding factors add to 0, when multiple values are sampled from \mathbf{G}_T , the last blinding value is set to the sum of all previously drawn elements multiplied by -1 (in \mathcal{R}). We provide a detailed description of `MulPub` in Protocol 2. \mathbf{G}_T and S_p are defined as in multiplication.

In this algorithm, each party draws the same number of elements from each \mathbf{G}_T in its possession to ensure that after single execution of this algorithm all parties are in the same state (by any given party may discard some of the computed values). Similar to the computation in multiplication, we order the parties based on the values of their IDs. Because any given share T is stored at $t + 1$ parties, there are t calls to each \mathbf{G}_T per invocation of this operation. Then the participant with the lowest ID among the parties with access to T ($j = 0$) uses the first element of \mathbf{G}_T to protect its value $v^{(p)}$ and disregards the $t - 1$ other elements, the participant with the next lowest ID uses the second element, etc. The participant with the highest ID among those with access to T ($j = t$) computes the sum of all t elements drawn from \mathbf{G}_T and subtracts the sum from its $v^{(p)}$. Correctness follows from the fact that the sum of all blinding values over all parties and all shares is equal to 0, i.e., $c = \sum_p c^{(p)} = \sum_p v^{(p)}$ (in \mathcal{R}).

To show security, we prove the following result:

Protocol 2 $c \leftarrow \text{MulPub}([a], [b])$

```
// pre-distributed values are  $[k]$  and public map  $\rho$ 
1: each  $p \in [1, n]$  does the following:
2:    $v^{(p)} = c^{(p)} = \sum_{T_1, T_2 \in \mathcal{T}, \rho(T_1, T_2) = p} a_{T_1} b_{T_2}$  (in  $\mathcal{R}$ );
3:   for  $T \in S_p$  do
4:     let  $j$  be the number of parties  $p' < p$  such that  $p' \notin T$ ;
5:     for  $i = 0$  to  $t - 1$  do
6:        $z = G_T.\text{next}$ ;
7:       if  $j = t$  then  $c^{(p)} = c^{(p)} - z$  (in  $\mathcal{R}$ );
8:       else if  $i = j$  then  $c^{(p)} = c^{(p)} + z$  (in  $\mathcal{R}$ );
9:       end if
10:    end for
11:  end for
12:  send  $c^{(p)}$  to all other parties;
13:   $c = c^{(p)}$ ;
14:  for  $i = 1$  to  $n - 1$  do
15:    upon receiving  $c^{(p')}$  from distinct  $p'$ , set  $c = c + c^{(p')}$  (in  $\mathcal{R}$ );
16:  end for
17: return  $c$ ;
```

Theorem 2. *The protocol $\text{MulPub}([a], [b])$ is secure according to definition 1 in the (n, t) setting with $t = (n - 1)/2$ assuming PRG is a pseudo-random generator.*

Before proceeding with the proof, we demonstrate intuition behind it on the three-party example in Figure 3. Let z_T denote the output of $G_T.\text{next}$. Then party 1 transmits $c^{(1)} = v^{(1)} + z_{\{2\}} + z_{\{3\}}$, party 2 transmits $c^{(2)} = v^{(2)} + z_{\{1\}} - z_{\{3\}}$, and party 3 transmits $c^{(3)} = v^{(3)} - z_{\{1\}} - z_{\{2\}}$, where $c = v^{(1)} + v^{(2)} + v^{(3)}$ and each $v^{(i)}$ needs to be protected (arithmetic is in \mathcal{R}). Without loss of generality, let party 3 be corrupt. Then party 3 (with access to $z_{\{1\}}$ and $z_{\{2\}}$) can compute $v^{(1)} + z_{\{3\}}$, $v^{(2)} - z_{\{3\}}$, and the output of the computation c , but no information about $v^{(1)}$ or $v^{(2)}$ (assuming security of the PRG) other than their sum $v^{(1)} + v^{(2)}$. The latter, however, is already computable by party 3 using the output c and its share $v^{(3)}$, which reveals no extra information about a and b beyond their product. The full proof is given in Appendix B.

Similar to multiplication, MulPub can be generalized to evaluate any (multi-variate) polynomial of degree 2 and open the result.

4.4 Inputting Private Values

There will be a need to enter private values into the computation and we discuss the corresponding protocols in this section. We start with a general case when a participant who is not a computational party supplies their input into the computation and consequently discuss an optimized version when the input owner is one of the computational parties.

The input owner holds a private value a which will be represented as an element of ring \mathcal{R} . The input owner will need to generate replicated shares that correspond to a and send them to the computational parties. This will be the easiest way to proceed when there is only one element to share. However, when someone is sharing a vector of elements, we can save on communication by using pseudo-random shares. All shares except one for any element can be pseudo-random and computed locally by computational parties after obtaining a PRG seed. This means that among all

Protocol 3 $[a_1], \dots, [a_m] \leftarrow \text{Input}(a_1, \dots, a_m)$

```

1: for  $T \in \mathcal{T} \setminus \{T^*\}$  do
2:   input owner generates random  $k_T$  and sends it to each  $p \in T$ ;
3: end for
4: for  $i \in [1, m]$  do
5:   for  $T \in \mathcal{T} \setminus \{T^*\}$  do
6:     each  $p \notin T$  sets share  $a_{i,T} = \text{PRG}(k_T).\text{next}$ ;
7:   end for
8:   input owner computes  $a_{i,T^*} = a_i - \sum_{T \in \mathcal{T} \setminus \{T^*\}} \text{PRG}(k_T).\text{next}$  (in  $\mathcal{R}$ ) and sends it to  $p \notin T^*$ ;
9:   each  $p \notin T^*$  sets share  $a_{i,T^*}$  to the value received from input owner;
10: end for
11: return  $[a_1], \dots, [a_m]$ ;

```

shares $T \in \mathcal{T}$, one is marked as special and is denoted as T^* . The corresponding share is computed by the input owner and is communicated to all parties with access to that share. The construction is given as Protocol 3.

When the input owner is one of the computational parties, we can capitalize on the fact that the parties already have pre-distributed PRG seeds. We denote the input party as p^* . Note that p^* has access to a subset of the PRG seeds corresponding to the shares it is entitled to have access to, but not to all seeds. While we could generate new seeds for each T such that $p^* \in T$ and make it available to all $p \notin T$ and p^* , these seeds will be accessible to more than t parties and do not contribute to security. Therefore, we instead choose to set such shares to 0 and use only shares accessible to p^* . As a result, T^* will be such that $p^* \notin T^*$, the parties will set shares $a_T = \text{PRG}(k_T).\text{next}$ for each T such that $p \notin T$ and $T \neq T^*$, share T^* will be computed as $a_{T^*} = a - \sum_{T \text{ s.t. } p \notin T \wedge T \neq T^*} a_T$ (in \mathcal{R}) by p^* and communicated to all $p \notin T^*$, and all remaining shares a_T are set to 0.

All variants use a single round. When a single input is shared by an external party, the input owner simply generates all $\binom{n}{t}$ shares and communicates them to the computational parties (each share is stored by $t+1$ participants). This cost (which becomes sharing of a PRG seed) is amortized among all inputs when sharing multiple inputs. The additional cost per input for the input owner becomes generation $\binom{n}{t} - 1$ pseudorandom ring elements and communicating the last, computed share to $t+1$ computational parties, i.e., the total communication is $t+1$ ring elements. Each computational party needs to generate $\binom{n-1}{t}$ or $\binom{n-1}{t} - 1$ pseudo-random ring elements. When the input is shared by a computational party, there is no setup cost. The input owner need to generate $\binom{n-1}{t} - 1$ pseudo-random elements (i.e., similar to the number of shares it stores per shared value) and communicate the computed share to t other parties. Each other party computes $\binom{n-2}{t}$ (i.e., the number of shares it has in common with the data owner) or $\binom{n-2}{t} - 1$ pseudo-random ring elements. As will be relevant later, when a computational party is sharing a ring element in the (3,1) setting, the input owner communicates a single ring element to another party (and only one pseudo-random element is computed by the input owner and the remaining computational party).

Security can be shown as before (see Appendix B for the proof):

Theorem 3. *Input is secure according to definition 1 in the (n, t) setting with $t = (n-1)/2$ in the presence of secure communication channels and assuming PRG is a pseudo-random generator.*

5 Composite Protocols

While the previous operations can be instantiated to work with any finite ring, the techniques in this section work only in a ring \mathbb{Z}_{2^k} for some k . Ring \mathbb{Z}_{2^k} (for an appropriate choice of k) is the primary reason for supporting secure computation over rings because it allows us to utilize native CPU instructions for ring operations.

The goal of this work is to enable efficient general-purpose computation over rings \mathbb{Z}_{2^k} , we therefore focus on major building blocks which can be consequentially used to compose a protocol for arbitrary functionalities including machine learning tasks. Of central importance to this effort is the development of comparison protocols (for both less-than comparison and equality testing functionalities), which are known to be difficult to design in a framework where the elementary techniques are based on arithmetic gates. Others include bit decomposition and truncation (i.e., division by a power of 2). Combined together, these techniques can enable Boolean, integer, fixed-point, and even floating-point arithmetic, as well as array and related operations, giving the ability to compose general-purpose protocols.

Because a number of protocols for common operations over \mathbb{Z}_{2^k} have already been developed, some of the constructions that we mention in this sections are adaptations of prior protocols to our setting and we defer their specification to the appendix. In particular, Appendix A provides specification of random bit generation protocol, **RandBit**, that produces a bit shared in \mathbb{Z}_{2^k} and a more recent version from [31], **edaBit**, that generates a number (k in our case) of random bits r_i shared in \mathbb{Z}_2 together with a representation of the bits as an integer $r = \sum_{i=1}^k 2^i r_i$ shared in \mathbb{Z}_{2^k} . The former can be computed in a single round, while the latter uses noticeably lower communication per bit, but the round complexity is logarithmic in k and t .

We also describe a comparison algorithm for computing $[a] \leq [b]$, which is commonly implemented by determining the most significant bit of the difference between a and b and denoted by **MSB**. Performance of these protocols is summarized in Tables 2 and 11.

Truncation is a necessary building block when working with fixed-point values or simulating fixed-point computation using integer arithmetic and permits us to minimize the ring size. Starting from [16], probabilistic truncation of input a by m bits that produces $\lfloor a/2^m \rfloor + u$, where u is a bit, is significantly faster than precise truncation that rounds down. It is biased towards rounding to the nearest integer to $a/2^m$ and is sufficient for our purpose. The protocol we present, **TruncPr**($[a], m$), is a constant-round solution that combines the approach from [?, 21] with **edaBits** from [31] and inherits from [31] the requirement that input a is 1 bit shorter than the ring size, i.e., $\text{MSB}(a) = 0$. We use notation $[x]_\ell$ to denote that SS is over \mathbb{Z}_{2^ℓ} .

