

Tetrad: Actively Secure 4PC for Secure Training and Inference

Nishat Koti*, Arpita Patra*, Rahul Rachuri†, Ajith Suresh*

*Indian Institute of Science, Bangalore, Email: {kotis, arpita, ajith}@iisc.ac.in

†Aarhus University, Denmark, Email: rachuri@cs.au.dk

Abstract—Mixing arithmetic and boolean circuits to perform privacy-preserving machine learning has become increasingly popular. Towards this, we propose a framework for the case of four parties with at most one active corruption called Tetrad.

Tetrad works over rings and supports two levels of security, fairness and robustness. The fair multiplication protocol costs 5 ring elements, improving over the state-of-the-art Trident (Chaudhari et al. NDSS’20). A key feature of Tetrad is that robustness comes for free over fair protocols. Other highlights across the two variants include (a) probabilistic truncation without overhead, (b) multi-input multiplication protocols, and (c) conversion protocols to switch between the computational domains, along with a tailor-made garbled circuit approach.

Benchmarking of Tetrad for both training and inference is conducted over deep neural networks such as LeNet and VGG16. We found that Tetrad is up to 4 times faster in ML training and up to 5 times faster in ML inference. Tetrad is also lightweight in terms of deployment cost, costing up to 6 times less than Trident.

I. INTRODUCTION

Increased concerns about privacy coupled with policies such as European Union General Data Protection Regulation (GDPR) make it harder for multiple parties to collaborate on machine learning computations. The emerging field of privacy-preserving machine learning (PPML) addresses this issue by offering tools to let parties perform computations without sacrificing privacy of the underlying data. PPML can be deployed across various domains such as healthcare, recommendation systems, text translation, etc., with works like [1] demonstrating practicality.

One of the main ways in which PPML is realised is through the paradigm of secure outsourced computation (SOC). Clients can outsource the training/prediction computation to powerful servers available on a ‘pay-per-use’ basis from cloud service providers. Of late, secure multiparty computation (MPC) based techniques [2]–[10] have been gaining interest, where a server enacts the role of a party in the MPC protocol. MPC [11], [12] allows mutually distrusting parties to compute a function in a distributed fashion while guaranteeing *privacy* of the parties’ inputs and *correctness* of their outputs against any coalition of t parties.

The goal of PPML is practical deployment, making *efficiency* a primary consideration. Functions such as comparison, activation functions (e.g., ReLU), are heavily used in machine learning. Instantiating these functions via MPC naively turns out to be prohibitively inefficient due to their non-linearity. Hence, there is motivation to design specialised protocols that can compute these functions efficiently. We work towards this goal in the 4-party (4PC) setting, assuming honest majority [2], [4], [13], [14]. 4PC is interesting because it buys us the following over 3PC (which is threshold optimal): (1) *independence from broadcast*: broadcast can be achieved by a simple protocol in which the sender sends to everyone and residual parties exchange and apply a majority rule (2) *efficient dot-product*: 4PC offers a more efficient dot-product protocol (which is an important building block for several ML algorithms) with communication complexity independent of feature size (3) *simplicity and efficiency*: protocols are vastly more efficient and simpler in terms of design (as shown in this and prior works). To enhance practical efficiency, many recent works [4], [8], [15], [16] resort to the preprocessing paradigm, which splits the computation into two phases; a preprocessing phase where input-independent (but function-dependent) computationally heavy tasks can be computed, followed by a fast online phase. Since the same functions in ML are evaluated several times, this paradigm naturally fits the case of PPML, where the ML algorithm is known beforehand. Further, recent works [15], [17], [18] propose MPC protocols over 32 or 64 bit rings to leverage CPU optimizations.

MPC protocols can be categorized as high-throughput [3], [4], [6], [8], [14], [19]–[23] and low-latency [24], [25], where the former, based on secret-sharing, requires less communication compared to the latter (garbled circuits). High-throughput protocols typically work over the boolean ring \mathbb{Z}_2 or an arithmetic ring \mathbb{Z}_{2^e} and aim to minimize communication overhead (bandwidth) at the expense of non-constant rounds. While high-throughput protocols enable efficient computation of functions such as addition, multiplication and dot-product, other functions such as division are best performed using garbled circuits. Activation functions such as ReLU used in neural networks (NN) alternate between multiplication and comparison, wherein multiplication is better suited to the arithmetic world and comparison to the boolean world. Hence, MPC protocols working over different representations (arithmetic/boolean/garbled circuit based) can be mixed to achieve better efficiency. This provided motivation for mixed protocols where each subprotocol is executed in a world where it performs best. Mixed-protocol frameworks [4], [6], [7], [10], [17], [23], [26], [27] have support for efficient ways to switch between the worlds, thereby getting the best from each of them.

# Parties	Reference ^a	#Active Parties ^b	Security	Dot Product ^c		Dot Product with Truncation		Conversions ^d		
				Comm _{pre} ^e	Comm _{on}	Comm _{pre}	Comm _{on}	A	B	G
3	ABY3 [6]	3	Abort	12dℓ	9dℓ	12dℓ + 84ℓ	9dℓ + 3ℓ	✓	✓	✓
	BLAZE [8]	2	Fair	3ℓ	3ℓ	15ℓ	3ℓ	✓	✓	✗
	SWIFT (3PC) [14]	2	Robust	3ℓ	3ℓ	15ℓ	3ℓ	✓	✓	✗
4	Mazloom et al. [28]	4	Abort	2ℓ	4ℓ	2ℓ	4ℓ	✓	✓	✗
	Trident [4]	3	Fair	3ℓ	3ℓ	6ℓ	3ℓ	✓	✓	✓
	Tetrad	2	Fair	2ℓ	3ℓ	2ℓ	3ℓ	✓	✓	✓
	SWIFT (4PC) [14]	2	Robust	3ℓ	3ℓ	4ℓ	3ℓ	✓	✓	✗
	Fantastic Four [29] (Best) ^f	4	Robust	-	6ℓ	ℓ	9ℓ	✓	✓	✗
	Fantastic Four [29] (Worst)	3	Robust	-	6(ℓ + κ)	≈ 80ℓ + 76κ	9ℓ + 6κ	✓	✓	✗
	Tetrad-R	2	Robust	2ℓ	3ℓ	2ℓ	3ℓ	✓	✓	✓

^aAmortized costs are reported for 1 million operations ^bparties that carry out most of the computation during online phase ^cℓ - size of ring in bits, κ - security parameter, d - length of the vectors. ^dA, B, G indicate support for arithmetic, boolean, and garbled worlds respectively ^e‘Comm’ - communication, ‘pre’ - preprocessing, ‘on’ - online ^fcf. §A-D for details

Table I: Comparison of actively-secure MPC frameworks (3PC and 4PC) for PPML.

This work proposes a mixed-protocol PPML framework via MPC with four parties and honest majority with active security.

Works such as [6], [9], [28] typically go for active security with abort, where the adversary can act maliciously to obtain the output and make honest parties abort. The stronger notion of fairness guarantees that either all or none of the parties obtain the output. This incentivizes the adversary to behave honestly in resource-expensive tasks such as PPML, as causing an abort will waste its resources. Trident [4] showed that fairness can be achieved at the cost of security with abort. In cases where the risk of failure of the system is too high, for instance, when deploying PPML for healthcare applications, participants might want to avoid the case when none of them receive the output. The way to tackle this issue is to modify protocols to guarantee that the correct output is always delivered to the participants irrespective of an adversary’s misbehaviour. This is provided by guaranteed output delivery (GOD) or robustness. A robust protocol prevents the adversary from repeatedly causing the computations to rerun, thereby upholding the trust in the system. We propose two variants of the framework – one with fairness and the other with robustness. We detail the related work in §A and continue with our contributions.

A. Our Contributions

We make several contributions towards designing a practically efficient 4PC mixed-protocol framework, tolerating at most one active corruption. It operates over the ring \mathbb{Z}_{2^ℓ} and provides *end-to-end* conversions to switch between arithmetic, boolean and garbled worlds. We assume a one-time key setup phase and work in the (function-dependent) preprocessing model which paves the way for a fast online phase.

Depending on the sensitivity of the application and the underlying data, one might want different levels of security. For this, we propose two variants of the framework, covering fairness (Tetrad) and robustness (Tetrad-R) guarantees. The fair variant improves upon the state-of-the-art *fair* framework of Trident [4]. Tetrad-R improves communication over the best robust protocols [14], [29], while offering support for secure training of neural networks, which was not supported in previous works.

1) *Improved Arithmetic/Boolean 4PC*: In Tetrad, the multiplication protocol has a communication cost of only 5 ring elements as opposed to 6 in the state-of-the-art framework of Trident [4]. Security is elevated to robustness via Tetrad-R, which has a minimal overhead over the fair one, in the preprocessing. Concretely, for a 64-bit ring with 40-bit statistical security, the overhead per multiplication is 0.027 bits for a circuit containing 2^{20} multiplications. This means robustness essentially comes free in the case of large circuits. A notable contribution is the design of the multiplication protocol. It gives the following benefits – i) support for on-demand applications, ii) probabilistic truncation without overhead and iii) multi-input multiplication gates.

On-demand applications: The design of the multiplication protocol allows Tetrad to support on-demand applications where a preprocessing phase is not available. This variant of the protocols (cf. §B) has a round complexity that is the same as that of the online phases of the protocols in the preprocessing model and retains the same overall communication. It takes advantage of parallelization, which is often not possible in the *function-dependent* preprocessing model where the preprocessing and the online phases must be executed sequentially.

Probabilistic truncation without any overhead: Multiplication (and dot product) with truncation forms an essential component while working with fixed-point values. Techniques for probabilistic truncation were proposed by [6], [7]. Recently, [28] gave an efficient instantiation of truncation for 4PC with abort, based on the technique of ABY3. Using that as a baseline, we demonstrate for the *first time* in the fair and robust settings, how multiplication (and dot-product) with truncation can be performed without any additional cost over a multiplication.

Multi-input multiplication: Inspired by [23], [30], we propose new protocols for 3 and 4-input multiplication, allowing multiplication of 3 and 4 inputs in one online round. Naively, performing a 4-input multiplication follows a tree-based approach, and the required communication is that of three 2-input multiplications and 2 online rounds. Our contribution lies in keeping the communication and the round of the online phase the same as that of 2-input multiplication (i.e. invariant of the number of inputs). To achieve this, we trade off the

preprocessing cost. Looking ahead, multi-input multiplication, when coupled with the optimized parallel prefix adder circuit from [23], brings in a $2\times$ improvement in online rounds. It also cuts down the online communication of secure comparison, impacting PPML applications.

2) *4PC Mixed-Protocol Framework*: In addition to relying on the improved arithmetic/boolean world, we observe that a large portion of the computation in most MPC-based PPML frameworks is done over the arithmetic and boolean worlds. The garbled world is used only to perform the non-linear operations (e.g. softmax) that are expensive in the arithmetic/boolean world and switch back immediately after. Leveraging this observation we propose tailor-made GC-based protocols with *end-to-end* conversion techniques.

The tailor-made GC for the fair protocols, has the following advantages over Trident – i) no use of commitments for the inputs, and ii) no requirement of an explicit input sharing and output reconstruction phase, as explained later. The overall communication cost remains the same as Trident with 1 GC and 2 online rounds. In addition, for time-constrained applications we offer a variant that trades off 1 GC at the expense of 1 lesser online round. When it comes to robustness, the state-of-the-art for GC protocols are [31], costing 12 GC and 2 rounds, and [24], costing 2 GC and 4 rounds. We propose robust GC conversions for the first time, and they cost 2 GC and have an amortized round complexity of 1.

As mentioned earlier, the framework operates over three domains - arithmetic, boolean, and garbled (§IV). For an operation that required computing over the garbled domain, the standard approach is to first switch from *Arithmetic to Garbled* and evaluate the garbled circuit to obtain a garbled-shared output. These shares are brought back to the arithmetic domain using a *Garbled to Arithmetic* conversion. Our approach instead is to modify the garbled circuit such that the output is in the arithmetic domain. This eliminates the need for an explicit *Garbled to Arithmetic* conversion, saving in both communication and rounds in the online phase. More generally, end-to-end conversions are of the form “ \times -Garbled- \times ” where \times can be either arithmetic or boolean, and need a single round for the garbled world (cf. §IV).

Comparison of Tetrad with actively secure PPML frameworks in 3PC and 4PC is presented in Table I. The dot product is chosen as a parameter as it is one of the most crucial building blocks in PPML applications.

3) *Benchmarking and PPML*: We demonstrate the practicality of the framework, which combines the arithmetic, boolean, garbled worlds via benchmarking. The training and inference phases of deep neural networks such as LeNet [32] and VGG16 [33] and the inference phase of Support Vector Machines are benchmarked.

The implementation section is presented through the lens of deployment scenarios with two different goals. Participants in the first scenario are interested in the shortest online runtime for the computation, whereas participants in the second one want to minimize the deployment cost. Correspondingly, there are variants of our framework that cater to both scenarios.

Considering online runtime as the metric, Tetrad_T is the time-optimized (T) variant with the fastest online

phase. Tetrad_C is the cost-optimized (C) variant, minimizing deployment cost. This is measured via *monetary cost* [34], which helps to capture the effect of the total runtime of the parties, and communication together. Both variants are compared against Trident [4], and their relative performance is indicated in Table II. The comparison is with respect to run time, communication, monetary cost, and throughput (Table V).

Protocol	Training & Inference ^a			Training	Inference
	Time _{on} ^b	Com _{tot}	CT _{tot}	Cost	TP _{on}
Tetrad _T	●	◐	●	◐	●
Tetrad _C	◐	●	◐	●	◐
Trident	○	○	○	○	○

^a ‘Com’ - Communication, ‘Time’ - Runtime, ‘CT’ - Cumulative Runtime, ‘Cost’ - Monetary Cost, ‘TP_{on}’ - Online Throughput, on - online, tot - total
^b ○ - good, ◐ - better, ● - best, (w.r.t parameter considered)

Table II: Comparison of Trident [4] with the versions of Tetrad for deep neural networks (cf. NN-4 in §VI).

Trident requires 3 parties to be active for most of the online phase, the 4th party coming in only towards the end of the computation. In Tetrad, it is brought down to 2, having a significant impact on the monetary cost.

Table II shows that Tetrad is better when compared to Trident across all the parameters considered. Within Tetrad, Tetrad_T fares better when it comes to online run time for both training and inference, while Tetrad_C does better in terms of communication. When it comes to inference, throughput is more relevant than the cost, and here, the time-optimized variant fares the best. Robust variants follow the same trends, and the reasons behind them are elaborated in §VI.

II. PRELIMINARIES AND DEFINITIONS

We consider 4 parties denoted by $\mathcal{P} = \{P_0, P_1, P_2, P_3\}$ that are connected by pair-wise private and authentic channels in a synchronous network, and a static, active adversary that can corrupt at most 1 party. In the secure outsourced computation (SOC) setting, the 4 servers hired to carry out the computation enact the role of the 4 parties mentioned above. In this setting inputs, intermediate values, and outputs exist in a secret-shared form. For ML training, data owners secret-share their data to the servers, which train the model using MPC. The trained model can then be reconstructed towards the data owners. Our framework is secure even if the corrupt server colludes with an arbitrary number of data owners. For ML inference, the model owner secret-shares a pre-trained model among the servers. A client secret-shares its query amongst the servers, who carry out the inference via MPC. The output is reconstructed towards the client. Security is guaranteed against a corrupt server that colludes either with the model owner or with the client. We do not guarantee the privacy of the training data against attacks such as attribute inference, membership inference, or model inversion [35]–[37]. This is an orthogonal problem, and we consider it as out-of-scope of this work.

