

Megha Byali, Harsh Chaudhari*, Arpita Patra, and Ajith Suresh

FLASH: Fast and Robust Framework for Privacy-preserving Machine Learning

Abstract: Privacy-preserving machine learning (PPML) via Secure Multi-party Computation (MPC) has gained momentum in the recent past. Assuming a minimal network of pair-wise private channels, we propose an efficient four-party PPML framework over rings \mathbb{Z}_{2^t} , FLASH, the first of its kind in the regime of PPML framework, that achieves the strongest security notion of Guaranteed Output Delivery (all parties obtain the output irrespective of adversary’s behaviour). The state of the art ML frameworks such as ABY3 by *Mohassel et.al* (ACM CCS’18) and SecureNN by *Wagh et.al* (PETS’19) operate in the setting of 3 parties with one malicious corruption but achieve the *weaker* security guarantee of *abort*. We demonstrate PPML with real-time efficiency, using the following custom-made tools that overcome the limitations of the aforementioned state-of-the-art– (a) *dot product*, which is independent of the vector size unlike the state-of-the-art ABY3, SecureNN and ASTRA by *Chaudhari et.al* (ACM CCSW’19), all of which have linear dependence on the vector size. (b) *Truncation* and *MSB Extraction*, which are constant round and free of circuits like Parallel Prefix Adder (PPA) and Ripple Carry Adder (RCA), unlike ABY3 which uses these circuits and has round complexity of the order of depth of these circuits. We then exhibit the application of our FLASH framework in the secure server-aided prediction of vital algorithms– Linear Regression, Logistic Regression, Deep Neural Networks, and Binarized Neural Networks. We substantiate our theoretical claims through improvement in benchmarks of the aforementioned algorithms when compared with the current best framework ABY3. All the protocols are implemented over a 64-bit ring in LAN and WAN. Our experiments demonstrate that, for MNIST dataset, the improvement (in terms of throughput) ranges from $24\times$ to $1390\times$ over LAN and WAN together.

Keywords: Privacy, Machine Learning, Robust 4PC

DOI 10.2478/popets-2020-0036

Received 2019-08-31; revised 2019-12-15; accepted 2019-12-16.

Megha Byali: Indian Institute of Science, E-mail: megha@iisc.ac.in

1 Introduction

Secure Multi-party Computation (MPC) [8, 21, 28, 33, 54] has evolved over the years in its pursuit of enabling a set of n mutually distrusting parties to compute a joint function f , in a way that no coalition of t parties can disrupt the true output of computation (correctness) or learn any information beyond what is revealed by the output of the computation (privacy). The area of secure MPC can be broadly categorized into honest majority [4, 8, 13, 44] and dishonest majority [19, 21, 28, 42, 54]. Over the years, MPC has progressed from being simply of theoretical interest to providing real-time practical efficiency. In terms of efficient constructions, the special case of dishonest-majority setting, namely two-party computation (2PC) [39, 40, 47, 54] has been in limelight over the last decade. However lately, the setting of three parties (3PC) [3, 4, 13, 44] and four parties (4PC) [13, 29, 34] have drawn phenomenal attention due to the customization in techniques and efficiency that the constructions have to offer. In this direction, the area of MPC in a small domain with an honest majority is quite fascinating due to variety of reasons mentioned below.

First, the most widely known real-time applications such as Danish Sugar-Beet Auction [10], Distributed Credential Encryption [44], Fair-play MPC [7], VIFF [27], Sharemind [9] explore MPC with 3 parties. Second, the expensive public-key primitives such as Oblivious Transfer (OT) known to be necessary for 2PC can be eliminated in the honest majority. Thus, the resulting constructions use only light-weight primitives and can even be information-theoretically secure. Third, the recent advances in secure Machine Learning (ML) have indicated real-time applications involving a small num-

***Corresponding Author: Harsh Chaudhari:** Indian Institute of Science, E-mail: chaudharim@iisc.ac.in

Arpita Patra: Indian Institute of Science. This author is supported by SERB Women Excellence Award 2017 (DSTO 1706). E-mail: arpita@iisc.ac.in

Ajith Suresh: Indian Institute of Science, This author is supported by Google Phd Fellowship 2019. E-mail: ajith@iisc.ac.in



ber of parties [2, 6, 14, 41, 43, 45, 53]. Furthermore, the stronger security notions of fairness (the adversary gets the output if and only if the honest parties do) and robustness aka guaranteed output delivery (GOD) (all parties obtain the output irrespective of adversary’s behaviour) are guaranteed only in the honest majority setting [16].

In this work, we strongly motivate the need for robustness in privacy-preserving machine learning as a service (MLaaS) and then go on to explore the setting of 4PC and demonstrate that our constructions are highly efficient compared to the existing state of the art 3PC ML frameworks. The guarantee of robustness is of utmost importance in the area of MLaaS. Consider the following scenario where an entity owns a trained ML model and wants to provide prediction as a service. The model owner outsources her trained model parameters to a set of three servers, which uses one of the aforementioned 3PC ML frameworks for secure prediction. These frameworks keep the privacy of the model parameters and the queries of the clients intact even when one of the servers is maliciously corrupted, but cannot guarantee an output to a given client’s query as the adversary can cause the protocol to abort. Thus in the practical setting, one simple strategy of the adversary would be to make the protocol abort for all the client queries. Eventually, this would steer the entity towards loss of monetary value and trust of the clients.

Motivation for 4PC Framework: The specific problem of MPC with 4-parties tolerating one corruption is of special interest to us. There are three primary motivations for us to consider this setting for achieving GOD– (a) avoid theoretical necessity of broadcast channel; (b) avoid expansive public-key primitives and (c) communication efficiency. We elaborate these points below. The popular setting of 3PC, when considered to achieve robustness, suffers from the necessity of an expensive robust broadcast channel as proven in the result of [17]. By moving to 4PC from 3PC, the need for a broadcast channel is removed, which results in highly efficient constructions [13] when compared to 3PC [14, 43, 53]. Additionally in 4PC, for any message sent by a party that needs an agreement, a simple honest majority rule over the residual three parties suffices. Such a property cannot be counted on in 3PC which leads to the use of costly workarounds than 4PC. [29] was the most recent work to propose guaranteed output delivery (robustness) in the 4PC setting. A major concern with GOD variant of multiplication protocol in [29] was utilizing Digital Signatures and expensive public-key primitives: Broadcast and a PKI Setup. Since our

end goal is an efficient and robust framework for ML, we let go their approach and propose a simple primitive coupled with a new secret sharing scheme which requires only symmetric-key primitives to achieve robustness.

Moreover, the state-of-the-art 3PC ML frameworks, like ABY3 and ASTRA, focused on highly efficient frameworks for machine learning in the semi-honest setting but suffered from efficiency loss for the primitives dot product, MSB extraction, and truncation in the malicious setting. For example, many of the widely used ML algorithms like Linear Regression, Logistic Regression, and Neural Networks use dot product computation as its building block. While the above frameworks incur a communication cost which is linearly dependent on the underlying size of the feature vector, we are able to eliminate this limitation and provide a dot product protocol whose communication is independent of the vector size. Additionally, we also make all our building blocks constant round and free of any circuits, unlike ABY3 which uses expensive non-constant round circuits like Parallel Prefix Adder (PPA) and Ripple Carry Adder (RCA) in their protocols.

Lastly, we choose build our framework over rings. Most of the computer architectures, Intel x64 for example, have their primitive data-types over rings. These architectures have specially designed hardware which can support fast and efficient arithmetic operations over rings. This led the way for efficient protocols over rings [3, 9, 11, 14, 20, 24] as opposed to fields, which are usually 10-20x slower since they have to rely on external libraries. Thus, our protocols over rings give the additional advantage of faster performance when implemented in the real-world architectures.

1.1 Our Contribution

We propose FLASH, the first robust framework for privacy-preserving machine learning in the four party (4PC) honest majority setting over a ring \mathbb{Z}_{2^ℓ} . We summarize our contributions below:

Robust 4PC protocol: We present an efficient and robust MPC protocol for four parties tolerating one malicious corruption. Concretely, for the multiplication operation, we require an overall communication of just 12 elements in the amortized sense. This is $\approx 2\times$ improvement in terms of communication over the state-of-the-art protocol of [29]. Moreover, our solution forgoes the need for Digital Signatures and expensive primitives like Broadcast and Public-Key Setup, unlike [29].

Protocol	Equation	ABY3		ASTRA		FLASH	
		Rounds	Comm.	Rounds	Comm.	Rounds	Comm.
Multiplication	$[[x], [y]] \rightarrow [x \cdot y]$	5	21ℓ	7	25ℓ	5	12ℓ
Dot Product	$[[\vec{x} \odot \vec{y}]] = [[\sum_{i=1}^d x_i y_i]]$	5	$21m\ell$	7	$23m\ell + 2\ell$	5	12ℓ
MSB Extraction	$[[x]] \rightarrow [[\text{msb}(x)]]^B$	$\log \ell + 4$	42ℓ	10	$52\ell + 4$	6	$16\ell + 4$
Truncation	$[[x], [y]] \rightarrow [((xy)^t)]$	$2\ell - 1$	$\approx 108\ell$	—	—	5	14ℓ
Bit Conversion	$[[b]]^B \rightarrow [b]$	6	42ℓ	—	—	5	14ℓ
Bit Insertion	$[[b]]^B [[x]] \rightarrow [bx]$	7	63ℓ	—	—	5	18ℓ

Table 1. Comparison of FLASH framework with ABY3 and ASTRA; ℓ and m denote the ring size and number of features respectively.

The removal of this additional setup of Digital Signatures, PKI and Broadcast primarily comes from two factors – i) a new secret sharing scheme which we call as *mirrored*-sharing, enables two disjoint sets of parties to perform the computation and perform an effective validation in a single execution, and ii) a simple yet novel *bi-convey* primitive, which enables two designated parties, say S_1, S_2 , to send a value to a designated party R with the help of a fourth party T .

The bi-convey primitive guarantees that if both S_1 and S_2 are honest, then party R will receive the value x for sure. If not, either the party R will be able to obtain x or both the parties R and T identify that one among S_1, S_2 is corrupt. Our construction for the bi-convey primitive requires a commitment scheme as the only cryptographic tool, which is considered inexpensive. Moreover, the commitments can be clubbed together for several instances and thus the cost of commitment gets amortized as well. Looking ahead, most of our constructions are designed in such a way that every message to be communicated will be made available to at least two parties and thus we can use the bi-convey primitive for the same.

Building Blocks for Machine Learning: We propose practically efficient building blocks that form the base for secure prediction. While ABY3 and SecureNN propose building blocks for security with abort, ASTRA elevates the security of these blocks from abort to fairness. We further strengthen the security and make all the building blocks robust. Additionally, we achieve significant efficiency improvements in all the building blocks due to the aid provided by an additional honest party in our setting. The improvements for each block are summarized as follows:

i) Dot Product: The aforementioned 3PC frameworks involve communication, linear in the order of vector size, we overcome this limitation with an efficient technique, independent of the vector size. This indepen-

dence stems from the peculiar structure of our mirrored sharing alongside the multiplication protocol in 4PC.

ii) Truncation: Overflow caused by repeated multiplications may cause accuracy loss which can be prevented with truncation. Truncation has been expensive in the 3PC framework, especially ABY3 uses a Ripple Carry Adder (RCA) circuit which consumes around 108 ring elements to achieve MSB Extraction. We propose a simple yet efficient technique with a total of just 14 ring elements and does not require any circuits. The technical novelty comes from the specific roles played by the parties, in conjunction with the multiplication protocol of 4PC. We defer the detailed analysis of our truncation protocol and the corresponding roles of the parties to Section 5.4.

iii) MSB Extraction: Comparing two arithmetic values in a privacy-preserving manner is one of the major hurdles in realizing efficient privacy-preserving ML algorithms. The state of the art SecureML[45] and ABY3 made an effort in this direction with the use of a garbled circuit technique and parallel prefix adder (PPA) respectively. Yet, these techniques still involve significant computation and communication which are a bottleneck to efficiency. We propose a technique free of any circuit computation and instead relies on the multiplication protocol of our 4PC.

iv) Bit Conversion and Insertion: Operating interchangeably in the arithmetic and boolean worlds often demand conversion of a boolean bit to its arithmetic equivalent (*bit conversion*) or the multiplication of a boolean bit with an arithmetic value (*bit insertion*). We propose efficient techniques to achieve the same with innovations coming from our mirrored secret sharing and its linearity property. Ours is the first work in 4PC that proposes these transformations and is even superior to the state-of-the-art 3PC ML frameworks ABY3 and ASTRA, in terms of both efficiency and security guarantee. Table 1 provides a detailed comparison in terms of com-

munication (Comm.) and rounds with ABY3 and ASTRA, where ℓ and m denote the ring size and number of features respectively.

Secure Prediction: We aim at secure prediction in a server-aided setting. Here, the model owner (M) holds a set of *trained* model parameters which are used to predict output to client’s (C) input query, while preserving the privacy of the inputs of both the parties. The servers perform computation and reconstruct the output towards the client. Security is provided against a malicious adversary corrupting one server along with either model owner or client. We extend our techniques for vital machine learning algorithms namely: i) Linear Regression, ii) Logistic Regression, iii) Deep Neural Network (DNN) and iv) Binarized Neural Network (BNN). While Linear Regression is extensively used in Market Analytics, Logistic Regression is used in a variety of applications like customer segmentation, insurance fraud detection and so on. Despite being computationally cheap and smaller in size, the performance accuracy of BNNs is comparable to that of deep neural networks. They are the go-to networks for running neural networks on low-end devices. These use cases exhibit the importance of these algorithms in real-time and we make an effort to efficiently perform the secure evaluation for these algorithms.

ML Algorithm	Setting	
	LAN	WAN
Linear Regression	1390×	125.4×
Logistic Regression	601×	48.5×
Deep Neural Network	344×	29.1×
Binarized Neural Network	277×	23.8×

Table 2. Improvement over ABY3 in terms of throughput for MNIST dataset

We provide implementation results for all our protocols over a ring $\mathbb{Z}_{2^{64}}$. We summarize the efficiency gain of our protocols over the state-of-the-art ABY3 and ASTRA, albeit more elaborate details follow in Section 6. The latency and throughput (the number of operations per unit time) of the protocols are measured in the LAN (1Gbps) and WAN (20Mbps) setting while communication complexity is measured independent of the network. We compare the most crucial building blocks, namely i) Dot Product, ii) MSB Extraction and iii) Truncation of our framework with state-of-the-art and show the practical improvement observed in each of the building blocks. We also provide throughput compar-

isons (# queries per sec in LAN and # queries per min in WAN) for the aforementioned algorithms, over multiple real-world datasets. Table 2 below shows the improvement over ABY3 for MNIST dataset [38] which has 784 features. We omit comparison with ASTRA as ABY3 outperforms ASTRA in terms of total communication (ref. Table 1). The improvements for DNN and BNN stated in Table 2 are for a network having 2 hidden layers, each layer consisting of 128 nodes.

4PC Abort: As an extension, we also propose protocols for the weaker abort setting. The abort variant for the protocols are achieved by tweaking the bi-convey primitive present in the robust protocols. We give a detailed analysis and comparison with state-of-the-art works in Appendix A.

1.2 Related Work:

In the regime of MPC over a small domain, interesting works that achieve guaranteed output delivery have been carried out mainly in the class of low-latency (consisting of small constant number of rounds) protocols [12, 13, 50]. However, in the view of practical efficiency, high throughput (light in communication and computation complexity) is desirable. Yet the literature of high throughput protocols has seen limited work [29] in guaranteeing security notions stronger than abort. The existing state-of-the-art includes notable works that are highly efficient, but trade security for efficiency [3–5, 15, 26, 49]. In this work, we attempt to bridge the gap between the security achieved and the corresponding efficiency, thus providing highly efficient PPML framework using robust 4PC as the backbone. Below we summarize the contributions closest to our setting.