The truncation protocol, given as Protocol 4, uses related random values r and \hat{r} , bit decomposition of which are known, where $r = \sum_{i=0}^{k-1} 2^i r_i$ is a full-size random value and $\hat{r} = \sum_{i=m}^{k-1} 2^i r_i$ is the portion remaining after truncating m bits. We thus modify the **edaBit** protocol to produce those values simultaneously. Each $[r]$ and $[\hat{r}]$ is computed as a sum of $t + 1$ integers, so we must compensate for two types of carries: (i) addition of m least significant bits in r will produce carry bits into the next bits which are not accounted for in \hat{r} and (ii) while the carry bits past the k bits are automatically removed in the ring when computing r , these bits remain in \hat{r} due to its shorter length. Because we compute the bitwise representation of r using bitwise addition protocol **BitAdd**, we can also extract the carry bit into any desired position which is already computed during the addition. The logic of the truncation protocol necessitates the removal of the $(k - 1)$ th bit. For this reason, we capture carries into the m th and $(k - 1)$ th positions and denote those bits from the i th call to **BitAdd** as $\text{cr}_{i,m}$ and $\text{cr}_{i,k-1}$, respectively (line 10). We subsequently convert the $2 \log(t + 1)$ carry bits and the most significant bit of r , denoted as b_{k-1} , from shares over \mathbb{Z}_2 to \mathbb{Z}_{2^k} using binary-to-arithmetic sharing protocol **B2A** (from [23]). All interactive operations except the

Protocol 4 $[a/2^m]_k \leftarrow \text{TruncPr}([a]_k, m)$

```
1: for  $p = 1, \dots, t + 1$  in parallel do
2:   party  $p$  samples  $r_0^{(p)}, \dots, r_{k-1}^{(p)} \in \mathbb{Z}_2$  and computes  $r^{(p)} = \sum_{j=0}^{k-1} r_j^{(p)} 2^j$  and  $\hat{r}^{(p)} = \sum_{j=m}^{k-1} r_j^{(p)} 2^j$ ;
3:   simultaneously execute  $[r^{(p)}]_k \leftarrow \text{Input}(r^{(p)}, k)$ ,  $[\hat{r}^{(p)}]_k \leftarrow \text{Input}(\hat{r}^{(p)}, k)$ , and
    $[r_i^{(p)}]_1 \leftarrow \text{Input}(r_i^{(p)}, 1)$  for  $i=1, \dots, k$  with  $p$  being the input owner;
4: end for
5:  $[r]_k = \sum_{p=1}^{t+1} [r^{(p)}]_k$ ;  $[\hat{r}]_k = \sum_{p=1}^{t+1} [\hat{r}^{(p)}]_k$ ;
6:  $s = t + 1$ ;
7: for  $i = 1, \dots, \lceil \log(t + 1) \rceil$  do
8:   for  $j = 1, \dots, \lfloor s/2 \rfloor$  in parallel do
9:      $\ell = j + s \cdot (i - 1)$ ;
10:     $\langle [r_1^{(j)}]_1, \dots, [r_{k-1}^{(j)}]_1 \rangle, [\text{cr}_{\ell, m-1}]_1, [\text{cr}_{\ell, k}]_1$ 
     $\leftarrow \text{BitAdd}(\langle [r_1^{(2j-1)}]_1, \dots, [r_{k-1}^{(2j-1)}]_1 \rangle,$ 
       $\langle [r_1^{(2j)}]_1, \dots, [r_{k-1}^{(2j)}]_1 \rangle)$ ;
11:    if  $s \bmod 2 = 0$  then  $s = s/2$ ;
    else
12:     $\langle [r_1^{(\frac{s+1}{2})}]_1, \dots, [r_{k-1}^{(\frac{s+1}{2})}]_1 \rangle = \langle [r_1^{(s)}]_1, \dots, [r_{k-1}^{(s)}]_1 \rangle$ ;
13:    end if
14:  end for
15: end for
16: end for
17:  $[b_0]_1, \dots, [b_{k-1}]_1 = [r_0^{(1)}]_1, \dots, [r_{k-1}^{(1)}]_1$ ;
18:  $[b_{k-1}]_k, \langle [\text{cr}_{\ell, m}]_k \rangle, \langle [\text{cr}_{\ell, k-1}]_k \rangle$ 
    $\leftarrow \text{B2A}([b_{k-1}]_1, \langle [\text{cr}_{\ell, m}]_1 \rangle, \langle [\text{cr}_{\ell, k-1}]_1 \rangle)$  for  $\ell=1, \dots, t$ ;
19:  $[\hat{r}]_k = [\hat{r}]_k - [b_{k-1}]_k \cdot 2^{k-m-1} + \sum_{\ell=1}^t ([\text{cr}_{\ell, m}]_k - [\text{cr}_{\ell, k-1}]_k \cdot 2^{k-m-1})$ ;
20:  $c \leftarrow \text{Open}([a]_k + [r]_k)$ ;
21:  $c' = (c/2^m) \bmod 2^{k-m-1}$ ;
22:  $[b]_k = (c/2^{k-1}) + [b_{k-1}]_k - 2(c/2^{k-1})[b_{k-1}]_k$ ;
23: return  $c' - [\hat{r}]_k + [b]_k \cdot 2^{k-m-1}$ ;
```

last one (line 20) can be precomputed. Security follows from the protocol logic as specified in prior work and from security of the building blocks.

It is also possible to use the above protocol to truncate an input $[a]$ by a private number of bits $[m]$ as outlined in [21]: Let M be some public upper bound on m . Protocol $\text{TruncPriv}([a], [m], M)$ then needs to securely compute $[2^{M-m}] \cdot [a]$ and can subsequently call $\text{TruncPr}([2^{M-m} \cdot a], M)$. A performance summary is given in Table 2.

6 Neural Network Applications

Today it is typical to benchmark secure multi-party frameworks on machine learning applications, e.g., NN inference. We briefly introduce NN basics and describe two mechanisms for improving efficiency of secure NN inference.

A *neural network* is a series of interconnected layers consisting of neurons. Each neuron has an associated weight and bias used for computation on some input data and outputs a prediction based on that data. A NN network layer takes the form $\mathbf{y} = g(\mathbf{x}\mathbf{W} + \mathbf{b})$, where \mathbf{x} is the input vector from the previous layer, \mathbf{W} is the weight tensor, \mathbf{b} is the bias vector, and g is some activation function.

Protocol	Rand. Protocol	Rounds	Communication
MSB($[a]_k$)	RandBit	$\log(k - 1) + 3$	$2t(k + 3)$
	edaBit	$\log(t + 1)(\log(k) + 1) + \log(k - 1) + 4$	$t^2(\log(k) + 1) + 7t + 1/2$
TruncPr($[a]_k, m$)	RandBit	2	$t(2k + 1)$
	edaBit	$\log(t + 1)(\log(k) + 1) + 4$	$t^2(\log(k) + 2/k + 4) + t(1/k + 4) + 1/2$
Convert($[a]_k, k, k'$)	RandBit	$\log(k) + 4$	$2t(k + k') + t(\log(k) + 2)$
	edaBit	$\log(t + 1)(\log(k) + 1) + \log(k) + 3$	$t^2(\log(k) + 1) + t(2k' + \log(k) + 3)$

Table 2: Performance of composite protocols with communication measured in the number of ring elements sent per party over $\mathbb{Z}_{2^{k+2}}$ for RandBit and \mathbb{Z}_{2^k} for edaBit(k).

Protocol 5 $[a]_{k'} \leftarrow \text{Convert}([a]_k, k, k')$, where $k' > k$

- 1: $[r]_k, [r_0]_1, \dots, [r_{k-1}]_1 \leftarrow \text{edaBit}(k)$;
 - 2: $c \leftarrow \text{Open}([a]_k - [r]_k)$;
 - 3: $[a_0]_1, \dots, [a_{k-1}]_1 \leftarrow \text{BitAdd}(c, [r_0]_1, \dots, [r_{k-1}]_1)$;
 - 4: for $i = 0$ to $k - 1$ in parallel $[a_i]_{k'} \leftarrow \text{B2A}([a_i]_1, k')$;
 - 5: **return** $[a]_{k'} = \sum_{i=0}^{k-1} [a_i]_{k'} 2^i$;
-

Sample activation functions are Rectified Linear Unit (ReLU), which on input $\mathbf{x} = (x_1, \dots, x_N)$ computes $\mathbf{y} = (y_1, \dots, y_N)$ where each $y_i = \max(0, x_i)$, and its variant ReLU6 which computes $y_i = \min(\max(0, x_i), 6)$.

6.1 Share Conversion

Conventional NN evaluation uses floating-point arithmetic, while secure evaluations for performance reasons typically employ fixed-point computation or emulate it on integers. If inputs are represented in the form of fixed-length integers, the values will grow with each layer that performs matrix multiplication. This can impact on performance because comparison-based activation and pooling operations have cost linear in the bitlength of ring elements. For this reason, it can be advantageous to start with a smaller ring size and increase it mid-computation to accommodate longer values.

This approach involves converting secret-shared $[a]_k$ over \mathbb{Z}_{2^k} to a different representation $[a]_{k'}$ over $\mathbb{Z}_{2^{k'}}$ for $k' > k$. Conversion techniques between certain types of fields are known [27], but they do not apply to our case. Simply casting k -bit shares to k' -bit shares for $k' > k$ affects correctness because the overflow due to share addition is not reduced modulo 2^k . Thus, the task is to leave k least significant bits of the value and erase the remaining bits in a longer share representation. One way to achieve this is to invoke truncation as $([a] \cdot 2^{k'-k}) \gg 2^{k'-k}$ or $[a] - ([a] \gg k) 2^k$. However, because computing precise truncation is costlier for rings than fields, we design a more efficient version based on bit decomposition. In particular, we perform bit decomposition of $[a]_k$ into shares of bits in \mathbb{Z}_2 , convert the bit shares to $\mathbb{Z}_{2^{k'}}$, and reassemble $[a]_{k'}$.

This procedure is denoted by Convert and given as Protocol 5 using edaBits. An equivalent version can be constructed using RandBit. It is based on bit decomposition from [23] and uses Boolean to arithmetic conversion, B2A, from \mathbb{Z}_2 to $\mathbb{Z}_{2^{k'}}$ and bitwise integer addition, BitAdd. Performance is summarized in Table 2.