In Tetrad, parties rely on a one-time shared key setup (cf. §A for the ideal functionality) [2]–[4], [6], [8] to facilitate generation of correlated randomness non-interactively. Our

protocols work over the arithmetic ring \mathbb{Z}_{2^ℓ} or boolean ring \mathbb{Z}_2 . We use fixed-point arithmetic (FPA) [2]–[4], [6], [8] representation to deal with floating-point values where a decimal value is represented as an ℓ -bit integer in signed 2's complement representation. The most significant bit (MSB) represents the sign bit and x least significant bits are reserved for the fractional part. The ℓ -bit integer is then treated as an element of \mathbb{Z}_{2^ℓ} and operations are performed modulo 2^ℓ . We set $\ell = 64$, $x = 13$, with $\ell - x - 1$ bits for the integral part.

Notation II.1. For a vector \vec{a} , a_i denotes the i^{th} element in the vector. For two vectors \vec{a} and \vec{b} of length d , the dot product is given by, $\vec{a} \odot \vec{b} = \sum_{i=1}^d a_i b_i$. Given two matrices \mathbf{A}, \mathbf{B} , the operation $\mathbf{A} \circ \mathbf{B}$ denotes the matrix multiplication.

Notation II.2. For a bit $b \in \{0, 1\}$, b^R denotes the representation of the bit value b over the arithmetic ring \mathbb{Z}_{2^ℓ} . In detail, all the bits of b^R will be zero except for the least significant bit, which is set to b .

Primitives: For our constructs we use two standard primitives (cf. §A) (a) a collision-resistant hash function, denoted as $H(\cdot)$; (b) a garbling scheme $\mathcal{G} = (\text{Gb}, \text{En}, \text{Ev}, \text{De})$.

a) Sharing Semantics: To enforce security, we perform computation on secret-shared data. For the arithmetic and boolean sharing, we follow a $(4, 1)$ replicated secret sharing (RSS) [4], denoted by $[\![\cdot]\!]_R$. To leverage the benefits of the preprocessing paradigm, we associate meaning to the shares and demarcate the parties in terms of their roles. Three of the shares of a $(4, 1)$ RSS for a value v can be generated in the preprocessing phase independent of the value to be shared, and their sum can be interpreted as a mask. The fourth share, dependent on v , can be computed in the online phase and can be treated as the masked value. We denote the three preprocessed shares as $\lambda_v^1, \lambda_v^2, \lambda_v^3$ and the mask as $\lambda_v = \lambda_v^1 + \lambda_v^2 + \lambda_v^3$. The masked value is denoted as m_v , and $m_v = v + \lambda_v$.

Type	P_0	P_1	P_2	P_3
$[\cdot]$ -sharing ^a	–	v^1	v^2	–
(\cdot) -sharing	–	v^1	v^2	v^3
$\langle \cdot \rangle$ -sharing	–	(v^1, v^3)	(v^2, v^3)	(v^1, v^2)
$[\![\cdot]\!]_R$ -sharing ^b	$(\lambda_v^1, \lambda_v^2, \lambda_v^3)$	$(m_v, \lambda_v^1, \lambda_v^3)$	$(m_v, \lambda_v^2, \lambda_v^3)$	$(m_v, \lambda_v^1, \lambda_v^2)$

^a $v = v_1 + v_2 (+v_3)$ ^b $m_v = v + \lambda_v$

Table III: Sharing semantics for a value $v \in \mathbb{Z}_{2^\ell}$ in Tetrad. All the shares are ℓ -bit ring elements.

Next, we distinguish the four parties into two sets; the *eval* set $\mathcal{E} = \{P_1, P_2\}$ which is assigned the task of carrying out the computation, and is active throughout the online phase. The *helper* set $\mathcal{D} = \{P_0, P_3\}$ is used to assist \mathcal{E} in verification, so it is only active towards the end of the computation. Complying with the roles and the RSS format, the distribution is done as follows: $P_0 : \{\lambda_v^1, \lambda_v^2, \lambda_v^3\}$, $P_1 : \{\lambda_v^1, \lambda_v^3, m_v\}$, $P_2 : \{\lambda_v^2, \lambda_v^3, m_v\}$, and $P_3 : \{\lambda_v^1, \lambda_v^2, m_v\}$. The shares are distributed among \mathcal{D} such that P_3 gets m_v whereas P_0 gets all the shares of λ_v . During preprocessing, P_0 computes a part of the data needed for verification (cf. Fig. 3) using its input independent shares, which is communicated to P_3 . This enables a verification in the online without P_0 , for the fair protocols.

Exploiting the asymmetry of the roles allows for minimal online participation, giving a huge improvement in the cumulative runtime (sum of uptime of all the parties), thereby saving in monetary costs (cf. §VI). The RSS sharing semantics are presented in Table III, denoted by $[\![\cdot]\!]_R$, in a modular way with the help of three intermediate sharing semantics $[\cdot]$, (\cdot) and $\langle \cdot \rangle$. All the schemes used are linear i.e. given shares of values v_1, \dots, v_m and public constants c_1, \dots, c_m , sharing of $\sum_{i=1}^m c_i v_i$ can be computed locally for an integer m .

Notation II.3. (a) For the $[\![\cdot]\!]_R$ -shares of n values a_1, \dots, a_n , $\gamma_{a_1 \dots a_n} = \prod_{i=1}^n \lambda_{a_i}$ and $m_{a_1 \dots a_n} = \prod_{i=1}^n m_{a_i}$ (b) We use superscripts \mathbf{B} , and \mathbf{G} to denote sharing semantics in boolean, and garbled world, respectively– $[\![\cdot]\!]_R^{\mathbf{B}}$, $[\![\cdot]\!]_R^{\mathbf{G}}$. We omit the superscript for arithmetic world.

Sharing semantics for boolean sharing over \mathbb{Z}_2 is similar to arithmetic sharing except that addition is replaced with XOR. The semantics for garbled sharing are described in §IV with the relevant context.

III. 4PC PROTOCOL

This section covers the details of our 4PC protocol over an arithmetic ring \mathbb{Z}_{2^ℓ} . We begin by explaining the relevant primitives in §III-A. The multiplication protocol with abort is presented in §III-B, followed by details on elevating the security to fairness in §III-C. Lastly, in §III-D, we show how to improve the security to robustness¹. Formal details along with the cost analysis for the protocols are provided in the full version of the paper [38].

A. Primitives

a) Joint-Send (jsnd): The Joint-Send (jsnd) primitive allows two parties P_i, P_j to relay a message v to a third party P_k ensuring either the delivery of the message or abort in case of inconsistency. Towards this, P_i sends v to P_k , while P_j sends a hash of the same, $H(v)$, to P_k . Party P_k accepts the message if the hash values are consistent and aborts otherwise. Note that the communication of the hash can be clubbed together for several instances and be deferred to the end of the protocol, amortizing the cost.

b) Joint-Send (jsnd) for robust protocols: To achieve robustness, we instantiate jsnd using the joint-message passing (jmp) primitive of [14]. The jsnd primitive (Fig. 11) allows two senders P_i, P_j to relay a common message, $v \in \mathbb{Z}_{2^\ell}$, to a recipient P_k , either by ensuring successful delivery of v , or by establishing a conflicting pair of parties, one among which is guaranteed to be corrupt. This implies the residual two parties are honest, one of which is then entrusted to take the computation to completion by enacting the role of a trusted party (P_{TP}). The instantiation of jsnd can be viewed as consisting of two phases (*send, verify*), where the *send* phase consists of P_i sending v to P_k and the rest of the protocol steps go to *verify* phase (which ensures correct *send* or P_{TP} identification). This requires 1 round of interaction and ℓ bits of communication. To leverage amortization, *verify* will be executed only once, at the end of the computation, and requires 2 rounds.

¹The classical notion of robustness is achieved

The jsnd primitive is instantiated depending on the desired security guarantee. For simplicity, we give common constructions for fair and robust variants of the protocols, when they only differ in the instantiation of jsnd.

Notation III.1. Protocol Π_{jsnd} denotes the instantiation of Joint-Send (jsnd) primitive. We say that P_i, P_j jsnd v to P_k when they invoke $\Pi_{\text{jsnd}}(P_i, P_j, v, P_k)$.

c) *Sharing:* Protocol Π_{Sh} (Fig. 1) enables P_i to generate $\llbracket \cdot \rrbracket$ -share of a value v . During the preprocessing phase, λ -shares are sampled non-interactively using the pre-shared keys (cf. §A-B) in a way that P_i will get the entire mask λ_v . During the online phase, P_i computes $m_v = v + \lambda_v$ and sends to P_1, P_2, P_3 , which exchange the hash values to check for consistency. Parties abort in the fair protocol in case of inconsistency, whereas for robust security, parties proceed with a default value.

Protocol $\Pi_{\text{Sh}}(P_i, v)$

Preprocessing: Sample the following:

$$P_i, P_0, P_1, P_3 : \lambda_v^1 \quad \Bigg| \quad P_i, P_0, P_2, P_3 : \lambda_v^2 \quad \Bigg| \quad P_i, P_0, P_1, P_2 : \lambda_v^3$$

Online:

- 1) P_i computes $m_v = v + \lambda_v$ and sends to P_1, P_2, P_3 .
- 2) P_1, P_2, P_3 mutually exchange $H(m_v)$ and accept the sharing if there exists a majority. Else parties abort for the case of fairness and accept a default value for the case of robust security.

Figure 1: $\llbracket \cdot \rrbracket$ -sharing of a value v by party P_i .

d) *Joint Sharing:* Protocol Π_{JSh} enables parties P_i, P_j to generate $\llbracket \cdot \rrbracket$ -share of a value v . The protocol is similar to Π_{Sh} except that P_j ensures the correctness of the sharing performed by P_i . During the preprocessing, λ -shares are sampled such that both P_i, P_j will get the entire mask λ_v . During the online phase, P_i, P_j compute and jsnd $m_v = v + \lambda_v$ to parties P_1, P_2, P_3 .

For joint-sharing a value v possessed by P_0 along with another party in the preprocessing, the communication can be optimized further. The protocol steps based on the (P_i, P_j) pair are summarised below:

- $(P_0, P_1) : \mathcal{P} \setminus \{P_2\}$ sample $\lambda_v^1 \in_R \mathbb{Z}_{2^\ell}$; Set $\lambda_v^2 = m_v = 0$; P_0, P_1 jsnd $\lambda_v^3 = -v - \lambda_v^1$ to P_2 .
- $(P_0, P_2) : \mathcal{P} \setminus \{P_3\}$ sample $\lambda_v^2 \in_R \mathbb{Z}_{2^\ell}$; Set $\lambda_v^1 = m_v = 0$; P_0, P_2 jsnd $\lambda_v^3 = -v - \lambda_v^2$ to P_3 .
- $(P_0, P_3) : \mathcal{P} \setminus \{P_1\}$ sample $\lambda_v^3 \in_R \mathbb{Z}_{2^\ell}$; Set $\lambda_v^1 = m_v = 0$; P_0, P_3 jsnd $\lambda_v^2 = -v - \lambda_v^3$ to P_1 .

e) *Reconstruction:* Protocol $\Pi_{\text{Rec}}(\mathcal{P}, v)$ (Fig. 13) enables parties in \mathcal{P} to compute v , given its $\llbracket \cdot \rrbracket$ -share. Note that each party misses one share to reconstruct the output, and the other 3 parties hold this share. 2 out of the 3 parties will jsnd the missing share to the party that lacks it. Reconstruction towards a single party can be viewed as a special case.

f) $\mathcal{F}_{\text{zero}}$ - *Generating additive shares of zero:* In Tetrad, we make use of a functionality $\mathcal{F}_{\text{zero}}$ to enable parties P_0, P_i obtain Z_i for $i \in \{1, 2, 3\}$ such that $Z_1 + Z_2 + Z_3 = 0$.

We observe that the functionality can be instantiated non-interactively using the pre-shared keys (cf. §A-B). For this, parties in $\mathcal{P} \setminus \{P_j\}$ sample random value r_j for $j \in \{1, 2, 3\}$. The shares are then defined as $Z_1 = r_3 - r_2, Z_2 = r_1 - r_3$ and $Z_3 = r_2 - r_1$.

g) *Multiplication of $\langle a \rangle, \langle b \rangle$, held in clear by P_0 :* To multiply $\langle a \rangle, \langle b \rangle$, where $a, b \in \mathbb{Z}_{2^\ell}$ are held in clear by P_0 , and generate $\langle z \rangle$ such that $z = ab$, Π_{MulR} (Fig. 2) proceed as follows. Parties locally generate a (\cdot) -sharing of z , where P_0 knows all three (\cdot) -shares. To complete the generation of $\langle z \rangle$, P_0, P_i for $i \in \{1, 2, 3\}$, randomize their (\cdot) -share of z using (\cdot) -share of 0, and jsnd $\langle z \rangle^i$, to one other party.

Protocol $\Pi_{\text{MulR}}(\langle a \rangle, \langle b \rangle)$

- 1) Invoke $\mathcal{F}_{\text{zero}}$ to enable P_0, P_j obtain Z_j for $j \in \{1, 2, 3\}$ such that $Z_1 + Z_2 + Z_3 = 0$.

$$P_0, P_1 \text{ jsnd } \langle z \rangle^1 = a^1 b^3 + a^3 b^1 + a^3 b^3 + Z_1 \text{ to } P_2.$$

$$P_0, P_2 \text{ jsnd } \langle z \rangle^2 = a^2 b^3 + a^3 b^2 + a^2 b^2 + Z_2 \text{ to } P_3.$$

$$P_0, P_3 \text{ jsnd } \langle z \rangle^3 = a^1 b^2 + a^2 b^1 + a^1 b^1 + Z_3 \text{ to } P_1.$$

- 2) Set $\langle z \rangle$ as $z^1 = \langle z \rangle^3, z^2 = \langle z \rangle^2, z^3 = \langle z \rangle^1$.

Figure 2: Multiplication of $\langle \cdot \rangle$ -shared values, held on clear by P_0 .

B. Multiplication in Tetrad

Given the shares of a, b , the goal of the multiplication protocol is to generate shares of $z = ab$. The protocol is designed such that parties P_1, P_2 obtain a masked version of the output z , say $z - r$ in the online phase, and P_0, P_3 obtain the mask r in the preprocessing phase. Parties then generate $\llbracket \cdot \rrbracket$ -sharing of these values by executing Π_{JSh} , and locally compute $\llbracket z - r \rrbracket + \llbracket r \rrbracket$ to obtain the final output.

a) *Online:* Note that,

$$\begin{aligned} z - r &= ab - r = (m_a - \lambda_a)(m_b - \lambda_b) - r \\ &= m_{ab} - m_a \lambda_b - m_b \lambda_a + \gamma_{ab} - r \quad (\text{cf. notation II.3}) \end{aligned} \quad (1)$$

In Eq 1, P_1, P_2 can compute m_{ab} locally, and hence we are interested in computing $y = (z - r) - m_{ab}$. Let us view y as $y = y_1 + y_2 + y_3$, where y_1 and y_2 can be computed respectively by P_1 and P_2 , and y_3 consists of terms that can be computed by both.

$$\begin{aligned} P_1 : y_1 &= -\lambda_a^1 m_b - \lambda_b^1 m_a + [\gamma_{ab} - r]_1 \\ P_2 : y_2 &= -\lambda_a^2 m_b - \lambda_b^2 m_a + [\gamma_{ab} - r]_2 \\ P_1, P_2 : y_3 &= -\lambda_a^3 m_b - \lambda_b^3 m_a \end{aligned} \quad (2)$$

The preprocessing is set up such that P_1, P_2 receive additive shares $([\cdot])$ of $\gamma_{ab} - r$. P_1, P_2 then mutually exchange the missing share to reconstruct y and subsequently $z - r$.

b) *Verification:* To ensure correctness of the values exchanged in the online phase, we use the assistance of P_3 . Concretely, P_3 obtains $y_1 + y_2 + s$, where s is a random mask known to P_0, P_1, P_2 . For this, P_3 needs $\gamma_{ab} + s$, which it obtains from the preprocessing phase. The mask s is used to prevent the leakage from γ_{ab} to P_3 . P_3 computes a hash of $y_1 + y_2 + s$ and sends it to P_1, P_2 , which abort if it is inconsistent.

c) *Preprocessing*: Parties should obtain the following values from the preprocessing phase:

$$i) P_1, P_2 : [\gamma_{ab} - r] \quad | \quad ii) P_0, P_3 : r \quad | \quad iii) P_3 : \gamma_{ab} + s$$

For i) and ii), let $\gamma_{ab} = \gamma_{ab}^1 + \gamma_{ab}^2 + \gamma_{ab}^3$, where P_0 along with P_i can compute γ_{ab}^i for $i \in \{1, 2, 3\}$. For P_1, P_2 , to form an additive sharing of $(\gamma_{ab} - r)$, it suffices for them to define their share as $\gamma_{ab}^i + [\gamma_{ab}^3 - r]$. Instead of sampling a fresh random value for r , P_0, P_3 , along with P_i , sample the share for $\gamma_{ab}^3 - r$ as u^i for $i \in \{1, 2\}$. P_0, P_3 compute r as $\gamma_{ab}^3 - u^1 - u^2$. Note that r computed this way is still uniformly random, as u^1, u^2 are sampled uniformly at random.