The study of MPC in high-throughput networks accelerated with the celebrated work of [23]. The works of [3–5, 15, 26, 49] swiftly followed. These works focus on the evaluation of arithmetic circuits over rings or finite fields. [5] is semi-honest and operates over both rings and fields. The works of [3, 20, 26] achieve abort security over rings with one malicious corruption. A compiler to transform semi-honest security to malicious-security was proposed by [15]. This conversion is obtained at twice the cost of the semi-honest protocol. The work of [29] explores 4PC and the security notions of fairness and guaranteed output delivery. However, [29] is dual execution based and relies on expensive public-key primitives and broadcast channel to achieve guaranteed output delivery. [49] improvises over [15] by presenting a batch multiplication technique and additionally explores the notion of fairness.

The influence of ML has found its way in a broad range of areas such as facial recognition [52], banking, medicine [25], recommendation systems and so on. Consequently, technology giants such as Amazon, Google are providing ML as a service (MLaaS) for both training and prediction purposes, where the parties outsource their computation to a set of servers. However, for confidential purposes, government regulations and competitive edge, such data cannot be made publicly available. Thus, there is a need for privacy of data while still enabling customers to perform training and prediction. This need for privacy has given rise to the culmination of MPC and ML. Recent works [14, 41, 43, 45, 51, 53] have shown the need of MPC in achieving efficient techniques for privacy-preserving machine learning in server aided setting, where parties outsource their data to a set of servers and the servers compute for purposes of training or classification. There have been works dedicated to linear regression [14, 43, 45], logistic regression [14, 43, 45] and neural networks [36, 43, 45, 51, 53] for both training and inference. Recent works have dived into variants of neural networks like Deep Neural Networks (DNNs) [43, 46, 51], Convolutional Neural Networks (CNNs) [36, 51, 53], Binarized Neural Networks (BNNs) [37] and Quantized Neural Networks (QNNs) [1, 35]. DNNs and CNNs have become one of the most powerful machine learning models in recent history with amount of data available to train them and are one of the most widely considered models for training and prediction tasks for low power devices. MOBIUS [37] was the first to explore secure prediction in BNNs for semi-honest 2PC.

2 Preliminaries and Definitions

We consider a set of four parties $\mathcal{P} = \{V_1, V_2, E_1, E_2\}$ connected by pair-wise private and authentic channels in a synchronous network. E_1, E_2 define the role of the parties as *evaluators* in the computation while parties V_1, V_2 enact the role of *verifiers* in the computation. We use \mathbf{E} and \mathbf{V} to denote the set of evaluators $\{E_1, E_2\}$ and verifiers $\{V_1, V_2\}$ respectively. The function f to be evaluated is expressed as a circuit ckt , with a publicly known topology and is evaluated over either an arithmetic ring \mathbb{Z}_{2^ℓ} or a Boolean ring \mathbb{Z}_{2^1} , consisting of 2-input addition and multiplication gates. d denotes the multiplicative depth of ckt .

We use a collision-resistant hash function, denoted by $H()$ and a commitment scheme, denoted by $\text{com}()$, in

our protocols for practical efficiency. The details of the same can be found in Section B.1.

Security Model: For MPC, each party is modelled as a non-uniform probabilistic polynomial time (PPT) interactive Turing Machine. We operate in a static security model with an honest majority, where a PPT adversary \mathcal{A} can corrupt a party at the onset of the protocol. \mathcal{A} can be malicious in our setting i.e., the corrupt parties can arbitrarily deviate from the protocol specification. The computational security parameter is denoted by κ .

Robustness or Guaranteed Output Delivery: A protocol is said to be *robust* if all the parties can compute the output of the protocol irrespective of the behaviour of the adversary. The security of our protocols is proved in the standard real/ideal world paradigm. The details for the ideal world functionality $\mathcal{F}_{\text{robust}}$ that realizes the same in the 4PC setting is presented in Fig 19.

Shared Key Setup: We adopt a one-time key setup to minimize the overall communication of the protocol. We use three types of key setup namely, between i) a pair of parties, ii) a committee of three parties and iii) all the four parties. In each type, the parties in consideration can run an MPC protocol to agree on a randomness and use it as the key for pseudo-random function (PRF) to derive any subsequent co-related randomness. We model the protocol for the shared key setup as functionality $\mathcal{F}_{\text{setup}}$ and is presented in Fig 17.

3 Sharing Semantics

We use additive secret sharing of secrets over either an arithmetic ring \mathbb{Z}_{2^ℓ} or a Boolean ring \mathbb{Z}_{2^1} . We define two variants of secret sharing that are used in this work.

- **Additive sharing ([·]-sharing):** A value x is additively shared between two parties if $x = x^1 + x^2$, where one party holds the first share x^1 while the other party holds x^2 . We use $[x] = (x^1, x^2)$ to denote [·]-sharing of x .

- **Mirrored sharing ([·]-sharing):** A value x is said to be [·]-shared among the parties in \mathcal{P} if:

- There exist values σ_x, μ_x such that $\mu_x = x + \sigma_x$.
- σ_x is [·]-shared among parties in \mathbf{E} as $[\sigma_x]_{E_1} = \sigma_x^1$ and $[\sigma_x]_{E_2} = \sigma_x^2$, while parties in \mathbf{V} hold both σ_x^1 and σ_x^2 .
- μ_x is [·]-shared among parties in \mathbf{V} as $[\mu_x]_{V_1} = \mu_x^1$ and $[\mu_x]_{V_2} = \mu_x^2$, while parties in \mathbf{E} hold both μ_x^1 and μ_x^2 .

The shares of each party can be summarized as:

$$\begin{array}{l} \overline{E_1 : \llbracket x \rrbracket_{E_1} = (\sigma_x^1, \mu_x^1, \mu_x^2) \quad V_1 : \llbracket x \rrbracket_{V_1} = (\sigma_x^1, \sigma_x^2, \mu_x^1)} \\ \overline{E_2 : \llbracket x \rrbracket_{E_2} = (\sigma_x^2, \mu_x^1, \mu_x^2) \quad V_2 : \llbracket x \rrbracket_{V_2} = (\sigma_x^1, \sigma_x^2, \mu_x^2)} \end{array}$$

We use the notation $\llbracket x \rrbracket = ([\sigma_x], [\mu_x])$ to denote $\llbracket \cdot \rrbracket$ -sharing of value x . Sharing techniques and protocols for the boolean variant (\mathbb{Z}_{2^1}) are identical to their arithmetic counterparts apart from addition and subtraction operations being replaced with XOR and multiplication with AND. We use $\llbracket \cdot \rrbracket^{\mathbf{B}}$ to denote the sharing over a boolean ring.

• **Linearity of $\llbracket \cdot \rrbracket$ -sharing and $\llbracket \cdot \rrbracket^{\mathbf{B}}$ -sharing:** Given $[x] = (x^1, x^2)$, $[y] = (y^1, y^2)$ and public constants $c_1, c_2 \in \mathbb{Z}_{2^e}$, we have

$$[c_1x + c_2y] = (c_1x^1 + c_2y^1, c_1x^2 + c_2y^2) = c_1[x] + c_2[y]$$

Thus, $[c_1x + c_2y]$ and $c_1[x] + c_2[y]$ are equivalent and implies that parties can compute shares of any linear function of $\llbracket \cdot \rrbracket$ -shared values locally. It is easy to see that the linearity property extends to our $\llbracket \cdot \rrbracket$ -sharing as well.

4 Robust 4PC

In this section, we present a robust and efficient 4PC protocol with security against one malicious adversary. Our protocol incurs 12 ring elements per multiplication and removes the need for any additional setup of Broadcast, Digital Signatures, and Public-Key Setup, unlike [29]. We begin this section by introducing a "bi-convey primitive", which forms the core for the majority of our constructions. As mentioned in the introduction, bi-convey primitive enables two designated parties to send a value x to the third party with the aid of the fourth party. The remainder of the section describes a high-level overview of our protocol which is divided into 3 stages– i) input sharing, ii) circuit evaluation and iii) output computation. We elaborate on our primitive and each of the stages below:

4.1 Bi-Convey Primitive

Bi-convey primitive enables either i) two parties, say S_1, S_2 , to convey a value $x \in \mathbb{Z}_{2^e}$ to a designated party R or ii) allows party R to identify that one among S_1, S_2 is corrupt. The technical innovation of our construction for the 4 party case lies in using the fourth party available, say T , in an efficient manner. To elaborate, the protocol proceeds as follows. Parties S_1, S_2 both send the value x to R . In parallel, they send a commitment of the same

($\text{com}(x)$) to the fourth party T . Note that the randomness used to prepare the commitment is picked from the common source of the randomness of S_1, S_2 and R . If the received copies of x match, party R accepts the value and sends continue to T , and discards any message received from T . If not, R will identify that one among (S_1, S_2) is corrupt and thus T is honest. She then sends her internal randomness to T and waits for a message from T . Note that, the internal randomness of R which is forwarded to T , in our setting are all the keys of R (established during the shared key setup phase) that are not available with T . Party T , on the other hand, first checks if the commitments received from S_1, S_2 match or not. If they match, she will forward $\text{com}(x)$ to R else, she will identify that one among (S_1, S_2) is corrupt and thus sends her internal randomness to R . Now, if R receives $\text{com}(x)$ from T , then she will accept the version of x that matches with the received $\text{com}(x)$ and stops. If not, then both R and T have identified that one among (S_1, S_2) is corrupt.

\mathcal{F}_{bic} receives x, x', l_R and l_T from the parties S_1, S_2, R and T respectively. Here l_R and l_T denote the internal randomness of parties R and T respectively. \mathcal{F}_{bic} sets $\text{msg}_{S_1} = \text{msg}_{S_2} = \perp$.

- If $x = x'$, then \mathcal{F}_{bic} sets $\text{msg}_T = \perp$ and $\text{msg}_R = x$. Else it sets $\text{msg}_T = l_R$ and $\text{msg}_R = l_T$.
- \mathcal{F}_{bic} sends $\text{msg}_{S_1}, \text{msg}_{S_2}, \text{msg}_R$ and msg_T to parties S_1, S_2, R and T respectively.

Fig. 1. Functionality \mathcal{F}_{bic} : Ideal Functionality for party R to receive value x from S_1 and S_2 .

The formal protocol appears in Fig 2 and the details for corresponding ideal world functionality \mathcal{F}_{bic} appears in Fig 1.

- **Input:** Parties S_1, S_2, R and T input x, x, l_R and l_T respectively.
- **Output:** Parties S_1, S_2 receive \perp . Parties R and T receive x and \perp as outputs respectively, when S_1, S_2 are honest. For the case when one among S_1, S_2 is corrupt, party R obtains either x or l_T , while party T obtains either l_R or \perp , depending on the adversary's strategy.
- Parties S_1, S_2 send the value x to party R . In parallel, S_1, S_2 compute commitment of x , $\text{com}(x)$, using shared randomness known to R as well (sampled from the key shared amongst S_1, S_2 and R established during the shared key setup phase) and send it to T .
- If the received values match, party R sets $\text{msg}_R = \text{continue}$, accept the value x and discard any further message from T . Else, he sets $\text{msg}_R = l_R$, where l_R denotes the internal randomness of R .

- If the received commitments match, party T sets $\text{msg}_T = \text{com}(x)$, else sets $\text{msg}_T = l_T$, where l_T denotes the internal randomness of T .
- Parties R and T mutually exchange the msg values.
- If $\text{msg}_R = l_R$ and $\text{msg}_T = \text{com}(x)$, then R accepts the value x that is consistent with $\text{com}(x)$.

Fig. 2. $\Pi_{\text{bic}}(S_1, S_2, x, R, T)$: Protocol for S_1, S_2 to convey a value x to R with the help of T

We now provide a brief motivation for the need of bi-convey primitive in our framework. Looking ahead, the bi-convey primitive is used as a black-box in almost all of our subsequent protocol constructions. Consider the case where a call to this primitive from the outer protocol results in exchange of internal randomness among two parties. This implies both the parties conclude one among the remaining parties is corrupt and can safely trust each other. Thus both the honest parties combined, act as a single trusted party and use the received randomness to compute the inputs of all the parties in clear. Note that, both the honest parties together are able to compute the inputs in clear primarily because of the specific design of our mirrored sharing format (Section 3) where two parties together possess all the shares to reconstruct the inputs of the circuit. The honest parties then compute the final circuit output and send it to the remaining two parties ensuring guaranteed output delivery. We give a more detailed explanation of a use case of bi-convey primitive fitting in a larger protocol in Section 4.3.

4.2 Input Sharing

The goal is to robustly generate a $[\cdot]$ -sharing of a party's input. We call a party who wants to share the input as a Dealer. On a high level, if a dealer D wants to share a value x , parties start by locally sampling σ_x^1, σ_x^2 and μ_x^1 , according to the defined sharing semantics. The dealer then sets the last share as $\mu_x^2 = x + \sigma_x - \mu_x^1$. In case when the dealer is a verifier (say V_1), we enforce V_1 to send μ_x^2 to both the evaluators and $\text{com}(\mu_x^2)$ to V_2 . Now, all parties except V_1 , exchange $\text{com}(\mu_x^2)$ and compute the majority. If there exists no majority then V_1 is known to be corrupt and eliminated from the computation. The remaining parties can then run a semi-honest three-party protocol to compute the output. A similar idea follows for the case when the dealer is an evaluator. We provide the formal details of our Π_{sh} in Fig 3 below.

- **Input:** Party D inputs value x while others input \perp .
- **Output:** Parties obtain $[\![x]\!]$ as the output.
- **If $D = E_1$:** Parties in \mathbf{V} and E_1 locally sample σ_x^1 , while all the parties in \mathcal{P} locally sample σ_x^2 . Parties in \mathbf{V} and E_1 locally compute $\sigma_x = \sigma_x^1 + \sigma_x^2$. Similar steps are done for $D = E_2$.
- **If $D = V_i$ for $i \in \{1, 2\}$:** Parties in \mathbf{V} and E_1 locally sample σ_x^1 , while parties in \mathbf{V} and E_2 locally sample σ_x^2 . Parties in \mathbf{V} locally compute $\sigma_x = \sigma_x^1 + \sigma_x^2$.
- **If $D = V_1$:** Party V_1 computes $\mu_x = x + \sigma_x$. Parties in \mathbf{E} and V_1 locally sample μ_x^1 . Party V_1 computes and sends $\mu_x^2 = \mu_x - \mu_x^1$ to parties in \mathbf{E} and V_2 . Parties in \mathbf{E} and V_2 exchange the received copy of μ_x^2 . If there exists no majority, then they identify V_1 to be corrupt and engage in semi-honest 3PC excluding V_1 (with default input for V_1). Else, they set μ_x^2 to the computed majority. Similar steps are done for $D = V_2$.
- **If $D = E_i$ for $i \in \{1, 2\}$:** Party E_i computes $\mu_x = x + \sigma_x$. Parties in \mathbf{E} and V_1 locally sample μ_x^1 . Party E_i computes and sends $\mu_x^2 = \mu_x - \mu_x^1$ to V_2 and the co-evaluator. E_i sends $\text{com}(\mu_x^2)$ to V_1 . Parties other than the dealer exchange the commitment of μ_x^2 to compute majority (the co-evaluator and V_2 also exchange their copies of μ_x^2). If no majority exists, then they identify E_i to be corrupt and engage in semi-honest 3PC excluding E_i (with default input for E_i). Else, they set μ_x^2 to the computed majority.