6.2 Quantized Neural Networks

To improve efficiency of NN inference, it is common to use quantized weights and activation values, which makes the resulting models suitable for deployment in constrained environments and is a

well-studied field (see, e.g., [34]). We outline the standard quantization approach from [37] and its privacy-preserving realization from [21] for quantized TFLite models and consequently describe our optimizations.

For a vector \mathbf{x} , each real-valued x_i is represented as $x_i = m(\bar{x}_i - z)$, where $m \in \mathbb{R}$ is the scale and z and \bar{x}_i are 8-bit integers with z being the zero point. Given an input column vector $\mathbf{x} = (x_1, \dots, x_N)$ and a row vector $\mathbf{w} = (w_1, \dots, w_N)$ of \mathbf{W} with quantization parameters (m_1, z_1) and (m_2, z_2) , respectively, the dot product of \mathbf{x} and \mathbf{w} , $y = \sum_{i=1}^N x_i w_i$, is specified with quantization parameters (m_3, z_3) . Since $y \approx m_3 \cdot (\bar{y} - z_3)$, $x_i \approx m_1 \cdot (\bar{x}_i - z_1)$, and $w_i \approx m_2 \cdot (\bar{w}_i - z_2)$, quantized \bar{y} is computed as

$$\bar{y} \approx z_3 + m_1 m_2 / m_3 \cdot \sum_{i=1}^N (\bar{x}_i + z_1) \cdot (\bar{w}_i - z_2) = z_3 + m \cdot s.$$

Computing s requires integer-only arithmetic and is guaranteed to fit in $16 + \log N$ bits. The scale $m = m_1 m_2 / m_3$ is a small real number. It can be written as $m = 2^{-e} m'$ with normalized $m' \in [0.5, 1)$ which informs the value of e and represented as a 32-bit integer m'' , where $m' \approx 2^{-31} m''$.

Two-dimensional convolutions typically add a quantized bias \bar{b} once the dot product is computed. This is handled by setting the scale of the bias to $m_1 m_2$ and the zero-point to 0, such that the the bias can be added to s prior to scaling. The last step of a convolution layer is to apply an activation function such as ReLU6. In a quantized NN, this functions as a clamping operation which eliminates values outside of range $[0, 255]$ and uses $m_3 = 6/255$ and $z_3 = 0$. This guarantees correct range while maximizing precision with 8-bit quantized values. Going forward, m_3 becomes m_1 for the next layer and thus all intermediate layers share the same $m_1 = m_3 = 6/255$. Other activation functions such as sigmoid would be handled differently, but we only consider clamping-based functions since they are often sufficient.

To compute a convolution layer securely, the model owner needs to enter private quantization parameters into the computation. This includes all zero points z_i , modified scale m'' , and integer scale adjustment 2^{M-e-31} , where M is an upper bound set to 63. After privately computing the dot product $[s] + [\bar{b}]$, the result is multiplied by $[m'']$ and need to be truncated by private amount $31 + e$. The truncation is accomplished by multiplying the scaled dot product by $[2^{M-n-31}]$ and $[m \cdot s]$ and consequently truncating by M bits. Lastly, after adding $[z_3]$ locally, clamping the result to the interval $[0, 255]$ is performed using two comparisons.

A limitation of [21]’s approach is that it required large scaling factors and consequently a large ring size of $k = 72$ for working with real numbers, using M -bit truncation with $M = 63$. We propose a modified approach where we fold the scales into other aspects of the layer computation and conduct smaller truncation at the end of each layer, which guarantees compact representation of intermediate results.

Let superscript $\langle i \rangle$ denote the layer number. Starting from layer 0, the entire layer computation (dot product, scaling, and clamping) can be interpreted as computing $0 \leq \bar{y}^{(0)} \leq 255$, where

$$\bar{y}^{(0)} = \frac{m_1^{(0)} m_2^{(0)}}{m_3^{(0)}} \cdot \left(\sum_{i=1}^N (\bar{x}_i^{(0)} - z_1^{(0)}) \cdot (\bar{w}_i^{(0)} - z_2^{(0)}) \right) + \bar{b}^{(0)},$$

and $z_3^{(i)}$ was observed to be 0 for all layers except the last one. Because $m_3^{(0)} = 6/255$, we can scale the equation to redefine $\bar{y}^{(0)}$ as

$$\bar{y}^{(0)} = \sum_{i=1}^N (\bar{x}_i^{(0)} - z_1^{(0)}) \cdot (\bar{w}_i^{(0)} - z_2^{(0)}) + \bar{b}^{(0)},$$

where $0 \leq \bar{y}^{(0)} \leq 6/m_1^{(0)}m_2^{(0)}$. Now, our clamping operation can use these bounds, with the upper bound being privately entered by the model owner to avoid division. As before, the output of this layer becomes the input for the subsequent layer, i.e., $\bar{x}^{(i)} = \bar{y}^{(i-1)}$. Our modified incoming vector, denoted $\hat{x}^{(1)}$, is coupled with an additional scaling factor of $(255m_1^{(0)}m_2^{(0)})/6$, such that

$$\bar{x}^{(1)} = 255m_1^{(0)}m_2^{(0)}\hat{x}^{(1)}/6 = \delta^{(1)}\hat{x}^{(1)}.$$

Using $\bar{x}^{(1)} = \delta^{(1)}\hat{x}^{(1)}$ gives us

$$\bar{y}^{(1)} = \left(\sum_{i=1}^N (\hat{x}_i^{(1)} - z_1^{(1)}/\delta^{(1)}) \cdot (\bar{w}_i^{(1)} - z_2^{(1)}) \right) + \bar{b}^{(1)}/\delta^{(1)}$$

with $0 \leq \bar{y}^{(1)} \leq 6/(\delta^{(1)}m_1^{(1)}m_2^{(1)})$. This expression can be evaluated securely without needing fixed-point multiplication or large truncation, and all bounds are computed by the model owner prior to privately entering them in the computation.

Evaluating subsequent layers in this fashion causes the outputs to grow by factor $\delta^{(i+1)}$, which can be computed as

$$\delta^{(i+1)} = \delta^{(i)}255m_1^{(i)}m_2^{(i)}/6$$

with $\delta^{(0)} = 1$. However, we can ensure values remain small by truncating the output $\bar{y}^{(i+1)}$ by $\ell^{(i)}$ bits. With the right choice of $\ell^{(i)}$ we are able to maintain the necessary accuracy, and the value of $\delta^{(i+1)}$ consequently becomes

$$\delta^{(i+1)} = \delta^{(i)} \cdot 255m_1^{(i)}m_2^{(i)}/(6 \cdot 2^{\ell^{(i)}}).$$

The maximum number of bits we can truncate in a layer needs to comply with constraint

$$\delta^{(i)} \cdot 255m_1^{(i)}m_2^{(i)}/(6 \cdot 2^{\ell^{(i)}}) \geq 1,$$

which leads to $\ell^{(i)} \leq \left\lfloor \log_2(255\delta^{(i)}m_1^{(i)}m_2^{(i)}/6) \right\rfloor$. Once again, these values are independent of the input data and become a part of the model. We thus can use TruncPriv outlined in Section 5 for truncation by a private amount. The net result is that we are able to use a significantly smaller bound M and consequently substantially shorter ring size k . In practice, the coefficients introduced in our methodology can reasonably be folded into the scaling factors m themselves.

Other layers such as average pooling can be approximated by substituting the division by some integer d with truncation by $\lfloor \log d \rfloor$ bits, and softmax can be replaced with argmax when computing the final prediction. These changes can slightly impact the scaling factors, but have no impact on the accuracy.

7 Performance Evaluation

We implemented the protocols described in this work and evaluate their performance. We run both micro-benchmarks to evaluate the individual operations as well as offer evaluation of machine learning applications.

The implementation was done in C++. We use AES from the OpenSSL cryptographic library [1] to instantiate the PRF and also to implement secure communication channels between each pair of

	Protocol	Batch Size						
		1	10	10 ²	10 ³	10 ⁴	10 ⁵	10 ⁶
LAN	Field (30 bits)	0.113	0.165	0.426	1.27	10.5	105	1,118
	Field (60 bits)	0.117	0.167	0.436	1.28	11.8	120	1,271
	Ring (30 bits)	0.083	0.084	0.113	0.188	1.16	10.1	105
	Ring (60 bits)	0.086	0.086	0.109	0.251	2.11	19.7	192
WAN	Field (30 bits)	14.8	14.9	15.1	17.4	33.6	256	2,517
	Field (60 bits)	14.9	14.9	15.2	17.8	50.7	335	3,752
	Ring (30 bits)	14.4	14.4	14.4	14.8	16.9	103	928
	Ring (60 bits)	14.4	14.4	14.4	15.1	32.1	184	1,739

Table 3: Performance of multiplication protocols in milliseconds.

	Protocol	Matrix Size			
		10 × 10	100 × 100	500 × 500	1000 × 1000
LAN	Field (30 bits)	0.556	191	18,219	1.68·10 ⁶
	Field (60 bits)	0.547	207	19,138	1.75·10 ⁶
	Ring (30 bits)	0.0987	2.45	284	3,682
	Ring (60 bits)	0.510	3.87	333	4,903
WAN	Field (30 bits)	22.7	249	18,387	–
	Field (60 bits)	22.9	365	19,285	–
	Ring (30 bits)	14.5	20.1	634	17,675
	Ring (60 bits)	14.7	34.3	917	18,915

Table 4: Performance of matrix multiplication in milliseconds.

the computational parties. We report the average execution time of 1000 executions for the micro-benchmark experiments and the average time of 5 executions for the application experiments. The runtimes are also averaged across the computation parties.

All experiments were run in a three-party setting. For LAN experiments, we used three 8-core 2.1GHz machines with 64GB of RAM. They were connected via a 1 Gbps Ethernet link with one-way latency of 0.15ms. For WAN benchmarks, we used two of the machines above and one remote 2.4GHz machine. One-way latency between the remote and local machines was 23ms. Although the machine configurations are slightly different, this should not introduce inconsistencies in the experiments because our protocols are interactive and symmetric, and the computation time depends on the communication links and performance of the slowest machine. All experiments use a single core.

7.1 Micro-benchmarks

In this section we report performance of individual operations such as multiplication, matrix multiplication, random bit generation (RandBit and edaBit) and comparison (MSB). The experiments used two bitlengths, $k = 30$ and $k = 60$, which allows us to use the `unsigned` and `unsigned long` integer types, respectively, to implement ring operations.