For iii), P_3 needs $w = \gamma_{ab}^1 + \gamma_{ab}^2 + s$. To tackle this, P_0, P_1, P_2 sample s_1, s_2 , and set $s = s_1 + s_2$. P_0, P_i , for $i \in \{1, 2\}$, send $\gamma_{ab}^i + s_i$ to P_3 . This requires a communication of 2 elements. As an optimization, P_0 sends w to P_3 . If P_0 is malicious, it might send a wrong value to P_3 . However, in this case, every party in the online phase would be honest. And since P_1, P_2 do not use w in their computation, any error in w is bound to get caught in the verification phase.

Protocol $\Pi_{\text{Mult}}(a, b, \text{isTr})$

Let isTr be a bit that denotes whether truncation is required ($\text{isTr} = 1$) or not ($\text{isTr} = 0$).

Preprocessing:

1) Locally compute:

$$\begin{aligned} P_0, P_1 : \gamma_{ab}^1 &= \lambda_a^1 \lambda_b^3 + \lambda_a^3 \lambda_b^1 + \lambda_a^3 \lambda_b^3 \\ P_0, P_2 : \gamma_{ab}^2 &= \lambda_a^2 \lambda_b^3 + \lambda_a^3 \lambda_b^2 + \lambda_a^2 \lambda_b^2 \\ P_0, P_3 : \gamma_{ab}^3 &= \lambda_a^1 \lambda_b^2 + \lambda_a^2 \lambda_b^1 + \lambda_a^1 \lambda_b^1 \end{aligned}$$

2) P_0, P_3 and P_j sample random $u^j \in_R \mathbb{Z}_{2^\ell}$ for $j \in \{1, 2\}$. Let $u^1 + u^2 = \gamma_{ab}^3 - r$ for a random $r \in_R \mathbb{Z}_{2^\ell}$.

3) P_0, P_3 compute $r = \gamma_{ab}^3 - u^1 - u^2$ and set $q = r^t$ if $\text{isTr} = 1$, else set $q = r$. P_0, P_3 execute $\Pi_{\text{JSh}}(P_0, P_3, q)$ to generate $\llbracket q \rrbracket$.

4) P_0, P_1, P_2 sample random $s_1, s_2 \in_R \mathbb{Z}_{2^\ell}$ and set $s = s_1 + s_2$. P_0 sends $w = \gamma_{ab}^1 + \gamma_{ab}^2 + s$ to P_3 .

Online: Let $y = (z - r) - m_a m_b$.

1) Locally compute:

$$\begin{aligned} P_1 : y_1 &= -\lambda_a^1 m_b - \lambda_b^1 m_a + \gamma_{ab}^1 + u^1 \\ P_2 : y_2 &= -\lambda_a^2 m_b - \lambda_b^2 m_a + \gamma_{ab}^2 + u^2 \\ P_1, P_2 : y_3 &= -\lambda_a^3 m_b - \lambda_b^3 m_a \end{aligned}$$

2) P_1 sends y_1 to P_2 , while P_2 sends y_2 to P_1 , and they locally compute $z - r = (y_1 + y_2 + y_3) + m_a m_b$.

3) If $\text{isTr} = 1$, P_1, P_2 set $p = (z - r)^t$, else $p = z - r$. P_1, P_2 execute $\Pi_{\text{JSh}}(P_1, P_2, p)$ to generate $\llbracket p \rrbracket$.

4) Locally compute $\llbracket o \rrbracket = \llbracket p \rrbracket + \llbracket q \rrbracket$. Here $o = z^t$ if $\text{isTr} = 1$ and z otherwise.

5) *Verification*: P_3 computes $v = -(\lambda_a^1 + \lambda_a^2) m_b - (\lambda_b^1 + \lambda_b^2) m_a + u^1 + u^2 + w$ and sends $H(v)$ to P_1 and P_2 . Parties P_1, P_2 abort iff $H(v) \neq H(y_1 + y_2 + s)$.

^aFor the fair protocol, it is enough for P_0, P_1, P_2 to sample s directly.

d) *Truncation*: For a value $v = v_1 + v_2$, SecureML [7] showed that the truncated value $v/2^x$, denoted by v^t , can be computed as $v_1^t + v_2^t$. With high probability, a truncated value having at most one bit error in the least significant position is generated. It was shown in SecureML that accuracy drop for ML algorithms due to the one bit error is minimal. However, the method cannot be generalized to more than two parties. ABY3 [6] demonstrated the extension to 3-party setting with a generic design that uses a truncation pair of the form (r, r^t) . Here, r is a random value and r^t denotes its truncated version. Given this pair, z can be truncated by opening $z - r$ towards all, and computing z^t as $z^t = (z - r)^t + r^t$. Note that all operations are carried out on shares. The design of our multiplication allows for truncation to be carried out this way without any additional overhead in communication. Towards this, P_1, P_2 locally truncate $(z - r)$ and generate $\llbracket \cdot \rrbracket$ -shares of it in the online phase. Similarly, P_0, P_3 truncate r in the preprocessing phase and generate its $\llbracket \cdot \rrbracket$ -shares. Then $\llbracket z^t \rrbracket = \llbracket (z - r)^t \rrbracket + \llbracket r^t \rrbracket$

e) *Multiplication by constant*: This operation in MPC is typically local: given constant α and $\llbracket v \rrbracket$, the product can be written as $\alpha v = \beta^1 + \beta^2$ where $\beta^1 = \alpha \cdot (m_v - \lambda_v^3)$ and $\beta^2 = \alpha \cdot (-\lambda_v^1 - \lambda_v^2)$. However, in FPA, we need to perform a truncation on the output. For this P_1, P_2 truncate β^1 and execute Π_{JSh} , while P_0, P_3 do the same with β^2 .

C. Achieving Fairness

Here we show how to extend the security of Tetrads from abort to fairness using techniques from Trident [4]. Before proceeding with the output reconstruction, we need to ensure that all the honest parties are alive after the verification phase. For this, all the parties maintain an *aliveness* bit, say b , which is initialized to *continue*. If the verification phase is not successful for a party, it sets $b = \text{abort}$. In the first round of reconstruction, the parties mutually exchange their b bit and accept the value that forms the majority. Since we have only one corruption, it is guaranteed that all the honest parties will be in agreement on b . If $b = \text{continue}$, then the parties exchange their missing shares and accept the majority. As per the sharing semantics, every missing share is possessed by three parties, out of which there can be at most one corruption. As an optimization, for instances where many values are reconstructed, two out of the three parties can send the share while the third can send a hash of the same.

D. Achieving Robustness

Here we show how to extend the security of Tetrads to provide robustness while retaining the same amortized communication complexity. The robust variant, denoted by Tetrads-R, additionally requires a verification check in the preprocessing phase of multiplication as compared to Tetrads. Moreover, the reconstruction protocol is similar to the fair counterpart, except that aliveness check is not required since a cheating would result in identifying an honest party (P_{TP}).

The multiplication protocol Π_{Mult} (Fig. 3) is modified as follows. First, the robust variant of Π_{JSh} is used instead of the fair one. This ensures correctness of messages to be communicated or identifies a conflicting pair of parties, one among which is guaranteed to be corrupt. Next, to ensure the correctness of w sent by P_0 alone in the preprocessing

Figure 3: Multiplication with / without truncation in Tetrads.

phase, we introduce $\Pi_{\text{Vrfy}P_0}$ (Fig. 4). If $\Pi_{\text{Vrfy}P_0}$ fails, parties identify a P_{TP} in the preprocessing phase itself. Finally, in case of an abort in the online phase (which proceeds similar to the that of Π_{Mult}), P_0 is assigned as the P_{TP} . Since P_0 does not participate in the online phase of multiplication, and its communication in the preprocessing has been verified via $\Pi_{\text{Vrfy}P_0}$, this assignment is safe.

Verifying the communication by P_0 . In Π_{Mult} (Fig. 3), P_0 computes and sends $w = \gamma_{ab}^1 + \gamma_{ab}^2 + s_1 + s_2$ to P_3 , where P_0, P_1, P_2 know s_1, s_2 in clear. Note that $w = w^1 + w^2$ for $w^1 = \gamma_{ab}^1 + s_1$ and $w^2 = \gamma_{ab}^2 + s_2$. Also, P_0 along with P_1, P_2 and P_3 possess the values w^1, w^2 and w respectively. Checking the correctness of w thus reduces to verifying if $w = w^1 + w^2$.

To verify this relation for all M multiplication gates in the circuit, i.e. $\{w_j \stackrel{?}{=} w_j^1 + w_j^2\}_{j \in [M]}$, one approach is to compute a random linear combination and verify the relation on the sum. While working over a field \mathbb{F}_p , this solution has an error probability $1/|\mathbb{F}_p|$, where $|\mathbb{F}_p|$ denotes the size of \mathbb{F}_p . However, this solution does not work naively over rings since not every element in the ring has an inverse, as opposed to fields. Concretely, the check can still pass with a probability of at most $1/2$ [39], [40]. To reduce the cheating probability, the check is repeated κ times, thereby bounding the cheating probability by $1/2^\kappa$. As an optimization, it is sufficient to choose the random combiners from $\{0, 1\}$. Thus, for one check, parties need to sample only a binary string of M bits using the shared-key. The formal verification protocol appears in Fig. 4.

Protocol $\Pi_{\text{Vrfy}P_0}(\{\{w_j\}_{j=1}^M\})$

Repeat the following κ times, in parallel.

- 1) Sample random values $\tau_1, \dots, \tau_M \in \mathbb{Z}_{2^\ell}$.
- 2) Locally compute: $P_0, P_1 : e^1 = \sum_{j=1}^M \tau_j w_j^1$; $P_0, P_2 : e^2 = \sum_{j=1}^M \tau_j w_j^2$; $P_0, P_3 : e = \sum_{j=1}^M \tau_j w_j$.
- 3) (P_0, P_1) , (P_0, P_2) and (P_0, P_3) generate $[\cdot]$ -shares of e^1, e^2 and e respectively using Π_{Jsh} .
- 4) Locally compute $[g] = [e] - [e^1] - [e^2]$.
- 5) Robustly reconstruct g and check if $g \stackrel{?}{=} 0$.

If for all κ repetitions, $g = 0$, then continue with rest of the computation. Else, P_0 is identified to be corrupt and $P_{\text{TP}} = P_1$.

Figure 4: Verification of P_0 's communication in the multiplication protocol of Tetrad-R

The robust protocol can be optimized further if cheating is detected (abort signal is generated) in the preprocessing phase. Concretely, this can be identified in the preprocessing phase either from the verification of jsnd instances or output of $\Pi_{\text{Vrfy}P_0}$. When such a cheating is detected, the corrupt party is identified as follows. Parties first broadcast their shared keys established in the key-setup phase. They recompute all the preprocessing data and verify against the data that was communicated to identify the corrupt party. Note that disclosing the shared keys does not violate input privacy because the preprocessing data is input independent. On identifying the corrupt party, it is eliminated from the computation, and a semi-honest 3-party computation is performed from this point onwards.

E. The complete 4PC

The above primitives can be compiled to compute an arithmetic circuit over \mathbb{Z}_{2^ℓ} as follows.

Parties first invoke the key-setup functionality $\mathcal{F}_{\text{Setup}}$ (Fig. 8) for key distribution, and preprocessing of input sharing (Π_{Sh}) and multiplication (Π_{Mult}), as per the given circuit. This generates the masks (λ) for all the wires in the circuit as per the sharing semantics. The preprocessing for linear gates can be performed non-interactively. The verification of all the protocols is executed before moving on to the online phase.

During the online phase, $P_i \in \mathcal{P}$ shares its input x_i by executing online steps of Π_{Sh} (Fig. 1). Parties then evaluate the gates in the circuit in the topological order, with linear gates being computed locally, and multiplication gates being computed via online phase of Π_{Mult} (Fig. 3). Finally, Π_{Rec} (Fig. 13) is executed for the output wires to reconstruct the output.

F. Supporting on-demand computations

For on-demand applications where the underlying function to be computed is not known in advance, the preprocessing model is not desirable. We observe that the Tetrad protocol can be modified by executing the preprocessing phase in the online phase itself, keeping the same overall communication cost. The formal protocol appears in Fig. 12.

IV. MIXED PROTOCOL FRAMEWORK

In the applications we consider, the garbled circuit is used as an intermediary to evaluate certain functions where the input to the function as well as the output are in $[\cdot]$ -shared (or $[\cdot]^B$ -shared) form. For this, we design end-to-end conversions which are of the form “x-Garbled-x” where x can be either arithmetic or boolean.

Similar to Trident [4], we design a fair GC world, using techniques from [41], that requires communicating 1 GC and 2 rounds for end-to-end conversions. We further extend it to provide robustness without inflating the cost. Due to its close resemblance to Trident, the details are provided in the full version of the paper [38]. We observe that the online rounds for end-to-end conversions can be further reduced to 1 at the expense of communicating one more GC in a parallel execution. Note that a similar approach of using 2 parallel executions in Trident does not lead to obtaining a 1-round conversion due to their protocol design and reliance on piecewise conversions. A high-level comparison is provided in Table IV, and more details appear in the full version [38].

Protocol ^a	Reference	Communication ^b (Preprocessing)	Rounds (Online)	Communication (Online)
2 GC variant	Trident	$6\ell\kappa + \ell$	2	$4\ell\kappa + 2\ell$
	Tetrad		1	$4\ell\kappa + \ell$
1 GC variant	Trident	$3\ell\kappa + \ell$	2	$2\ell\kappa + 3\ell$
	Tetrad		2	$2\ell\kappa + 2\ell$

^a Notations: ℓ - size of ring in bits, κ - computational security parameter.

^b Cost of GC is omitted, see [38] for more details.

Table IV: End-to-end conversions in Trident [4] and Tetrad.

When compared to the standalone protocol of [41], the customized fair GC protocol for mixed framework eliminates

the need for commitments to ensure input consistency and explicit input sharing and output reconstruction phases. For robustness, the standalone GC protocols of [31] requires communicating 12 GCs in 2 rounds while [24] communicates 2 GCs in 4 rounds. On the other hand, the robust variant in this work requires communicating 2 GC in 1 round. Moreover, these protocols leverage the benefit of amortization which comes from using jsnd.

Leveraging an honest majority among the garblers and using jsnd, we only need semi-honest GC computation to get active security. Moreover, the state-of-the-art GC optimizations of free-XOR [42], [43], half gates [44], [45], and fixed AES-key [46] are deployed in our protocol.

A. GC for mixed protocol framework

The 2 GC variant has two parallel executions, each comprising of 3 garblers and 1 evaluator. P_1, P_2 act as evaluators in two independent executions and the parties in $\Phi_1 = \{P_0, P_2, P_3\}$, $\Phi_2 = \{P_0, P_1, P_3\}$ act as garblers, respectively. Note that it suffices for only P_0, P_3 to generate and jsnd the GC to the evaluator.