Fig. 3. $\Pi_{\text{sh}}(D, x)$: Protocol to generate $[\![x]\!]$ by dealer D .

4.3 Circuit Evaluation

The circuit is evaluated in topological order where for every gate g the following invariant is maintained: given the $[\cdot]$ -sharing of the inputs, the output is generated in the $[\cdot]$ -shared format. When g is an addition gate ($z = x + y$), the linearity of $[\cdot]$ -sharing suffices to maintain this invariant.

- **Input:** Parties input their $[\![x]\!]$ and $[\![y]\!]$ shares.
- **Output:** Parties obtain $[\![z]\!]$ as the output, where $z = xy$.
- Parties in \mathbf{V} and E_1 collectively sample σ_z^1 and δ_{xy}^1 , while parties in \mathbf{V} and E_2 together sample σ_z^2 .
- Verifiers V_1, V_2 compute $\delta_{xy} = \sigma_x \sigma_y$, set $\delta_{xy}^2 = \delta_{xy} - \delta_{xy}^1$ and invoke $\Pi_{\text{bic}}(V_1, V_2, \delta_{xy}^2, E_2, E_1)$, which makes sure that E_2 receives δ_{xy}^2 .
- Parties in \mathbf{V} and E_1 collectively sample Δ_1 . Parties V_1 and E_1 compute $A_1 = -\mu_x^1 \sigma_y^1 - \mu_y^1 \sigma_x^1 + \delta_{xy}^1 + \sigma_z^1 + \Delta_1$ and invoke $\Pi_{\text{bic}}(V_1, E_1, A_1, E_2, V_2)$, such that E_2 receives A_1 .
- Similarly, parties in \mathbf{V} and E_2 collectively sample Δ_2 . Parties V_1 and E_2 compute $A_2 = -\mu_x^1 \sigma_y^2 - \mu_y^1 \sigma_x^2 + \delta_{xy}^2 + \sigma_z^2 + \Delta_2$ and invoke $\Pi_{\text{bic}}(V_1, E_2, A_2, E_1, V_2)$, such that E_1 receives A_2 .
- Parties V_2 and E_1 compute $B_1 = -\mu_x^2 \sigma_y^1 - \mu_y^2 \sigma_x^1 - \Delta_1$ and invoke $\Pi_{\text{bic}}(V_2, E_1, B_1, E_2, V_1)$. Similarly, V_2 and

E_2 compute $B_2 = -\mu_x^2\sigma_y^2 - \mu_y^2\sigma_x^2 - \Delta_2$ and invoke $\Pi_{\text{bic}}(V_2, E_2, B_2, E_1, V_1)$.

– Evaluators compute $\mu_z = A_1 + A_2 + B_1 + B_2 + \mu_x\mu_y$ locally. Parties in \mathbf{E} and V_1 collectively sample μ_z^1 followed by evaluators setting $\mu_z^2 = \mu_z - \mu_z^1$ and invoking $\Pi_{\text{bic}}(E_1, E_2, \mu_z^2, V_2, V_1)$ for V_2 to receive μ_z^2 .

Fig. 4. $\Pi_{\text{mult}}(x, y, z)$: Multiplication Protocol

For a multiplication gate g ($z = xy$), the goal is for the evaluators to robustly compute μ_z where

$$\begin{aligned}\mu_z &= xy + \sigma_z = (\mu_x - \sigma_x)(\mu_y - \sigma_y) + \sigma_z \\ &= \mu_x\mu_y - \mu_x\sigma_y - \mu_y\sigma_x + \sigma_x\sigma_y + \sigma_z\end{aligned}$$

followed by evaluators setting μ_z^2 share and robustly sending it to V_2 . On a high level, we view the aforementioned equation of μ_z as: $\mu_z = \mu_x\mu_y + A + B$, where $A = -\mu_x^1\sigma_y - \mu_y^1\sigma_x + \delta_{xy} + \sigma_z + \Delta$ is solely possessed by V_1 and $B = -\mu_x^2\sigma_y - \mu_y^2\sigma_x - \Delta$ is possessed by V_2 . In order for evaluators to compute μ_z , E_1 and E_2 need to robustly receive $A + B$. Note that $\mu_x\mu_y$ is already available with the evaluators. Thus A is further split into $A_1 + A_2$, such that each $A_j \in \{1, 2\}$ is possessed by V_1 and E_j . Similarly, B is split such that each $B_j \in \{1, 2\}$ is possessed by V_2 and E_j . Now parties need to simply invoke Π_{bic} protocol, one for each A_j and B_j with the co-evaluator acting as the receiving party. Thus evaluators are able to compute $A + B$ correctly. After computing μ_z , the evaluators set $\mu_z^2 = \mu_z - \mu_z^1$ and call Π_{bic} protocol to send μ_z^2 to V_2 , where μ_z^1 is collectively sampled by parties in \mathbf{E} and V_1 . We provide the formal details of our $\Pi_{\text{mult}}(x, y, z)$ in Fig 4. For correctness of μ_z ,

$$\begin{aligned}\mu_z &= xy + \sigma_z = (\mu_x - \sigma_x)(\mu_y - \sigma_y) + \sigma_z \\ &= \mu_x\mu_y - \mu_x\sigma_y - \mu_y\sigma_x + \sigma_x\sigma_y + \sigma_z \\ &= (-\mu_x^1\sigma_y - \mu_y^1\sigma_x + \delta_{xy}^1 + \sigma_z^1 + \Delta_1 + \Delta_2) \\ &\quad + (-\mu_x^2\sigma_y - \mu_y^2\sigma_x + \delta_{xy}^2 + \sigma_z^2 - \Delta_1 - \Delta_2) \\ &= \mu_x\mu_y + (A_1 + A_2) + (B_1 + B_2)\end{aligned}$$

where $A_j = -\mu_x^1\sigma_y^j - \mu_y^1\sigma_x^j + \delta_{xy}^j + \sigma_z^j + \Delta_j$ and $B_j = -\mu_x^2\sigma_y^j - \mu_y^2\sigma_x^j - \Delta_j$ for $j \in \{1, 2\}$. The evaluators receive A_1, A_2, B_1 and B_2 , whose correctness is guaranteed by Π_{bic} protocol. Thus the evaluators can correctly compute $\mu_z = \mu_x\mu_y + (A_1 + A_2) + (B_1 + B_2)$. Verifier V_2 also correctly receives μ_z^2 share from the evaluators, by the underlying correctness guarantee of Π_{bic} protocol.

We now analyze how Π_{bic} primitive fits into the larger Π_{mult} protocol to make it robust. Consider Step 2 of the protocol Π_{mult} where parties invoke $\Pi_{\text{bic}}(V_1, V_2, \delta_{xy}^2, E_2, E_1)$. As mentioned in Section 4.1, primitive Π_{bic} guarantees that either i) party E_2 receives the correct value δ_{xy}^2 or ii) both E_1 and E_2 identify that one among (V_1, V_2) is corrupt. In the first case, parties

can proceed with the execution of the protocol. For the second case, parties E_1 and E_2 mutually exchange their internal randomness (this includes the keys established during the shared key setup phase). Using the received randomness, both E_1 and E_2 can compute the missing part of her share corresponding to the $[\cdot]$ -sharing of the inputs and hence obtain all the inputs in clear. Given the inputs in clear, both E_1 and E_2 can compute the function output in clear and send it to the remaining two parties.

4.4 Output Computation

The output computation phase is comparatively simple. The missing share of the output with respect to each party is possessed by the remaining three parties. Thus two out of the three parties send the missing share and the third party sends the corresponding hash. Thus each party sets the missing share as the majority among the received values and reconstruct the output. The formal details of our robust output computation protocol Π_{oc} is given in Fig 5.

- **Input:** Parties input their $[\mathbf{z}]$ shares.
- **Output:** Parties obtain \mathbf{z} as the output.
 - For $i, j \in \{1, 2\}$ and $i \neq j$, E_i receives σ_z^j from parties in \mathbf{V} and $H(\sigma_z^j)$ from E_j .
 - V_2 receives μ_z^1 from parties in \mathbf{E} and $H(\mu_z^1)$ from V_1 .
 - V_1 receives μ_z^2 from parties in \mathbf{E} and $H(\mu_z^2)$ from V_2 .
 - Each party sets the missing share as the majority among the received values and outputs $\mathbf{z} = \mu_z^1 + \mu_z^2 - \sigma_z^1 - \sigma_z^2$.

Fig. 5. Π_{oc} : Protocol for Robust Reconstruction

5 ML Building Blocks

In this section, we provide constructions for our crucial building blocks necessary to achieve secure training and prediction for algorithms namely– i) Linear Regression, ii) Logistic Regression, iii) Deep Neural Network (DNN) and iv) Binarized Neural Network (BNN). We provide the formal details of the corresponding lemmas and proofs to Appendix C.

5.1 Arithmetic/Boolean Couple Sharing Primitive

Two parties, either $\{V_1, V_2\}$ (set \mathbf{V}) or $\{E_1, E_2\}$ (set \mathbf{E}) own a common value x and want to create a $[\![\cdot]\!]$ -sharing of x . We abstract out this procedure (Fig 6) and define it as *couple sharing* of a value.

<p>Case 1: ($\mathbf{S} = \mathbf{E}$)</p> <ul style="list-style-type: none"> • Input: E_1 and E_2 input x while others input \perp. • Output: Parties obtain $[\![x]\!]$ as the output. <p>– Parties set $\sigma_x^1 = 0$ and $\sigma_x^2 = 0$. Parties in \mathbf{E} and V_1 collectively sample random $\mu_x^1 \in \mathbb{Z}_{2^\ell}$.</p> <p>– E_1 and E_2 set $\mu_x^2 = x - \mu_x^1$. Parties then execute $\Pi_{\text{bic}}(E_1, E_2, \mu_x^2, V_2, V_1)$, such that V_2 receives μ_x^2.</p> <p>Case 2: ($\mathbf{S} = \mathbf{V}$)</p> <ul style="list-style-type: none"> • Input: V_1 and V_2 input x while others input \perp. • Output: Parties obtain $[\![x]\!]$ as the output. <p>– Parties set $\mu_x^1 = 0$ and $\mu_x^2 = 0$. Parties in \mathbf{V} and E_1 collectively sample random $\sigma_x^1 \in \mathbb{Z}_{2^\ell}$.</p> <p>– V_1 and V_2 set $\sigma_x^2 = x - \sigma_x^1$. Parties then execute $\Pi_{\text{bic}}(V_1, V_2, \sigma_x^2, E_2, E_1)$, such that E_2 receives σ_x^2.</p>
--

Fig. 6. $\Pi_{\text{csh}}(\mathbf{S}, x)$: Protocol to generate couple sharing of x

On a high level when set $\mathbf{S} = \mathbf{E}$, in order to share a value x , parties set $\sigma_x^1 = \sigma_x^2 = 0$. A random μ_x^1 is collectively sampled and the owners of the value set μ_x^2 such that $\mu_x^1 + \mu_x^2 = x$ and send μ_x^2 to V_2 using Π_{bic} protocol. The shares of parties can be viewed as:

$$\begin{array}{ll} E_1 : [\![x]\!]_{E_1} = (0, \mu_x^1, \mu_x^2) & V_1 : [\![x]\!]_{V_1} = (0, 0, \mu_x^1) \\ E_2 : [\![x]\!]_{E_2} = (0, \mu_x^1, \mu_x^2) & V_2 : [\![x]\!]_{V_2} = (0, 0, \mu_x^2) \end{array}$$

For the case when set $\mathbf{S} = \mathbf{V}$ and value x , parties in \mathbf{V} and E_1 collectively sample random σ_x^1 followed by \mathbf{V} setting $\sigma_x^2 = -x - \sigma_x^1$ and robustly sending it to E_2 .

$$\begin{array}{ll} E_1 : [\![x]\!]_{E_1} = (\sigma_x^1, 0, 0) & V_1 : [\![x]\!]_{V_1} = (\sigma_x^1, \sigma_x^2, 0) \\ E_2 : [\![x]\!]_{E_2} = (\sigma_x^2, 0, 0) & V_2 : [\![x]\!]_{V_2} = (\sigma_x^1, \sigma_x^2, 0) \end{array}$$

5.2 Dot Product

Given vectors \vec{x} and \vec{y} , each of size d , the goal is to compute the dot product $z = \vec{x} \odot \vec{y} = \sum_{i=1}^d x_i y_i$. The recent works of ABY3 and ASTRA have tackled dot product computation in the semi-honest setting with cost equal to that of a single multiplication thus, making the total cost independent of the vector size. However, in the malicious setting, their techniques become expensive, with cost dependent on the vector size. In this work, we remove this dependency and retain the cost to be the

same as that of a single multiplication. This independence stems from the peculiar structure of our sharing and our robust multiplication method. On a high level, instead of calling Π_{bic} protocol for A_{1i}, A_{2i}, B_{1i} and B_{2i} corresponding to each product $z_i = x_i y_i$, the parties add up their shares and then invoke Π_{bic} once for each of the summed up share. To facilitate this modification, verifiers also adjust $\delta_{xy}^2 = \sum_{i=1}^d \delta_{x_i y_i} - \delta_{xy}^1$ before sending to E_2 . Formal details are presented in Fig 7 below.

<ul style="list-style-type: none"> • Input: Parties input their $[\![\vec{x}]\!]$ and $[\![\vec{y}]\!]$ shares. • Output: Parties obtain $[\![z]\!]$ as output, where $z = \vec{x} \odot \vec{y}$. <p>– Parties in \mathbf{V} and E_1 collectively sample σ_z^1 and δ_{xy}^1, while parties in \mathbf{V} and E_2 together sample σ_z^2.</p> <p>– Verifiers V_1, V_2 compute $\delta_{xy} = \sum_{i=1}^d \sigma_{x_i} \sigma_{y_i}$, set $\delta_{xy}^2 = \delta_{xy} - \delta_{xy}^1$ and invoke $\Pi_{\text{bic}}(V_1, V_2, \delta_{xy}^2, E_2, E_1)$, such that E_2 receives δ_{xy}^2.</p> <p>– Parties in \mathbf{V} and E_1 collectively sample Δ_1. Parties V_1 and E_1 compute $A_1 = \sum_{i=1}^d (-\mu_{x_i}^1 \sigma_{y_i}^1 - \mu_{y_i}^1 \sigma_{x_i}^1) + \sigma_z^1 + \delta_{xy}^1 + \Delta_1$ and invoke $\Pi_{\text{bic}}(V_1, E_1, A_1, E_2, V_2)$, such that E_2 receives A_1.</p> <p>– Similarly, parties in \mathbf{V} and E_2 collectively sample Δ_2. Parties V_1 and E_2 compute $A_2 = \sum_{i=1}^d (-\mu_{x_i}^1 \sigma_{y_i}^2 - \mu_{y_i}^1 \sigma_{x_i}^2) + \sigma_z^2 + \delta_{xy}^2 + \Delta_2$ and invoke $\Pi_{\text{bic}}(V_1, E_2, A_2, E_1, V_2)$, such that E_1 receives A_2.</p> <p>– V_2 and E_1 compute $B_1 = \sum_{i=1}^d (-\mu_{x_i}^2 \sigma_{y_i}^1 - \mu_{y_i}^2 \sigma_{x_i}^1) - \Delta_1$ and invoke $\Pi_{\text{bic}}(V_2, E_1, B_1, E_2, V_1)$. Similarly, V_2 and E_2 compute $B_2 = \sum_{i=1}^d (-\mu_{x_i}^2 \sigma_{y_i}^2 - \mu_{y_i}^2 \sigma_{x_i}^2) - \Delta_2$ and execute $\Pi_{\text{bic}}(V_2, E_2, B_2, E_1, V_1)$.</p> <p>– Evaluators compute $\mu_z = \mu_x \mu_y + A_1 + A_2 + B_1 + B_2$ locally. Parties in \mathbf{E} and V_1 collectively sample μ_z^1 followed by evaluators setting $\mu_z^2 = \mu_z - \mu_z^1$ and execute $\Pi_{\text{bic}}(E_1, E_2, \mu_z^2, V_2, V_1)$ for V_2 to receive μ_z^2.</p>
--

Fig. 7. $\Pi_{\text{dp}}([\![\vec{x}]\!], [\![\vec{y}]\!])$: Dot Product of two vectors

5.3 MSB Extraction

The goal is to check if $u < v$, given two elements $[\![u]\!]$ and $[\![v]\!]$. Most state-of-the-art protocols [43, 45] adopt a circuit based approach to perform comparison which is a bottleneck for efficiency. However recently, ASTRA [14] proposed a solution for 3 parties that is free of any circuits. Inspired from their idea, we present a solution that consumes only $16\ell + 4$ bits in total. The problem of comparison can be reduced to checking the MSB of the value, represented as $\text{msb}(a)$, where $a = u - v$. If $u < v$, then $\text{msb}(a) = 1$, else 0. Evaluators collectively sample random $r \in \mathbb{Z}_{2^\ell}$, compute $\text{msb}(r)$ and generate $[\![\text{msb}(r)]\!]^{\mathbf{B}}$. Parties then execute Π_{mult} on r and a , followed by reconstruction of ra towards V_1 and V_2 . Verifiers then compute $\text{msb}(ra)$ and generate $[\![\text{msb}(ra)]\!]^{\mathbf{B}}$.