Tables 3 and 4 report performance of multiplication and matrix multiplication, respectively. As we strive to measure performance improvement when we switch computation from a field to a ring, we compare performance of our protocols to those using Shamir SS in the same setting (i.e., semi-honest security with honest majority) using PICCO implementation [65] with recent improvements from [11]. The field size is set to accommodate 30- and 60-bit integers. Batch size denotes how many operations were executed at the same time in a single batch.

	Protocol	Batch Size						
		1	10	10^2	10^3	10^4	10^5	10^6
LAN	Field (30 bits)	0.143	0.203	0.787	5.34	48.7	456	4,694
	Field (60 bits)	0.144	0.228	0.95	6.34	57.1	569	6,012
	Ring (30 bits)	0.103	0.105	0.144	0.404	3.20	32.1	337
	Ring (60 bits)	0.100	0.109	0.170	0.723	6.57	66.2	656
WAN	Field (30 bits)	21.8	21.9	22.9	32.9	113	1,056	11,089
	Field (60 bits)	21.8	21.9	23.4	35.4	126	1,414	15,870
	Ring (30 bits)	21.8	21.8	21.8	22.2	26.8	170	1,534
	Ring (60 bits)	21.8	21.8	22.6	22.7	53.3	308	3,040

Table 5: Performance of RandBit protocols in milliseconds.

	Protocol	Batch Size						
		1	10	10^2	10^3	10^4	10^5	10^6
LAN	Ours (30 bits)	0.60	0.65	1.23	3.53	19.0	179	1,762
	Ours (60 bits)	0.73	1.08	1.35	6.33	38.7	484	4,173
	MP-SPDZ (32 bits)	21.6	19.2	20.1	20.9	24.1	194	1,884
	MP-SPDZ (64 bits)	28.0	29.0	28.8	29.0	35.2	287	2,736
WAN	Ours (30 bits)	102	102	102	105	148	935	8,751
	Ours (60 bits)	117	117	117	123	259	1,881	18,688
	MP-SPDZ (32 bits)	871	878	885	889	907	8,762	87,485
	MP-SPDZ (64 bits)	1,331	1,339	1,352	1,365	1,382	13,219	131,350

Table 6: Performance of edaBit protocols in milliseconds, compared to MP-SPDZ [3].

We observe that ring realization of multiplication in Table 3 is up to 10 times faster on the LAN for a sufficiently large batch size compared to field, despite using the same amount of communication and the need to repeat the computation twice (once for each share) with RSS. This indicates that using native CPU instructions for secure arithmetic has remarkable advantage. On WAN, as expected, performance is heavily dominated by the network latency for small batches and we gain up to 2.7 times performance improvement for large batch sizes.

Matrix multiplication in Table 4 is performed in a single round using the necessary number of dot-products. Because local work is the bottleneck, we see performance improvement by up to a factor of 450 when we switch to a ring! Furthermore, performance improvement is well-pronounced even for matrices of small size, as we have a 5.6-fold improvement for 10×10 matrices on the LAN. (The largest matrix size could not be handled by the remote machine, but the times on LAN and WAN are expected to be similar.)

Tables 5 and 6 provide random bit generation results. To support k -bit integers, ring-based RandBit requires a $(k + 2)$ -bit ring. Field-based RandBit from [16] does not increase the field size; however, all uses of RandBit we are aware of are for operations such as comparisons that utilize statistical hiding and, as a result, increase the field size by a statistical security parameter κ (typically set to 48 in implementations). For this reason, our field-based RandBit and MSB benchmarks utilize 79- and 109-bit fields. Both versions of RandBit in Table 5 communicate the same number of field or ring elements; however, the performance gain of the ring version grows as we increase the batch size, reaching 15-fold improvement and indicating that local field-based computation is the bottleneck. This is in large part due to the need to perform modulo exponentiation (see [16]).

The concept of edaBit is recent and for that reason in Table 6 we compare our implementation to that reported in the original publication [31], available through MP-SPDZ repository [3]. Note

	Protocol	Batch Size						
		1	10	10 ²	10 ³	10 ⁴	10 ⁵	10 ⁶
LAN	Field (30 bits)	1.43	3.77	22.5	211	2,351	24,349	256,785
	Field (60 bits)	2.19	6.52	52.2	506	5,333	54,286	572,439
	Ring (RandBit) (30 bits)	0.73	0.88	2.53	11.9	99.7	1,140	10,492
	Ring (RandBit) (60 bits)	0.87	1.31	5.70	44.5	416	4,039	38,373
	Ring (edaBit) (30 bits)	1.29	1.39	1.85	4.55	29.7	276	2,804
	Ring (edaBit) (60 bits)	1.49	1.72	2.43	9.03	60.9	672	5,989
	MP-SPDZ (32 bits)	27.7	31.9	31.6	30.0	37.3	268	3,113
	MP-SPDZ (64 bits)	34.3	38.1	36.4	39.0	44.6	366	3,769
WAN	Field (30 bits)	153	158	198	615	5,265	54,798	573,482
	Field (60 bits)	155	165	287	1,390	12,513	132,874	1,389,257
	Ring (RandBit) (30 bits)	131	131	133	169	616	5,633	53,040
	Ring (RandBit) (60 bits)	131	131	138	314	2,008	19,653	190,160
	Ring (edaBit) (30 bits)	218	218	218	223	311	1,750	16,114
	Ring (edaBit) (60 bits)	261	261	262	271	489	3,372	31,540
	MP-SPDZ (32 bits)	1,204	1,206	1,206	1,310	1,445	10,350	88,930
	MP-SPDZ (64 bits)	1,678	1,679	1,760	1,725	2,067	16,045	151,570

Table 7: Performance of MSB protocols in milliseconds.

that each edaBit corresponds to generating k random bits together with the corresponding k -bit random integer. It is clear from the table that MP-SPDZ’s implementation is optimized for large sizes and fast networks. In particular, it gives comparable runtime for batches of size 1 and 1,000. Furthermore, our implementation is several times faster on a WAN for all sizes and becomes slightly slower on LAN for the largest size. Note that the times we measured for MP-SPDZ are very different from those originally provided in [31], which reported the ability to generate 7.18 million 64-bit edaBits per second. This is 20 times faster than the fastest time per operation we record and stems from the differences in hardware. In particular, experiments in [31] were run multi-threaded on powerful AWS c5.9xlarge instances with 36 cores and a 10 Gbps link. This difference highlights the need to reproduce experiments on similar hardware to draw meaningful comparisons about performance of different algorithms.

Table 7 reports performance of multiple MSB protocols: field-based implementation from PICCO, our ring implementations using RandBit and using edaBit, and edaBit-based implementation from MP-SPDZ. The gap between the first two shows the gain after switching from field-based to ring-based arithmetic. Both of them make a linear in k number of calls to RandBit, but our implementation executes BitLT over \mathbb{Z}_2 , while field-based uses a fixed field for all operations. As a result, our ring RandBit-based MSB is up to 24.5 times faster than the field version on LAN and up to 10.8 on WAN. If we compare our RandBit and edaBit MSB implementations, the advantage of the edaBit version is well pronounced starting from batch sizes of 100 on LAN and 1000 on WAN. This is due to its substantially lower communication, but higher round complexity makes it slower for small sizes. MP-SPDZ’s edaBit-based implementation in the same setting generally took longer to run than our edaBit-based implementation, especially on WAN, but the gap between smaller and larger ring sizes is smaller.

We also visualize time per operation with variable batch sizes on LAN in Figure 4. Multiplication and RandBit sub-figures compare ring vs. field protocols, indicating a substantial gap as expected; edaBit sub-figure compares our and MP-SPDZ implementations in the same setting; and MSB sub-figure compares RandBit and edaBit variants.

It is also informative to compare our field vs. ring results with those of SPDZ. While SPDZ [26]

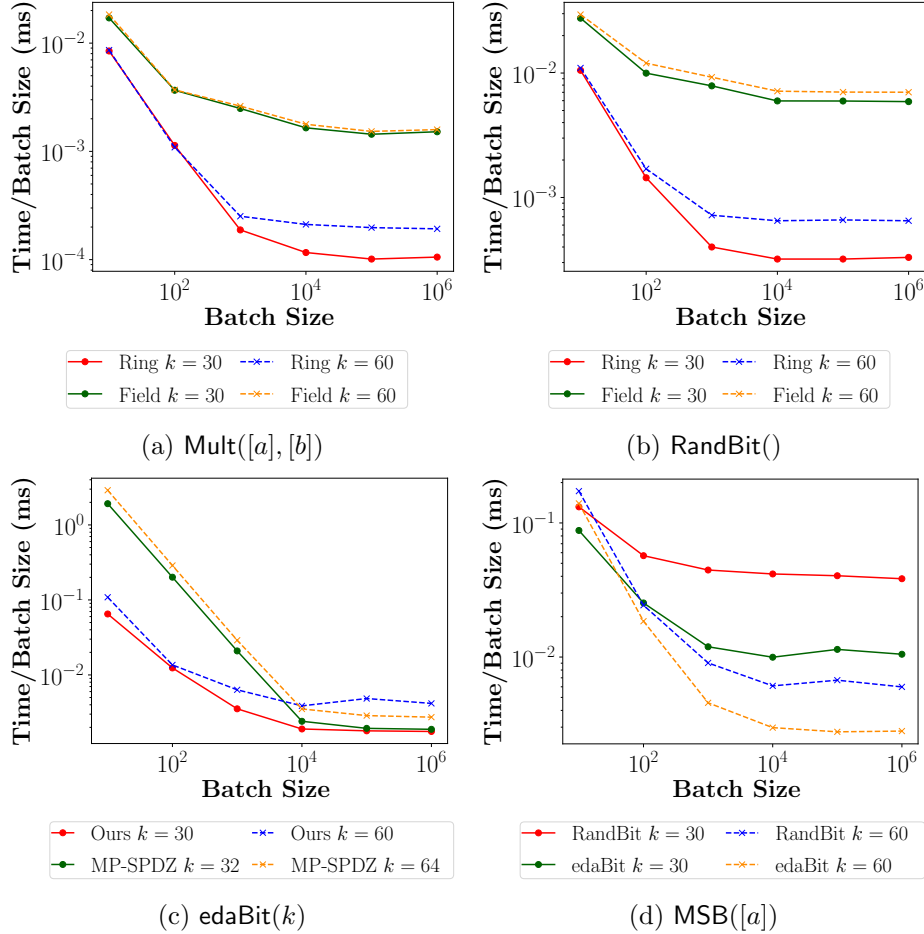


Figure 4: LAN microbenchmarks results.

and its ring version $\text{SPDZ}_{\mathbb{Z}_2^k}$ [19, 23] use a much stronger adversarial model and different type of SS, we would like to know whether similar savings are achievable in different settings. [23] reports that performance improved by a factor of 4.6–4.9 for multiplication and by a factor of 5.2–6.0 for RandBit-based comparison on a 1Gbps LAN. The results are only provided as throughput improvement and do not report different batch sizes. In our experiments we observed greater improvements, up to 10 times for multiplication and up to 24.5 improvement for MSB. This may be explained by the fact that our techniques are more lightweight and perhaps switching to faster arithmetic makes less of an impact in the SPDZ setting.