Garbled evaluation proceeds in three phases– i) Input phase, ii) Evaluation, and iii) Output phase. The input phase involves transferring the keys to the evaluators for every input to the GC. Note here that the function (to be evaluated via the GC) input is already $[\cdot]^B$ -shared. Since each share of the function input is available with two garblers in each garbling instance, the correct key transfer is ensured via jsnd. The evaluation consists of GC transfer followed by GC evaluation. Lastly, in the output phase, evaluators obtain the encoded output. Preliminary details about the garbling scheme and additional details of the GC protocol are given in full version of the paper [38].

a) Input Phase: Given that the function input x is already available as $[x]^B$, the boolean values $m_x, \alpha_x, \lambda_x^3$, where $\alpha_x = \lambda_x^1 \oplus \lambda_x^2$ and $x = m_x \oplus \alpha_x \oplus \lambda_x^3$, act as the *new* inputs for the garbled computation, and garbled sharing ($[\cdot]^G$) is generated for each of these values. The semantics of $[\cdot]^B$ -sharing ensures that each of these shares ($m_x, \alpha_x, \lambda_x^3$) is available with two garblers in each garbling instance. The keys for the shares can either be sent (using jsnd) correctly to the evaluators or the inconsistency is detected. This key delivery essentially generates $[\cdot]^G$ -sharing for each of these three values which enables GC evaluation. Thus, the goal of our input phase is to create the compound sharing, $[x]^G = ([m_x]^G, [\alpha_x]^G, [\lambda_x^3]^G)$ for every input x to the function to be evaluated via the GC. We first discuss the semantics for $[\cdot]^G$ -sharing followed by steps for generating $[\cdot]^C$ -sharing.

b) Garbled sharing semantics: A value $v \in \mathbb{Z}_2$ is $[\cdot]^G$ -shared (garbled shared) amongst \mathcal{P} if $P_i \in \{P_0, P_3\}$ holds $[v]_i^G = (K_v^{0,1}, K_v^{0,2})$, P_1 holds $[v]_1^G = (K_v^{v,1}, K_v^{v,2})$ and P_2 holds $[v]_2^G = (K_v^{0,1}, K_v^{v,2})$. Here, $K_v^{v,j} = K_v^{0,j} \oplus v\Delta^j$ for $j \in \{1, 2\}$, and Δ^j , which is known only to the garblers in Φ_j , denotes the global offset with its least significant bit set to 1 and is same for every wire in the circuit. A value $x \in \mathbb{Z}_2$ is said to be $[\cdot]^C$ -shared (compound shared) if each value from $(m_x, \alpha_x, \lambda_x^3)$, which are as defined above, is $[\cdot]^G$ -shared. We write $[x]^G = ([m_x]^G, [\alpha_x]^G, [\lambda_x^3]^G)$.

c) Generation of $[v]^G$ and $[x]^C$: Protocol $\Pi_{\text{Sh}}^G(\mathcal{P}, v)$ (Fig. 5) enables generation of $[v]^G$ where two garblers in each garbling instance hold v , and proceeds as follows. Consider the first garbling instance with evaluator P_1 where garblers P_k, P_l hold v . Garblers in Φ_1 generate $\{K_v^{b,1}\}_{b \in \{0,1\}}$ which denotes the key for value b on wire v , following the free-XOR technique [42], [43]. P_k, P_l jsnd $K_v^{v,1}$ to evaluator P_1 . Similar steps carried out with respect to the second garbling instance, at the end of which, garblers in Φ_2 possess $\{K_v^{b,2}\}_{b \in \{0,1\}}$ while the evaluator P_2 holds $K_v^{v,2}$. Following this, the shares $[v]_s^G$ held by $P_s \in \mathcal{P}$ are defined as $[v]_0^G = [v]_3^G = (K_v^{0,1}, K_v^{0,2})$, $[v]_1^G = (K_v^{v,1}, K_v^{0,2})$, $[v]_2^G = (K_v^{0,1}, K_v^{v,2})$.

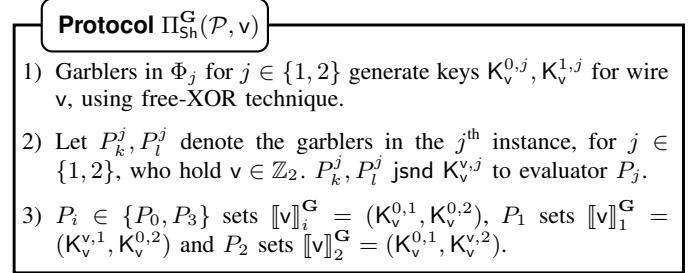


Figure 5: Generation of $[v]^G$

To generate $[x]^C$, we need a way to generate $([m_x]^G, [\alpha_x]^G, [\lambda_x^3]^G)$, given $[x]^B$. For this, Π_{Sh}^G is invoked for each of $m_x, \alpha_x, \lambda_x^3$.

B. Conversions involving Garbled World

Assume the GC is required to compute a function f on inputs $x, y \in \mathbb{Z}_{2^\ell}$ and let the output be $f(x, y)$. All the conversions described are for the 2 GC variant. Conversions for the 1 GC variant are straightforward, hence we omit the details. The conversions are generic for fair and robust variants, where the security follows from that of the underlying primitives.

Case I: Boolean-Garbled-Boolean. Since the inputs to the GC are available in boolean form, say $[x]^B, [y]^B$, parties generate $[x]^C, [y]^C$ by invoking the garbled sharing protocol Π_{Sh}^G . Additionally, parties P_0, P_3 sample $R \in \mathbb{Z}_{2^\ell}$ to mask the function output, $f(x, y)$, and generate $[R]^B$ (using the joint sharing protocol) and $[R]^G$. Garblers $P_g \in \{P_0, P_2, P_3\}$ garble the circuit which computes $z = f(x, y) \oplus R$, and send the GC along with the decoding information to evaluator P_1 . Analogous steps are performed for evaluator P_2 . Upon GC evaluation and output decoding, evaluators obtain $z = f(x, y) \oplus R$, and jointly boolean share z to generate $[z]^B$. Parties then compute $[f(x, y)]^B = [z]^B \oplus [R]^B$.

Case II: Boolean-Garbled-Arithmetic. This is similar to *Case I* except that the circuit which computes $z = f(x, y) + R$ is garbled instead. Boolean sharing of z is replaced with arithmetic, followed by computing $[f(x, y)] = [z] - [R]$.

Cases III & IV: Input in Arithmetic Sharing. The function to be computed $f(x, y)$, is modified as $f'(m_x, \alpha_x, \lambda_x^3, m_y, \alpha_y, \lambda_y^3) = f(m_x - \alpha_x - \lambda_x^3, m_y - \alpha_y - \lambda_y^3)$ where inputs x, y are replaced by the triples $\{m_x, \alpha_x, \lambda_x^3\}, \{m_y, \alpha_y, \lambda_y^3\}$ and $\alpha_x = \lambda_x^1 + \lambda_x^2$ and $\alpha_y = \lambda_y^1 + \lambda_y^2$.

The circuit to be garbled thus, corresponds to the function f' . Parties generate $[[m_x]]^G, [[\alpha_x]]^G, [[\lambda_x^3]]^G, [[m_y]]^G, [[\alpha_y]]^G, [[\lambda_y^3]]^G$ via Π_{Sh}^G , following which, parties proceed with the rest of the computation whose steps are similar to *Case I*, and *II*, depending on the requirement on the output sharing.

C. Other Conversions

a) Arithmetic to Boolean: To convert arithmetic sharing of $v \in \mathbb{Z}_{2^\ell}$ to boolean sharing, observe that $v = v_1 + v_2$ where $v_1 = m_v - \lambda_v^3$ is possessed by parties P_1, P_2 , while $v_2 = -(\lambda_v^1 + \lambda_v^2)$ is possessed by parties P_0, P_3 . Thus, $[[v]]^B$ can be computed as $[[v]]^B = [[v_1]]^B + [[v_2]]^B$, where $[[v_2]]^B$ can be generated in the preprocessing phase, and $[[v_1]]^B$ can be generated in the online phase by the respective parties executing joint boolean sharing protocol. Boolean addition, when instantiated using the adder of ABY2.0 [23], requires $\log_4(\ell)$ rounds.

b) Boolean to Arithmetic: To convert a boolean sharing of v into an arithmetic sharing, we use techniques from [4], [14]. For a value $v \in \mathbb{Z}_{2^\ell}$, note that

$$\begin{aligned} v &= \sum_{i=0}^{\ell-1} 2^i v_i = \sum_{i=0}^{\ell-1} 2^i (\lambda_{v_i} \oplus m_{v_i}) \\ &= \sum_{i=0}^{\ell-1} 2^i \left(m_{v_i}^R + \lambda_{v_i}^R (1 - 2m_{v_i}^R) \right) \end{aligned}$$

where $\lambda_{v_i}^R, m_{v_i}^R$ denote the arithmetic value of bits λ_{v_i}, m_{v_i} over the ring \mathbb{Z}_{2^ℓ} . For each bit v_i of v , parties generate the arithmetic sharing of $\lambda_{v_i}^R$ in the preprocessing, using techniques from bit to arithmetic protocol (cf. full version [38]). During the online phase, additive shares for each bit v_i is locally computed similar to bit to arithmetic protocol. Parties then multiply the i th share with 2^i and locally add up to obtain an additive sharing of v . The rest of the steps are similar to the bit to arithmetic protocol, and the formal protocol appears in full version [38].

V. BUILDING BLOCKS

This section covers the primitives needed for realising privacy-preserving variants of the applications considered, and elaborate details appear in the full version [38]. The building blocks can be combined to construct different layers in a neural network, as shown in [10] (Fig. 3).

a) Dot Product (Scalar Product): Given $[[\vec{a}]], [[\vec{b}]]$ with $|\vec{a}| = |\vec{b}| = d$, protocol Π_{dotp} computes $[[z]]$ such that $z = (\vec{a} \odot \vec{b})^t$ if truncation is enabled, else $z = \vec{a} \odot \vec{b}$. Following [4], [14], we combine the partial products from the multiplication protocol across d multiplications and communicate them in a single shot. This makes the communication cost of the dot product independent of the vector size. Due to the similarity of Π_{dotp} protocol with multiplication (Fig. 3), we provide the formal details in the full version [38]. The protocol for robust setting follows similarly.

Matrix multiplication is an extension of the dot product protocol. We abuse notation and follow the $[[\cdot]]$ -sharing semantics (ref. §II) for matrices as well. For $\mathbf{X}^{u \times v}$, we have

$m_{\mathbf{X}} = \mathbf{X} \oplus [\lambda_{\mathbf{X}}^1] \oplus [\lambda_{\mathbf{X}}^2] \oplus [\lambda_{\mathbf{X}}^3]$. Here $m_{\mathbf{X}}, [\lambda_{\mathbf{X}}^1], [\lambda_{\mathbf{X}}^2]$, and $[\lambda_{\mathbf{X}}^3]$ are matrices of dimension $u \times v$, and \oplus denote the matrix addition operation. Looking ahead \ominus, \odot will be used to denote matrix subtraction and multiplication operation, respectively. Multiplication of two matrices, $\mathbf{X}^{u \times v}, \mathbf{Y}^{v \times w}$ is a collection of uw independent dot product operations over vectors of length v .

In a convolutional neural network, a convolution operation can be reduced to matrix multiplications [14], [47] as follows. Consider an $f \times f$ kernel over a $w \times h$ input with $p \times p$ padding using $s \times s$ stride having i input channels and o output channels. A convolution can be computed as a matrix multiplication on matrices of dimension $(w' \cdot h') \times (i \cdot f \cdot f)$ and $(i \cdot f \cdot f) \times (o)$ where $w' = \frac{w - f + 2p}{s} + 1$ and $h' = \frac{h - f + 2p}{s} + 1$.

b) Multi-input Multiplication: Inspired from ABY2.0 [23], we design 3-input and 4-input multiplication protocols for our setting. We remark that the multi-input multiplication, when coupled with the optimized PPA circuit from [23], improves the rounds as well as communication in the online phase.

The goal of 3-input multiplication is to generate $[[\cdot]]$ -sharing of $z = abc$ given $[[a]], [[b]], [[c]]$, without the need for performing two sequential multiplications (i.e. first ab then abc). For this parties proceed similar to the multiplication protocol (see §III-B), where they compute $[[z]] = [[z - r]] + [[r]]$. Observe that

$$\begin{aligned} z - r &= abc - r = (m_a - \lambda_a)(m_b - \lambda_b)(m_c - \lambda_c) - r \\ &= m_{abc} - m_{ac}\lambda_b - m_{bc}\lambda_a - m_{ab}\lambda_c + m_a\gamma_{bc} + m_b\gamma_{ac} \\ &\quad + m_c\gamma_{ab} - \gamma_{abc} - r \end{aligned}$$

Similar to the 2-input fair multiplication Π_{Mult} (Fig. 3), the goal of the preprocessing phase is to generate additive shares of $\gamma_{ab}, \gamma_{ac}, \gamma_{bc}, \gamma_{abc}$ among P_1, P_2 . Informally, the terms that P_1, P_2 cannot compute locally for the aforementioned γ values, can be computed by P_0, P_3 , as evident from our sharing semantics. P_0, P_3 compute the missing terms and share them among P_1, P_2 in the preprocessing phase. P_1, P_2 proceed with online phase similar to Π_{Mult} , to compute $z - r$. Thus the online complexity is retained as that of Π_{Mult} while the preprocessing communication is increased to 9 elements. For the 4-input case, the goal is to compute $z = abcd$ for which the additive shares of $\gamma_{ab}, \gamma_{ac}, \gamma_{ad}, \gamma_{bc}, \gamma_{bd}, \gamma_{cd}, \gamma_{abc}, \gamma_{acd}, \gamma_{bcd}, \gamma_{abcd}$ needs to be generated in the preprocessing.

c) Secure Comparison: To compute $a > b$ in the FPA representation, given its $[[\cdot]]$ -sharing, Π_{bitext} uses the technique of extracting the most significant bit (msb) of the value $v = a - b$ [6], [8], [14]. To compute the msb, we use two variants - i) the communication optimized parallel prefix adder (PPA) circuit from ABY3 [6] ($2(\ell - 1)$ AND gates, $\log \ell$ depth), and ii) the round optimized bit extraction circuit from ABY2 [23]. The circuit of ABY2 uses multi-input AND gates and has a multiplicative depth of $\log_4(\ell)$. Both these circuits take two ℓ -bit values in boolean sharing as the input and outputs the result in boolean sharing form. Note that $v = (m_v - \lambda_v^3) + (-\lambda_v^1 - \lambda_v^2)$ as per the sharing semantics (cf. Table III). P_0, P_3 execute Π_{JSn}^B on $(-\lambda_v^1 - \lambda_v^2)$ during the preprocessing, while P_0, P_3 execute Π_{JSn}^B on $(m_v - \lambda_v^3)$ during the online phase to generate the respective boolean sharing.

d) *Bit to Arithmetic*: Protocol Π_{bit2A} enables computing $\llbracket b \rrbracket$ of a bit b given its boolean sharing $\llbracket b \rrbracket^{\text{B}}$. Let b^{R} denotes the value of $b \in \{0, 1\}$ over the arithmetic ring \mathbb{Z}_{2^ℓ} . Then for $b = b_1 \oplus b_2$, note that $b^{\text{R}} = (b_1^{\text{R}} - b_2^{\text{R}})^2$.

Let $b_1 = m_b \oplus \lambda_v^3$ and $b_2 = \lambda_v^1 \oplus \lambda_v^2$. To compute $\llbracket b \rrbracket$, a pair of parties can generate the arithmetic sharing corresponding to b_1^{R} and b_2^{R} by executing Π_{JSh} . $\llbracket b \rrbracket$ can be computed by invoking Π_{Mult} once with inputs $x = y = b_1^{\text{R}} - b_2^{\text{R}}$.