Parties locally XOR their boolean shares of $\text{msb}(r)$ and $\text{msb}(ra)$ to obtain $\llbracket \text{msb}(a) \rrbracket^{\mathbf{B}}$. The formal protocol is presented in Figure 8 below.

- **Input:** Parties input their $\llbracket a \rrbracket$ shares.
- **Output:** Parties obtain $\llbracket \text{msb}(a) \rrbracket^{\mathbf{B}}$ as the output.
 - Parties in \mathbf{E} sample random $r \in \mathbb{Z}_{2^\ell}$ and set $p = \text{msb}(r)$.
 - Parties execute $\Pi_{\text{cSh}}(\mathbf{E}, r)$ and $\Pi_{\text{cSh}}^{\mathbf{B}}(\mathbf{E}, p)$ to generate $\llbracket r \rrbracket$ and $\llbracket p \rrbracket^{\mathbf{B}}$ respectively.
 - Parties execute $\Pi_{\text{mult}}(\llbracket r \rrbracket, \llbracket a \rrbracket)$ to generate $\llbracket ra \rrbracket$. Parties also execute $\Pi_{\text{bic}}(\mathbf{E}_1, \mathbf{E}_2, \mu_{ra}^2, \mathbf{V}_1, \mathbf{V}_2)$ and $\Pi_{\text{bic}}(\mathbf{E}_1, \mathbf{E}_2, \mu_{ra}^1, \mathbf{V}_2, \mathbf{V}_1)$ to reconstruct ra towards \mathbf{V}_1 and \mathbf{V}_2 respectively. Verifiers then set $q = \text{msb}(ra)$.
 - Parties execute $\Pi_{\text{cSh}}^{\mathbf{B}}(\mathbf{V}, q)$ to generate $\llbracket q \rrbracket^{\mathbf{B}}$ followed by locally computing $\llbracket \text{msb}(a) \rrbracket^{\mathbf{B}} = \llbracket p \rrbracket^{\mathbf{B}} \oplus \llbracket q \rrbracket^{\mathbf{B}}$.

Fig. 8. $\Pi_{\text{msb}}(\llbracket a \rrbracket)$: Extraction of MSB from a value

5.4 Truncation

We use ℓ -bit integers in signed 2's complement form to represent a decimal value where the sign of the decimal value is represented by the most significant bit (MSB). Consider a decimal value z represented in the signed 2's complement form. We use d_z to denote the least significant bits that represent its fractional part and $i_z = \ell - d_z$ to represent its integral part. It is observed that in the face of repeated multiplications, d_z and i_z needed to represent the output z keeps doubling with every multiplication and can eventually lead to an overflow. To avoid this multiplication overflow while preserving the accuracy and correctness, truncation is performed at the output of a multiplication gate. Truncation of a value z is defined as $z^t = z/2^{d_z}$, where the value z is *right arithmetic shifted* by d_z bits.

SecureML [45] proposed an efficient truncation method for the two-party setting, where the parties locally truncate the shares after a multiplication. They showed that this technique introduces at most 1 bit error in the least significant bit (LSB) position and thus causes a minor reduction in the accuracy. Later ABY3 [43] showed that this idea cannot be trivially extended to three party setting and proposed an alternative technique to achieve truncation. Their main idea revolves around generating $(\llbracket r \rrbracket, \llbracket r^t \rrbracket)$ pair, where r is a random ring element and $r^t = r/2^d$. Parties then compute $z - r$ in clear and locally truncate it to obtain $(z - r)^t$. This is followed by generating $\llbracket (z - r)^t \rrbracket$ and adding it to $\llbracket r^t \rrbracket$ to obtain $\llbracket z^t \rrbracket$. Similar to SecureML, this technique may also incur a one-bit error in the LSB position of z^t . To generate $(\llbracket r \rrbracket, \llbracket r^t \rrbracket)$, ABY3 requires two expensive cir-

cuit evaluations and leading to a total cost of more than 100 ring elements per multiplication. While we adopt ABY3's idea of using (r, r^t) pair in our Π_{mult} protocol to achieve truncation, we remove the need of expensive circuits and maintain the total cost to 14 ring elements.

We begin with the generation of (r, r^t) pair. Parties in \mathbf{V} and \mathbf{E}_1 sample random $r_1 \in \mathbb{Z}_{2^\ell}$, while parties in \mathbf{V} and \mathbf{E}_2 sample r_2 . Verifiers \mathbf{V}_1 and \mathbf{V}_2 set $r = r_1 + r_2$. Then parties \mathbf{V}_1 and \mathbf{V}_2 locally truncate r to obtain r^t and execute Π_{cSh} to generate $\llbracket r^t \rrbracket$. Thus, the pair $(\llbracket r \rrbracket, \llbracket r^t \rrbracket)$ is generated. Unlike Π_{mult} (Fig 4), evaluators instead reconstruct $(z - r)$, followed by locally truncating it to obtain $(z - r)^t$. Evaluators execute Π_{cSh} to generate $\llbracket (z - r)^t \rrbracket$ followed by locally adding to $\llbracket r^t \rrbracket$ to obtain $\llbracket z^t \rrbracket$. The formal details of our protocol Π_{mulTr} appears in Fig 9 below.

- **Input:** Parties input their $\llbracket x \rrbracket$ and $\llbracket y \rrbracket$ shares.
- **Output:** Parties obtain $\llbracket z^t \rrbracket$ as output, where $z^t = (xy)^t$.
 - Parties in \mathbf{V} and \mathbf{E}_1 collectively sample σ_z^1 and r_1 , while parties in \mathbf{V} and \mathbf{E}_2 together sample σ_z^2 and r_2 .
 - Verifiers set $r = r_1 + r_2$ and truncate r by d bits to obtain r^t . Parties execute $\Pi_{\text{cSh}}(\mathbf{V}, r^t)$ to generate $\llbracket r^t \rrbracket$ sharing.
 - Verifiers locally set $\delta_{xy} = \sigma_x \cdot \sigma_y$ and compute $\delta_{xy}^2 = \delta_{xy} - \delta_{xy}^1$, where δ_{xy}^1 is collectively sampled by parties in \mathbf{V} and \mathbf{E}_1 . Parties then execute $\Pi_{\text{bic}}(\mathbf{V}_1, \mathbf{V}_2, \delta_{xy}^2, \mathbf{E}_2, \mathbf{E}_1)$, such that \mathbf{E}_2 receives δ_{xy}^2 .
 - Parties in \mathbf{V} and \mathbf{E}_1 collectively sample Δ_1 . Parties \mathbf{V}_1 and \mathbf{E}_1 compute $A_1 = -\mu_x^1 \sigma_y^1 - \mu_y^1 \sigma_x^1 + \delta_{xy}^1 - r_1 + \Delta_1$ and execute $\Pi_{\text{bic}}(\mathbf{V}_1, \mathbf{E}_1, A_1, \mathbf{E}_2, \mathbf{V}_2)$, such that \mathbf{E}_2 receives A_1 .
 - Similarly, parties in \mathbf{V} and \mathbf{E}_2 collectively sample Δ_2 . Parties \mathbf{V}_1 and \mathbf{E}_2 compute $A_2 = -\mu_x^1 \sigma_y^2 - \mu_y^1 \sigma_x^2 + \delta_{xy}^2 - r_2 + \Delta_2$ and execute $\Pi_{\text{bic}}(\mathbf{V}_1, \mathbf{E}_2, A_2, \mathbf{E}_1, \mathbf{V}_2)$, such that \mathbf{E}_1 receives A_2 .
 - Parties \mathbf{V}_2 and \mathbf{E}_1 compute $B_1 = -\mu_x^2 \sigma_y^1 - \mu_y^2 \sigma_x^1 - \Delta_1$ and execute $\Pi_{\text{bic}}(\mathbf{V}_2, \mathbf{E}_1, B_1, \mathbf{E}_2, \mathbf{V}_1)$. Similarly, \mathbf{V}_2 and \mathbf{E}_2 compute $B_2 = -\mu_x^2 \sigma_y^2 - \mu_y^2 \sigma_x^2 - \Delta_2$ and execute $\Pi_{\text{bic}}(\mathbf{V}_2, \mathbf{E}_2, B_2, \mathbf{E}_1, \mathbf{V}_1)$.
 - Evaluators compute $z - r = \mu_x \mu_y + A_1 + A_2 + B_1 + B_2$ and truncate it by d bits to obtain $(z - r)^t$.
 - Parties execute $\Pi_{\text{cSh}}(\mathbf{E}, (z - r)^t)$ to generate $\llbracket (z - r)^t \rrbracket$ sharing and locally add to obtain $\llbracket z^t \rrbracket = \llbracket (z - r)^t \rrbracket + \llbracket r^t \rrbracket$

Fig. 9. $\Pi_{\text{mulTr}}^A(x, y)$: Truncation Protocol

5.5 Bit Conversion

Here, we describe a protocol to transform $\llbracket \cdot \rrbracket^{\mathbf{B}}$ -sharing of bit b to its arithmetic equivalent. For this transformation, we use the following equivalence relation:

$$b = \sigma_b \oplus \mu_b = \mu_{b'} + \sigma_{b'} - 2\mu_{b'}\sigma_{b'}$$

where $\mu_{b'}$ and $\sigma_{b'}$ denote the bits μ_b and σ_b respectively over \mathbb{Z}_{2^ℓ} . Parties who hold μ_b and σ_b in clear convert them to $\mu_{b'}$ and $\sigma_{b'}$ respectively. Parties generate $[\![\cdot]\!]$ -sharing of $\sigma_{b'}$ and $\mu_{b'}$ by executing Π_{cSh} followed by multiplication of $[\![\mu_{b'}]\!]$ and $[\![\sigma_{b'}]\!]$. We call the resultant protocol as Π_{btr} and the formal details are given below.

- **Input:** Parties input their $[\![b]\!]^{\mathbf{B}}$ shares.
 - **Output:** Parties obtain $[\![b]\!]$ as the output.
- Parties execute $\Pi_{\text{cSh}}(\mathbf{V}, \sigma_{b'})$ and $\Pi_{\text{cSh}}(\mathbf{E}, \mu_{b'})$ to generate $[\![\sigma_{b'}]\!]$ and $[\![\mu_{b'}]\!]$ respectively.
 - Parties execute $\Pi_{\text{mult}}([\![\mu_{b'}]\!], [\![\sigma_{b'}]\!])$ to generate $[\![\mu_{b'}\sigma_{b'}]\!]$, followed by locally computing $[\![b]\!] = [\![\mu_{b'}]\!] + [\![\sigma_{b'}]\!] - 2[\![\mu_{b'}\sigma_{b'}]\!]$.

Fig. 10. $\Pi_{\text{btr}}([\![b]\!]^{\mathbf{B}})$: Conversion of a bit to arithmetic equivalent

We observe that cost of multiplication in Π_{btr} can be reduced from 12ℓ to 10ℓ bits. Note that the value $\sigma_{\mu_{b'}}$ is set to zero, when Π_{cSh} is executed to generate $[\![\mu_{b'}]\!]$. This implies $\delta_{\mu_{b'}\sigma_{b'}} = 0$ and thus removes the extra call to Π_{bic} protocol.

5.6 Bit Insertion

Given a bit $b \in \{0, 1\}$ in $[\![\cdot]\!]^{\mathbf{B}}$ -shared form and $x \in \mathbb{Z}_{2^\ell}$ in $[\![\cdot]\!]$ -shared form, we have to compute $[\![bx]\!]$. A trivial solution is to convert $[\![b]\!]^{\mathbf{B}}$ to $[\![b]\!]$ using Π_{btr} followed by a multiplication with $[\![x]\!]$, which requires a total of 26 ring elements and 10 rounds. Instead, we propose a better solution that requires 18ℓ ring elements and 5 rounds in total. We can view the equation for bit insertion as follows:

$$\begin{aligned}
\mu_{bx} &= (\mu_b \oplus \sigma_b) \cdot (\mu_x - \sigma_x) + \sigma_{bx} \\
&= (\mu_{b'} + \sigma_{b'} - 2\mu_{b'}\sigma_{b'}) \cdot (\mu_x - \sigma_x) + \sigma_{bx} \\
&= \gamma_{b'x} - \mu_{b'}\sigma_x + (\mu_x - 2\gamma_{b'x})\sigma_{b'} + (2\mu_{b'} - 1)\delta_{b'x} + \sigma_{bx} \\
&= \gamma_{b'x} + (-\mu_{b'}^1\sigma_x + (\mu_x^1 - 2\gamma_{b'x}^1)\sigma_{b'} + (2\mu_{b'}^1 - 1)\delta_{b'x} + \sigma_{bx}) \\
&\quad + (-\mu_{b'}^2\sigma_x + (\mu_x^2 - 2\gamma_{b'x}^2)\sigma_{b'} + (2\mu_{b'}^2 - 1)\delta_{b'x}) \\
&= \gamma_{b'x} + (\mathbf{A}_1 + \mathbf{A}_2) + (\mathbf{B}_1 + \mathbf{B}_2)
\end{aligned}$$

where $\gamma_{b'x} = \mu_{b'}\mu_x$, $\delta_{b'x} = \sigma_{b'}\sigma_x$ and $\mu_{b'}$, $\sigma_{b'}$ represent μ_b and σ_b over \mathbb{Z}_{2^ℓ} respectively. In the above equation, we observe that, given the $[\![\cdot]\!]$ -shares of $\mu_{b'}$, $\sigma_{b'}$, $\gamma_{b'x}$ and $\delta_{b'x}$, parties can robustly compute $[\![\cdot]\!]$ -sharing of μ_{bx} . The protocol proceeds as follows: Parties begin by generating $[\![\cdot]\!]$ -shares of $\mu_{b'}$, $\gamma_{b'x}$ towards set \mathbf{V} and $\sigma_{b'}$, $\delta_{b'x}$ towards set \mathbf{E} , so that parties can compute \mathbf{A}_1 , \mathbf{A}_2 , \mathbf{B}_1 and \mathbf{B}_2 . This is followed by parties executing Π_{bic} protocol for each \mathbf{A}_i and \mathbf{B}_i , so that \mathbf{E}_1 and \mathbf{E}_2 are able to compute μ_{bx} . The formal details appear in Fig 11.