7.2 Machine Learning Applications

We next evaluate our protocols on machine learning applications and show that they exhibit good performance. We consider NNs, quantized NNs, and support vector machine (SVM) evaluation, in part to facilitate comparison to prior work.

7.2.1 Neural Networks

There are many different types of NNs, and for our standard benchmarking we chose the NN from MiniONN [46] for the MNIST dataset [43] (Figure 12 in [46], Network B in [62], and Network C in

	Field	MiniONN	GAZELLE	SecureNN	FALCON [44]	FALCON [63]	[21]	Ours			
Batch Size	1	1	1	1	1	1	1	1	5	10	50
Latency	1798	9320	810*	130*	820	42*	67*, 18*	99.8	77.4	74.4	74.6

Table 8: Performance of MNIST NN prediction in milliseconds. (*) denotes experiments on more powerful hardware.

[63]), because it is a popular choice for evaluating privacy-preserving NN inference. The MNIST NN evaluation (given as Protocol 9 in Appendix A for completeness) uses convolution, fully-connected layers, an ReLU activation function, and max pooling of a window 2×2 to compute the maximum element in that window.

We use MiniONN’s implementation choices and, in particular, run the computation on integer inputs. To avoid using floating-point arithmetic, [46] scaled inputs by a factor of 1000 and rounded to the nearest integer. To compensate for the bitlength of the intermediate results growing with each multiplication, [46]’s implementation ran the computation using a 37-bit modulus and avoided the use of truncation. However, we determined that this size is too small and 49 bits are needed to correctly evaluate the model, which we subsequently use. Our implementation achieves the same 99.0% precision as reported for this model in [47] (which corrects [46]).

While it is possible to perform the entire computation in $\mathbb{Z}_{2^{49}}$, we observe that the initial steps are of the largest size and use significantly shorter integers than 49 bits. Because the cost of comparisons is linear in the bitlength of the ring elements, we can substantially improve performance by starting computation on shorter values and converting the intermediate results to a larger ring prior to multiplication that increases the size of the intermediate results. Therefore, we start computation with 20-bit integers and increase the ring size by 10 bits at different points of the computation(see Protocol 9 for detail).

Performance of MNIST NN inference with three parties (total time) is presented in Table 8. We also ran the same computation over a field (using [65, 11]), which required a 89-bit modulus. To closely mimic our ring-based implementation, this implementation computes with integers of increasing sizes, but uses the same modulus throughout the computation. We also include runtimes of two-party MiniONN [46], two-party Gazelle [38], two-party FALCON [44], SecureNN with custom three-party arithmetic [62], three-party FALCON [63], and three-party Dalskov et al. [21] with two types of truncation (`TruncPr` and `TruncPrSp`, respectively). All runtimes except our and field implementations were taken from the respective publications and many of them do not result in an accurate comparison due to the differences in the computing environment and implementations. In particular, SecureNN and GAZELLE used more powerful hardware (AWS `c4.8xlarge`), the use of multi-threading was not specified; FALCON [63] uses a similar configuration and multi-threads all data-independent computation, which constitutes the bulk of the work; Dalskov et al. [21] used even more powerful setup (AWS `c5.9xlarge`, multi-threaded, 10Gps connection), which in the context of `edaBit` generation was more than an order of magnitude faster environment than ours. This makes our solution attractive compared to the state of the art and the rapid progress with PPML makes it clear that results a couple to a few years old become obsolete.

Several other publications benchmarked NN predictions [52, 58, 6, 18, 17, 55, 51, 42, 57]. However, because they do not support or do not run MiniONN’s MNIST NN evaluation, we cannot directly compare our performance. For example, while ABY³ [52] is said to use MiniONN’s MNIST NN, evaluation is actually based on a different, simpler model used in Chameleon [58].

α	Ours				MP-SPDZ, [21]				
	0.25	0.5	0.75	1.0	0.25	0.5	0.75	1.0	
ρ	128	3.17	6.36	9.91	13.4	3.77	8.22	12.7	16.5
	160	5.04	10.2	15.4	20.6	5.26	11.1	17.5	25.1
	192	6.80	14.8	22.4	30.3	6.82	15.6	24.4	34.2
	224	9.88	20.0	30.2	40.8	9.92	20.0	30.6	45.2

Table 9: Performance of quantized MobileNets prediction in seconds.

Protocol	Integer size	Batch Size					
		1	5	10	50	100	250
Ours	30	15.0	4.80	3.39	2.08	1.77	1.59
	60	19.8	7.03	5.55	3.60	3.31	3.20
SPD \mathbb{Z}_{2^k} [23]	32	242 + 3055	162 + 3055	–	–	–	–
	64	362 + 10006	244 + 10006	–	–	–	–
Field	30	276	265	–	–	–	–
	60	635	628	–	–	–	–

Table 10: Performance of ALOI SVM classification in milliseconds.

7.2.2 Quantized Neural Networks

Benchmarks for quantized NNs were based on the MobileNets [35] architecture, which consists of 28 layers and 1000 output classes. The network alternates between 3×3 depthwise convolutions and 1×1 pointwise convolutions. A resolution multiplier ρ scales the dimensions of the input image, and a width multiplier α scales the size of the input and output channels. The models we used are hosted on TensorFlow’s online repository [2] and are trained on the ImageNet [29] dataset. We experimentally determined that an upper bound of $M = 16$ is sufficient for truncation by a private value, since all computed $\ell^{(i)}$ s are ≤ 9 for all model configurations.

Performance of quantized MobileNet inference is presented in Table 9. For accurate comparison, we executed [21]’s implementation on our machines using the same setting. Our methodology from Section 6.2 allowed us to reduce the ring size from $k = 72$ to $k = 30$ or less, potentially reducing the time by a factor of 2. The improvement that we observe in Table 9 is not drastic and can be explained by the differences in the algorithms. In particular, Escudero et al. [32] experimentally determined that [21]’s daBit implementation was superior to edaBits we use only in a single setting, namely, for the semi-honest, honest majority setting over \mathbb{Z}_{2^k} . In addition, MP-SPDZ’s optimization for large computation also aids its efficiency. This demonstrates that our quantized NN solution can aid efficiency.

7.2.3 Support Vector Machines

A *support vector machine* (SVM) is a type of a supervised learning classifier, where the computation is parameterized by the number of classes q and features m . We choose to do SVM classification and specifically run the computation for the ALOI dataset [33] to be able to compare performance of our framework to that of SPD \mathbb{Z}_{2^k} reported in [23]. In particular, SPD \mathbb{Z}_{2^k} achieves security in the malicious model with no honest majority over a ring – a much stronger security model than ours – and we are interested in knowing the computational price of the differences in the security assumptions.

SVM computation consists of a series of parallel dot products of the feature vector and support matrix, followed by argmax computation of the resulting values. This computation is given as

Protocol 10 in Appendix A, where argmax is computed in a hierarchical manner. The results of our experiments for the ALOI dataset with 463 classes and 128 features are presented in Table 10.

We see that $\text{SPD}\mathbb{Z}_{2^k}$'s performance (combined offline and online) in the two-party setting is about 200 times slower due to their significantly stronger security model. We also see that field-based implementation is an order of magnitude slower and the times do not reduce as significantly with the increased batch size as in our ring-based setting. This tells us that a single invocation of SVM evaluation is network-bound over a ring, but this is not the case with field-based computation.

8 Conclusions

In this work we study multi-party threshold secret sharing over a ring in the semi-honest model with honest majority with the goal of improving performance compared to field-based computation. We design low-level operations for n -party replicated secret sharing over any ring and consequentially build on them to enable general-purpose protocols over ring \mathbb{Z}_{2^k} . Our implementation results demonstrate that ring-based implementations of different operations are significantly faster than their field-based equivalents with $n = 3$. This allows us to improve performance of different applications including privacy-preserving machine learning tasks. We specifically test performance of neural network, quantized neural network and support vector machine classification and determine that performance of our techniques is on par with the best custom three-party protocols for those functions.

Acknowledgments

The authors would like to thank Jian Liu for help with understanding and reproducing the computation associated with the MNIST neural network evaluation in MiniONN [46]. The authors also would like to thank Marcel Keller for help with MP-SPDZ experiments and several constructions included in MP-SPDZ. In addition, the authors acknowledge support from the Emulab project [64] for utilizing an Emulab machine for WAN experiments. This work was supported in part by a Google Faculty Research Award and Buffalo Blue Sky Initiative. Any opinions, findings, and conclusions or recommendations expressed in this publication are those of the authors and do not necessarily reflect the views of the funding sources.

References

- [1] GMP – The GNU multiple precision arithmetic library. <http://www.gmpilib.org>.
- [2] Tensorflow repository. https://tensorflow.org/lite/guide/hosted_models.
- [3] MP-SPDZ repository. <https://github.com/data61/MP-SPDZ>, 2021.
- [4] M. Abspoel, A. Dalskov, D. Escudero, and A. Nof. An efficient passive-to-active compiler for honest-majority MPC over rings. In *International Conference on Applied Cryptography and Network Security (ACNS)*, pages 122–152, 2021.
- [5] M. Abspoel, D. Escudero, and N. Volgushev. Secure training of decision trees with continuous attributes. In *PoPETS*, pages 167–187, 2021.
- [6] N. Agrawal, A. Shahin Shamsabadi, M. J. Kusner, and A. Gascón. QUOTIENT: Two-party secure neural network training and prediction. In *CCS*, pages 1231–1247, 2019.