Using the techniques from [4], [14], we obtain a communication-optimized variant by trading off computation in the preprocessing. For this, note that

$$b^{\text{R}} = (m_b \oplus \lambda_b)^{\text{R}} = m_b^{\text{R}} + (\lambda_b)^{\text{R}}(1 - 2m_b^{\text{R}}) \quad (3)$$

Let $v = m_b^{\text{R}}$ and $u = (\lambda_b)^{\text{R}}$. During the preprocessing, P_0 generates $\langle \cdot \rangle$ -sharing of u and a check is executed to verify its correctness. The online phase consists of each pair of parties (P_1, P_3) , (P_2, P_3) and (P_1, P_2) locally computing an additive sharing of b^{R} , generating the corresponding $\llbracket \cdot \rrbracket$ -sharing using Π_{JSh} , and locally adding the shares to obtain $\llbracket b \rrbracket$.

e) *Bit Injection*: Protocol Π_{bitInj} enables computing $\llbracket bv \rrbracket$, given the boolean sharing $\llbracket b \rrbracket^{\text{B}}$ of a bit b and the arithmetic sharing $\llbracket v \rrbracket$ of a value $v \in \mathbb{Z}_{2^\ell}$. Similar to Π_{bit2A} ,

$$\begin{aligned} (bv)^{\text{R}} &= (m_b \oplus \lambda_b)^{\text{R}}(m_v - \lambda_v) \\ &= (m_b^{\text{R}} + (\lambda_b)^{\text{R}}(1 - 2m_b^{\text{R}}))(m_v - \lambda_v) \\ &= m_b^{\text{R}}m_v - m_b^{\text{R}}\lambda_v + (2m_b^{\text{R}} - 1)((\lambda_b)^{\text{R}}\lambda_v - m_v(\lambda_b)^{\text{R}}) \end{aligned}$$

During preprocessing, P_0 generates $\langle \cdot \rangle$ -sharing of λ_b^{R} , followed by verifying its correctness, similar to Π_{bit2A} . $\langle \cdot \rangle$ -shares of $(\lambda_b)^{\text{R}}\lambda_v$ are generated by multiplying $\langle (\lambda_b)^{\text{R}} \rangle$ and $\langle \lambda_v \rangle$ using Π_{Mult} (Fig. 2). In the online phase, each pair of parties (P_1, P_3) , (P_2, P_3) and (P_1, P_2) locally compute an additive sharing of $(bv)^{\text{R}}$, generate its $\llbracket \cdot \rrbracket$ -sharing using Π_{JSh} , and locally add these shares to generate $\llbracket (bv)^{\text{R}} \rrbracket$.

f) *Oblivious Selection*: Given $\llbracket \cdot \rrbracket$ -shares of $x_0, x_1 \in \mathbb{Z}_{2^\ell}$ and $\llbracket b \rrbracket^{\text{B}}$ where $b \in \{0, 1\}$, oblivious selection (Π_{obv}) enables parties to generate re-randomized $\llbracket \cdot \rrbracket$ -shares of $z = x_b$. The protocol is similar in spirit to Oblivious Transfer primitive. Note that z can be written as $z = b(x_1 - x_0) + x_0$. Parties invoke Π_{bitInj} to compute $\llbracket b(x_1 - x_0) \rrbracket$, and sum it with $\llbracket x_0 \rrbracket$ to generate $\llbracket z \rrbracket$.

g) *Piece-wise Polynomials*: Piece-wise polynomial functions are constructed as a series of constant public polynomials f_1, \dots, f_m and $c_1 < \dots < c_m$ such that,

$$f(y) = \begin{cases} 0, & y < c_1 \\ f_1, & c_1 \leq y < c_2 \\ \dots & \\ f_m, & c_m \leq y \end{cases}$$

f can be computed as, $f(y) = \sum_{i=1}^m b_i \cdot (f_i - f_{i-1})$, where $f_0 = 0$, $f_m = 1$, and $b_i = 1$ if $y \geq c_i$ and 0 otherwise, for $i \in \{1, \dots, m\}$. Given the $\llbracket \cdot \rrbracket$ -shares of y , one can obtain the $\llbracket \cdot \rrbracket^{\text{B}}$ -shares of the bits b_1, \dots, b_m using secure comparison. Shares of the product terms, $b_i \cdot (f_i - f_{i-1})$, can thus be generated by invoking m Π_{bitInj} , followed by a local addition. A naive application of Π_{bitInj} involves sharing (via Π_{JSh}) additive shares of $b_i \cdot (f_i - f_{i-1})$, thereby requiring m

Π_{JSh} in the online phase. Instead, it can be made independent of m by first computing additive shares of $f(y)$, and then invoking one Π_{JSh} .

Non-linear activation functions, such as Rectified Linear Unit and Sigmoid, can be viewed as instantiations of piece-wise polynomial functions as shown in ABY3 [6].

h) *ArgMin/ ArgMax*: Protocol Π_{argmin} computes the index of the smallest element in a vector $\vec{x} = (x_1, \dots, x_m)$ of m elements, where \vec{x} is $\llbracket \cdot \rrbracket$ -shared, i.e. each element $x_i \in \mathbb{Z}_{2^\ell}$ of \vec{x} is $\llbracket \cdot \rrbracket$ -shared. The protocol outputs a $\llbracket \cdot \rrbracket^{\text{B}}$ -shared bit vector \vec{b} of size m which has a 1 at the index associated with the minimum value in \vec{x} , and 0 elsewhere. We follow the standard tree-based approach [18] to recursively find the minimum value in \vec{x} while also updating \vec{b} to reflect the index of this smallest element. Each bit of \vec{b} is initialized to 1. The elements of \vec{x} are grouped into pairs and securely compared to find their pairwise minimum. Using this information, \vec{b} is updated such that b_j 's are reset to 0 for x_j 's $\in \vec{x}$ which do not form the minimum in their respective pair; the other bits in \vec{b} still equal 1. The protocol recurses on the remaining elements $x_j \in \vec{x}$, which were the pairwise minimums. Eventually, only one $b_j \in \vec{b}$ equals 1, indicating that x_j is the minimum, with index j . Computing Π_{argmax} can be done similarly.

VI. IMPLEMENTATION AND BENCHMARKING

We benchmark training and inference phases for deep NNs with varying parameter sizes and the inference phase for Support Vector Machines (SVM) using MNIST [48] and CIFAR-10 [49] dataset. Training phase of SVM requires additional tools and primitives, and is out of scope of this work. Benchmarks of the protocols are against the state-of-the-art 4PC of Trident [4] and SWIFT [14] 4PC (supports only inference).

a) *Benchmarking Environment Details*: The protocols are benchmarked over a Wide Area Network (WAN), instantiated using n1-standard-64 instances of Google Cloud², with machines located in East Australia (P_0), South Asia (P_1), South East Asia (P_2), and West Europe (P_3). The machines are equipped with 2.0 GHz Intel (R) Xeon (R) (Skylake) processors supporting hyper-threading, with 64 vCPUs, and 240 GB of RAM Memory. Parties are connected by pairwise authenticated bidirectional synchronous channels (e.g., instantiated via TLS over TCP/IP). We use a bandwidth of 40 MBps between every pair of parties and the average round-trip time (rtt)³ values among P_0 - P_1 , P_0 - P_2 , P_0 - P_3 , P_1 - P_2 , P_1 - P_3 , and P_2 - P_3 are 153.74ms, 93.39ms, 274.84ms, 62.01ms, 174.15ms, and 219.46ms respectively.

For a fair comparison, we implemented and benchmarked all the protocols, including the protocols of Trident and SWIFT, building on the ENCRYPTO library [50] in C++17. Primitives such as maxpool, which Trident and SWIFT do not support, have been run using our building blocks. We would like to clarify that our code is developed for benchmarking, is not optimized for industry-grade use, and optimizations like GPU support can further enhance performance. Our protocols

²<https://cloud.google.com/>

³Time for communicating 1 KB of data between a pair of parties

are instantiated over a 64-bit ring ($\mathbb{Z}_{2^{64}}$), and the collision-resistant hash function is instantiated using SHA-256. We use multi-threading, and our machines are capable of handling a total of 64 threads. Each experiment is run 10 times, and the average values are reported. We use 1 KB = 8192 bits and use a batch size of $B = 128$ for training.

b) Benchmarking Parameters: We evaluate the protocols across a variety of parameters as given in Table V. In addition to parameters such as runtime, communication, and *online throughput* (TP) [4], [6], [19], [21], the cumulative runtime (sum of the up-time of all the hired servers) is also reported. This is because when deployed over third-party cloud servers, one pays for them by the communication and the uptime of the hired servers. To analyze the cost of deployment of the framework, *monetary cost* (Cost) [51] is reported. This is done using the pricing of Google Cloud Platform⁴, where for 1 GB and 1 hour of usage, the costs are USD 0.08 and USD 3.04, respectively. For protocols with an asymmetric communication graph, communication load is unevenly distributed among all the servers, leaving several communication channels underutilized. Load balancing improves the performance by running several execution threads in parallel, each with the roles of the servers changed. Load balancing has been performed in all the protocols benchmarked.

Notation	Description
$T_{on,i}$	Online runtime of party P_i .
$T_{tot,i}$	Total runtime of party P_i .
PT_{on}	Protocol online runtime; $\max_i \{T_{on,i}\}$.
PT_{tot}	Protocol total runtime; $\max_i \{T_{tot,i}\}$.
CT_{on}	Cumulative online runtime; $\sum_i T_{on,i}$.
CT_{tot}	Cumulative total runtime; $\sum_i T_{tot,i}$.
$Comm_{on}$	Online communication.
$Comm_{tot}$	Total communication.
Cost	Total monetary cost.
TP	Online throughput (higher = better) (#iterations / #queries per minute in online)

Table V: Benchmarking parameters (lower is better, except for TP)

c) Network Architectures: We consider the following networks for benchmarking. These were chosen based on the different range of model parameters and types of layers used in the networks. We refer readers to [7], [52] for the architecture and a detailed description of the training and inference steps for the ML algorithms.

- *SVM:* Consists of 10 categories for classification [18].
- *NN-1:* Fully connected network with 3 layers and around 118K parameters [6], [8].
- *NN-2:* Convolutional neural network comprising of 2 hidden layers, with 100 and 10 nodes [4], [6], [10].
- *NN-3:* LeNet [32], comprises of 2 convolutional and fully connected layers, followed by maxpool for convolutional layers. This has approximately 431K parameters.
- *NN-4:* VGG16 [33] has 16 layers in total and contains fully-connected, convolutional, ReLU activation and maxpool layers. This has ≈ 138 million parameters.

⁴See <https://cloud.google.com/vpc/network-pricing> for network cost and <https://cloud.google.com/compute/vm-instance-pricing> for computation cost.

d) Datasets: We use the following datasets:

- MNIST [48] is a collection of 28×28 pixel, handwritten digit images with a label between 0 and 9 for each. It has 60,000 and respectively, 10,000 images in training and test set. We evaluate NN-1, NN-3, SVM on this dataset.
- CIFAR-10 [49] has 32×32 pixel images of 10 different classes such as dogs, horses, etc. It has 50,000 images for training and 10,000 for testing, with 6,000 images in each class. NN-2, NN-4 are evaluated on this dataset.

e) Discussion: Broadly speaking, we consider two deployment scenarios – optimized for time (T), and for cost (C). In the first one, participants want the result of the output as soon as possible while maximizing the online throughput. In the second one, they want the overall monetary cost of the system to be minimal and are willing to tolerate an overhead in the execution time. Using multi-input multiplication gates and the 2 GC variant of the garbled makes the online phase faster but incur an increase in monetary cost. This is because they cause an overhead in communication in the preprocessing phase, and communication affects monetary cost more than uptime (in our setting).

$Tetrad_T$ makes use of multi-input multiplication gates and the 2 GC variant of the garbled world and is the fastest variants of the framework. On the other hand, $Tetrad_C$ is the variant with minimal monetary cost. We only report the numbers for the fair variant of Tetrad and not the robust variant. The overhead for the robust variant over the fair one is minimal, and is primarily due to (i) the use of *robust* joint-send primitive and (ii) the augmented one-time verification check at the end of the preprocessing phase. The overhead amortises for deep networks, like the ones considered in this work.

A. ML Training

For training we consider NN-1, NN-2, NN-3 and NN-4 networks. We report values corresponding to one iteration, that comprises of a forward propagation followed by a backward propagation. More details are provided in §C.

Starting with the time-optimized variant, $Tetrad_T$ is 3–4× faster than Trident in online runtime. The primary factor is the reduction in online rounds of our protocol due to multi-input gates. More precisely, we use the depth-optimized bit extraction circuit while instantiating the ReLU activation function using multi-input AND gates (cf. §V). Looking at the total communication ($Comm_{tot}$) in Table VI, we observe that the gap in $Comm_{tot}$ between $Tetrad_T$ vs. Trident decreases as the networks get deeper. This is justified as the improvement in communication of our dot product with truncation outpaces the overhead in communication caused by multi-input gates. The impact of this is more pronounced with NN-4, as observed by the lower monetary cost of $Tetrad_T$ over Trident. Another reason is that there are two active parties (P_1, P_2) in our framework, whereas Trident has three. Given the allocation of servers, the best rtt Trident can get with three parties (P_0, P_1, P_2) is $153.74ms$, as compared to $62.01ms$ of Tetrad, contributing to Tetrad being faster. However, if the rtt among all the parties were similar, this gap would be closed. Concretely, the online runtime (PT_{on}) of Trident will be similar to that of $Tetrad_C$.

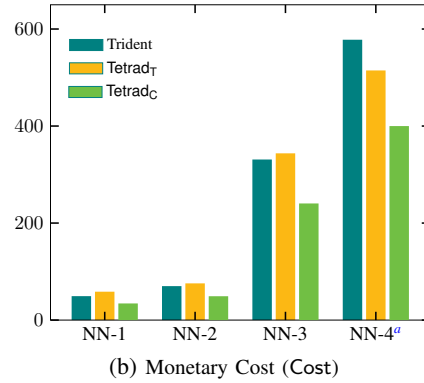
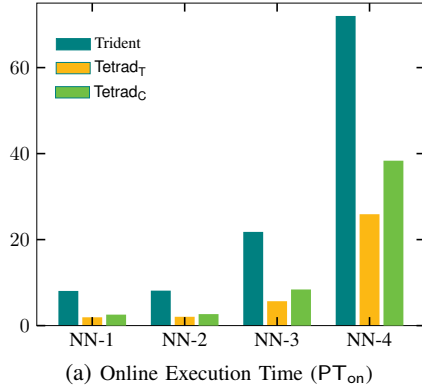


Figure 6: Training of Neural Networks: in terms of PT_{on} and Cost (lower is better) (cf. Table V)

Algorithm	Parameter	Trident	Tetrad _T	Tetrad _C
NN-1	PT_{on}	8.06	1.93	2.55
	PT_{tot}	10.76	5.05	5.27
	CT_{tot}	27.90	12.69	11.22
	$Comm_{tot}$	0.16	0.30	0.16
	Cost	49.33	58.51	34.29
	TP	1904.79	3792.64	3725.49
NN-2	PT_{on}	8.13	2.05	2.67
	PT_{tot}	11.47	5.79	6.14
	CT_{tot}	30.88	14.82	13.40
	$Comm_{tot}$	0.28	0.39	0.24
	Cost	70.00	75.67	49.16
	TP	428.16	652.75	644.69
NN-3	PT_{on}	21.79	5.67	8.40
	PT_{tot}	30.66	15.14	17.87
	CT_{tot}	91.68	40.01	42.76
	$Comm_{tot}$	1.59	1.94	1.28
	Cost	331.01	343.73	240.41
	TP	53.62	55.71	54.13
NN-4	PT_{on}	72.01	25.90	38.35
	PT_{tot}	283.89	182.13	194.58
	CT_{tot}	859.09	500.13	522.32
	$Comm_{tot}$	31.59	29.52	22.24
	Cost	5779.27	5146.10	3999.30
	TP	2.55	2.61	2.56

Table VI: Benchmarking of the training phase of ML algorithms. Time (in seconds) and communication (in GB) are reported for 1 iteration. Monetary cost (USD) is reported for 1000 iterations.