- **Input:** Parties input their $[\![b]\!]^{\mathbf{B}}$ and $[\![x]\!]$ shares.
 - **Output:** Parties obtain $[\![bx]\!]$ as the output.
- Parties in \mathbf{V} and \mathbf{E}_1 collectively sample random $\sigma_{bx}^1 \in \mathbb{Z}_{2^\ell}$, while parties in \mathbf{V} and \mathbf{E}_2 together sample random σ_{bx}^2 .
 - Parties in \mathbf{V} and \mathbf{E}_1 collectively sample random $\sigma_{b'}^1$ followed by \mathbf{V}_1 and \mathbf{V}_2 setting $\sigma_{b'}^2 = \sigma_{b'} - \sigma_{b'}^1$. Parties then execute $\Pi_{\text{bic}}(\mathbf{V}_1, \mathbf{V}_2, \sigma_{b'}^2, \mathbf{E}_2, \mathbf{E}_1)$, such that \mathbf{E}_2 receives $\sigma_{b'}^2$. The same procedure is used for \mathbf{E}_2 to receive $\delta_{b'x}^2$.
 - Parties in \mathbf{E} and \mathbf{V}_1 collectively sample random $\mu_{b'}^1$ followed by \mathbf{E}_1 and \mathbf{E}_2 setting $\mu_{b'}^2 = \mu_{b'} - \mu_{b'}^1$. Parties then execute $\Pi_{\text{bic}}(\mathbf{E}_1, \mathbf{E}_2, \mu_{b'}^2, \mathbf{V}_2, \mathbf{V}_1)$, such that \mathbf{V}_2 receives $\mu_{b'}^2$. The same procedure is used for \mathbf{V}_2 to receive $\gamma_{b'x}^2$.
 - Parties in \mathbf{V} and \mathbf{E}_1 collectively sample Δ_1 . Parties \mathbf{V}_1 and \mathbf{E}_1 compute $\mathbf{A}_1 = -\mu_{b'}^1\sigma_x^1 + (\mu_x^1 - 2\gamma_{b'x}^1)\sigma_{b'}^1 + (2\mu_{b'}^1 - 1)\delta_{b'x}^1 + \sigma_{bx}^1 + \Delta_1$ and invoke $\Pi_{\text{bic}}(\mathbf{V}_1, \mathbf{E}_1, \mathbf{A}_1, \mathbf{E}_2, \mathbf{V}_2)$.
 - Similarly, parties in \mathbf{V} and \mathbf{E}_2 collectively sample Δ_2 . Parties \mathbf{V}_1 and \mathbf{E}_2 compute $\mathbf{A}_2 = -\mu_{b'}^1\sigma_x^2 + (\mu_x^1 - 2\gamma_{b'x}^1)\sigma_{b'}^1 + (2\mu_{b'}^1 - 1)\delta_{b'x}^1 + \sigma_{bx}^2 + \Delta_2$ and invoke $\Pi_{\text{bic}}(\mathbf{V}_1, \mathbf{E}_2, \mathbf{A}_2, \mathbf{E}_1, \mathbf{V}_2)$.
 - Parties \mathbf{V}_2 and \mathbf{E}_1 compute $\mathbf{B}_1 = -\mu_{b'}^2\sigma_x^1 + (\mu_x^2 - 2\gamma_{b'x}^2)\sigma_{b'}^1 + (2\mu_{b'}^2 - 1)\delta_{b'x}^1 - \Delta_1$ and invoke $\Pi_{\text{bic}}(\mathbf{V}_2, \mathbf{E}_1, \mathbf{B}_1, \mathbf{E}_2, \mathbf{V}_1)$. Similarly, \mathbf{V}_2 and \mathbf{E}_2 compute $\mathbf{B}_2 = -\mu_{b'}^2\sigma_x^2 + (\mu_x^2 - 2\gamma_{b'x}^2)\sigma_{b'}^2 + (2\mu_{b'}^2 - 1)\delta_{b'x}^2 - \Delta_2$ and invoke $\Pi_{\text{bic}}(\mathbf{V}_2, \mathbf{E}_2, \mathbf{B}_2, \mathbf{E}_1, \mathbf{V}_1)$.
 - Evaluators compute $\mu_{b'x} = \mathbf{A}_1 + \mathbf{A}_2 + \mathbf{B}_1 + \mathbf{B}_2 + \gamma_{b'x}$ locally. Parties in \mathbf{E} and \mathbf{V}_1 collectively sample $\mu_{b'x}^1$ followed by evaluators setting $\mu_{b'x}^2 = \mu_{b'x} - \mu_{b'x}^1$ and invoking $\Pi_{\text{bic}}(\mathbf{E}_1, \mathbf{E}_2, \mu_{b'x}^2, \mathbf{V}_2, \mathbf{V}_1)$.

Fig. 11. $\Pi_{\text{bin}}([\![b]\!]^{\mathbf{B}}, [\![x]\!])$: Insertion of bit b in a value

6 Secure Prediction

In this section, we provide detailed protocols for the prediction phase of the following ML algorithms – i) Linear Regression, ii) Logistic Regression, iii) Deep Neural Network and iv) Binarized Neural Network, using the building blocks constructed earlier in Section 5.

6.1 Our Model

We consider a server-aided setting where both model owner M and client C outsource their trained model parameters and query to a set of four non-colluding servers $\{\mathbf{V}_1, \mathbf{V}_2, \mathbf{E}_1, \mathbf{E}_2\}$, in a $[\![\cdot]\!]$ -shared fashion. The servers then compute the function using our 4PC protocol and finally reconstruct the result towards C . We assume the existence of a malicious adversary \mathcal{A} , who can corrupt either M or C and at most one among $\{\mathbf{V}_1, \mathbf{V}_2, \mathbf{E}_1, \mathbf{E}_2\}$. Recall that \mathbf{E} and \mathbf{V} denote the set of servers $\{\mathbf{E}_1, \mathbf{E}_2\}$ and $\{\mathbf{V}_1, \mathbf{V}_2\}$ respectively. We begin with the assump-

tion that both M and C have already outsourced their input vectors to $\{V_1, V_2, E_1, E_2\}$.

Notations: We use bold smalls to denote a vector. Given a vector $\vec{\mathbf{a}}$, the i^{th} element in the vector is denoted by \mathbf{a}_i . Model Owner M holds a vector of *trained* model parameters denoted by $\vec{\mathbf{w}}$. C 's query is denoted by $\vec{\mathbf{z}}$. Both $\vec{\mathbf{w}}$ and $\vec{\mathbf{z}}$ are vectors of size d , where d denotes the number of features.

6.2 Linear Regression

In case of linear regression model, the output of the prediction phase for a query $\vec{\mathbf{z}}$ is given by $\vec{\mathbf{w}} \odot \vec{\mathbf{z}} = \sum_{i=1}^d \mathbf{w}_i z_i$. Thus the prediction phase boils down to servers executing Π_{dp} protocol with inputs as $\llbracket \vec{\mathbf{w}} \rrbracket$ and $\llbracket \vec{\mathbf{z}} \rrbracket$, to obtain $\llbracket \cdot \rrbracket$ shares of $\vec{\mathbf{w}} \odot \vec{\mathbf{z}}$.

6.3 Logistic Regression

The prediction phase of logistic regression model for a query $\vec{\mathbf{z}}$ is given by $\text{sig}(\vec{\mathbf{w}} \odot \vec{\mathbf{z}})$, where $\text{sig}(\cdot)$ denotes the sigmoid function. The sigmoid function is defined as $\text{sig}(u) = \frac{1}{1+e^{-u}}$. SecureML [45] showed the drawbacks of using sigmoid function for a general MPC setting and proposed a MPC friendly approximation, defined as follows :

$$\text{sigx}(u) = \begin{cases} 0 & u < -\frac{1}{2} \\ u + \frac{1}{2} & -\frac{1}{2} \leq u \leq \frac{1}{2} \\ 1 & u > \frac{1}{2} \end{cases}$$

The above equation can also be viewed as, $\text{sigx}(u) = \overline{b_1} b_2 (u + 1/2) + \overline{b_2}$, where bit $b_1 = 1$ if $u + 1/2 < 0$, bit $b_2 = 1$ if $u - 1/2 < 0$. Servers execute $\Pi_{\text{msb}}(u + 1/2)$ and $\Pi_{\text{msb}}(u - 1/2)$ to generate $\llbracket b_1 \rrbracket^{\mathbf{B}}$ and $\llbracket b_2 \rrbracket^{\mathbf{B}}$ respectively. Servers can locally compute $\llbracket \overline{b_i} \rrbracket^{\mathbf{B}}$ from $\llbracket b_i \rrbracket^{\mathbf{B}}$. After this, $\Pi_{\text{mult}}^{\mathbf{B}}(\llbracket \overline{b_1} \rrbracket, \llbracket b_2 \rrbracket)$ is executed to generate $\llbracket b \rrbracket^{\mathbf{B}}$, where $b = \overline{b_1} b_2$. Servers then invoke Π_{bin} on $\llbracket b \rrbracket^{\mathbf{B}}$ and $\llbracket (u + 1/2) \rrbracket$ to generate $\llbracket \overline{b_1} b_2 (u + 1/2) \rrbracket$, and $\Pi_{\text{btr}}(\llbracket \overline{b_2} \rrbracket^{\mathbf{B}})$ to generate $\llbracket \overline{b_2} \rrbracket$. Servers then locally add their shares to obtain $\llbracket \text{sigx}(u) \rrbracket$. Thus the cost for one query prediction in a logistic regression model is the same as the cost of linear regression, plus the additional overhead of computing $\text{sigx}(\vec{\mathbf{w}} \odot \vec{\mathbf{z}})$.

6.4 Deep Neural Networks (DNN)

All the techniques used to tackle the above models can be easily extended to support neural network prediction. We follow a similar procedure as ABY3, where

each node across all layers, use ReLU ($\text{rel}(\cdot)$) as its activation function. It comprises of computation of activation vectors for all the layers of the network. The activation vector for a given layer i of the network is defined as $\vec{\mathbf{a}}_i = \text{rel}(\vec{\mathbf{u}}_i)$, where $\vec{\mathbf{u}}_i = \mathbf{W}_i \times \vec{\mathbf{a}}_{i-1}$ is a matrix multiplication of weight matrix \mathbf{W}_i with the activation vector of the previous layer. Weight matrix $\mathbf{W}_i \in \mathbb{R}^{n_i \times n_{i-1}}$ contains all the weights connecting the nodes between layers i and $i - 1$, where n_i represents the number nodes in layer i . We set matrix $\vec{\mathbf{a}}_0 = \vec{\mathbf{z}}$, where $\vec{\mathbf{z}}$ is the input query of the client. All the above operations, that are needed for prediction, are simply a composition of several multiplications, dot products along with the evaluation of many ReLU functions. We now define the ReLU function below and also explain how to tackle it in our setting.

ReLU: The ReLU function is given as $\max(0, u)$. We view it as $\text{rel}(u) = \overline{b}u$, where bit $b = 1$ if $u < 0$, and \overline{b} is the complement of b . Servers execute $\Pi_{\text{msb}}(u)$ to generate $\llbracket b \rrbracket^{\mathbf{B}}$. Servers locally compute $\llbracket \overline{b} \rrbracket^{\mathbf{B}}$ from $\llbracket b \rrbracket^{\mathbf{B}}$, followed by executing Π_{bin} on $\llbracket \overline{b} \rrbracket^{\mathbf{B}}$ and $\llbracket u \rrbracket$ to generate $\llbracket \overline{b}u \rrbracket$.

6.5 Binarized Neural Network (BNN)

MOBIUS [37] proposed a secure prediction protocol for BNN in two party setting with one semi-honest corruption over \mathbb{Z}_{2^t} . In the original work of BNN [32], a batch normalization operation is performed at the output of every hidden layer of the binarized network, which requires bit-shifting mechanism. Performing bit-shifting in two party setting is very expensive. As a countermeasure, MOBIUS proposed an alternate solution for batch normalization with cost equal to that of one multiplication. The alternate solution is as follows: Suppose x_l^i be the output of node i in the l^{th} hidden layer, instead of using bit-shifting to normalize x_l^i , they perform $x_l^i = p_l^i x_l^i + q_l^i$, where x_l^i is the normalized output and p_l^i, q_l^i are the normalization batch parameters for node i of hidden layer l , which are provided by M . MOBIUS also showed that this method drops the accuracy by a negligible amount. Inspired from the ideas of MOBIUS, we now provide a secure prediction protocol for our setting. Note that, $\llbracket \cdot \rrbracket$ -shares of the weight matrices $\mathbf{W}_l \in \{-1, 1\}^{n_l \times n_{l-1}}$, batch normalization parameters $\vec{p}_l, \vec{q}_l, \forall l \in \{1, \dots, l_{\text{final}}\}$ and the query $\vec{\mathbf{z}}$ are already available among the servers.

We describe our protocol layer by layer. We use n_l to denote the number of nodes in layer l . The computation

in each layer l consists of three stages: i) The first stage comprises of matrix multiplication $\bar{\mathbf{x}}_l = \mathbf{W}_l \times f(\bar{\mathbf{x}}'_{l-1})$, where $\bar{\mathbf{x}}'_{l-1}$ denotes an n_{l-1} -sized vector and $f(\bar{\mathbf{x}}'_{l-1})$ denotes the vector obtained by applying activation function f on it. The activation function for a given value a is defined as

$$f(a) = \begin{cases} -1 & a < 0 \\ 1 & a \geq 0 \end{cases}$$

The matrix multiplication can be viewed as n_l dot product (protocol Π_{dp}) computations. ii) Servers, then perform batch normalization process on vector $\bar{\mathbf{x}}_l$ to obtain $\bar{\mathbf{x}}'_l = \bar{\mathbf{p}}_l \circ \bar{\mathbf{x}}_l + \bar{\mathbf{q}}_l$, where \circ denotes element wise multiplication. As evident, we use n_l multiplications and additions to compute the $\llbracket \cdot \rrbracket$ -sharing of $\bar{\mathbf{x}}'_l$. iii) This stage consists of passing the $\bar{\mathbf{x}}'_l$ through the activation function f to obtain $f(\bar{\mathbf{x}}'_l)$. To compute the activation function $f(a)$ in a $\llbracket \cdot \rrbracket$ -shared fashion, servers execute Π_{msb} on $\llbracket a \rrbracket$ to extract the MSB $msb(a)$, followed by executing Π_{btr} on $\llbracket msb(a) \rrbracket^B$ to generate $\llbracket msb(a) \rrbracket$. Finally, the servers locally compute $\llbracket f(a) \rrbracket = 2\llbracket msb(a) \rrbracket - 1$. For the input layer ($l = 0$), servers set $f(\bar{\mathbf{x}}'_0) = \bar{\mathbf{z}}$. Note that stage three is not required at the output layer.