- [7] T. Araki, J. Furukawa, Y. Lindell, A. Nof, and K. Ohara. High-throughput semi-honest secure three-party computation with an honest majority. In *CCS*, pages 805–817, 2016.
- [8] D. Beaver. Efficient multiparty protocols using circuit randomization. In *CRYPTO*, pages 420–432, 1991.
- [9] D. Beaver and A. Wool. Quorum-based secure multi-party computation. In *EUROCRYPT*, pages 375–390, 1998.
- [10] G. R. Blakley. Safeguarding cryptographic keys. In *International Workshop on Managing Requirements Knowledge (MARK)*, pages 313–318, 1979.
- [11] M. Blanton, A. Kang, and C. Yuan. Improved building blocks for secure multi-party computation based on secret sharing with honest majority. In *International Conference on Applied Cryptography and Network Security (ACNS)*, pages 377–397, 2020.
- [12] F. Boemer, R. Cammarota, D. Demmler, T. Schneider, and H. Yalame. MP2ML: A mixed-protocol machine learning framework for private inference. In *International Conference on Availability, Reliability and Security (ARES)*, pages 1–10, 2020.
- [13] F. Boemer, Y. Lao, R. Cammarota, and C. Wierzynski. nGraph-HE: A graph compiler for deep learning on homomorphically encrypted data. In *ACM International Conference on Computing Frontiers*, pages 3–13, 2019.
- [14] D. Bogdanov, S. Laur, and J. Willemsen. Sharemind: A framework for fast privacy-preserving computations. In *European Symposium on Research in Computer Security (ESORICS)*, pages 192–206, 2008.
- [15] S. Carpov, K. Deforth, N. Gama, M. Georgieva, D. Jetchev, J. Katz, I. Leontiadis, M. Mohammadi, A. Sae-Tang, and M. Vuille. Manticore: Efficient framework for scalable secure multiparty computation protocols. IACR Cryptology ePrint Archive Report 2021/200, 2021.
- [16] O. Catrina and S. De Hoogh. Improved primitives for secure multiparty integer computation. In *International Conference on Security and Cryptography for Networks (SCN)*, pages 182–199, 2010.
- [17] H. Chaudhari, A. Choudhury, A. Patra, and A. Suresh. ASTRA: High throughput 3PC over rings with application to secure prediction. In *ACM Workshop on Cloud Computing Security (CCSW)*, pages 81–92, 2019.
- [18] H. Chaudhari, R. Rachuri, and A. Suresh. Trident: Efficient 4PC framework for privacy preserving machine learning. In *NDSS*, 2020.
- [19] R. Cramer, I. Damgård, D. Escudero, P. Scholl, and C. Xing. SPDZ_{2^k}: Efficient MPC mod 2^k for dishonest majority. In *CRYPTO*, pages 769–798, 2018.
- [20] R. Cramer, I. Damgård, and Y. Ishai. Share conversion, pseudorandom secret-sharing and applications to secure computation. In *Theory of Cryptography Conference (TCC)*, pages 342–362, 2005.
- [21] A. Dalskov, D. Escudero, and M. Keller. Secure evaluation of quantized neural networks. In *PoPETS*, pages 355–375, 2020.

- [22] A. Dalskov, D. Escudero, and M. Keller. Fantastic four: Honest-majority four-party secure computation with malicious security. In *USENIX Security Symposium*, pages 2183–2200, 2021.
- [23] I. Damgård, D. Escudero, T. Frederiksen, M. Keller, P. Scholl, and N. Volgushev. New primitives for actively-secure MPC over rings with applications to private machine learning. In *IEEE Symposium on Security and Privacy*, pages 1102–1120, 2019.
- [24] I. Damgård and J. Nielsen. Scalable and unconditionally secure multiparty computation. In *CRYPTO*, pages 572–590, 2007.
- [25] I. Damgård, C. Orlandi, and M. Simkin. Yet another compiler for active security or: Efficient MPC over arbitrary rings. In *CRYPTO*, pages 799–829, 2018.
- [26] I. Damgård, V. Pastro, N. P. Smart, and S. Zakarias. Multiparty computation from somewhat homomorphic encryption. In *CRYPTO*, pages 643–662, 2012.
- [27] I. Damgård and R. Thorbek. Efficient conversion of secret-shared values between different fields. IACR Cryptology ePrint Archive Report 2008/221, 2008.
- [28] Daniel Demmler, Thomas Schneider, and Michael Zohner. ABY - a framework for efficient mixed-protocol secure two-party computation. In *NDSS*, 2015.
- [29] J. Deng, W. Dong, R. Socher, L. Li, K Li, and L. Fei-Fei. ImageNet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 248–255, 2009.
- [30] H. Eerikson, M. Keller, C. Orlandi, P. Pullonen, J. Puura, and M. Simkin. Use your brain! arithmetic 3PC for any modulus with active security. In *Conference on Information-Theoretic Cryptography (ITC)*, pages 5:1–5:24, 2020.
- [31] D. Escudero, S. Ghosh, M. Keller, R. Rachuri, and P. Scholl. Improved primitives for MPC over mixed arithmetic-binary circuits. In *CRYPTO*, pages 823–852, 2020.
- [32] D. Escudero, S. Ghosh, M. Keller, R. Rachuri, and P. Scholl. Improved primitives for MPC over mixed arithmetic-binary circuits. IACR Cryptology ePrint Archive Report 2020/338, 2020.
- [33] J. Geusebroek, G. J. Burghouts, and A. Smeulders. The Amsterdam Library of Object Images. *International Journal of Computer Vision*, 61(1):103–112, 2005.
- [34] A. Gholami, S. Kim, Z. Dong, Z. Yao, M. W. Mahoney, and K. Keutzer. A survey of quantization methods for efficient neural network inference. arXiv preprint arXiv:2103.13630, 2021.
- [35] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam. MobileNets: Efficient convolutional neural networks for mobile vision applications. arXiv preprint arXiv:1704.04861, 2017.
- [36] M. Ito, A. Saito, and T. Nishizeki. Secret sharing schemes realizing general access structures. In *IEEE Global Telecommunication Conference (Globecom)*, pages 99–102, 1987.
- [37] B. Jacob, S. Kligys, B. Chen, M. Zhu, M. Tang, A. Howard, H. Adam, and D. Kalenichenko. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2704–2713, 2018.

- [38] C. Juvekar, V. Vaikuntanathan, and A. Chandrakasan. GAZELLE: A low latency framework for secure neural network inference. In *USENIX Security Symposium*, pages 1651–1669, 2018.
- [39] M. Keller. MP-SPDZ: A versatile framework for multi-party computation. In *ACM Conference on Computer and Communications Security (CCS)*, pages 1575–1590, 2020.
- [40] Marcel Keller and Ke Sun. Secure quantized training for deep learning. arXiv preprint arXiv:2107.00501, 2021.
- [41] L. Kerik, P. Laud, and J. Randmets. Optimizing MPC for robust and scalable integer and floating-point arithmetic. In *International Conference on Financial Cryptography and Data Security Workshops*, pages 271–287, 2016.
- [42] N. Kumar, M. Rathee, N. Chandran, D. Gupta, A. Rastogi, and R. Sharma. CrypTFlow: Secure tensorflow inference. In *IEEE Symposium on Security and Privacy*, pages 336–353, 2020.
- [43] Y. LeCun and C. Cortes. MNIST handwritten digit database. <http://yann.lecun.com/exdb/mnist/>, 2010.
- [44] S. Li, K. Xue, B. Zhu, C. Ding, X. Gao, D. Wei, and T. Wan. FALCON: A fourier transform based approach for fast and secure convolutional neural network predictions. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8705–8714, 2020.
- [45] Y. Lindell and A. Nof. A framework for constructing fast MPC over arithmetic circuits with malicious adversaries and an honest majority. In *ACM Conference on Computer and Communications Security (CCS)*, pages 259–276, 2017.
- [46] J. Liu, M. Juuti, Y. Lu, and N. Asokan. Oblivious neural network predictions via MiniONN transformations. In *ACM Conference on Computer and Communications Security (CCS)*, pages 619–631, 2017.
- [47] J. Liu, M. Juuti, Y. Lu, and N. Asokan. Oblivious neural network predictions via MiniONN transformations. IACR Cryptology ePrint Archive Report 2017/452, 2017.
- [48] E. Makri, D. Rotaru, F. Vercauteren, and S. Wagh. Rabbit: Efficient comparison for secure multi-party computation. IACR Cryptology ePrint Archive Report 2021/119, 2021.
- [49] U. Maurer. Secure multi-party computation made simple. In *Security in Communication Networks (SCN)*, pages 14–28, 2002.
- [50] S. Micali, O. Goldreich, and A. Wigderson. How to play any mental game. In *ACM Symposium on Theory of Computing (STOC)*, pages 218–229, 1987.
- [51] P. Mishra, R. Lehmkuhl, A. Srinivasan, W. Zheng, and R. A. Popa. DELPHI: A cryptographic inference service for neural networks. In *USENIX Security Symposium*, pages 2505–2522, 2020.
- [52] P. Mohassel and P. Rindal. ABY³: A mixed protocol framework for machine learning. In *ACM Conference on Computer and Communications Security (CCS)*, pages 35–52, 2018.
- [53] P. Mohassel and Y. Zhang. SecureML: A system for scalable privacy-preserving machine learning. In *IEEE Symposium on Security and Privacy*, pages 19–38, 2017.

- [54] A. Patra, T. Schneider, A. Suresh, and H. Yalame. ABY2.0: Improved mixed-protocol secure two-party computation. In *USENIX Security Symposium*, pages 2165–2182, 2021.
- [55] A. Patra and A. Suresh. BLAZE: Blazing fast privacy-preserving machine learning. In *NDSS*, 2020.
- [56] D. Rathee, M. Rathee, R. K. K. Goli, D. Gupta, R. Sharma, N. Chandran, and A. Rastogi. SIRNN: A math library for secure rnn inference. arXiv preprint arXiv:2105.04236, 2021.
- [57] D. Rathee, M. Rathee, N. Kumar, N. Chandran, D. Gupta, A. Rastogi, and R. Sharma. CryptTFlow2: Practical 2-party secure inference. In *ACM Conference on Computer and Communications Security (CCS)*, pages 325–342, 2020.
- [58] M. S. Riazi, C. Weinert, O. Tkachenko, E. M. Songhori, T. Schneider, and F. Koushanfar. Chameleon: A hybrid secure computation framework for machine learning applications. In *Asia Conference on Computer and Communications Security (ASIACCS)*, pages 707–721, 2018.
- [59] D. Rotaru and T. Wood. Marbled circuits: Mixing arithmetic and boolean circuits with active security. In *INDOCRYPT*, pages 227–249, 2019.
- [60] SecureSCM. Deliverable D9.2, EU FP7 Project Secure Supply Chain Management (SecureSCM). <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.221.393&rep=rep1&type=pdf>, 2009.
- [61] A. Shamir. How to share a secret. *Communications of the ACM*, 22(11):612–613, 1979.
- [62] S. Wagh, D. Gupta, and N. Chandran. SecureNN: Efficient and private neural network training. *Proceedings of Privacy Enhancing Technologies Symposium (PoPETS)*, 2019(3):26–49, 2019.
- [63] S. Wagh, S. Tople, F. Benhamouda, E. Kushilevitz, P. Mittal, and T. Rabin. Falcon: Honest-majority maliciously secure framework for private deep learning. In *PoPETS*, pages 188–208, 2021.
- [64] B. White, J. Lepreau, L. Stoller, R. Ricci, S. Guruprasad, M. Newbold, M. Hibler, C. Barb, and A. Joglekar. An integrated experimental environment for distributed systems and networks. In *Symposium on Operating Systems Design and Implementation (OSDI)*, pages 255–270, 2002.
- [65] Y. Zhang, A. Steele, and M. Blanton. PICCO: A general-purpose compiler for private distributed computation. In *ACM Conference on Computer and Communications Security (CCS)*, pages 813–826, 2013.