The cost-optimized variant $Tetrad_C$ on the other hand, is $1.5\times$ slower in the online phase compared to $Tetrad_T$. However, it is still faster than Trident owing to the rtt setup, as discussed above. When it comes to monetary cost, this variant is up to 20 – 40% cheaper than its time-optimized counterpart and cheaper by around 30% over Trident.

These trends can be better captured with a pictorial representation as given in Figure 6.

a) *Varying batch sizes and feature sizes:* Table VII shows the online throughput (TP) of neural network (NN-1) training over varying batch sizes and feature sizes using synthetic datasets. We find that both $Tetrad_T$, $Tetrad_C$ are up to $1.8\times$ higher in TP. However, as the batch size and feature size increase, both Trident and Tetrad experience a bandwidth bottleneck. The effect of the bandwidth limitation is higher for Tetrad; hence the gain in TP over Trident decreases a bit.

^ascaled down by a factor of 10 for better visibility

Batch Size	Features	Trident	Tetrad _T	Tetrad _C
128	10	1905.58	5407.35	5271.88
	100	1905.58	5152.29	5029.14
	1000	1904.4	3500.89	3443.6
256	10	1905.58	2818.4	2744.87
	100	1905.58	2747.5	2677.58
	1000	1849.78	2195.3	2150.43

Table VII: Online throughput (TP) of NN-1 training (iterations per minute) over various batch sizes and features.

B. ML Inference

We benchmark the inference phase of SVM and the aforementioned NNs. In addition to Trident [4], we also benchmark against the 4PC robust protocol of SWIFT [14] since it supports NN inference. Note that the best case performance of Fantastic Four [29] when cast in the preprocessing model resembles that of SWIFT, while their worst case execution (3PC malicious) is an order of magnitude slower (cf. §A-D), as demonstrated in their paper (cf. Table 2 of [29]).

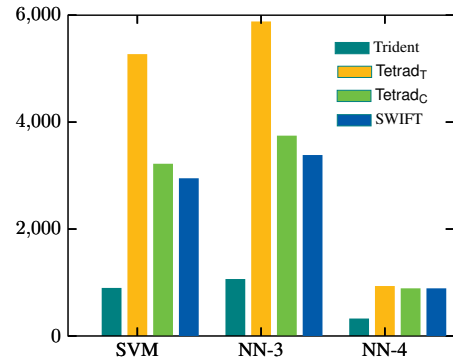


Figure 7: Inference of SVM, NN-3 and NN-4: in terms of TP (higher is better)

Similar to training, the time-optimized variant for inference is faster when it comes to PT_{on} , by 4 – $6\times$ over Trident. This is also reflected in the TP, where the improvement is about 2.8 – $5.5\times$, as evident from Figure 7. In inference, the communication is in the order of megabytes, while run time is in the order of a few seconds. The key observation is that communication is well suited for the bandwidth used (40

MBps). So unlike training, the monetary cost in inference depends more on run time rather than on communication. This is evident from Table VIII which shows that Tetrad_T saves on monetary cost up to a factor of 6 over Trident.

Algorithm	Parameter	Trident	Tetrad_T	Tetrad_C	SWIFT
SVM	PT_{on}	17.09	2.91	4.77	5.21
	PT_{tot}	17.37	3.19	5.05	6.04
	CT_{tot}	47.02	6.99	10.70	14.47
	Comm_{tot}	1.36	2.34	1.25	1.36
	Cost	39.92	6.26	9.23	12.43
	TP	898.80	5271.74	3221.29	2949.76
NN-1	PT_{on}	5.87	1.31	1.87	2.31
	PT_{tot}	6.15	1.58	2.14	3.13
	CT_{tot}	16.75	3.76	4.88	8.65
	Comm_{tot}	0.06	0.09	0.05	0.06
	Cost	14.15	3.19	4.13	7.32
	TP	2615.35	11734.60	8226.93	6661.00
NN-2	PT_{on}	5.87	1.31	1.87	2.31
	PT_{tot}	6.15	1.58	2.14	3.13
	CT_{tot}	16.75	3.77	4.88	8.66
	Comm_{tot}	0.26	0.37	0.22	0.25
	Cost	14.19	3.24	4.16	7.35
	TP	2615.35	11734.60	8226.93	6661.00
NN-3	PT_{on}	14.42	2.61	4.10	4.54
	PT_{tot}	14.71	2.91	4.39	5.39
	CT_{tot}	39.92	6.43	9.40	13.18
	Comm_{tot}	5.62	8.42	4.76	5.39
	Cost	34.59	6.74	8.68	11.97
	TP	1065.35	5882.44	3746.89	3384.51
NN-4	PT_{on}	47.05	7.85	12.69	13.13
	PT_{tot}	47.61	8.44	13.28	14.33
	CT_{tot}	129.41	17.77	27.46	31.35
	Comm_{tot}	85.69	124.09	71.27	81.33
	Cost	122.66	34.40	34.32	39.18
	TP	326.46	934.34	891.19	891.19

Table VIII: Benchmarking of the inference phase of ML algorithms. Time (in seconds) and communication (in MB) are reported for 1 query. Monetary cost (USD) is reported for 1000 queries.

Note that the cost-optimized variant under performs in terms of monetary cost compared to Tetrad_T . This is because, as mentioned earlier, run time plays a bigger role in monetary cost than communication. Hence for inference, the time-optimized variant becomes the optimal choice.

C. Comparison operations

Table IX compares the performance of the frameworks for circuits of varying depth. At each layer of the circuits, we perform 128 comparisons where the comparison results are generated in arithmetic shared form. The idea is that each layer emulates a comparison layer in an NN with a batch size of 128.

Interestingly, beyond a depth of roughly 100, the time-optimized variant (Tetrad_T) starts outperforming in every metric, especially monetary cost, over the cost-optimized one (Tetrad_C). This is because as the depth increases, runtime (CT) grows at a much higher rate than the total communication. What we can infer from Table IX is that if one were to use a DNN with a depth of over 100, Tetrad_T becomes the optimal choice.

FUTURE WORK

Tetrad requires the preprocessing to be function-dependent. Decoupling the preprocessing from the function to

Depth	Parameter	Trident	Tetrad_T	Tetrad_C
128	PT_{on}	3.55	0.53	0.93
	CT_{tot}	9.6	1.06	1.85
	Cost	0.49	0.05	0.09
1024	PT_{on}	28.42	4.23	7.41
	CT_{tot}	76.79	8.47	14.82
	Cost	3.89	0.43	0.75
8192	PT_{on}	227.34	33.87	59.27
	CT_{tot}	614.3	67.76	118.56
	Cost	31.27	3.48	6.03

Table IX: Benchmarking of comparisons over various depths. Each of the layer has 128 comparisons. Time is reported in minutes, and monetary cost in USD.

be computed in the online phase will make the framework more generic and is left as an interesting direction to pursue. Even though fixed-point arithmetic is efficient for the applications considered, in some cases, other representations such as floating-point and posit arithmetic might be desirable. Supporting alternative representations may require rethinking parts of the framework; hence it is left as an open problem.

The following are some of the challenges to be addressed while extending Tetrad to support training of other ML algorithms such as SVM, ResNet and LSTMs. In SVM training, the choice of kernel function plays an important role in determining the efficiency, especially for the non-linear classifiers. Some of the most widely used non-linear kernels include i) Polynomial: $(\vec{x} \odot \vec{y})^d$, ii) Gaussian: $\exp(-\gamma \|\vec{x} - \vec{y}\|^2)$ for $\gamma > 0$, and iii) Hyperbolic: $\tanh(\mu \vec{x} \odot \vec{y} + c)$ for some $\mu > 0$ and $c < 0$, where \vec{x}, \vec{y} denote the input vectors. These kernels are expensive to compute (computation and communication) using standard MPC approaches such as circuit garbling, and hence, demand new MPC-friendly protocols which guarantee efficiency without losing out on accuracy (e.g., Sigmoid approximation of [7]). Further, note that using the naive MPC protocols for training would demand a non-linear increase in bit-size of fixed-point arithmetic to accommodate for an increased dataset size [53]. Concretely, for a dataset with only 212 entries and 14 features, the ring size should be at least 246 bits. Thus, it is necessary to redesign the protocols to enable computation within the standard ring sizes. For deep networks such as ResNet and LSTMs, they require performing batch normalization multiple times, each of which involves division and square-root operations [52]. Since the latter is expensive to perform over rings, designing efficient protocols for these operations is an interesting question.

Finally, although it is known how to instantiate the required primitives securely using standard MPC techniques, they are far from being practically efficient. Moreover, since the secure variant is known to have an overhead over the plaintext computation, sophisticated techniques are required to handle the large amount of intermediate data generated while training very deep networks. Existing PPML frameworks lack support for training the above ML algorithms to the best of our knowledge. We believe that accounting for the points above can bring the existing PPML frameworks, including Tetrad , one step closer to the efficient realization of these algorithms.

ACKNOWLEDGEMENTS

The authors would like to acknowledge support from Google PhD Fellowship 2019, Centre for Networked Intelligence (a Cisco CSR initiative) 2021, SERB MATRICS (Theoretical Sciences) Grant 2020 and Google India AI/ML Research Award 2020. The authors would also like to acknowledge the financial support from Google Cloud to perform the benchmarking.

This project has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation program (grant agreements No. 850990 (PSOTI) and No. 803096 (SPEC)) and from the Digital Research Centre Denmark (DIREC). This work was co-funded by the Deutsche Forschungsgemeinschaft (DFG) – SFB 1119 CROSSING/236615297.

REFERENCES

- [1] J. Alvarez-Valle, P. Bhatu, N. Chandran, D. Gupta, A. V. Nori, A. Rastogi, M. Rathee, R. Sharma, and S. Ugare, “Secure medical image analysis with cryptflow,” *CoRR*, vol. abs/2012.05064, 2020. [Online]. Available: <https://arxiv.org/abs/2012.05064> 1
- [2] M. Byali, H. Chaudhari, A. Patra, and A. Suresh, “FLASH: Fast and robust framework for privacy-preserving machine learning,” *PoPETS*, vol. 2020, no. 2, pp. 459–480, Apr. 2020. 1, 3, 4
- [3] H. Chaudhari, A. Choudhury, A. Patra, and A. Suresh, “ASTRA: High Throughput 3PC over Rings with Application to Secure Prediction,” in *ACM CCSW@CCS*, 2019. [Online]. Available: <https://eprint.iacr.org/2019/429> 1, 3, 4, 15
- [4] H. Chaudhari, R. Rachuri, and A. Suresh, “Trident: Efficient 4PC framework for privacy preserving machine learning,” in *NDSS 2020*. The Internet Society, Feb. 2020. 1, 2, 3, 4, 6, 7, 9, 10, 11, 12, 16, 18
- [5] E. Makri, D. Rotaru, N. P. Smart, and F. Vercauteren, “EPIC: Efficient private image classification (or: Learning from the masters),” in *CT-RSA 2019*, ser. LNCS, M. Matsui, Ed., vol. 11405. Springer, Heidelberg, Mar. 2019, pp. 473–492. 1
- [6] P. Mohassel and P. Rindal, “ABY³: A mixed protocol framework for machine learning,” in *ACM CCS 2018*, D. Lie, M. Mannan, M. Backes, and X. Wang, Eds. ACM Press, Oct. 2018, pp. 35–52. 1, 2, 3, 4, 6, 9, 10, 11, 15, 16, 18
- [7] P. Mohassel and Y. Zhang, “SecureML: A system for scalable privacy-preserving machine learning,” in *2017 IEEE Symposium on Security and Privacy*. IEEE Computer Society Press, May 2017, pp. 19–38. 1, 2, 6, 11, 13, 18
- [8] A. Patra and A. Suresh, “BLAZE: Blazing fast privacy-preserving machine learning,” in *NDSS 2020*. The Internet Society, Feb. 2020. 1, 2, 3, 4, 9, 11, 15, 18
- [9] S. Wagh, D. Gupta, and N. Chandran, “SecureNN: 3-party secure computation for neural network training,” *PoPETS*, vol. 2019, no. 3, pp. 26–49, Jul. 2019. 1, 2
- [10] 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 18*, J. Kim, G.-J. Ahn, S. Kim, Y. Kim, J. López, and T. Kim, Eds. ACM Press, Apr. 2018, pp. 707–721. 1, 9, 11
- [11] A. C.-C. Yao, “Protocols for secure computations (extended abstract),” in *23rd FOCS*. IEEE Computer Society Press, Nov. 1982, pp. 160–164. 1
- [12] O. Goldreich, S. Micali, and A. Wigderson, “How to play any mental game or A completeness theorem for protocols with honest majority,” in *19th ACM STOC*, A. Aho, Ed. ACM Press, May 1987, pp. 218–229. 1
- [13] S. D. Gordon, S. Ranellucci, and X. Wang, “Secure computation with low communication from cross-checking,” in *ASIACRYPT 2018, Part III*, ser. LNCS, T. Peyrin and S. Galbraith, Eds., vol. 11274. Springer, Heidelberg, Dec. 2018, pp. 59–85. 1, 15, 16
- [14] N. Koti, M. Pancholi, A. Patra, and A. Suresh, “SWIFT: Super-fast and Robust Privacy-Preserving Machine Learning,” in *USENIX Security’21*, 2021, <https://eprint.iacr.org/2020/592>. 1, 2, 4, 9, 10, 12, 16, 17
- [15] I. Damgård, C. Orlandi, and M. Simkin, “Yet another compiler for active security or: Efficient MPC over arbitrary rings,” in *CRYPTO 2018, Part II*, ser. LNCS, H. Shacham and A. Boldyreva, Eds., vol. 10992. Springer, Heidelberg, Aug. 2018, pp. 799–829. 1
- [16] M. Keller, V. Pastro, and D. Rotaru, “Overdrive: Making SPDZ great again,” in *EUROCRYPT 2018, Part III*, ser. LNCS, J. B. Nielsen and V. Rijmen, Eds., vol. 10822. Springer, Heidelberg, Apr. / May 2018, pp. 158–189. 1
- [17] D. Demmler, T. Schneider, and M. Zohner, “ABY - A framework for efficient mixed-protocol secure two-party computation,” in *NDSS 2015*. The Internet Society, Feb. 2015. 1, 16
- [18] I. Damgård, D. Escudero, T. K. Frederiksen, M. Keller, P. Scholl, and N. Volgushev, “New primitives for actively-secure MPC over rings with applications to private machine learning,” in *2019 IEEE Symposium on Security and Privacy*. IEEE Computer Society Press, May 2019, pp. 1102–1120. 1, 10, 11, 18
- [19] 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*, E. R. Weippl, S. Katzenbeisser, C. Kruegel, A. C. Myers, and S. Halevi, Eds. ACM Press, Oct. 2016, pp. 805–817. 1, 11, 15
- [20] 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, Part II*, ser. LNCS, J.-S. Coron and J. B. Nielsen, Eds., vol. 10211. Springer, Heidelberg, Apr. / May 2017, pp. 225–255. 1, 15
- [21] 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 *2017 IEEE Symposium on Security and Privacy*. IEEE Computer Society Press, May 2017, pp. 843–862. 1, 11, 15
- [22] M. Abspoel, A. Dalskov, D. Escudero, and A. Nof, “An efficient passive-to-active compiler for honest-majority MPC over rings,” *Cryptology ePrint Archive*, Report 2019/1298, 2019, <https://eprint.iacr.org/2019/1298>. 1
- [23] A. Patra, T. Schneider, A. Suresh, and H. Yalame, “ABY2.0: Improved Mixed-Protocol Secure Two-Party Computation,” in *USENIX Security’21*, 2021, <https://eprint.iacr.org/2020/1225>. 1, 2, 3, 9, 16
- [24] M. Byali, A. Joseph, A. Patra, and D. Ravi, “Fast secure computation for small population over the internet,” in *ACM CCS 2018*, D. Lie, M. Mannan, M. Backes, and X. Wang, Eds. ACM Press, Oct. 2018, pp. 677–694. 1, 3, 8, 16
- [25] M. Byali, C. Hazay, A. Patra, and S. Singla, “Fast actively secure five-party computation with security beyond abort,” in *ACM CCS 2019*, L. Cavallaro, J. Kinder, X. Wang, and J. Katz, Eds. ACM Press, Nov. 2019, pp. 1573–1590. 1
- [26] D. Rotaru and T. Wood, “MARbled circuits: Mixing arithmetic and Boolean circuits with active security,” in *INDOCRYPT 2019*, ser. LNCS, F. Hao, S. Ruj, and S. Sen Gupta, Eds., vol. 11898. Springer, Heidelberg, Dec. 2019, pp. 227–249. 1, 16
- [27] D. Escudero, S. Ghosh, M. Keller, R. Rachuri, and P. Scholl, “Improved primitives for MPC over mixed arithmetic-binary circuits,” in *CRYPTO 2020, Part II*, ser. LNCS, D. Micciancio and T. Ristenpart, Eds., vol. 12171. Springer, Heidelberg, Aug. 2020, pp. 823–852. 1, 16
- [28] S. Mazloom, P. H. Le, S. Ranellucci, and S. D. Gordon, “Secure parallel computation on national scale volumes of data,” in *USENIX Security 2020*, S. Capkun and F. Roesner, Eds. USENIX Association, Aug. 2020, pp. 2487–2504. 2, 16
- [29] A. Dalskov, D. Escudero, and M. Keller, “Fantastic Four: Honest-Majority Four-Party Secure Computation With Malicious Security,” in *USENIX Security’21*, 2021, <https://eprint.iacr.org/2020/1330>. 2, 12, 16, 17
- [30] S. Ohata and K. Nuida, “Communication-efficient (client-aided) secure two-party protocols and its application,” in *FC 2020*, ser. LNCS, J. Bonneau and N. Heninger, Eds., vol. 12059. Springer, Heidelberg, Feb. 2020, pp. 369–385. 2