7 Implementation & Benchmarks

We show the practicality of our framework by providing implementation results and compare with ABY3, in their respective settings over a ring of $\mathbb{Z}_{2^{64}}$.

i) Experimental Setup: Our experiments have been carried out both in the LAN and WAN setting. In the LAN setting, our machines are equipped with Intel Core i7-7790 CPU with 3.6 GHz processor speed and 32 GB RAM. Each of the four cores were able to handle eight threads, resulting in a total of 32 threads. We had a bandwidth of 1Gbps and an average round-trip time (rtt) of $\approx 0.26ms$. In the WAN setting, we use Microsoft Azure Cloud Services (Standard D8s v3, 2.4 GHz Intel Xeon® E5-2673 v3 (Haswell), 32GB RAM, 8 vcpus) with machines located in North Central US (S_1), South East Asia (S_2), Australia East (S_3) and West Europe (S_4). Each of the eight cores was capable of handling 16 threads resulting in a total of 128 threads. The bandwidth was limited to 20Mbps and the average rtt times are as follows:

S_1-S_2	S_1-S_3	S_1-S_4	S_2-S_3	S_2-S_4	S_3-S_4
161.76ms	197.03ms	97.32ms	116.36ms	225.34ms	236.56ms

We build on the ENCRYPTO library [18], following the standards of C++11. Due to the unavailability of the code of ABY3 [43], we implement their framework for comparison. For our executions, we report the average values over a run of 15 times.

ii) Parameters for Comparison: We consider three parameters for comparison– a) Latency (calculated as the maximum runtime of the servers), b) Communication complexity and c) Throughput (number of operations per unit time). The latency and throughput are evaluated over both LAN and WAN settings. The communication complexity is measured independent of the network. For the aforementioned algorithms, the throughput is calculated as the number of queries that can be computed per second and min in LAN and WAN respectively.

iii) Server Assignment: We assign the roles to the servers to maximize the performance of each of the frameworks, that we use for benchmarking. The table below provides the assignment of roles to the corresponding servers. For the 4PC setting, V_1, V_2 represent the set of verifiers while E_1, E_2 represent the set of evaluators. P_0, P_1, P_2 represent the parties, in the 3PC setting. we omit comparison with ASTRA framework as ABY3 outperforms ASTRA in terms of total communication (ref. Table 1).

Work	S_1	S_2	S_3	S_4
FLASH	E_1	E_2	V_1	V_2
ABY3	P_1	P_2	P_3	–

Table 3. Server Assignment for FLASH and ABY3 frameworks

iv) Datasets: We pick real-world datasets to measure the throughput for the prediction phase. The datasets we pick have features ranging from 13 to 784, which cover a range of feature sizes for a wide span of commonly used datasets.

ML Algorithm	Dataset	#features	#samples
Linear Reg.	Boston Housing Prices [30]	14	≈ 500
	Weather Conditions [48]	31	≈ 119000
Logistic Reg.	Candy Power Ranking [31]	13	≈ 85
	Food Recipes [22]	680	≈ 20000
DNN & BNN	MNIST [38]	784	≈ 70000

Table 4. Real World datasets for Comparison

For Linear Regression, we use Boston Housing Prices Dataset (Boston) [30] and the dataset obtained from [48] about the Weather Conditions in World War Two (Weather). The Boston dataset has ≈ 500 samples,

each with 14 features, while the Weather dataset has $\approx 119,000$ samples with 31 features.

For Logistic Regression we use the dataset from [22] which categorizes and gives the rating for recipes (Recipes) and Candy Power Ranking (Candy) dataset from [31] which predicts the most popular Halloween candy. The Candy dataset is small with only 13 features and ≈ 85 samples whereas the Recipe dataset is large with 680 features and $\approx 20,000$ samples.

For Deep Neural Network and Binarized Neural Network, we use MNIST [38] which contains 784 pixel images of handwritten numbers, each of size 28×28 . We also use synthetic datasets as it provides freedom to tune the number of features parameter and showcase the improvement with increasing feature size.

7.1 ML Building Blocks

We begin by comparing our protocols for some of the crucial ML building blocks, namely i) Dot Product, ii) MSB Extraction and iii) Truncation, against the state of the art protocols of ABY3 [43]. The comparison is mainly to show the substantial improvement we achieve in each building block when we shift from 3PC to 4PC setting, along with robustness guarantee. Later in Section 7.2 and 7.3 we show how the improvement in these blocks help us achieve massive improvements (Table.2) for our ML algorithms.

i) Dot Product: Dot Product is one of the vital building blocks for many machine learning algorithms like Linear Regression, Logistic Regression and Neural Network to name a few.

Work	LAN Latency (ms)	WAN Latency (s)
ABY3	3.55	1.10
FLASH	1.51	1.08

Table 5. Latency of 1 dot product computation for 784 features

Table 5 gives the comparison of our work with ABY3 with respect to the completion of one dot product computation for $d = 784$ features. We observe that for the LAN setting, even though the number of rounds required for completion of one dot product execution for both frameworks is 5 rounds, the latency of ABY3 is still twice of our FLASH. This discrepancy happens because the rtt of the network varies drastically with increase in the size of communication. In case of ABY3, due to their dot product protocol being dependent on the number of features the per party communication turns out to be 42.8KB, whereas our protocol incurs a

tiny cost of 0.09KB. Such a discrepancy is not observed in WAN as the communication threshold to vary the rtt is very high, under which all our protocols operate. We also plot the number of dot product computations that can be performed per sec, for varying feature sizes.

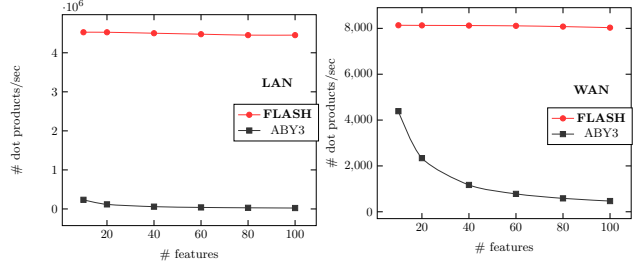


Fig. 12. # of dot product computations with increasing features.

It is clear from Figure.12 that varying the number of features has minimal impact on our throughput, since the communication cost of ours is independent of the feature size, while ABY3 suffers with increase in number of features. Thus for any machine learning algorithm which is heavily dependent on dot product computations, our protocol outperforms ABY3.

ii) MSB Extraction: MSB Extraction is the crux for many classification algorithms. Deep Neural Network and Binarized Neural Network where a large number of sequential comparisons are required. Table 6 gives the comparison of our work with ABY3, with respect to the completion of one MSB Extraction.

Work	LAN Latency (ms)	WAN Latency (s)
ABY3	3.53	2.29
FLASH	1.77	1.31

Table 6. Latency for single execution of MSB Extraction protocol

We also provide a latency graph with respect to the number of sequential comparisons to emphasize the effect of having a constant round protocol as opposed to a non-constant one.

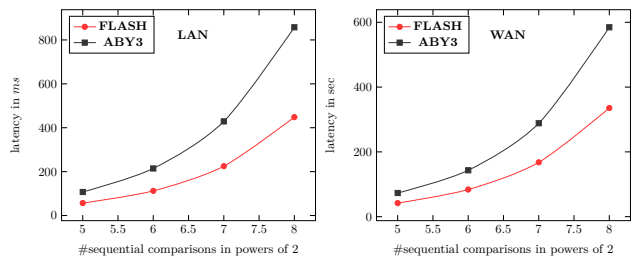


Fig. 13. Latency with increasing sequential comparisons

We observe from Figure 13 that our protocol outperforms ABY3 by a large margin with the increase in sequential comparisons. This is because our protocol requires only 6 rounds per comparison as opposed to 11 rounds for ABY3, when instantiated over a 64 bit ring. For the prediction phase of an ML algorithm like Deep Neural Network, the gap will keep growing bigger with the increase in depth (# hidden layers) of the neural network.

iii) Truncation: To showcase the effect of our efficient truncation protocol, we compare our protocol with that of ABY3. Table 7 gives the comparison with respect to the completion of a single execution of the protocol.

Work	LAN Latency (ms)	WAN Latency (s)
ABY3	1.52	1.11
FLASH	1.51	1.07

Table 7. Latency for a single execution of Truncation protocol

In the case of ABY3, though the truncation protocol takes $2\ell - 1$ rounds, the latency of both the frameworks in Table.7 are almost identical. This is because the goal of ABY3 was to have a high throughput framework, thus they compute $\approx 2^{20}$ parallel instances of $([r], \llbracket r^t \rrbracket)$ pairs so that the amortized time for a single execution of truncation protocol reduces. On the flip side, we do not have any such restriction on the number of $([r], \llbracket r^t \rrbracket)$ pair instances and the latency remains the same even if only one pair is required. Table 8 provides the throughput, measured as the number of multiplications with truncation performed, over both LAN (#mult/sec) and WAN (#mult/min) settings.

Work	LAN		WAN	
	#mult/sec	Improv.	#mult/min	Improv.
ABY3	0.45M	8.8×	4.76M	8.81×
FLASH	3.97M		0.54M	

Table 8. Throughput Comparison wrt # multiplications with truncation

We observe a minimum improvement of $8.8\times$ over ABY3. The improvement comes from the fact that ABY3 requires ≈ 6300 bits per truncation as compared to 896 bits for our case, when instantiated over a 64 bit ring. Our protocol will outperform ABY3 for all the ML algorithms that require repeated multiplications in the prediction phase.

7.2 Linear and Logistic Regression

In this section, we compare the concrete improvement of our framework against ABY3, for Linear and Logistic Regression. The performance is reported in terms of throughput of the protocol, the units being # queries/sec over LAN and # queries/min over WAN. We begin by comparing our framework with ABY3 over synthesized datasets as it provides us the freedom to tune the number of features parameter and showcase the improvement with the increase in #features. Table 9 provides a throughput comparison for #features $d = 10, 100$ and 1000 .

Setting	# Features	Ref.	Linear Reg.	Logistic Reg.
LAN (ms)	10	ABY3	1.68	5.59
		FLASH	1.51	3.26
	100	ABY3	2.03	5.94
		FLASH	1.51	3.26
	1000	ABY3	3.63	7.54
		FLASH	1.52	3.27
WAN (sec)	10/100/1000	ABY3	1.11	3.78
		FLASH	1.08	2.46

Table 9. Latency of frameworks for Linear and Logistic Reg.

As mentioned earlier in Section.7.1, the increase in feature size changes the LAN latency for ABY3 from 1.68ms to 3.63ms and 5.59ms to 7.54ms for Linear and Logistic regression respectively, whereas our latency stays stable to ≈ 1.5 ms and ≈ 3.26 ms for the same. The reason for the stability in our latency is the underlying dot product which is independent of the feature size. We now test on real-world datasets as mentioned in Table 4 for Linear and Logistic Regression. Figures 14 and 15 provide a comparison with ABY3 in terms of the number of queries computed per second and minute in LAN and WAN setting respectively. For Linear Regression, we observe a minimum throughput gain of $\approx 38\times$. The improvement primarily comes from the underlying Π_{dp} protocol and its independence of feature size property. Similarly, for Logistic Regression, we observe a throughput gain of around $\approx 48\times$, where protocols Π_{dp} and Π_{msb} become the prime contributors for the improvements in Logistic Regression.

7.3 Deep and Binarized Neural Network

In this section, we compare our framework with ABY3, for DNN and BNN. The accuracy of our predictions has the same bit-error that ABY3 mentions due to

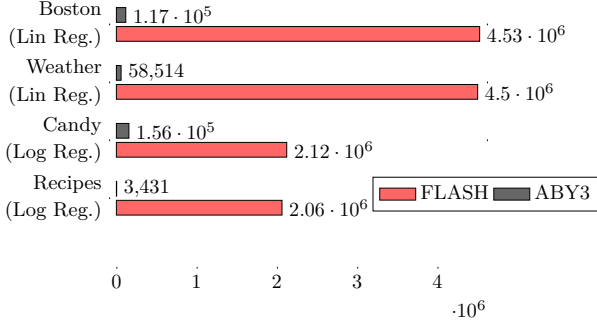


Fig. 14. Throughput Comparison (# queries/sec) for Linear and Logistic Regression in LAN setting

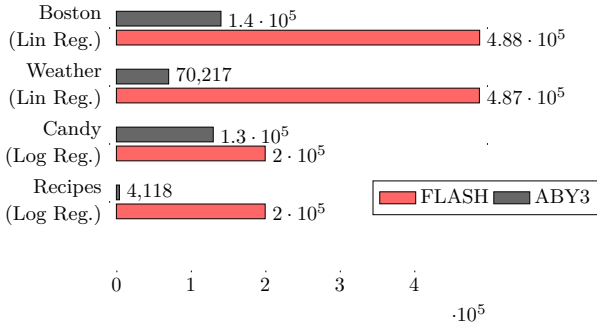


Fig. 15. Throughput Comparison (# queries/min) for Linear and Logistic Regression in WAN setting

the similarity in the approach to truncation. We begin by comparing (Table 10) over synthesized datasets and show the improvement in terms of latency for #features $d = 10, 100$ and 1000 .

Setting	# Features	Ref.	DNN	BNN
LAN (ms)	10	ABY3	59.71	59.73
		FLASH	18.65	23.37
	100	ABY3	67.78	67.77
		FLASH	18.74	23.69
	1000	ABY3	146.37	147.36
		FLASH	19.06	23.80
WAN (sec)	10/100/1000	ABY3	13.56	13.56
		FLASH	11.24	13.68

Table 10. Latency of frameworks for DNN and BNN

Figure 16 also shows how the depth of the neural network affects the throughput of the two frameworks. We consider a neural network with each hidden layer having 128 nodes and the final output layer having 10 nodes. The network is tested on MNIST dataset with $d = 784$ features.

It is clear from Figure 16, that we achieve impressive throughput gains of $\approx 200\times$ and $\approx 18\times$ for LAN and

WAN setting, even when the depth of the neural network goes up to 8 hidden layers. Such massive improvements primarily come from amalgamation of the improvements observed in the underlying building blocks (Section 7.1). Similar to DNN, we also achieve similar massive improvements for the case of BNN due to the aforementioned reasons. When tested on MNIST dataset ($d = 784$ features) for a BNN having 2 hidden layers, we observed throughput gains of $\approx 277\times$ in LAN and $\approx 23.8\times$ in WAN setting.

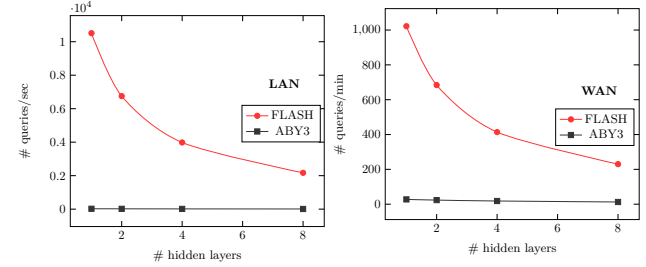


Fig. 16. Throughput Comparison for DNN with increasing number of hidden layers.

Acknowledgment: We thank Ananth Raghunathan, Yupeng Zhang and the anonymous reviewers of PETS 2020 for their valuable comments, which helped us improve the paper. Arpita Patra is supported by SERB Women Excellence Award 2017 (DSTO 1706). Ajith Suresh is supported by Google Phd Fellowship 2019.