A Additional Protocols

A.1 Random Bit Generation

Random bit generation is a crucial component of a variety of protocols including different types of comparisons, bit decomposition, division, etc. Therefore, it is of paramount importance to support this functionality for general-purpose computation. In this work we examine two variants: (i) generating shares of a single bit as full-size ring elements and (ii) generating shares of k -bit random r as full-size ring elements together with generating shares of individual bits of r in \mathbb{Z}_2 .

Protocol 6 $[b] \leftarrow \text{RandBit}()$

- 1: $[u]_{k+2} \leftarrow \text{PRandR}(k+2)$;
 - 2: $[a]_{k+2} = 2[u]_{k+2} + 1$;
 - 3: $e \leftarrow \text{MulPub}([a]_{k+2}, [a]_{k+2})$;
 - 4: compute the smallest root of e modulo 2^{k+2} and denote it by c ; compute the inverse of c modulo 2^{k+2} and denote it by c^{-1} ;
 - 5: $[d]_{k+2} = c^{-1}[a]_{k+2} + 1$;
 - 6: for each $T \in \mathcal{T}$, let share $b_T = d_T/2$;
 - 7: **return** k least significant bits of each b_T as $[b]_k$;
-

Protocol	Rounds	Communication
$\text{RandBit}()$	1	$n - 1$
$\text{edaBit}(k)$	$\log(t+1)(\log(k)+1)+1$	$t^2(\log(k)+1)+t+1/2$

Table 11: Performance of random bit generation protocols with communication measured in the number of ring elements sent per party over $\mathbb{Z}_{2^{k+2}}$ for RandBit and \mathbb{Z}_{2^k} for $\text{edaBit}(k)$.

The first variant, denoted RandBit , originated in [16] for field-based SS and was modified in [23] to work in \mathbb{Z}_{2^k} . We use the logic of [23] and adjust the algorithm to work in our setting. The result is shown as Protocol 6.

To achieve 50% probability of each outcome of the output bit, the computation uses a larger ring $\mathbb{Z}_{2^{k+2}}$ for most steps of the protocol when the remaining computation uses ring \mathbb{Z}_{2^k} . Consequently, we use notation $[x]_\ell$ with variable ℓ to denote that shares and computation are over ring \mathbb{Z}_{2^ℓ} . We also parameterize function PRandR by the desired bitlength and $\text{PRandR}(\ell)$ denotes that the function returns a random ring element from \mathbb{Z}_{2^ℓ} .

Correctness of Protocol 6 follows from [23] and security follows from the logic. That is, because the protocol only discloses random e and otherwise uses secure building blocks, no information about private values can be leaked. The protocol runs in one round using the same communication as MulPub over $\mathbb{Z}_{2^{k+2}}$. To improve performance, in our implementation we compute the square root and inverse operations on line 4 simultaneously.

The second variant of random bit generation is based on the computation described in [31] and is denoted as $\text{edaBit}(k)$, where the parameter k specifies the number of generated random bits as well as the bitlength of their representation as integer r . It produces secret-shared k -bit integer r together with shares of the individual bits of r in \mathbb{Z}_2 . We use a simplified version with k being equal to the bitlength of the ring elements (i.e., the ring is \mathbb{Z}_{2^k}), which eliminates certain operations for dealing with carry after addition. The construction is given as Protocol 7. The idea consists of $t+1$ parties (without loss of generality, we chose the first $t+1$ parties for this role) each locally generating k random bits and computing representation of those bits as a k -bit integer (line 2). The bits are input into the computation using SS over \mathbb{Z}_2 , while the integers are entered using shares in \mathbb{Z}_{2^k} (line 3). Because we use Input to generate shares over different rings, we specify the second argument ℓ , which indicates that the shares need to be produced in \mathbb{Z}_{2^ℓ} . The output that the protocol produces is the sum of the $t+1$ random integers (without the carry bits) and its bit decomposition is computed using bitwise addition BitAdd from [60] of the $t+1$ integers represented as bits in a tree-like manner.

Protocol 7 ($[r]_k, [b_0]_1, \dots, [b_{k-1}]_1 \leftarrow \text{edaBit}(k)$)

```
1: for  $p = 1, \dots, t + 1$  in parallel do
2:   party  $p$  samples  $r_0^{(p)}, \dots, r_{k-1}^{(p)} \in \mathbb{Z}_2$  and computes  $r^{(p)} = \sum_{j=0}^{k-1} r_j^{(p)} 2^j$ ;
3:   simultaneously execute  $[r^{(p)}]_k \leftarrow \text{Input}(r^{(p)}, k)$  and  $[r_i^{(p)}]_1 \leftarrow \text{Input}(r_i^{(p)}, 1)$  for  $i = 1, \dots, k$ 
   with  $p$  being the input owner;
4: end for
5:  $[r]_k = \sum_{p=1}^{t+1} [r^{(p)}]_k$ ;
6:  $s = t + 1$ ;
7: for  $i = 1, \dots, \lceil \log(t + 1) \rceil$  do
8:   for  $j = 1, \dots, \lfloor s/2 \rfloor$  in parallel do
9:      $\langle [r_1^{(j)}]_1, \dots, [r_{k-1}^{(j)}]_1 \rangle \leftarrow \text{BitAdd}(\langle [r_1^{(2j-1)}]_1, \dots, [r_{k-1}^{(2j-1)}]_1 \rangle, \langle [r_1^{(2j)}]_1, \dots, [r_{k-1}^{(2j)}]_1 \rangle)$ ;
10:    if  $s \bmod 2 = 0$  then
11:       $s = s/2$ ;
12:    else
13:       $\langle [r_1^{((s+1)/2)}]_1, \dots, [r_{k-1}^{((s+1)/2)}]_1 \rangle = \langle [r_1^{(s)}]_1, \dots, [r_{k-1}^{(s)}]_1 \rangle$ ;
14:       $s = (s + 1)/2$ ;
15:    end if
16:  end for
17: end for
18:  $[b_0]_1, \dots, [b_{k-1}]_1 = [r_0^{(1)}]_1, \dots, [r_{k-1}^{(1)}]_1$ 
19: return ( $[r]_k, [b_0]_1, \dots, [b_{k-1}]_1$ )
```

A.2 Comparisons

Less-than comparisons, $[a] < [b]$, are traditionally computed using SS by determining the most significant bit of the difference between a and b . Starting from [16], comparison protocols blind the difference by adding a random integer bit decomposition of which is known, open the sum, truncate all but one bit, and compensate for any carry caused by the addition. This logic was adapted to the ring setting in [23] by using building blocks that work over \mathbb{Z}_{2^k} . In the solution that we present as Protocol 8, we incorporate the edaBit protocol from [31] for efficient random bit generation into the construction of [23] adopted to the semi-honest setting. The presence of carry is determined using sub-protocol BitLT which performs comparison of two bit-decomposed values, one of which is given in the clear, using binary computation over \mathbb{Z}_2 .

Security of the algorithm follows from prior work and the fact that we use a composition of secure building blocks. In particular, the only values revealed in the protocol (in steps 4 and 9) are information-theoretically protected using freshly generated randomness. The complexity of this protocol and its prior version that makes calls to RandBit is given in Table 2.

To correctly implement comparison of two k -bit integers over ring \mathbb{Z}_{2^k} , one would need to invoke the MSB protocol 3 times. However, correctness is also guaranteed if we compare two $(k - 1)$ -bit integers over ring \mathbb{Z}_{2^k} using a single call to MSB. We use the latter approach in our implementation of machine learning algorithms.

There are noteworthy differences in the design of protocols developed for a ring as opposed to original protocols for a field. Certain operations such as prefix multiplication are not available in a ring and we resort to logarithmic round building blocks when protocols over a field achieve constant round complexity. In the context of comparison, a typical tool for realizing them was truncation (i.e., right shift), the cost of which was linear in the number of bits truncated, but the modulus

Protocol 8 $[a_{k-1}]_k \leftarrow \text{MSB}([a]_k)$, where $a = \sum_{i=0}^{k-1} a_i 2^i \in \mathbb{Z}_{2^k}$

- 1: $[r]_k, [r_0]_1, \dots, [r_{k-1}]_1 \leftarrow \text{edaBits}(k)$;
 - 2: $[b]_k \leftarrow \text{RandBit}()$;
 - 3: $[r']_k = [r]_k - [r_{k-1}]_1 2^{k-1}$;
 - 4: $c \leftarrow \text{Open}([a]_k + [r]_k)$;
 - 5: $c' = c \bmod 2^{k-1}$;
 - 6: $[u]_1 \leftarrow \text{BitLT}(c', [r_0]_1, \dots, [r_{k-2}]_1)$;
 - 7: $[a']_k = c' - [r']_k + 2^{k-1}[u]_1$;
 - 8: $[d]_k = [a]_k - [a']_k$;
 - 9: $e \leftarrow \text{Open}([d]_k + 2^{k-1}[b]_k)$ and let e_{k-1} be the most significant bit of e ;
 - 10: $[a_{k-1}]_k = e_{k-1} + [b]_k - 2e_{k-1}[b]_k$;
 - 11: **return** $[a_{k-1}]_k$;
-

had to be increased by a statistical security analysis to support such operations. In a ring, on the other hand, there is no significant increase in the ring size, but the communication cost is linear in the bitlength of the ring and not in the bitlength of the truncated portion. This brings different trade-offs, but the availability of faster arithmetic in a ring will still lead to significant savings.