- [31] Y. Ishai, R. Kumaresan, E. Kushilevitz, and A. Paskin-Cherniavsky, “Secure computation with minimal interaction, revisited,” in *CRYPTO 2015, Part II*, ser. LNCS, R. Gennaro and M. J. B. Robshaw, Eds., vol. 9216. Springer, Heidelberg, Aug. 2015, pp. 359–378. 3, 8, 16
- [32] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-based learning applied to document recognition,” *Proceedings of the IEEE*, pp. 2278–2324, 1998. 3, 11
- [33] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014. 3, 11
- [34] B. Pinkas, M. Rosulek, N. Trieu, and A. Yanai, “SpOT-light: Lightweight private set intersection from sparse OT extension,” in *CRYPTO 2019, Part III*, ser. LNCS, A. Boldyreva and D. Micciancio, Eds., vol. 11694. Springer, Heidelberg, Aug. 2019, pp. 401–431. 3
- [35] M. Fredrikson, S. Jha, and T. Ristenpart, “Model inversion attacks that exploit confidence information and basic countermeasures,” in *ACM CCS 2015*, I. Ray, N. Li, and C. Kruegel, Eds. ACM Press, Oct. 2015, pp. 1322–1333. 3
- [36] F. Tramèr, F. Zhang, A. Juels, M. K. Reiter, and T. Ristenpart, “Stealing machine learning models via prediction APIs,” in *USENIX Security 2016*, T. Holz and S. Savage, Eds. USENIX Association, Aug. 2016, pp. 601–618. 3
- [37] R. Shokri, M. Stronati, C. Song, and V. Shmatikov, “Membership inference attacks against machine learning models,” in *2017 IEEE Symposium on Security and Privacy*. IEEE Computer Society Press, May 2017, pp. 3–18. 3
- [38] N. Koti, A. Patra, R. Rachuri, and A. Suresh, “Tetrad: Actively Secure 4PC for Secure Training and Inference,” *IACR Cryptology ePrint Archive*, 2021, <https://eprint.iacr.org/2021/755>. 4, 7, 8, 9
- [39] M. Abspoel, R. Cramer, I. Damgård, D. Escudero, and C. Yuan, “Efficient information-theoretic secure multiparty computation over $\mathbb{Z}/p^k\mathbb{Z}$ via galois rings,” in *TCC 2019, Part I*, ser. LNCS, D. Hofheinz and A. Rosen, Eds., vol. 11891. Springer, Heidelberg, Dec. 2019, pp. 471–501. 7
- [40] E. Boyle, N. Gilboa, Y. Ishai, and A. Nof, “Practical fully secure three-party computation via sublinear distributed zero-knowledge proofs,” in *ACM CCS 2019*, L. Cavallaro, J. Kinder, X. Wang, and J. Katz, Eds. ACM Press, Nov. 2019, pp. 869–886. 7, 15
- [41] P. Mohassel, M. Rosulek, and Y. Zhang, “Fast and secure three-party computation: The garbled circuit approach,” in *ACM CCS 2015*, I. Ray, N. Li, and C. Kruegel, Eds. ACM Press, Oct. 2015, pp. 591–602. 7, 16
- [42] V. Kolesnikov and T. Schneider, “Improved garbled circuit: Free XOR gates and applications,” in *ICALP 2008, Part II*, ser. LNCS, L. Aceto, I. Damgård, L. A. Goldberg, M. M. Halldórsson, A. Ingólfssdóttir, and I. Walukiewicz, Eds., vol. 5126. Springer, Heidelberg, Jul. 2008, pp. 486–498. 8
- [43] V. Kolesnikov, P. Mohassel, and M. Rosulek, “FlexXOR: Flexible garbling for XOR gates that beats free-XOR,” in *CRYPTO 2014, Part II*, ser. LNCS, J. A. Garay and R. Gennaro, Eds., vol. 8617. Springer, Heidelberg, Aug. 2014, pp. 440–457. 8
- [44] S. Zahur, M. Rosulek, and D. Evans, “Two halves make a whole - reducing data transfer in garbled circuits using half gates,” in *EUROCRYPT 2015, Part II*, ser. LNCS, E. Oswald and M. Fischlin, Eds., vol. 9057. Springer, Heidelberg, Apr. 2015, pp. 220–250. 8
- [45] S. Gueron, Y. Lindell, A. Nof, and B. Pinkas, “Fast garbling of circuits under standard assumptions,” *Journal of Cryptology*, vol. 31, no. 3, pp. 798–844, Jul. 2018. 8
- [46] M. Bellare, V. T. Hoang, S. Keelveedhi, and P. Rogaway, “Efficient garbling from a fixed-key blockcipher,” in *2013 IEEE Symposium on Security and Privacy*. IEEE Computer Society Press, May 2013, pp. 478–492. 8
- [47] Stanford, “CS231n: Convolutional Neural Networks for Visual Recognition,” <https://cs231n.github.io/convolutional-networks/>. 9
- [48] Y. LeCun and C. Cortes, “MNIST handwritten digit database,” 2010. [Online]. Available: <http://yann.lecun.com/exdb/mnist/> 10, 11
- [49] A. Krizhevsky, V. Nair, and G. Hinton, “The CIFAR-10 dataset,” 2014, <https://www.cs.toronto.edu/~kriz/cifar.html>. 10, 11
- [50] Cryptography and P. E. G. at TU Darmstadt, “ENCRYPTO Utils,” https://github.com/encryptogroup/ENCRYPTO_utils, 2017. 10
- [51] P. Miao, S. Patel, M. Raykova, K. Seth, and M. Yung, “Two-sided malicious security for private intersection-sum with cardinality,” in *CRYPTO 2020, Part III*, ser. LNCS, D. Micciancio and T. Ristenpart, Eds., vol. 12172. Springer, Heidelberg, Aug. 2020, pp. 3–33. 11
- [52] S. Wagh, S. Tople, F. Benhamouda, E. Kushilevitz, P. Mittal, and T. Rabin, “Falcon: Honest-majority maliciously secure framework for private deep learning,” *PoPETs*, vol. 2021, no. 1, pp. 188–208, Jan. 2021. 11, 13, 15, 18
- [53] D. Cabarcas, H. D. Vanegas, and D. E. Escudero, “Privacy-preserving machine learning for support vector machines,” Privacy-Preserving Machine Learning Workshop (PPML@CRYPTO’21), 2021. 13
- [54] D. Boneh, E. Boyle, H. Corrigan-Gibbs, N. Gilboa, and Y. Ishai, “Zero-knowledge proofs on secret-shared data via fully linear PCPs,” in *CRYPTO 2019, Part III*, ser. LNCS, A. Boldyreva and D. Micciancio, Eds., vol. 11694. Springer, Heidelberg, Aug. 2019, pp. 67–97. 15
- [55] W. Henecka, S. Kögl, A.-R. Sadeghi, T. Schneider, and I. Wehrenberg, “TASTY: tool for automating secure two-party computations,” in *ACM CCS 2010*, E. Al-Shaer, A. D. Keromytis, and V. Shmatikov, Eds. ACM Press, Oct. 2010, pp. 451–462. 16
- [56] P. Rogaway and T. Shrimpton, “Cryptographic hash-function basics: Definitions, implications, and separations for preimage resistance, second-preimage resistance, and collision resistance,” in *FSE 2004*, ser. LNCS, B. K. Roy and W. Meier, Eds., vol. 3017. Springer, Heidelberg, Feb. 2004, pp. 371–388. 16
- [57] O. Goldreich, *Foundations of Cryptography: Basic Applications*. Cambridge, UK: Cambridge University Press, 2004, vol. 2. 16
- [58] Y. Lindell, “How to simulate it - A tutorial on the simulation proof technique,” *Cryptology ePrint Archive*, Report 2016/046, 2016, <https://eprint.iacr.org/2016/046>. 16
- [59] B. Alon, E. Omri, and A. Paskin-Cherniavsky, “MPC with friends and foes,” in *CRYPTO 2020, Part II*, ser. LNCS, D. Micciancio and T. Ristenpart, Eds., vol. 12171. Springer, Heidelberg, Aug. 2020, pp. 677–706. 17
- [60] P. Pullonen and S. Siim, “Combining secret sharing and garbled circuits for efficient private IEEE 754 floating-point computations,” in *FC 2015 Workshops*, ser. LNCS, M. Brenner, N. Christin, B. Johnson, and K. Rohloff, Eds., vol. 8976. Springer, Heidelberg, Jan. 2015, pp. 172–183. 18

APPENDIX A PRELIMINARIES

A. Related Work

Related work covers MPC protocols with an honest majority for high-throughput and constant-round setting and mixed-protocol frameworks for the case of PPML.

ABY3 [6] was the first framework for the case of 3 parties, supporting both training and inference. It had variants for both passive and active security, with the former being based on [19] and the latter on [20], [21]. ASTRA [3] improved upon the 3PC of [19]–[21] by proposing faster protocols for the online phase with active security. As a result, secure inference of ASTRA is faster than ABY3. Building on [54], BLAZE [8] proposed an actively secure framework that supports inference of neural networks. BLAZE pushes the expensive zero-knowledge part of the computation to the preprocessing phase, making its online phase faster than that of [54]. SWIFT (3PC) improved upon BLAZE by using the distributed zero-knowledge protocol of [40], thereby achieving GOD. In an orthogonal line of work, FALCON [52] focused on enhancing the efficiency of actively secure protocols for large convolutional neural networks, supporting training and inference.

In the high-throughput setting for 4PC, [13] explores protocols for the security notions of abort. Inspired by

the theoretical GOD construction in [13], FLASH proposed practical protocols with GOD for secure inference. Trident [4] improved protocols (in terms of communication) compared to [13] with a focus on security with fairness. In addition, it was the first work to propose a mixed-protocol framework for the case of 4 parties. More recently, [28] improved over [13] to provide support for fixed-point arithmetic with applications to graph parallel computation, albeit with abort security.

Improving the security of Trident to GOD, SWIFT [14] presented an efficient, robust PPML framework with protocols as fast as Trident. SWIFT only supports the secure inference of neural networks and lacks conversions similar to the ones from Trident and the garbled world. Fantastic Four [29] also provides robust 4PC protocols which are on par with SWIFT. While they claim to provide a better security model called *private robustness* compared to SWIFT, it has been shown in SWIFT that the two security models are theoretically equivalent. Our security model is also similar to SWIFT, and we elaborate on its equivalence to private robustness in §A-C.

In the regime of constant-round protocols, [41] presents 3PC protocols in the honest majority setting satisfying security with abort, which require communicating one garbled circuit and three rounds of interaction. The work of [31] presents a robust 4-party computation protocol (4PC) with GOD in 2-rounds (which is optimal) at the expense of 12 garbled circuits. Further, [24] presents efficient 3PC and 4PC constructions providing security notions of fairness and GOD.

A mixed-protocol framework for MPC was first shown to be practical, in the 2-party dishonest majority setting, by TASTY [55]. TASTY was a passively secure compiler supporting generation of protocols based on homomorphic encryption and garbled circuits. This was followed by ABY [17], which proposed a mixed protocol framework, also with passive security, combining the arithmetic, boolean and garbled worlds. The work of ABY2 [23] improves upon ABY, providing a faster online phase with applications to PPML. The work of [26], [27] proposed efficient mixed world conversions for the case of n parties with a dishonest majority. Both works have active security, with [26] supporting the inference of SVMs, and [27] supporting neural network inference.

In the honest majority setting, ABY3 [6] extended the idea to 3 parties and provided specialised protocols for the case of PPML. ABY3 was the first work to support secure training in the case of 3 parties, while Trident [4] extended it to the 4-party setting.

B. Basic Primitives

a) Shared Key Setup: Let $F : \{0, 1\}^\kappa \times \{0, 1\}^\kappa \rightarrow X$ be a secure pseudo-random function (PRF), with co-domain X being \mathbb{Z}_{2^ℓ} . The following set of keys are established.

- One key between every pair – k_{ij} for P_i, P_j .
- One key between every set of three parties – k_{ijk} for P_i, P_j, P_k .
- One shared keys $k_{\mathcal{P}}$ known to all parties in \mathcal{P} .

Suppose P_0, P_1 wish to sample a random value $r \in \mathbb{Z}_{2^\ell}$ non-interactively. To do so they invoke $F_{k_{01}}(id_{01})$ and obtain r . Here, id_{01} denotes a counter maintained by the parties, and is updated after every PRF invocation. The appropriate keys

used to sample is implicit from the context, from the identities of the pair that sample or from the fact that it is sampled by all, and, hence, is omitted.

Functionality $\mathcal{F}_{\text{SETUP}}$

$\mathcal{F}_{\text{SETUP}}$ interacts with the parties in \mathcal{P} and the adversary \mathcal{S} . $\mathcal{F}_{\text{SETUP}}$ picks random keys k_{ij} and k_{ijk} for $i, j, k \in \{0, 1, 2, 3\}$ and $k_{\mathcal{P}}$. Let y_s denote the keys corresponding to party P_s . Then

- $y_s = (k_{01}, k_{02}, k_{03}, k_{012}, k_{013}, k_{023}$ and $k_{\mathcal{P}})$ when $P_s = P_0$.
- $y_s = (k_{01}, k_{12}, k_{13}, k_{012}, k_{013}, k_{123}$ and $k_{\mathcal{P}})$ when $P_s = P_1$.
- $y_s = (k_{02}, k_{12}, k_{23}, k_{012}, k_{023}, k_{123}$ and $k_{\mathcal{P}})$ when $P_s = P_2$.
- $y_s = (k_{03}, k_{13}, k_{23}, k_{013}, k_{023}, k_{123}$ and $k_{\mathcal{P}})$ when $P_s = P_3$.

Output: Send (Output, y_s) to every $P_s \in \mathcal{P}$.