References

- [1] A.Barak, D.Escudero, A.P.K.Dalskov, and M.Keller. Secure evaluation of quantized neural networks. *IACR Cryptology ePrint Archive*, 2019.
- [2] Á.Kiss, M.Naderpour, J.Liu, N. Asokan, and T.Schneider. Sok: Modular and efficient private decision tree evaluation. In *PoPETs*, 2018.
- [3] T. Araki, A. Barak, J. Furukawa, T. Lichter, Y. Lindell, A. Nof, K. Ohara, A. Watzman, and O. Weinstein. Optimized Honest-Majority MPC for Malicious Adversaries - Breaking the 1 Billion-Gate Per Second Barrier. In *IEEE S&P*, 2017.
- [4] T. Araki, A. Barak, J. Furukawa, Y. Lindell, A. Nof, and K. Ohara. DEMO: high-throughput secure three-party computation of kerberos ticket generation. In *ACM CCS*, 2016.
- [5] T. Araki, J. Furukawa, Y. Lindell, A. Nof, and K. Ohara. High-Throughput Semi-Honest Secure Three-Party Computation with an Honest Majority. In *ACM CCS*, 2016.
- [6] A.Tueno, F.Kerschbaum, and S.Katzenbeisser. Private evaluation of decision trees using sublinear cost. In *PoPETs*, 2019.
- [7] A. Ben-David, N. Nisan, and B. Pinkas. Fairplaymp: a system for secure multi-party computation. In *ACM CCS*, 2008.
- [8] M. Ben-Or, S. Goldwasser, and A. Wigderson. Completeness Theorems for Non-Cryptographic Fault-Tolerant Distributed Computation (Extended Abstract). In *ACM STOC*, 1988.
- [9] D. Bogdanov, S. Laur, and J. Willemson. Sharemind: A framework for fast privacy-preserving computations. In *ESORICS*, 2008.
- [10] P. Bogetoft, D. L. Christensen, I. Damgård, M. Geisler, T. P. Jakobsen, M. Krøigaard, J. D. Nielsen, J. B. Nielsen, K. Nielsen, J. Pagter, M. I. Schwartzbach, and T. Toft. Secure Multiparty Computation Goes Live. In *FC*, 2009.
- [11] D. Boneh, E. Boyle, H. Corrigan-Gibbs, N. Gilboa, and Y. Ishai. How to prove a secret: Zero-knowledge proofs on distributed data via fully linear pcps. *CRYPTO*, 2019.
- [12] M. Byali, C. Hazay, A. Patra, and S. Singla. Fast actively secure five-party computation with security beyond abort. In *ACM CCS*, 2019.
- [13] M. Byali, A. Joseph, A. Patra, and D. Ravi. Fast secure computation for small population over the internet. *ACM CCS*, 2018.
- [14] H. Chaudhari, A. Choudhury, A. Patra, and A. Suresh. AS-TRA: High-throughput 3PC over Rings with Application to Secure Prediction. In *ACM CCSW*, 2019.
- [15] K. Chida, D. Genkin, K. Hamada, D. Ikarashi, R. Kikuchi, Y. Lindell, and A. Nof. Fast large-scale honest-majority MPC for malicious adversaries. In *CRYPTO*, 2018.
- [16] R. Cleve. Limits on the security of coin flips when half the processors are faulty (extended abstract). In *ACM STOC*, 1986.
- [17] R. Cohen and Y. Lindell. Fairness versus guaranteed output delivery in secure multiparty computation. In *ASIACRYPT*, 2014.
- [18] Cryptography and Privacy Engineering Group at TU Darmstadt. ENCRYPTO Utils. https://github.com/encryptogroup/ENCRYPTO_utils, 2017.
- [19] I. Damgård, M. Keller, E. Larraia, V. Pastro, P. Scholl, and N. P. Smart. Practical covertly secure MPC for dishonest majority - or: Breaking the SPDZ limits. In *ESORICS*, 2013.
- [20] I. Damgård, C. Orlandi, and M. Simkin. Yet another compiler for active security or: Efficient MPC over arbitrary rings. *CRYPTO*, 2018.
- [21] I. Damgård, V. Pastro, N. P. Smart, and S. Zakarias. Multiparty Computation from Somewhat Homomorphic Encryption. In *CRYPTO*, 2012.
- [22] H. Darwood. Epicurious - recipes with rating and nutrition. 2017.
- [23] D. Demmler, T. Schneider, and M. Zohner. ABY - A Framework for Efficient Mixed-Protocol Secure Two-Party Computation. In *NDSS*, 2015.
- [24] H. Eerikson, C. Orlandi, P. Pullonen, J. Puura, and M. Simkin. Use your brain! arithmetic 3pc for any modulus with active security. *IACR Cryptology ePrint Archive*, 2019.
- [25] A. Esteva, B. Kuprel, R. A. Novoa, J. Ko, S. M. Swetter, H. M. Blau, and S. Thrun. Dermatologist-level classification of skin cancer with deep neural networks. *Nature*, 2017.
- [26] J. Furukawa, Y. Lindell, A. Nof, and O. Weinstein. High-Throughput Secure Three-Party Computation for Malicious Adversaries and an Honest Majority. In *EUROCRYPT*, 2017.
- [27] M. Geisler. Viff: Virtual ideal functionality framework, 2007.
- [28] O. Goldreich, S. Micali, and A. Wigderson. How to Play any Mental Game or A Completeness Theorem for Protocols with Honest Majority. In *STOC*, 1987.
- [29] S. D. Gordon, S. Ranellucci, and X. Wang. Secure computation with low communication from cross-checking. In *ASIACRYPT*, 2018.
- [30] D. Harrison and D. L. Rubinfeld. Hedonic housing prices and the demand for clean air. *Journal of Environmental Economics and Management*, 1978.
- [31] W. Hickey. The ultimate halloween candy power ranking. 2017.
- [32] I. Hubara, M. Courbariaux, D. Soudry, R. El-Yaniv, and Y. Bengio. Binarized neural networks. In *NIPS*, 2016.
- [33] Y. Ishai, J. Kilian, K. Nissim, and E. Petrank. Extending Oblivious Transfers Efficiently. In *CRYPTO*, 2003.
- [34] Y. Ishai, R. Kumaresan, E. Kushilevitz, and A. Paskin-Cherniavsky. Secure computation with minimal interaction, revisited. In *CRYPTO*, 2015.
- [35] J. So, B. Guler, A. S. Avestimehr, and P. Mohassel. Coded-privateml: A fast and privacy-preserving framework for distributed machine learning. *CoRR*, 2019.
- [36] C. Juvekar, V. Vaikuntanathan, and A. Chandrakasan. GAZELLE: A low latency framework for secure neural network inference. In *USENIX*, 2018.
- [37] H. Kitai, J. P. Cruz, N. Yanai, N. Nishida, T. Oba, Y. Unagami, T. Teruya, N. Attrapadung, T. Matsuda, and G. Hanaoka. MOBIUS: model-oblivious binarized neural networks. *CoRR*, 2018.
- [38] Yann LeCun and Corinna Cortes. MNIST handwritten digit database. 2010.
- [39] Y. Lindell. Fast cut-and-choose-based protocols for malicious and covert adversaries. *J. Cryptology*, 2016.
- [40] Y. Lindell and B. Pinkas. An efficient protocol for secure two-party computation in the presence of malicious adversaries. In *EUROCRYPT*, 2007.

- [41] E. Makri, D. Rotaru, N. P. Smart, and F. Vercauteren. EPIC: efficient private image classification (or: Learning from the masters). *CT-RSA*, 2018.
- [42] P. Mohassel and M. K. Franklin. Efficiency tradeoffs for malicious two-party computation. In *PKC*, 2006.
- [43] P. Mohassel and P. Rindal. ABY³: A Mixed Protocol Framework for Machine Learning. In *ACM CCS*, 2018.
- [44] P. Mohassel, M. Rosulek, and Y. Zhang. Fast and Secure Three-party Computation: Garbled Circuit Approach. In *CCS*, 2015.
- [45] P. Mohassel and Y. Zhang. Secureml: A system for scalable privacy-preserving machine learning. In *IEEE S&P*, 2017.
- [46] M.S.Riazi, M.Samragh, H.Chen, K.Laine, K.E.Lauter, and F.Koushanfar. XONN: xnor-based oblivious deep neural network inference. 2019.
- [47] J. B. Nielsen and C. Orlandi. Cross and clean: Amortized garbled circuits with constant overhead. In *TCC*, 2016.
- [48] NOAA. Weather conditions in world war two. 2017.
- [49] P. S. Nordholt and M. Veeningen. Minimising Communication in Honest-Majority MPC by Batchwise Multiplication Verification. In *ACNS*, 2018.
- [50] A. Patra and D. Ravi. On the exact round complexity of secure three-party computation. *CRYPTO*, 2018.
- [51] 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 *AsiaCCS*, 2018.
- [52] F. Schroff, D. Kalenichenko, and J. Philbin. Facenet: A unified embedding for face recognition and clustering. In *IEEE CVPR*, 2015.
- [53] S. Wagh, D. Gupta, and N. Chandran. Securenn: Efficient and private neural network training. *19th Privacy Enhancing Technologies Symposium*, 2019.
- [54] A. C. Yao. Protocols for Secure Computations. In *FOCS*, 1982.

A Comparison with [29]

In this section we compare our work with state-of-the-art 4PC protocol of [29] in detail for the Abort scenario.

A.1 4PC with Abort:

All the aforementioned robust protocols can be easily converted to the abort variant by tweaking the Bi-convey primitive (Section 4.1). In case of abort setting, parties S_1 and S_2 in the Bi-Convey primitive send x and $H(x)$ respectively to R , who accepts x if the hashes match else aborts. Thus by swapping with the abort variant of the primitive, all the building blocks achieve security with abort. Table 11 provides round and communication complexity comparison of both the variants of the protocols.

Protocol	Equation	FLASH (Abort)		FLASH (Robust)	
		Rounds	Comm.	Rounds	Comm.
Multiplication	$[[x]].[[y]] \rightarrow [[x.y]]$	2	6ℓ	5	12ℓ
Dot Product	$[[\vec{x} \odot \vec{y}]] = [[\sum_{i=1}^d x_i y_i]]$	2	6ℓ	5	12ℓ
MSB Extraction	$[[x]] \rightarrow [[\text{msb}(x)]]^B$	4	$8\ell + 2$	6	$16\ell + 4$
Truncation	$[[x]].[[y]] \rightarrow [[(xy)^t]]$	2	7ℓ	5	14ℓ
Bit Conversion	$[[b]]^B \rightarrow [[b]]$	2	7ℓ	5	14ℓ
Bit Insertion	$[[b]]^B[[x]] \rightarrow [[bx]]$	2	9ℓ	5	18ℓ

Table 11. Comparison of Abort and Robust variants in FLASH.

As observed in Table 11, for the abort setting our cost of multiplication protocol is 6 elements which turns out to be the same as [29]. But from a practical viewpoint, if we cast ours and GRW18 multiplication protocol into the offline-online paradigm, where the offline phase generates the necessary offline values in order for a fast online phase to be executed when the client query becomes available, our protocol requires only 3 parties to be active (V_2 , E_1 and E_2) in the online phase, whereas [29] needs all parties to be active throughout the entire execution. This is helpful, because now the server associated with party V_1 is only needed to generate offline values and can be shut down for the entirety of the online phase which will, in turn, save a lot in terms of monetary cost for running the server on the cloud (WAN) setting. Hence, even though the communication and round complexity of both the works turns out to be the same with respect to a single multiplication, our work has better *practical* efficiency in terms of the number of servers required in the online phase.

Table 12 provides a concrete comparison of our framework with [29] below.

Work	Equation	Offline Phase		Online Phase	
		Rounds	Comm.	Rounds	Comm.
[29]	$[[x]].[[y]] \rightarrow [[x.y]]$	1	2ℓ	1	4ℓ
Ours	$[[x]].[[y]] \rightarrow [[x.y]]$	1	3ℓ	1	3ℓ

Table 12. Comparison of FLASH with [29] for the Abort setting.

B Building Blocks and Security

B.1 Building Blocks

i) Collision Resistant Hash: Consider a hash function family $H = \mathcal{K} \times \mathcal{L} \rightarrow \mathcal{Y}$. The hash function H is said to be collision resistant if for all probabilistic polynomial-time adversaries \mathcal{A} , given the description of H_k where $k \in_R \mathcal{K}$, there exists a negligible function $\text{negl}(\cdot)$ such that $\Pr[(x_1, x_2) \leftarrow \mathcal{A}(k) : (x_1 \neq$

$x_2) \wedge H_k(x_1) = H_k(x_2)] \leq \text{negl}(\kappa)$, where $m = \text{poly}(\kappa)$ and $x_1, x_2 \in_R \{0, 1\}^m$.

ii) **Commitment Scheme:** We use $\text{com}(x)$ to denote commitment of a value x . The commitment scheme ($\text{com}()$) possess two properties, namely – i) *hiding*, which ensures the privacy of value x given just the commitment, and ii) *binding*, which prevents a corrupt party from opening the commitment to a different value $x' \neq x$. The commitment scheme can be implemented via a hash function $\mathcal{H}()$, whose security can be proved in the random-oracle model (ROM). For example, $(c, o) = (\mathcal{H}(x||r), x||r) = \text{Com}(x; r)$.

B.2 Ideal World Functionalities

We prove the security of our protocols in the standard real/ideal world paradigm where we compare the view of the adversary in the real world and ideal world. In an ideal world execution, each party sends its input to an incorruptible trusted third party (TTP), who computes the given function $f()$ using the inputs received and sends back the respective output to each party.

Fig 17 denotes the ideal functionality $\mathcal{F}_{\text{setup}}$ that establishes the shared randomness among the parties.

$\mathcal{F}_{\text{setup}}$ interacts with the parties in \mathcal{P} and the adversary \mathcal{S} . $\mathcal{F}_{\text{setup}}$ picks random keys $k_{\mathbf{E}}, k_{\mathbf{V}}, k_{\mathbf{E}, \mathbf{V}_1}, k_{\mathbf{E}, \mathbf{V}_2}, k_{\mathbf{V}, \mathbf{E}_1}, k_{\mathbf{V}, \mathbf{E}_2}, k_{\mathcal{P}} \in \{0, 1\}^\kappa$. Let y_i denote the keys corresponding to party P_i . Then

- $y_i = (k_{\mathbf{V}}, k_{\mathbf{E}, \mathbf{V}_1}, k_{\mathbf{V}, \mathbf{E}_1}, k_{\mathbf{V}, \mathbf{E}_2}, k_{\mathcal{P}})$ when $P_i = \mathbf{V}_1$.
- $y_i = (k_{\mathbf{V}}, k_{\mathbf{E}, \mathbf{V}_2}, k_{\mathbf{V}, \mathbf{E}_1}, k_{\mathbf{V}, \mathbf{E}_2}, k_{\mathcal{P}})$ when $P_i = \mathbf{V}_2$.
- $y_i = (k_{\mathbf{E}}, k_{\mathbf{V}, \mathbf{E}_1}, k_{\mathbf{E}, \mathbf{V}_1}, k_{\mathbf{E}, \mathbf{V}_2}, k_{\mathcal{P}})$ when $P_i = \mathbf{E}_1$.
- $y_i = (k_{\mathbf{E}}, k_{\mathbf{V}, \mathbf{E}_2}, k_{\mathbf{E}, \mathbf{V}_1}, k_{\mathbf{E}, \mathbf{V}_2}, k_{\mathcal{P}})$ when $P_i = \mathbf{E}_2$.

Output: $\mathcal{F}_{\text{setup}}$ sends the keys y_i to party P_i .

Fig. 17. Functionality $\mathcal{F}_{\text{setup}}$

B.3 4PC Protocol

We present the 4PC protocol in Fig 18 and the corresponding ideal functionality appears in Fig 19.

Input Sharing: For each input value x , parties execute $\Pi_{\text{sh}}(D, x)$, where D is the owner of value x .

Addition gate: For every addition gate $z = x + y$ in the ckt, parties execute $\Pi_{\text{add}}(x, y, z)$.