A.3 Machine Learning Applications

Protocol 9 shows the computation associated with evaluating the MNIST neural network model from [46]. Protocol 10 provides computation associated with support vector machine evaluation on

Protocol 9 MNIST neural network evaluation

- 1: (Convolution) 28×28 input image, 5×5 window size, $(1, 1)$ stride, 16 output channels: $\mathbb{Z}_{2^{20}}^{16 \times 576} \leftarrow \text{MatMult}(\mathbb{Z}_{2^{20}}^{16 \times 25}, \mathbb{Z}_{2^{20}}^{25 \times 576})$;
 - 2: (ReLU) calculates ReLU for each entry of $\mathbb{Z}_{2^{20}}^{16 \times 576}$;
 - 3: (Max Pooling) input $\mathbb{Z}_{2^{20}}^{16 \times 576}$ and $1 \times 2 \times 2$ window size, outputs $\mathbb{Z}_{2^{20}}^{16 \times 12 \times 12}$;
 - 4: (Conversion) $\mathbb{Z}_{2^{30}}^{16 \times 144} \leftarrow \text{Convert}(\mathbb{Z}_{2^{20}}^{16 \times 144}, 20, 30)$;
 - 5: (Convolution) 5×5 window size, $(1, 1)$ stride, 16 output channels: $\mathbb{Z}_{2^{30}}^{16 \times 64} \leftarrow \text{MatMult}(\mathbb{Z}_{2^{30}}^{16 \times 400}, \mathbb{Z}_{2^{30}}^{400 \times 64})$;
 - 6: (ReLU) calculates ReLU for each entry of $\mathbb{Z}_{2^{30}}^{16 \times 64}$;
 - 7: (Max Pooling) input $\mathbb{Z}_{2^{30}}^{16 \times 64}$ and $1 \times 2 \times 2$ window size, outputs $\mathbb{Z}_{2^{30}}^{16 \times 4 \times 4}$;
 - 8: (Conversion) $\mathbb{Z}_{2^{40}}^{16 \times 16} \leftarrow \text{Convert}(\mathbb{Z}_{2^{30}}^{16 \times 16}, 30, 40)$;
 - 9: (Fully Connected) Connects 256 incoming nodes to 100 outgoing nodes : $\mathbb{Z}_{2^{40}}^{100 \times 1} \leftarrow \text{MatMult}(\mathbb{Z}_{2^{40}}^{100 \times 256}, \mathbb{Z}_{2^{40}}^{256 \times 1})$;
 - 10: (ReLU) calculates ReLU for each entry of $\mathbb{Z}_{2^{40}}^{100 \times 1}$;
 - 11: (Conversion) $\mathbb{Z}_{2^{49}}^{100 \times 1} \leftarrow \text{Convert}(\mathbb{Z}_{2^{40}}^{100 \times 1}, 40, 49)$;
 - 12: (Fully Connected) Connects 100 incoming nodes to 10 outgoing nodes : $\mathbb{Z}_{2^{49}}^{10 \times 1} \leftarrow \text{MatMult}(\mathbb{Z}_{2^{49}}^{10 \times 100}, \mathbb{Z}_{2^{49}}^{100 \times 1})$;
 - 13: **return** $\mathbb{Z}_{2^{49}}^{10 \times 1}$;
-

q classes and m features. For the ALOI SVM from [23], we use $q = 463$ and $m = 128$.

Protocol 10 SVM classification, where q is the number of classes, m is the number of features, $f_{i,j}$ is the feature vector, and b_j are the biases

- 1: **for each** $j = 1$ to q **in parallel do** $[c_j] \leftarrow [b_j] + \sum_{i=1}^m \text{DotProd}([f_{i,j}], [x_i]);$
 - 2: **return** $([a_{\text{ind}}], [\text{ind}]) \leftarrow \text{ArgMax}([c_1], \dots, [c_q]);$
-

B Security Definitions and Proofs

We formulate (simulation-based) security in the presence of semi-honest participants as follows:

Definition 1. Let parties P_1, \dots, P_n engage in a protocol Π that computes function $f(\text{in}_1, \dots, \text{in}_n) = (\text{out}_1, \dots, \text{out}_n)$, where in_i and out_i denote the input and output of party P_i , respectively. Let $\text{VIEW}_\Pi(P_i)$ denote the view of participant P_i during the execution of protocol Π . More precisely, P_i 's view is formed by its input and internal random coin tosses r_i , as well as messages m_1, \dots, m_k passed between the parties during protocol execution: $\text{VIEW}_\Pi(P_i) = (\text{in}_i, r_i, m_1, \dots, m_k)$. Let $I = \{P_{i_1}, P_{i_2}, \dots, P_{i_t}\}$ denote a subset of the participants for $t < n$, $\text{VIEW}_\Pi(I)$ denote the combined view of participants in I during the execution of protocol Π (i.e., the union of the views of the participants in I), and $f_I(\text{in}_1, \dots, \text{in}_n)$ denote the projection of $f(\text{in}_1, \dots, \text{in}_n)$ on the coordinates in I (i.e., $f_I(\text{in}_1, \dots, \text{in}_n)$ consists of the i_1 th, \dots , i_t th element that $f(\text{in}_1, \dots, \text{in}_n)$ outputs). We say that protocol Π is t -private in the presence of semi-honest adversaries if for each coalition of size at most t there exists a probabilistic polynomial time simulator S_I such that $\{S_I(\text{in}_I, f_I(\text{in}_1, \dots, \text{in}_n)), f(\text{in}_1, \dots, \text{in}_n)\} \equiv \{\text{VIEW}_\Pi(I), (\text{out}_1, \dots, \text{out}_n)\}$, where $\text{in}_I = \bigcup_{P_i \in I} \{\text{in}_i\}$ and \equiv denotes computational or statistical indistinguishability.

Proof of Theorem 1. Let I denote the set of corrupt parties. We consider the maximal amount of corruption with $|I| = t$. Because the computation proceeds on secret shares and the parties do not learn the result, no information should be revealed to the computational parties as a result of protocol execution.

We build a simulator S_I that interacts with the parties in I as follows: when a party $p \in I$ expects to receive a value from another party $p' \notin I$ in step 5 of the computation according to function χ , S_I chooses a random element of \mathcal{R} and sends it to p . S_I preserves consistency of the view and ensures that when the same value is to be sent by p' to multiple parties in I , all of them receive the same random value. This is the only portion of the protocol where corrupt parties can receive values (that the simulator produces), and the only portion of the protocol when a corrupt party p may send a value to an honest party p' is step 4, which S_I receives on behalf of p' . All other computation is local, in which S_I does not participate.

We next argue that the simulated view is computationally indistinguishable from the real view. First, note that the corrupt parties in I collectively hold shares a_T, b_T and keys k_T (and thus can compute values $G_T.\text{next}$) for each $T \in \mathcal{T}$ such that $\exists p \in I$ and $p \notin T$. This entitles the corrupt parties to computing the corresponding shares c_T , but the rest of the shares must remain unknown, so that they are unable to compute c . Next, notice that when $|I| = t$, there is only one share $T^* = I$ such that all parties $p \in I$ have no access to k_{T^*} and c_{T^*} , while all parties $p' \notin I$ store those values. Then there are two cases to consider: (1) If one or more parties $p \in I$ receive $\chi(p')$'s share of $v^{p'}$ from another party $p' \notin I$ (it must be the case that $\chi(p') \neq T^*$), the received share has been masked by a fresh pseudo-random element from G_{T^*} , is therefore pseudo-random and indistinguishable from random by any $p \in I$. (2) If no party $p \in I$ receives a value from any given $p' \notin I$, indistinguishability is trivially maintained. \square

Proof of Theorem 2. As before, let I denote the set of corrupt parties with $|I| = t$. We build a simulator S_I that interacts with the parties in I as follows: after S_I extracts shares a_T, b_T, k_T

($T \in \mathcal{T}$ such that $\exists p \in I$ and $p \notin T$) from the corrupt parties and receives the output c from the trusted party, S_I computes $v^{(p)}$ as prescribed by the protocol for each $p \in I$ and also their sum $v_I = \sum_{p \in I} v^{(p)}$ (in \mathcal{R}). S_I sets $v^{(p)}$ values for the remaining $n - t$ parties to random elements of \mathcal{R} subject to $\sum_{p \notin I} v^{(p)} = c - v_I$ (in \mathcal{R}). S_I , acting on behalf of party $p \notin I$, sends the corresponding $v^{(p)}$ to each party in I .

To show that this simulation is indistinguishable from the real protocol execution, recall that there will be at least one T , denoted by $T^* = I$, to which the parties in I have no access (and thus correspondingly cannot distinguish the output G_{T^*} from random elements of the ring). During real protocol execution the parties in I receive $t + 1$ values $c^{(p)}$, one per $p \notin I$. With the knowledge that the corrupt parties collectively have, they can remove the effect of all randomization except the use of the output of G_{T^*} . If we let z_{i,T^*} denote the i th call to $G_{T^*}.\text{next}$ during the execution of `MulPub` in Protocol 2, then the corrupt parties can recover t values of the form $v^{(p)} + z_{i,T^*}$ with unique p and i and one value of the form $v^{(p)} - \sum_{i=1}^t z_{i,T^*}$ for another p . The next thing to notice is that any t (out of $t + 1$) of these values are pseudo-random and computationally protect the corresponding $v^{(p)}$ values. The introduction of the remaining value reveals the sum of all $v^{(p)}$ s, but not other information (i.e., the last value corresponds to the difference to make the sum equal to $c - v_I$). This means that substituting these values with random elements subject to $\sum_{p \notin I} v^{(p)} = c - v_I$ provides the same information to the corrupt parties and achieves computational indistinguishability of the views. \square

Proof of Theorem 3. It is straightforward to show security of the full version of `Input` when the input owner is different from the computational parties. That is, the input owner creates proper shares according to the SS scheme using a PRG. Thus, as long as security of the PRG holds, the real view is computationally indistinguishable from a simulated view created without the use of any secrets.

However, when the input owner is one of the computational parties, only a reduced set of shares is produced. Thus, we need to evaluate the combined view of each coalition of t corrupt participants. There are two important cases to consider: (i) input owner p^* is a part of the coalition and (ii) it is not.

When p^* is a corrupt participant, building a simulator is trivial: the simulator simply receives shares from the input owner on behalf of honest participants and terminates. Because inputs a_i are available to the corrupt parties, no information need to be protected and the real and simulated views use identical values.

When there are t corrupt participants who are different from p^* , we simulate the view by choosing a random value for a_{i,T^*} and sending it to each corrupt $p \notin T^*$. What remains to show is that the t corrupt parties do not possess enough shares to reconstruct the secret and, as a result, cannot learn any information about it. In more detail, p^* distributes its secrets using only shares T such that $T \in \mathcal{T} \setminus \{T^*\}$. However, because we use (n, t) threshold SS, there will be a share T possessed by p^* which is not available to any of the t corrupt parties I . Specifically that share is available to all participants except corrupt minority I . This means that the corrupt parties will not be able to reconstruct information about the private inputs and the real and simulated views are indistinguishable as long as PRG's security holds. \square