Figure 8: Ideal functionality for shared-key setup

The key setup is modelled via a functionality $\mathcal{F}_{\text{SETUP}}$ (Fig. 8) that can be realised using any secure MPC protocol. A simple instantiation of such an MPC protocol is as follows. P_i samples key k_{ij} and sends to P_j . P_i samples k_{ijk} and sends to P_j . P_i, P_j jsnd k_{ijk} to P_k . Similarly, P_0 samples $k_{\mathcal{P}}$ and sends to P_3 . P_0, P_3 jsnd $k_{\mathcal{P}}$ to P_1 and P_2 .

b) Collision-Resistant Hash Function [56]: . A family of hash functions $\{H : \mathcal{K} \times \mathcal{M} \rightarrow \mathcal{Y}\}$ is said to be collision resistant if for all PPT adversaries \mathcal{A} , given the hash function H_k for $k \in_R \mathcal{K}$, the following holds: $\Pr[(x, x') \leftarrow \mathcal{A}(k) : (x \neq x') \wedge H_k(x) = H_k(x')] = \text{negl}(\kappa)$, where $x, x' \in \{0, 1\}^m$ and $m = \text{poly}(\kappa)$.

C. Security Model

We prove security using the real-world/ ideal-world simulation paradigm [57], [58]. The security is analyzed by comparing what an adversary can do in the real world's execution of the protocol with what it can do in an ideal world execution where there is a trusted third party and is considered secure by definition. In the ideal world, parties send their inputs to the trusted third party over perfectly secure channels that carries out the computation and sends the output to the parties. Informally, a protocol is secure if whatever an adversary can do in the real world can also be done in the ideal world.

Functionality $\mathcal{F}_{\text{FAIR}}$

Every honest party $P_i \in \mathcal{P}$ sends its input x_i to the functionality. Corrupted parties may send arbitrary inputs as instructed by the adversary. While sending the inputs, the adversary is also allowed to send a special abort command.

Input: On message (Input, x_i) from P_i , do the following: if (Input, $*$) already received from P_i , then ignore the current message. Otherwise, record $x'_i = x_i$ internally. If x_i is outside P_i 's domain, consider $x'_i = \text{abort}$.

Output: If there exists an $i \in \{0, 1, 2, 3\}$ such that $x'_i = \text{abort}$, send (Output, \perp) to all the parties. Else, compute $y = f(x'_0, x'_1, x'_2, x'_3)$ and send (Output, y) to all parties.

Figure 9: Fair functionality for computing function f

Let \mathcal{A} denote the probabilistic polynomial time (PPT) real-world adversary corrupting at most one party in \mathcal{P} , \mathcal{S} denote the corresponding ideal world adversary, and \mathcal{F} denote the ideal functionality. Let $\text{IDEAL}_{\mathcal{F}, \mathcal{S}}(1^\kappa, z)$ denote the joint output of the honest parties and \mathcal{S} from the ideal execution

with respect to the security parameter κ and auxiliary input z . Similarly, let $\text{REAL}_{\Pi, \mathcal{A}}(1^\kappa, z)$ denote the joint output of the honest parties and \mathcal{A} from the real world execution. We say that the protocol Π securely realizes \mathcal{F} if for every PPT adversary \mathcal{A} there exists an ideal world adversary \mathcal{S} corrupting the same parties such that $\text{IDEAL}_{\mathcal{F}, \mathcal{S}}(1^\kappa, z)$ and $\text{REAL}_{\Pi, \mathcal{A}}(1^\kappa, z)$ are computationally indistinguishable. The ideal functionality for computing a function f with fairness and robustness appears in Fig. 9 and Fig. 10, respectively.

Functionality $\mathcal{F}_{\text{ROBUST}}$

Every honest party $P_i \in \mathcal{P}$ sends its input x_i to the functionality. Corrupted parties may send arbitrary inputs as instructed by the adversary.

Input: On message (Input, x_i) from P_i , do the following: if (Input, $*$) already received from P_i , then ignore the current message. Otherwise, record $x'_i = x_i$ internally. If x_i is outside P_i 's domain, consider x'_i to be some predetermined default value.

Output: Compute $y = f(x'_0, x'_1, x'_2, x'_3)$ and send (Output, y) to all parties.

Figure 10: Robust functionality for computing function f

a) *On the security of robust Tetrad:* We emphasize that we follow the standard traditional (real-world / ideal-world based) security definition of MPC, according to which, in the 4-party setting with one corruption, exactly one party is assumed to be corrupt, and the rest are *honest*. As per this definition, disclosing the honest parties' inputs to a selected *honest* party is *not* a breach of security. Indeed in Tetrad, the data sharing and the computation on the shared data are done so that any malicious behaviour leads to establishing a trusted party P_{TP} who is enabled to receive all the inputs and compute the output on the clear. There has been a recent study on the additional requirement of hiding the inputs from a quorum of honest parties (treating them as semi-honest), termed as Friends-and-Foes (FaF) security notion [59]. This is a stronger security goal than the standard one. Informally, designing secure 4PC FaF protocols requires security against two independent corruptions. Our sharing semantics, designed to handle only one corruption, does not suffice. Hence, we leave FaF-secure 4PC for future exploration.

Another security notion, called *private robustness*, was recently proposed in the work of Dalskov et al. [29], where the protocol does not demand the inputs be sent to a P_{TP} . Their work, however, considers a more restricted security model, where it is assumed that parties will discard messages which are *non-intended* and are not a part of the protocol. This involves assuming a *secure erasure*. Under this assumption, our model is equivalent to that of private robustness since the trusted party P_{TP} will erase the input of the honest parties after computing the function output.

D. Comparison with Fantastic Four [29]

We analyse the performance of Fantastic Four [29] where execution proceeds in segments (cf. §6.4, [29]). Elaborately, computation is carried out optimistically for each segment, followed by a verification phase before proceeding to the next segment. If verification fails, the current segment is recomputed via an active 3PC protocol. Subsequent segments

also proceed with a 3PC execution until the verification fails again. In this case, a semi-honest 2PC with a helper is carried out for the current and rest of the segments. For analysis, we consider their best and worst-case execution cost.

Protocol	Dot Product w/ Truncation		#Active Parties
	Preprocessing	Online	
Fantastic Four: Case I	ℓ	9ℓ	4
Fantastic Four: Case II	$76(\ell + \kappa) + 54x + 12$	$9\ell + 6\kappa$	3
Tetrad-R(on-demand)	-	5ℓ	3

Table X: Comparison with Fantastic Four [29]

Observe that the best case happens when the verification is always successful, which we call as *Case I*. In this case, the communication cost is that of the 4PC execution. Note that an adversary can *always* make the verification fail in the first segment itself. This results in executing the entire protocol (all segments) with their active 3PC, which accounts for their worst-case cost. We denote this as *Case II*. Their 3PC protocols are designed to work over the extended ring of size $\ell + \kappa$ bits. As evident from Tables 2, 3 of their paper, their 3PC is at least $10\times$ more expensive than their 4PC in terms of both runtime and communication. Thus, the higher cost of 3PC defeats the purpose of having an additional honest party in the system.

Observe that their protocols are designed to work with a function-independent preprocessing. Thus, for a fair comparison, we compare both cases against the on-demand variant of our robust protocols (Tetrad-R). The results are summarised in Table X. We remark that the values for their cases are obtained from Table 1 of their paper [29].

APPENDIX B 4PC PROTOCOL

Here we detail the additional information regarding the 4PC protocols.

A. Joint-send for robust protocols

The formal protocol for Π_{jsnd} in the robust setting [14] is given in Fig. 11.

Protocol $\Pi_{\text{jsnd}}(P_i, P_j, v, P_k)$

$P_s \in \mathcal{P}$ initializes an inconsistency bit $b_s = 0$. If P_s remains silent instead of sending b_s in any of the following rounds, the recipient sets b_s to 1.

- *Send:* P_i sends v to P_k .
- *Verify:* P_j sends $H(v)$ to P_k .
- P_k sets $b_k = 1$ if the received values are inconsistent or if the value is not received.
- P_k sends b_k to all parties. P_s for $s \in \{i, j, l\}$ sets $b_s = b_k$.
- P_s for $s \in \{i, j, l\}$ mutually exchange their bits. P_s resets $b_s = b'$ where b' denotes the bit which appears in majority among b_i, b_j, b_l .
- All parties set $P_{\text{TP}} = P_l$ if $b' = 1$, terminate otherwise.

Figure 11: Joint-Send for robust protocols

Lemma B.1 (Communication). *Protocol Π_{jsnd} (Fig. 11) requires an amortized communication of ℓ bits and 1 round.*

Proof: In the protocol $\Pi_{\text{jsnd}}(P_i, P_j, v, P_k)$ for the fair variant, P_i communicates v to P_k requiring communication of ℓ bits and one round. The hash value communication from P_j to P_k can be clubbed for multiple instances with the same set of parties and hence the cost gets amortized. The analysis is similar for the robust case as well. Here, though the verification consists of multiple steps, the cost gets amortized over multiple instances. ■

B. Function-independent preprocessing

We provide the fair multiplication, $\Pi_{\text{Mult}}^{\text{NoPre}}$, for *function-independent* preprocessing in Fig. 12. The protocol incurs no overhead over the fair multiplication (Π_{Mult}) in Tetrad. This is due to the design of Π_{Mult} where values u^1, u^2 are sampled non-interactively in the preprocessing. Thus the joint-sharing by P_0, P_3 (Step 5 (a) in Fig. 12) can be performed along with the communication among P_1, P_2 (Step 4 in Fig. 12) in the online. Moreover, the rest of the communication can be deferred till the verification stage and thus, the online round complexity is retained. The protocol for robust setting is similar.

Protocol $\Pi_{\text{Mult}}^{\text{NoPre}}(a, b, \text{isTr})$

Let isTr be a bit that denotes whether truncation is required ($\text{isTr} = 1$) or not ($\text{isTr} = 0$).

Online:

- 1) Locally compute the following:

$$\begin{aligned} P_0, P_1 : \gamma_{ab}^1 &= \lambda_a^1 \lambda_b^3 + \lambda_a^3 \lambda_b^1 + \lambda_a^3 \lambda_b^3 \\ P_0, P_2 : \gamma_{ab}^2 &= \lambda_a^2 \lambda_b^3 + \lambda_a^3 \lambda_b^2 + \lambda_a^2 \lambda_b^2 \\ P_0, P_3 : \gamma_{ab}^3 &= \lambda_a^1 \lambda_b^2 + \lambda_a^2 \lambda_b^1 + \lambda_a^1 \lambda_b^1 \end{aligned}$$

- 2) P_0, P_3 and P_j sample random $u^j \in_R \mathbb{Z}_{2^\ell}$ for $j \in \{1, 2\}$. Let $u^1 + u^2 = \gamma_{ab}^3 - r$ for a random $r \in_R \mathbb{Z}_{2^\ell}$.

- 3) Let $y = (z - r) - m_a m_b$. Locally compute the following:

$$\begin{aligned} P_1 : y_1 &= -\lambda_a^1 m_b - \lambda_b^1 m_a + \gamma_{ab}^1 + u^1 \\ P_2 : y_2 &= -\lambda_a^2 m_b - \lambda_b^2 m_a + \gamma_{ab}^2 + u^2 \\ P_1, P_2 : y_3 &= -\lambda_a^3 m_b - \lambda_b^3 m_a \end{aligned}$$

- 4) P_1 sends y_1 to P_2 , while P_2 sends y_2 to P_1 .

- 5) Parties proceed as follows:

- a) P_0, P_3 : $r = \gamma_{ab}^3 - u^1 - u^2$; $q = r^t$ if $\text{isTr} = 1$, else $q = r$; Execute $\Pi_{\text{JSh}}(P_0, P_3, q)$.
- b) P_1, P_2 : $z - r = (y_1 + y_2 + y_3) + m_a m_b$; $p = (z - r)^t$ if $\text{isTr} = 1$, else $p = z - r$; Execute $\Pi_{\text{JSh}}(P_1, P_2, p)$.

- 6) Locally compute $[[o]] = [[p]] + [[q]]$. Here $o = z^t$ if $\text{isTr} = 1$ and z otherwise.

Verification:

- 1) P_0, P_1, P_2 sample random $s_1, s_2 \in_R \mathbb{Z}_{2^\ell}$ and set $s = s_1 + s_2$. P_0 sends $w = \gamma_{ab}^1 + \gamma_{ab}^2 + s$ to P_3 .
- 2) P_3 computes $v = -(\lambda_a^1 + \lambda_a^2) m_b - (\lambda_b^1 + \lambda_b^2) m_a + u^1 + u^2 + w$ and sends $H(v)$ to P_1 and P_2 . Parties P_1, P_2 abort iff $H(v) \neq H(y_1 + y_2 + s)$.

Figure 12: Fair multiplication without preprocessing.

C. Input Sharing & Output Reconstruction

Lemma B.2 (Communication). *Protocol Π_{Sh} (Fig. 1) requires an amortized communication of at most 3ℓ bits and 1 round in the online phase.*

Protocol $\Pi_{\text{Rec}}(\mathcal{P}, [[v]])$

- 1) P_1, P_0 jsnd λ_v^1 to P_2 ; P_2, P_0 jsnd λ_v^3 to P_3 ; P_3, P_0 jsnd λ_v^2 to P_1 ; P_1, P_2 jsnd m_v to P_0 .
- 2) Compute $v = m_v - \lambda_v^1 - \lambda_v^2 - \lambda_v^3$.

Figure 13: Reconstruction (with abort) of v among \mathcal{P} .

Lemma B.3 (Communication). *Protocol Π_{Rec} (Fig. 13) requires an amortized communication of 4ℓ bits and 1 round in the online phase.*

APPENDIX C ML ALGORITHMS

a) Training and Inference of NN: An NN can be divided into various layers, where each layer contains a predefined number of nodes. These nodes are a linear function composed of a non-linear ‘‘activation’’ function. The nodes at the input layer are evaluated on the input features to evaluate a neural network. The outputs from these nodes are fed as inputs to the nodes in the next layer. This process is repeated for all the layers to obtain the output. The underlying operation involved is a computation of activation matrices for all the layers. This constitutes the forward propagation phase. The backward propagation involves adjusting model parameters according to the difference in the computed output and the actual output and comprises computing error matrices.

Concretely, each layer comprises matrix multiplications followed by an application of the ReLU function. The maxpool layer additionally follows convolutional layers after the ReLU layer. After evaluating the layers in a sequential manner, at the output layer, we use the MPC friendly variant of the softmax activation function, $\text{softmax}(u_i) = \frac{\text{ReLU}(u_i)}{\sum_{j=1}^n \text{ReLU}(u_j)}$, proposed by SecureML [7]. To perform the division, we switch from arithmetic to garbled world and then use a division garbled circuit [60] followed by a switch back to the arithmetic world. For training, we use Gradient Descent, where the forward propagation comprises computing activation matrices for all the layers in the network. The backward propagation comprises computing error matrices involving matrix multiplications with derivative of maxpool and derivative of ReLU, depending on the network architecture. We refer readers to [4], [6]–[8], [52] for formal details.

b) Inference of SVM: SVM is a function which takes as input an n -dimensional *feature vector*, \vec{x} , and outputs the *category* to which the feature vector belongs. SVM is implemented as a matrix \mathbf{F} , of dimension $q \times n$ where each row of \mathbf{F} is called the support vector and a vector $\vec{b} = (b_1, \dots, b_q)$, is called the *bias*. Each element of \mathbf{F} and \vec{b} lies in \mathbb{Z}_{2^ℓ} . Each support vector along with a scalar from the bias can classify the input \vec{x} into a specific category. More precisely, let \mathbf{F}_i denote the i^{th} row of matrix \mathbf{F} . Then, the value $\mathbf{F}_i \cdot \vec{x} + b_i$ specifies how likely \vec{x} is to be in category i . To find the most likely category, we compute argmax over these values, i.e. $\text{category}(\vec{x}) = \text{argmax}_{i \in \{1, \dots, q\}} \mathbf{F}_i \cdot \vec{x} + b_i$. We refer the readers to [18] for more details.