Multiplication gate: For every multiplication gate ($z = xy$) in the ckt, parties execute $\Pi_{\text{mult}}(x, y, z)$.

Output Computation: For every output value z , parties execute Π_{oc} .

Fig. 18. $\Pi_{4\text{PC}}$: A 4PC Robust Protocol

$\mathcal{F}_{\text{robust}}$ receives input (**Input**, x) from party $P \in \{\mathbf{V}_1, \mathbf{V}_2, \mathbf{E}_1, \mathbf{E}_2\}$. While honest parties send their input correctly, corrupt parties may send arbitrary inputs as instructed by the adversary \mathcal{A} .

- For every party P , $\mathcal{F}_{\text{robust}}$ sets x to some pre-determined value if either $x = *$ or x is outside the domain of values allowed for input of P .
- $\mathcal{F}_{\text{robust}}$ computes output $y = f(x'_1, x'_2, x'_3, x'_4)$ and sends (**Output**, y) to all the parties in $\{\mathbf{V}_1, \mathbf{V}_2, \mathbf{E}_1, \mathbf{E}_2\}$.

Fig. 19. Functionality $\mathcal{F}_{\text{robust}}$ for 4PC protocol

C Lemmas and Proofs

C.1 4PC

Lemma C.1. *The designated receiver R either receives a given value x correctly in Π_{bic} or receiver R and helper T mutually exchange all their internal randomness.*

Proof. The case of R and T (who act as pair of honest parties) mutually exchanging their internal randomness occurs when when one of the senders (S_1, S_2) are corrupt and copies of x received by R and the hashes $H(x)$ received by T mismatch. In all the other cases there always exists a majority among the copies of x received by R . Thus R is able to correctly obtain x in the remaining cases. \square

Lemma C.2. *Π_{bic} protocol requires a communication cost (amortized) of 2ℓ bits and at most 2 rounds.*

Proof. For a given value x , the communication cost is equal to 2ℓ bits as the senders S_1, S_2 send x to the designated party R . Round complexity wise, in case of a corrupt sender, he/she can delay party R from receiving x by atmost 2 rounds. This case occurs when in the first round the copies of x received by R mismatch and the hashes $H(x)$ received by party T match. The second round simply involves party T sending $H(x)$ to R who accepts the copy which matches with the received hash. The case when R or T is corrupt, Π_{bic} will take exactly 1 round as S_1 and S_2 will always send the correct copies. \square

Lemma C.3. *For a gate $g = (x, y, z)$, given the $\llbracket \cdot \rrbracket$ -shares of inputs x and y , protocols Π_{add} and Π_{mult} compute $\llbracket \cdot \rrbracket$ -share of the output wire z .*

Proof. By linearity property of $\llbracket \cdot \rrbracket$ -sharing, the addition gates preserve the $\llbracket \cdot \rrbracket$ -sharing of their inputs. For every multiplication gate $g = (z = xy)$, the evaluators robustly compute μ_z , after which they set $\mu_z^2 = \mu_z - \mu_z^1$ (μ_z^1 chosen non-interactively) for consistent $\llbracket \cdot \rrbracket$ -sharing of z to preserve the invariant. The share μ_z^2 for every multiplication gate is later robustly communicated to the verifier V_2 to maintain a consistent $\llbracket \cdot \rrbracket$ -sharing for the entire circuit. \square

Lemma C.4. Π_{mult} protocol requires a communication cost (amortized) of 12ℓ bits and at most 5 rounds.

Proof. Π_{bic} of $\delta_{xy}^2, A_1, A_2, B_1$ and B_2 takes 10ℓ bits followed by Π_{bic} of μ_z^2 takes another 2ℓ bits. Round complexity wise, in case of a corrupt verifier, Π_{bic} of δ_{xy}^2 takes at most 2 rounds. Π_{bic} of A_1, A_2, B_1 and B_2 also takes at most 2 rounds followed by evaluators executing Π_{bic} of μ_z^2 consumes 1 round. A similar argument can be made when one of the evaluator is corrupt. \square

Lemma C.5. Each party either commits to his/her input in Π_{sh} or is identified to be corrupt.

Proof. In Π_{sh} , the mirrored sharing of inputs by each party is as in Π_{sh} with an additional step of identifying the adversary in case of mismatch. The step of eliminating the adversary uses the computation of honest majority on the dispersed shares. Since only, one corruption can occur, an honest party's input always gets committed irrespective of the behaviour of the adversary. However, the case of no honest majority can occur only when the dealer is corrupt. Hence only a corrupt party is eliminated if she does not commit to her input and a default value is taken. The uniqueness of the share also follows from collision resistant hash. Else, the chosen input is committed. \square

Lemma C.6. The protocol Π_{oc} is correct.

Proof. The correctness for output computation follows from the fact that each party receives 2 copies and a corresponding hash for its missing share from the remaining parties. Thus each party correctly reconstructs the output as a majority always exists. \square

Lemma C.7. The protocol $\Pi_{4\text{PC}}$ is correct.

Proof. We argue that the computed z corresponds to unique set of inputs. By Lemma C.5, a corrupt party either commits to its input in which case, we proceed to evaluation or is identified to be corrupt and elimi-

nated in which case, the output is computed on default input of the corrupt party. In the evaluation step, the computation of addition gates is local by the linearity property. For a multiplication gate $\Pi_{\text{mult}}(x, y, z)$, the correctness of A_1, A_2, B_1, B_2 and δ_{xy}^2 sharing follows from the correctness of Π_{bic} protocol. Hence evaluators correctly compute μ_z , and set $\mu_z^2 = \mu_z - \mu_z^1$. Verifier V_2 also correctly receives μ_z^2 , from the underlying correctness of Π_{bic} protocol. The protocol $\Pi_{4\text{PC}}$, relies on the the routines $\Pi_{\text{sh}}, \Pi_{\text{mult}}$ and Π_{oc} and thus its correctness follows from their correctness. \square

C.2 Privacy-Preserving Machine Learning

C.2.1 Arithmetic/Boolean Couple Sharing

C.2.1.1 i) Parties in E couple share

Lemma C.8. Π_{cSh} protocol requires a communication cost (amortized) of 2ℓ bits and at most 2 rounds.

Proof. The communication cost of 2ℓ bits comes directly from the cost of Π_{bic} protocol as the rest of the steps are local, which includes collectively sampling μ_x^1 . Round complexity argument also follow from Π_{bic} protocol. \square

C.2.1.2 ii) Parties in V couple share

Lemma C.9. Π_{cSh} protocol requires a communication cost (amortized) of 2ℓ bits and at most 2 rounds.

Proof. The communication cost of 2ℓ bits comes directly from the cost of Π_{bic} protocol as the rest of the steps are local, which includes collectively sampling σ_x^1 . Round complexity argument also follow from Π_{bic} protocol. \square

C.2.2 4PC Truncation

Lemma C.10. Π_{mulTr} protocol requires a communication cost (amortized) of 14ℓ bits and at most 5 rounds.

Proof. Π_{cSh} of $\llbracket r^t \rrbracket$ and δ_{xy} takes 4ℓ bits in total. Π_{bic} of A_1, A_2, B_1 and B_2 takes 8ℓ bits followed by Π_{cSh} of $(z - r)^t$ takes another 2ℓ bits. Round complexity wise, in case of a corrupt verifier, Π_{cSh} of $\llbracket r^t \rrbracket$ and δ_{xy} takes at most 2 rounds. Π_{bic} of A_1, A_2, B_1 and B_2 also takes at most 2 rounds followed by Π_{cSh} of $(z - r)^t$ consumes 1 round. A similar argument can be made when one of the evaluator is corrupt. \square

C.2.3 Dot Product

Lemma C.11. Π_{dp} protocol requires a communication cost (amortized) of 12ℓ bits and atmost 5 rounds.

Proof. The communication cost of 12ℓ bits comes directly from the cost of Π_{mult} protocol as the rest of the steps are local. Round complexity argument also follow from Π_{mult} protocol. \square

C.2.4 Bit Conversion

Lemma C.12. Π_{btr} protocol requires a communication cost (amortized) of 14ℓ bits and atmost 5 rounds.

Proof. Firstly, the protocol Π_{cSh} used to generate the arithmetic equivalent $[\cdot]$ -sharing of bit σ_b and μ_b consumes 4ℓ bits in total. The optimized multiplication of $\mu_{b'} \cdot \sigma_{b'}$ consumes 10ℓ bits in total as $\delta_{\mu_{b'} \sigma_{b'}} = 0$ so Π_{cSh} is not required the same. In case of a corrupt verifier Π_{cSh} of σ_b can take atmost 2 rounds, followed by 3 rounds for optimized multiplication (as $\delta_{\mu_{b'} \sigma_{b'}} = 0$) making the total rounds equal to 5. A similar argument can be made for the case when one of the evaluator is corrupt. \square

C.2.5 Bit Insertion

Lemma C.13. Π_{bin} protocol requires a communication cost (amortized) of 18ℓ bits and atmost 5 rounds.

Proof. Four calls to Π_{bic} for $\sigma_{b'}^2, \mu_{b'}^2, \gamma_{b'x}$ and $\delta_{b'x}$ consumes 8ℓ bits in total. Again four calls to Π_{bic} each for A_1, A_2, B_1 and B_2 consumes another 8ℓ bits followed by evaluators invoking Π_{bic} of $\mu_{b'x}^2$ which consumes 2ℓ bits. Round complexity wise, in case of a corrupt verifier, Π_{bic} for $\sigma_{b'}^2, \mu_{b'}^2, \gamma_{b'x}$ and $\delta_{b'x}$ takes atmost 2 rounds, followed by Π_{bic} of A_1, A_2, B_1 and B_2 which consumes atmost 2 more rounds. Finally, Π_{bic} of $\mu_{b'x}^2$ which requires 1 round. A similar argument can be made when one of the evaluator is corrupt. \square

C.2.6 MSB Extraction

Lemma C.14. Π_{msb} protocol requires a communication cost (amortized) of $16\ell + 4$ bits and atmost 6 rounds.

Proof. The protocols $\Pi_{cSh}(\mathbf{E}, r)$ and $\Pi_{cSh}^B(\mathbf{E}, p)$ to generate $[\![r]\!]$ and $[\![p]\!]^B$ consume 2ℓ bits and 2 bits respectively.

As a consequence of $\Pi_{cSh}(\mathbf{E}, r)$, $\delta_{ra} = 0$ and thus Π_{mult} of ra consumes 10ℓ bits in total. The reconstruction of ra towards V_1, V_2 requires 4ℓ bits. As $\Pi_{cSh}^B(\mathbf{V}, q)$ consumes 2 bits, thus making a total communication equal to $16\ell + 4$ bits. Round complexity wise, in case of a corrupt evaluator, $\Pi_{cSh}(\mathbf{E}, r)$ can take atmost 2 rounds followed by Π_{mult} of ra which will take atmost 3 rounds. Reconstruction of ra towards V_1, V_2 can be clubbed with Π_{bic} of μ_{ra}^2 followed by $\Pi_{cSh}^B(\mathbf{V}, q)$ which consumes 1 round. A similar argument can be made when one of the verifier is corrupt. \square

D Security of Bi-Convey

In this section, we provide a detailed security proof for our Bi-Convey Primitive (Π_{bic}), which forms the backbone for most of our constructions, in the stand-alone model. $\mathcal{S}_{\Pi_{bic}}^P$ denotes the simulator for the case of a corrupt party $P \in \{V_1, V_2, E_1, E_2\}$.

We begin with case of a corrupt S_1 . Since party S_1 is not receiving any messages in the protocol Π_{bic} , there is no need for $\mathcal{S}_{\Pi_{bic}}^{S_1}$ to simulate any messages. Based on the messages received from S_1 , simulator prepares the input value of corrupt S_1 and invoke the ideal functionality \mathcal{F}_{bic} . A detailed description of $\mathcal{S}_{\Pi_{bic}}^{S_1}$ is given in Fig 20. Note that, $\mathcal{S}_{\Pi_{bic}}^{S_1}$ has the knowledge of input value x , since it plays the role of an honest S_2 .

- $\mathcal{S}_{\Pi_{bic}}^{S_1}$ receives x' and $\text{com}(x'')$ from S_1 on behalf of parties R and T respectively.
- If $x' \neq x$ or $\text{com}(x'') \neq \text{com}(x)$, $\mathcal{S}_{\Pi_{bic}}^{S_1}$ sets the input message of S_1 as $x_{S_1} = \perp$. Else it sets $x_{S_1} = x$.
- $\mathcal{S}_{\Pi_{bic}}^{S_1}$ invokes the ideal functionality \mathcal{F}_{bic} on behalf of S_1 with input x_{S_1} .

Fig. 20. $\mathcal{S}_{\Pi_{bic}}^{S_1}$: Simulator for the case of corrupt S_1

It is easy to see that the view of the adversary \mathcal{A} in the real and simulated worlds are indistinguishable. The case for a corrupt S_2 follows similarly.

We now consider the case of a corrupt R . For this, $\mathcal{S}_{\Pi_{bic}}^R$ (Fig 21) samples a random value x on behalf of S_1, S_2 and prepares the commitment of x honestly. This is followed by sending the values x, x and $\text{com}(x)$ to R on behalf of S_1, S_2 and T respectively.

- $\mathcal{S}_{\Pi_{bic}}^R$ samples a random value x on behalf of S_1, S_2 . It then prepares the commitment $\text{com}(x)$ using a randomness shared with R .

- $\mathcal{S}_{\Pi_{\text{bic}}}^R$ sends x, x and $\text{com}(x)$ to R on behalf of S_1, S_2 and T respectively.
- $\mathcal{S}_{\Pi_{\text{bic}}}^R$ invokes the simulator for ideal functionality $\mathcal{F}_{\text{setup}}$ and obtains the internal randomness of R, ι_R . $\mathcal{S}_{\Pi_{\text{bic}}}^R$ invokes the ideal functionality \mathcal{F}_{bic} on behalf of R with ι_R as the input.

Fig. 21. $\mathcal{S}_{\Pi_{\text{bic}}}^R$: Simulator for the case of corrupt R

For the case of a corrupt T , $\mathcal{S}_{\Pi_{\text{bic}}}^T$ (Fig 22) proceeds as follows: $\mathcal{S}_{\Pi_{\text{bic}}}^T$ samples a random value x on behalf of S_1, S_2 and prepares the commitment of x honestly. This is followed by sending the values $\text{com}(x), \text{com}(x)$ and \perp to T on behalf of S_1, S_2 and R respectively.

- $\mathcal{S}_{\Pi_{\text{bic}}}^T$ samples a random value x on behalf of S_1, S_2 . It then prepares the commitment $\text{com}(x)$.
- $\mathcal{S}_{\Pi_{\text{bic}}}^T$ sends $\text{com}(x), \text{com}(x)$ and continue to T on behalf of S_1, S_2 and R respectively.
- $\mathcal{S}_{\Pi_{\text{bic}}}^T$ invokes the simulator for ideal functionality $\mathcal{F}_{\text{setup}}$ and obtains the internal randomness of T, ι_T . $\mathcal{S}_{\Pi_{\text{bic}}}^T$ invokes the ideal functionality \mathcal{F}_{bic} on behalf of T with ι_T as the input.

Fig. 22. $\mathcal{S}_{\Pi_{\text{bic}}}^T$: Simulator for the case of corrupt T

In each of the cases, since the simulator behaves entirely as an honest party in the protocol simulation, the view of the adversary \mathcal{A} in the real and simulated worlds are indistinguishable in a very straightforward manner. This concludes the proof.